# Spectral and active learning for enhanced and computationally scalable quantum molecular dynamics

Dissertation

zur Erlangung des Doktorgrades

der Fakultät für Mathematik, Informatik und Naturwissenschaften

der Universität Hamburg

vorgelegt im Fachbereich Mathematik von

Yahya SALEH

Hamburg, 2023

Als Dissertation angenommen vom Fachbereich Mathematik der Universität Hamburg auf Grund der Gutachten von

Accepted as a dissertation by the department of Mathematics at the Universität Hamburg based on the review of

Prof. Dr. Armin Iske Prof. Dr. Jochen Küpper Prof. Dr. Stephan Wojtowytsch

Die Disputation fand am 12. Mai, 2023 im Center for Free-Electron Laser Science CFEL, Hamburg statt.

The defense took place on May 12, 2023, in the Center for Free-Electron Laser Science CFEL, Hamburg.

Die Mitglieder der Prüngskommission waren

The members of the examination committee were

Prof. Dr. Ralf Holtkamp Prof. Dr. Armin Iske Prof. Dr. Jochen Küpper Prof. Dr. Jörg Teschner Prof. Dr. Stephan Wojtowytsch

> Prof. Dr. Armin Iske Leiter des Fachbereichs Mathematik

#### Abstract

Major problems in quantum molecular physics, such as quantum modeling of chemical dynamics and studying atmospheres of hot exoplanets, pose prohibitive computational challenges that are at the forefront of research efforts in numerical analysis and scientific computing. In particular, it is often required in these applications to approximate a vast number of highly-oscillatory solutions of infinitedimensional eigenvalue problems and evolution equations in very high dimensions. Such tasks are rather prohibitive using standard linear numerical schemes due to the curse of dimensionality phenomena and the need for high resolution for a reliable modeling of the oscillatory behavior. While nonlinear approximation concepts, e.g., neural networks, promise to mitigate the curse of dimensionality and offer better approximation abilities, they are rather fragile and not straightforward to use for large-scale problems in quantum molecular physics. This work aims at reducing the computational costs of numerical quantum simulations of molecular systems and their scaling *via* novel *spectral* and *active learning* algorithms. The proposed algorithms expand the approximation capabilities of standard linear methods while maintaining a high robustness.

In particular, standard spectral methods for solving differential equations are extended to a *spectral learning* framework, where standard bases of innerproduct spaces are composed with invertible neural networks. I provide sufficient conditions on the neural networks to guarantee that the resulting sequence of functions is also a basis of the underlying inner-product space. This allows one to define well-posed numerical schemes using these augmented bases to solve various approximation problems. As application, I derive convergence guarantees of spectral learning for approximating Schwartz functions and eigenvalues to Schrödinger operators as the size of the truncated basis goes to infinity. Moreover, it is shown that the convergence of spectral learning is faster than that of spectral methods. This is achieved by showing that the total variation of compactlysupported smooth functions with respect to the push-forward measures induced by neural networks admit minima. Theoretical results are supported by numerical simulations to compute the spectra of linear infinite-dimensional operators that characterize nuclear motions in polyatomic molecules. Results show a two order increased accuracy over standard spectral methods upon the use of invertible neural networks.

In addition, I consider the construction of molecular potential energy surfaces, i. e., the solution of the parametric electronic Schrödinger equation in an active learning paradigm. I derive an upper bound on the generalization error in an active learning setting, which gives a theoretical insight into the construction of active learning algorithms. I propose and implement an algorithm that follows this insight and that allows one to infer the solution operator with minimial datasets. Simulations are performed on polyatomic molecules and results indicate a roughly two-times faster convergence than what is possible *via* common active learning algorithms.

### Zusammenfassung

Wichtige Probleme der molekularen Quantenphysik, wie die quantenmechanische Modellierung chemischer Dynamiken und die Untersuchung von Atmosphären heißer Exoplaneten, stellen unüberwindbare Rechenschwierigkeiten dar, die an der Spitze der Forschungsbemühungen im Bereich der numerischen Analyse und des wissenschafltichen Rechnens stehen. In diesen Anwendungen ist es oft erforderlich, eine große Anzahl stark oszillierender Lösungen von Eigenwertproblemen unendlicher Dimensionlaität und Evolutionsgleichungen in sehr hohen Dimensionen zu approximieren. Derartige Aufgaben sind mit linearen Standardverfahren aufgrund des Fluches der Dimensionalität und der Notwendigkeit einer hohen Auflösung zur zuverlässigen Modellierung des oszillierenden Verhaltens schwer zu lösen. Während nichtlineare Approximationskonzepte, z. B. neuronale Netze, versprechen, den Fluch der Dimensionalität zu mildern, und bessere Approximationsfähigkeiten bieten, sind sie eher fragil und nicht einfach für groß angelegte Probleme in der molekularen Qauntenphysik anzuwenden. Diese Arbeit zielt darauf, die Rechenkosten numerischer quantenmechanischer Simulationen von molekularen Systemen und deren Skalierung durch neuartige spektrale und aktive Lernalgorithmen zu reduzieren. Die vorgeschlagenen Algorithmen erweitern die Approximationsfähigkeiten von linearen Standardverfahren und zeigen gleichzeitig eine hohe Robustheit.

Insbesondere werden spektrale Standardverfahren zur Lösung von Differentialgleichungen in einem spektralen Lernrahmen erweitert, in dem Standardbasen von Innenprodukträumen mit invertierbaren neuronalen Netzen komponiert werden. Es werden ausreichende Bedingungen für die neuronalen Netze bereitgestellt, um sicherzustellen, dass die resultierende Funktionenfolge ebenfalls eine Basis des zugrunde liegenden Innenproduktraums ist. Durch die Verwendung dieser erweiterten Basen können gut gestellte numerische Verfahren definiert werden, um verschiedene Approximationsprobleme zu lösen. Als Anwendung werden Konvergenzgarantien des spektralen Lernens zur Approximation von Schwartz-Funktionen sowie Eigenwerten von Schrödinger-Operatoren abgeleitet, wenn die Größe der abgeschnittenen Basis gegen Unendlich geht. Darüber hinaus wird gezeigt, dass die Konvergenz des spektralen Lernens schneller ist als die von spektralen Verfahren. Dies wird dadurch gezeigt, dass die Variation von kompakten glatten Funktionen bezüglich der durch neuronale Netze induzierten Push-Forward-Masse Minima aufweist. Die theoretischen Ergebnisse werden durch numerische Simulationen zur Berechnung der Spektren linearer unendlichdimensionaler Operatoren gestützt, die nukleare Bewegungen in polyatomaren Molekülen charakterisieren. Die Ergebnisse zeigen eine um zwei Ordnungen höhere Genauigkeit gegenüber spektraler Standardverfahren durch die Verwendung invertierbarer neuronaler Netze.

Zusätzlich betrachte ich die Konstruktion von Potentialenergieflächen von Molekülen, d.h. die Lösung der parametrischen elektronischen Schrödingergleichung in einem aktiven Lernparadigma. Ich leite eine obere Schranke für den Generalisierungsfehler in einem aktiven Lernparadigma her, die Einblicke in die Konstruktion aktiver Lernalgorithmen gibt. Diesem Ansatz folgend schlage ich einen Algorithmus vor und setze ihn um. Er erlaubt, den Lösungsoperator mit minimalen Datensätzen zu bestimmen. Simulationen werden an polyatomaren Molekülen durchgeführt und die Ergebnisse deuten auf eine etwa doppelt so schnelle Konvergenz hin als bei herkömmlichen aktiven Lernalgorithmen möglich ist.

#### Danksagung

Ich möchte mich bei meinen Betreuern, Dr. Andrey Yachmenev, Prof. Dr. Armin Iske und Prof. Dr. Jochen Küpper, dafür bedanken, dass sie mich in praxisrelevante und herausfordernde Rechenprobleme der molekularen Quantenphysik, sowie in Approximationstheorien für maschinelles Lernen eingeführt haben. Besonderes schätze ich die Beharrlichkeit, mit der sie sowohl mathematische Rigorosität als auch praxis-/experimentnahe Berechnungen gefördert haben, ohne eines über dem anderen zu vernachlässigen. Ich bin sehr dankbar für die Zeit, die ich im Team vom Dr. Yachmenev verbracht habe und durch die ich meine Programmierkenntnisse im Bereich des wissenschaftlichen Rechnens wesentlich verbessern konnte. Ich bedanke mich sehr bei Prof. Dr. Stephan Wojtowytsch für das kritische Gutachten meiner Doktorarbeit und die nützlichen Kommentare, die zur Verbesserung ihrer Qualität beigetragen haben. Ich danke Prof. Dr. Jörg Teschner für das Übernehmen des Prüfungskommissionsvorsitzes und Prof. Dr. Ralph Holtkamp für die Protokollführung. I am also very thankful to my friends and colleagues Vishnu Sanjay and Jannik Eggers, who accompanied my initial ideas on spectral learning and enriched them with deep and interesting discussions.

Während der Doktorarbeit genoß ich die Präsenz und Unterstützung vieler Freunde, die zu zahlreich sind, um sie namentlich zu erwähnen. In particolare ringrazio Eleonora per la sua gioiosa compagnia nei miei primi passi del dottorato e per il suo continuo supporto, soprattutto durante la pandemia. Un altro grande ringraziamento va a Nicola e Adriana, due amici su cui potevo e posso sempre contare, per tutti i momenti divertenti di cucina, danza e sport. Mein herzlicher Dank gilt Martina, die immer ein offenes Ohr hatte und einen sicheren Raum für mich und meine Gedanken angeboten hat. J'apprécie beaucoup le soutien et la patience de Nathalie et sa joyeuse présence avec moi dans plusieurs conférences, et festivals Forró. I am grateful for the fun times with my friends Thea, Carlos, Sebastian and Majd and their continuous support. Ich bin Saskia für ihres immer Dasein und Unterstützung in meinen letzten Schritten der Doktorarbeit sehr dankbar. Mein Dank gilt meiner Familie in München, der Familie von Ammo Hussam für die unbegrenzte Unterstützung und die vielen schönen Zeiten zusammen seit dem Beginn meines Masterstudiums. Besonderes dankbar bin ich für die spannenden Diskussionen in den langen Weihnachtsnächten mit Ammo, die mich geprägt haben und immer noch prägen. (أشكر عائلتي الصغيرة، أي، فادية، وأخي، Ich zolle meinem Vater, Hani, meine Ehrerbietung, der mich, meine Ambitionen, und meine Leidenschaft für Wissenschaft geprägt hat.

# Contents

1	Intr	Introduction			
2	Quantum molecular dynamics				
	2.1	Modeling of quantum molecular mechanics		10	
		2.1.1	Molecular Hamiltonians	12	
		2.1.2	Born-Oppenheimer approximation	13	
	2.2	Electr	onic Schrödinger equations and potential energy surfaces	15	
2.3 Schrödinger equations for the nuclei			dinger equations for the nuclei	16	
		2.3.1	A variational formulation	17	
		2.3.2	Spectral discretization and the curse of dimensionality $\ldots$	19	
3	Acti	ve lear	ning for constructing potential energy surfaces	25	
	3.1	Forma	al setting and notation	27	
3.2 Empirical risk minimization principles and generalization e		rical risk minimization principles and generalization errors .	28		
		3.2.1	Upper bound to the true risk in passive learning	30	
		3.2.2	Upper bound to the true risk in active learning	30	
	3.3	Practi	cal pool-based active learning	34	
3.4 Simulations on $pyrrole(H_2O)$		ations on $pyrrole(H_2O)$	37		
		3.4.1	Performance	40	
		3.4.2	Distribution of queried data	44	
		3.4.3	Batch size and size of initially labeled dataset	45	
3.5 Summary and Conclusion				47	

4	Spe	ctral le	arning: augmenting bases with normalizing flows	53		
	4.1	Augn	nenting expressivity <i>via</i> normalizing flows	. 59		
		4.1.1	Augmenting the expressivity of base distributions	. 59		
		4.1.2	Augmenting the expressivity of bases of $L^2(\mu)$	. 60		
	4.2	From	spectral methods to spectral learning	. 66		
		4.2.1	Linear convergence analysis	. 68		
		4.2.2	Faster convergence rates <i>via</i> normalizing flows	. 75		
	4.3 Computing the spectra of polyatomic molecules					
		4.3.1	Convergence of the numerical schemes	. 83		
		4.3.2	Error quantification	. 84		
		4.3.3	Loss function	. 86		
		4.3.4	Accessing of the approximate eigenfunctions	. 86		
	4.4	Sumn	nary and outlook	. 87		
5	Con	clusio	ns and outlook	95		
A	Hilbert spaces and linear operators thereon					
	A.1	Bases	of Hilbert spaces	. 99		
	A.2	Linea	r operators on Hilbert spaces	. 101		
B	Fror	n Carte	esian to internal coordinates of molecules	107		
C	Ran	dom fo	orest regressors	109		
D	Mea	asure th	ieory	113		
	D.1	Push-	forward measures	. 114		
	D.2	The R	adon-Nikodym theorem	. 114		
	D.3	Rador	n measures and a Riesz representation theorem	. 116		
	D.4	Proba	bility theory	. 117		
E	Fun	ctional	analysis	121		
Bi	bliog	raphy		123		
Li	st of ]	Illustra	ations	143		

Acronyms	145
Notation and terminology	147
A statement on data availability and reproducibility	155
Publications and conference contributions	157
Eidesstattliche Versicherung	163

xi

Meinem Vater.

# Chapter 1

# Introduction

Chemical reactions underlie the mechanisms of life. For instance, chemical reactions between macromolecules enable the replication of genetic phenomena, and the interplay of proteins with light define fundamental processes of nature, such as photosynthesis. Observing ultrafast chemical reactions, i. e., molecules "in action", has been a longstanding dream in the molecular sciences [1–3], whereby the observation of the transition state<sup>1</sup> [4], and the recording of nuclear and electronic motions [5] during the breaking or formation of bonds is of particular interest. As in the production of movies from a number of sequential images, observing such reactions consists of taking several snapshots of the positions of the nuclei and the electron densities that constitute molecules. Piecing these snapshots together creates the so-called quantum-molecular movie. While cameras from our everyday life capture light bouncing off objects, they do not work for imaging nature in its smallest scales. Molecules are very tiny compared to everyday objects. Hence, one would need cameras that operate at much shorter wavelengths than the visible light, wavelengths that are comparable to atomic distances. For this purpose one could use, e.g., x-rays or the light corresponding to energetic beams of electrons. The technological developments of pulsed light and electron sources have, indeed, paved way for important steps toward this dream. In particular, the so-called pump-probe observation scheme [6] has emerged as a powerful tool to probe the structure and dynamics of matter at atomic scale. In this scheme, a laser pulse

<sup>&</sup>lt;sup>1</sup>A very short-lived configuration of atoms at a local energy maximum in a reaction coordinate.

(pump) is used to trigger a chemical reaction, i.e., to put molecules "at work", whereby defining the start of the reaction. Then another laser pulse (probe) is used to monitor the progress of the reaction.

Accurate simulations of molecular motions and interactions with laser fields provide crucial information for designing and elucidating ultrafast imaging experiments [7–9] and are, hence, highly-required. Moreover, such simulations are essential for interpreting observations and experiments in a variety of other applications, such as spectroscopy [10, 11], astrophysics [12, 13], and cold chemical reactions [14]. Underlying these simulations is quantum mechanics, a theory that explains nature at small scales. At its core, quantum mechanics models objects as elements of complex Hilbert spaces and postulates that their physical properties are characterized by unbounded linear operators on these spaces. The possible outcomes of a certain measurement are mathematically represented by the spectrum of the corresponding operator. The state of the object as a function of time is governed by the time-dependent Schrödinger equation, a fundamental law for describing non-relativistic particles in physics and chemistry [15]. Throughout the last hundred years, quantum mechanics and generalizations thereof have proved an unprecedented ability to describe nature at small scales ranging from that of quarks, the tiniest known particles, to that of atoms, electrons and relatively large molecular systems such as proteins. However, Dirac's observation<sup>2</sup> on the impracticality of the quantum theory and the difficulty of subjugating it to numerical approximations is as true today as it was in 1930. Nowadays, quantum dynamics of molecules poses a variety of computational challenges that are at the forefront of research efforts in the fields of approximation theory and numerical analysis [15, 16]. These challenges include the approximation of vast number of eigenfunctions of unbounded linear operators and the simulation of ultrafast dynamics for long times. Additionally, the target functions are often highly oscillatory, and lie in high-dimensions [17].

The need to approximate highly-oscillatory functions is prevalent in quantum

<sup>&</sup>lt;sup>2</sup>"For dealing with atoms involving many electrons the accurate quantum theory, involving a solution of the wave equation in many-dimensional space, is far too complicated to be practical. One must therefore resort to approximate methods." P.A.M. Dirac, 1930

physics applications. For example, when imaging chemical reactions, strong laser fields can induce dissociation dynamics in molecules, i. e., the fragmentation of molecules into several parts [18, 19]. An accurate quantum theoretical description of such reactions requires the calculations of many eigenfunctions, i. e., molecular wavefunctions, of a very oscillatory nature. This problem is also encountered in computations of spectra of hot exoplanets, where hundreds and thousands of eigenfunctions corresponding to different molecular motions are required [20]. A reliable numerical approximation of many highly-oscillatory functions is rather prohibitive [21]. Oscillation is, in a sense, an artifact of resolution, i. e., upon zooming enough all functions oscillate mildly [21]. Hence, increasing the resolution of the underlying numerical scheme can, in principle, improve the accuracy of the approximation. However, an increase in the oscillatory behavior necessitates an exponential increase in resolution for a dependable approximation.

Another computational challenge in quantum simulations follows from the intrinsic high dimensionality of quantum systems. For instance, describing static and dynamic properties of two quantum particles, each having three degrees of freedom, involves solving an infinite dimensional eigenvalue problem and an evolution equation in six and seven dimensions, respectively, and the inclusion of only one more particle increases the dimensionality of these problems by three. This is particularly problematic for molecules. For example, benzene is composed of 12 nuclei and 42 electrons, and the corresponding static Schrödinger equation is, hence, 162-dimensional. The prohibitive difficulties in performing simulations of such high-dimensional quantum systems prompted extensive research into deriving reduced models, i. e., effective fundamental laws. For example, noting that the electrons in a molecule are lighter than the nuclei by at least a factor of one thousand gives rise to the famous Born-Oppenheimer approximation [17, 22, 23]. The Born-Oppenheimer approximation breaks the full-dimensional Schrödinger equation of the nuclei into two lower-dimensional equations. The first equation, known as the electronic Schrödinger equation, describes the motion of the electrons in a field of static nuclei. Here, the spatial configuration, i.e., positions of the nuclei act as a parameter to the equation. The second equation describes the motion of the nuclei under an effective potential generated by the electrons. While

separating the scales reduces the complexity of a full numerical treatment, simulating the reduced dynamics remains difficult due to two reasons. First, solving a reduced Schrödinger equation for the nuclei requires the computation of an effective potential, known as the potential energy surface (PES). The PES of a molecule is a function that maps spatial configurations of the nuclei to a certain eigenfunction of the corresponding electronic Schrödinger equations. PESs of molecules can be rough and high-dimensional. Their construction is often cast as a statistical regression problem, a task that is rendered more complex by the high expenses associated with generating the training data. Second, the standard methods for simulating the reduced-order, yet still high-dimensional, dynamics, such as finite volume [24, 25], finite differences [26, 27] or spectral methods [28–30] suffer from the *curse of dimensionality*, i. e., their costs scale exponentially with the dimension of the system.

The aforementioned bottlenecks, i. e., the high-oscillatory nature of solutions and the curse of dimensionality can be formally illustrated by typical upper bounds on the approximation error of classical approximation methods. Denoting by f,  $f^*$  the function to be approximated, and the approximation, respectively. Such bounds look like

$$||f - f^*||_{L^2} \le C(d) ||f||_{H^2},$$

where *C* is often exponentially dependent on the dimension *d* of the problem and the right-hand side converges to zero with an increasing resolution. The Sobolev  $H^2$  norm of *f* has bigger values for highly-oscillatory functions, meaning that approximating such functions suffers from a slower convergence.

The advent of nonlinear approximation methods in general and neural networks in particular has profoundly advanced simulations of quantum molecular dynamics. One of the first applications in quantum molecular physics to witness a revolution driven by the use of neural networks is the construction of PESs [31, 32]. The construction of PESs can be straightforwardly cast as a standard regression problem, and the potential of using neural networks to this end was, hence, recognized as early as 1995 [33]. Indeed, the use of neural networks for building PESs has been recognized as a paradigm shift for constructing PESs, especially in high-dimensions [31]. It led to neural networks that can predict a range of chemical properties for molecules and materials [34], such as force-fields, and polarizability.

Another, less mature and more recent application of neural networks for quantum simulations is approximating solutions of Schrödinger equations. In many applications of practical importance, the curse of dimensionality can be avoided by selecting a nonlinear method of approximation rather than a linear method. Indeed, nonlinear methods, e.g., neural networks, have shown impressive approximation capabilities in high-dimensional modeling of problems ranging from image processing to natural language processing. This has prompted extensive investigations into the applicability of such methods for solving differential equations in general [35–40], and Schrödinger equations in particular [41–49]. In practice, it was shown that such models do, indeed, provide high-accuracy solutions for solving high-dimensional differential equations, while promising smaller scaling with the dimension of the problem than that of linear models. However, this comes at the cost of less efficiency since a straightforward utilization of neural networks to approximate solutions to differential equations is often not possible and extensive architectural engineering is required. Moreover, direct ways of increasing the accuracy of standard neural networks often do not exist. These difficulties render the use of neural networks fragile [50]. Moreover, such nonlinear models do not lend themselves straightforwardly to constructive convergence analysis, and results herein are mainly limited to Barron spaces [50–52], which are specifically tailored to neural networks. Indeed, several results on analyzing these nonlinear models for solving differential equations assume that the solutions and the data of the equation lie in Barron spaces [53–58]. As for approximating highly-oscillatory functions via neural networks, less is known. While some theoretical results show that the approximation error of neural networks decays exponentially in the number of non-zero weights in the network for approximating oscillatory textures [59], the analysis assumes no constraints on the learning algorithm or on the size of the dataset. In fact, it is demonstrated that standard neural networks with gradient descent optimization algorithms tend to fit training data by a low-frequency function [60].

### This thesis

The present work tackles two major computational challenges encountered in quantum simulations of molecules. Specifically, the challenges are related to constructing PESs of molecules and solving static Schrödinger equations for the nuclei. While the two problems being addressed are formally different, they are both associated with the approximation of rough functions in high dimensions.

Chapter 2 provides a formal introduction to quantum molecular dyanmics, the Born-Oppenheimer approximation, and the emergence, therefrom, of the parametric electronic Schrödinger equation and an effective Schrödinger equation for the nuclei. Afterwards, the problem of constructing PESs (Section 2.2) and the Schrödinger equation for the nuclei in a variational formulation (Section 2.3) are introduced.

Chapter 3 describes the first contribution of this work. It starts with an extensive survey of the state-of-the-art approaches for constructing PESs in a supervised learning paradigm and emphasizes limitations thereof. I then highlight the need for constructing PESs in an *active learning* paradigm, where the choice of the training set in the regression task is optimized. I provide a theoretical insight for a good choice of the dataset in terms of an empirical risk minimization principle (Theorem 3.2). An algorithm that follows this insight is proposed (Algorithm 3) and applied in a novel simulation to solve the electronic Schrödinger equation of pyrrole(H<sub>2</sub>O) cluster. The proposed algorithm led to, roughly, two times faster convergence than commonly used learning algorithms for the construction of the PES of pyrrole(H<sub>2</sub>O).

Chapter 4 describes the second contribution of this work. I start the chapter by introducing standard spectral methods for solving differential equations in general and Schrödinger equations in particular. State-of-the-art methods for solving Schrödinger equations and limitations thereof are surveyed. I propose and develop *spectral learning*, a natural nonlinear extension of standard spectral methods that is based on augmenting the expressivity of standard bases of inner-product spaces using machine learning models. This idea mimics the use of normalizing flows to augment the expressivity of base distributions for modeling of probability

distributions. I characterize sufficient conditions on the utilized machine learning models, for the resulting sequences of functions to define bases of the underlying inner-product space (Theorem 4.1). Furthermore, it is shown that normalizing flows, i.e., standard invertible neural networks, satisfy these conditions. As an application, I prove that spectral learning is well-posed, in the sense that convergence guarantees as the size of the truncated augmented basis increases to infinity can be obtained for approximating Schwartz functions (Lemma 4.2) and eigenvalues of Schrödinger operators (Theorem 4.2). I demonstrate that faster convergence rates than linear spectral methods can be achieved *via* spectral learning (Theorem 4.5). This is shown by proving that the total variation of the target function with respect to the push-forward measures induced by the augmenting neural networks admits minima. I performed numerical simulations using the proposed nonlinear framework to compute the spectra corresponding to nuclear motions of three-atomic molecules. Numerical results, reported in Section 4.3, agree with the theoretical observations and show a 2-order of magnitude increase of performance. Results are particularly relevant for approximating eigenfunctions corresponding to large eigenvalues, which are typically highly-oscillatory.

I would like to draw the attention of the reader to the fact that the present work is of an interdisciplinary nature. It was developed in the framework of DASHH, Data science in Hamburg HELMHOLTZ Graduate School for the Structure of Matter, whose aim is to employ formal sciences, i. e., mathematics and computer science, and advances in machine learning to solve problems in the natural sciences. The work uses terminology and results from three different research domains, that of mathematical/numerical analysis, machine learning and quantum molecular physics. However, I made a special effort to make the text accessible to mathematicians, physicists and computer scientists. In particular, I assumed little to no previous expertise of the reader in various domains and extended the main text with extensive appendices to provide the previous knowledge necessary to grasp the contributions of the present work. I admit, though, that the formal discussion of learning algorithms and of spectral learning may remain rather inaccessible to the non-expert formal scientist. To make contributions more accessible, I made efforts to present the simulations of the proposed methods (Section 3.4 and Section 4.3) in a way that does not fully require the understanding of the theoretical underpinnings.

Finally, I call on readers with a formal science background to note that the proposed methods, although developed to solve problems in quantum molecular physics, are applicable to a wide variety of other domains. In particular, the proposed learning Algorithm 3 is applicable to any statistical inference problem and the proposed spectral learning paradigm (Definition 4.3) can be used to solve general differential equations. To facilitate the application of these methods to other problems, the proposed tools were presented in an abstract formulation.

### Chapter 2

# Quantum molecular dynamics

The simulations one requires for interpreting observations and experiments in domains such as imaging ultrafast molecular dynamics, spectroscopy, and astrophysics are, essentially, based on solving the evolution equation of quantum systems

$$i\hbar\frac{\partial\psi}{\partial t} = (H+H_t)\psi \tag{2.1}$$

with an appropriate initial condition. Here,  $i = \sqrt{-1}$  is the imaginary unit, and  $\hbar = h/2\pi$  where *h* is Planck's constant. *H*, the *Hamiltonian*, is a differential operator characterizing the total energy of the molecular system, while  $H_t$  is a time-dependent operator that can model, e. g., external fields. (2.1) is called the *time-dependent Schrödinger equation*. Generally, solving (2.1) starts by describing static behavior of molecular objects, which is governed by the infinite dimensional eigenvalue problem

$$H\psi_k = E_k\psi_k$$
, where  $\int_{\Omega} |\psi_k|^2 d\mu = 1$  for all  $k \in \mathbb{N}_{\geq 0}$ , (2.2)

and  $\Omega \subseteq \mathbb{R}^d$  is open. (2.2) is called the *time-independent Schrödinger equation* (*TISE*). Each eigenpair ( $\psi_k$ ,  $E_k$ ) of (2.2) represents a *quantum state*,  $\psi_k$ , and its corresponding *energy*,  $E_k$ . Solutions of (2.2) are then used as a basis to solve (2.1).

The postulates of quantum mechanics represent quantum states as elements of complex Hilbert spaces. In molecular physics, the Hilbert space is set to be  $L^2$ . Since energies, as all other physically measurable quantities, are real values, these postulates restrict the choice of Hamiltonians to those that have real spectra. The statistical interpretation of quantum mechanics views  $|\psi|^2$  as the probability density for the position of the particle to be located in a certain volume. The restriction of this quantity to integrate to one corresponds to the certainty of finding the particle somewhere in  $\Omega$ . Denote by  $N_l$ ,  $N_n$  the number of electrons and nuclei in the molecule. Since each particle has three degrees of freedom, a quantum description of a molecular system would, then, mean solving (2.2) in  $3(N_n + N_l)$  dimensions, a task that is prohibitive even for small molecules.

The focus in this work is on solving (2.2). I start this chapter by introducing the generic Hamiltonian for molecular systems and the assumptions thereon that comply with the postulates of quantum mechanics. I then discuss the Born-Oppenheimer (BO) approximation that allows one to split (2.2) for molecules into two TISEs. This model reduction theorem gives rise to the first research problem tackled in this work, mainly constructing potential energy surfaces (PESs). This task is often posed as a supervised learning problem. I highlight numerical difficulties that render the construction of PESs prohibitive for bigger molecular systems. Finally, the framework of the second contribution of this work is outlined in Section 2.3, where spectral methods for solving the time-independent Schrödinger equation (TISE) *via* the Rayleigh-Ritz principle are discussed, and limitations thereof are highlighted.

The necessary mathematical definitions and results from the theory of Hilbert spaces are summarized in Appendix A.

### 2.1 Modeling of quantum molecular mechanics

The Hamiltonian of a molecular system is a linear operator that is often composed of two parts [15, 61], a kinetic energy operator

$$T=-\frac{1}{2}\Delta,$$

where  $\Delta$  denotes the Laplacian operator, and a multiplication operator *V* that is called a potential energy term, i. e.,

$$Vf(x) = V(x)f(x), \quad x \in \Omega,$$

which describes static forces in the quantum system. To comply with the statistical interpretation of quantum mechanics, the range of Hamiltonian operators is assumed to be  $L^2$  defined on an open  $\Omega$  in the Euclidean space  $\mathbb{R}^d$  with functions taking values in  $\mathbb{R}$ . In what follows the dependence of functional spaces on the domain is ignored for notational simplicity. Denote by  $\langle ., . \rangle$  the  $L^2$  inner product and set  $\|\cdot\| = \sqrt{\langle ., . \rangle}$ .

Since Hamiltonians correspond to real measurable quantities they should have real spectra. The right mathematical condition to impose this requirement is that of self-adjointness (see Theorem A.2). Denote by D(T) the domain of T. To guarantee that the operator  $T : D(T) \rightarrow L^2$  is self-adjoint set  $D(T) = H^2$ . The right conditions on the potential for H to remain self-adjoint was a subtle problem in the development of quantum mechanics and was satisfactorily addressed by Kato and Rellich [15, 62, 63].

**Theorem 2.1** ([15]). *Let T* be a self-adjoint operator on a Hilbert space, and *V* be a symmetric operator satisfying

$$\|V\psi\| \le a\|\psi\| + b\|T\psi\|$$
 for all  $\psi \in D(T)$ .

Then, H = T + V is self-adjoint with domain D(H) = D(T).

For example, any bounded potential satisfies the above condition and the Hamiltonian is, thus, self-adjoint. To accommodate the Coulomb potential  $V = |x|^{-1}$  one often adopts the following setting.

Assumption 2.1.  $V = V_{\infty} + V_2$  with  $V_2 \in L^2$ ,  $V_{\infty} \in L^{\infty}$ .

This potential satisfies the above condition and H = T + V is a self-adjoint operator with  $D(H) = H^2$ . Note, however, that the Hamiltonian in this setting

is an unbounded linear operator due to the  $V_{\infty}$  term. This makes the analysis of quantum systems nontrivial.

### 2.1.1 Molecular Hamiltonians

Consider molecules composed of  $N_n$  nuclei of masses  $M_j$  and electric charges  $Z_j e$ , with  $Z_j$  denoting the atomic number of the *j*th nuclei, and  $N_l$  electrons of masses *m* and charges *e*.

Let  $x_j, y_l$  denote the spatial coordinates of the *j*th nucleus and *l*th electron, respectively, and let  $(x, y) \in \Omega_{nuc} \times \Omega_{el} \subseteq \Omega \subseteq \mathbb{R}^{3(N_n+N_l)}$  denote the electronic and nuclear coordinates, respectively. I also call a certain *x* a *nuclear/molecular geometry*. The molecular Hamiltonian is the sum of a kinetic and potential energy parts [15, 61]

$$H_{\rm mol} = \underbrace{T_{\rm nuc} + T_{\rm el}}_{T} + \underbrace{V_{\rm nn} + V_{\rm ne} + V_{\rm ee}}_{V},\tag{2.3}$$

where  $T_{nuc}$  and  $T_{el}$  are the kinetic energy operators of the nuclei and the electrons, respectively

$$T_{
m nuc} = -\sum_{j=1}^{N_n} rac{\hbar^2}{2M_j} \Delta_{x_j},$$
  
 $T_{
m el} = -\sum_{l=1}^{N_l} rac{\hbar^2}{2m} \Delta_{y_l}.$ 

The potential is the sum of the nucleus-nucleus

$$V_{nn}(x) = \sum_{1 \le k < j \le N_n} \frac{Z_k Z_j e^2}{|x_k - x_j|},$$

nucleus-electron

$$V_{\rm ne}(x,y) = -\sum_{l=1}^{N_l} \sum_{j=1}^{N_n} \frac{Z_j e^2}{|y_l - x_j|},$$

and electron-electron interactions

$$V_{\rm ee}(y) = \sum_{1 \le j < l \le L} \frac{e^2}{|y_j - y_l|}.$$

The TISE for molecules is extremely high-dimensional even for small molecules. Consequently, molecular quantum mechanics makes extensive use of reducedorder modeling and approximation theorems. I present the Born-Oppenheimer (BO) approximation, a model-reduction theory that allows for reducing the difficulty of solving the TISE for molecules.

### 2.1.2 Born-Oppenheimer approximation

The BO approximation [15, 17, 22, 23, 61] rests on the observation that the mass of the electrons is negligible in comparison to that of the nuclei. Given the same amount of momentum, electrons would then move on a much faster timescale than that of the nuclei. One can, hence, split the molecular TISE into two lowerdimensional TISEs. The first, known as the *electronic Schrödinger equation*, is a parametric Schrödinger equations that describes the motion of the electrons in a field of static nuclei. The second is a TISE that describes the motion of the nuclei in an effective field generated by the electrons.

Let  $H_{el} = T_{el} + V_{ne} + V_{ee} + V_{nn}$  denote the *electronic Hamiltonian* and note that it parametrically depends on the positions of the nuclei through the potentials  $V_{ne}$ ,  $V_{nn}$ . The electronic TISE reads

$$H_{\mathrm{el}}\phi_{\mathrm{el},k}(y;x) = E_{\mathrm{el},k}(x)\phi_{\mathrm{el},k}(y;x) \quad \text{for all } k \in \mathbb{N}_{>0}.$$

$$(2.4)$$

Note that (2.4) is a parametric differential equation, where the parameter x take values in the uncountably infinite set  $\Omega_{nuc}$ , corresponding to different nuclear (molecular) geometries.

Methods to compute the solutions of (2.4) exist [23] and are covered in the wide field of quantum chemistry. I now discuss how to use these solutions to solve the remainder of the molecular Schrödinger equation. Assume that the *j*th eigenfunction  $\psi_j$  of (2.2) with the molecular Hamiltonian (2.3) can be written as

 $\psi_j(x, y) = \sum_k \phi_{\text{el},k}(y; x) \phi_{\text{nuc},j}(x)$ . Substituting that into the molecular Schrödinger equation one has

$$\sum_{k} (T_{\text{nuc}} + H_{\text{el}})(\phi_{\text{el},k}(y;x)\phi_{\text{nuc},j}(x)) = \sum_{k} E_{j}\phi_{\text{el},k}(y;x)\phi_{\text{nuc},j}(x).$$

Thus, using chain rule one has

$$\sum_{k} (T_{\text{nuc}}\phi_{\text{el},k}(y;x)\phi_{\text{nuc},j}(x) + \sum_{k} \phi_{\text{el},k}(y;x)(T_{\text{nuc}}\phi_{\text{nuc},j}(x)) - 2\sum_{k} \sum_{j=1}^{N_{n}} \frac{\hbar^{2}}{2M_{j}} (\nabla_{x_{j}}[\phi_{\text{el},k}(y;x)).(\nabla_{x_{j}}\phi_{\text{nuc},j}(x)) + \sum_{k} E_{\text{el},k}(x)\phi_{\text{el},k}(y;x)\phi_{\text{nuc},j}(x) = \sum_{k} E_{j}\phi_{\text{el},k}(y;x)\phi_{\text{nuc},j}(x).$$

Multiplying by  $\phi_{el,p}^*$ , the adjoint of  $\phi_{el,p}$ , and integrating over  $\Omega_{el}$  one has

$$\begin{split} &\sum_{k} \langle \phi_{\mathrm{el},p}, T_{\mathrm{nuc}} \phi_{\mathrm{el},k} \rangle \phi_{\mathrm{nuc},j}(x) + T_{\mathrm{nuc}} \phi_{\mathrm{nuc},j}(x) - \\ &- 2 \sum_{k} \sum_{j=1}^{N_n} \frac{\hbar^2}{2M_j} \langle \phi_{\mathrm{el},p}, \nabla_{x_j} \phi_{\mathrm{el},k} \rangle . \nabla_{x_j} [\phi_{\mathrm{nuc},j}(x)] + \\ &+ E_{\mathrm{el},p}(x) \phi_{\mathrm{nuc},j}(x) = E_j \phi_{\mathrm{nuc},j}(x), \end{split}$$

where I used the fact that eigenfunctions of (2.4) are orthonormal. Note that

$$\langle \phi_{\mathrm{el},p}, T_{\mathrm{nuc}}H_{\mathrm{el}}\phi_{\mathrm{el},k} \rangle = E_{\mathrm{el},k} \langle \phi_{\mathrm{el},p}, T_{\mathrm{nuc}}\phi_{\mathrm{el},k} \rangle$$
$$\langle H_{\mathrm{el}}\phi_{\mathrm{el},p}, T_{\mathrm{nuc}}\phi_{\mathrm{el},k} \rangle = E_{\mathrm{el},p} \langle \phi_{\mathrm{el},p}, T_{\mathrm{nuc}}\phi_{\mathrm{el},k} \rangle.$$

Thus,

$$\langle \phi_{\mathrm{el},p}, T_{\mathrm{nuc}}\phi_{\mathrm{el},k} \rangle = rac{\langle \phi_{\mathrm{el},p}, [H_{\mathrm{el}}, T_{\mathrm{nuc}}]\phi_{\mathrm{el},k} \rangle}{E_{\mathrm{el},p} - E_{\mathrm{el},k}},$$

where [., .] denotes the commutator between two operators. Assuming that the eigenvalues of (2.4) are well separated,  $\langle \phi_{\text{el},p}, T_{\text{nuc}}\phi_{\text{el},k} \rangle$  is very small and can be neglected. Same argumentation can be developed for  $\langle \phi_{\text{el},p}, \nabla_{x_i}\phi_{\text{el},k} \rangle$ . Thus, one

ends up with a TISE for the nuclei

$$\underbrace{(T_{\text{nuc}} + E_{\text{el},p})}_{H_{\text{nuc}}} \phi_{\text{nuc},j} = E_j \phi_{\text{nuc},j} \quad \text{for all } j \in \mathbb{N}_{\ge 0},$$
(2.5)

where  $E_{\text{nuc},p} : \Omega_{\text{nuc}} \to \mathbb{R}$  is a function, called the *potential energy surface* that, given any arrangement *x* of the nuclei returns the corresponding *p*th eigenvalue of the electronic Schrödinger equation (2.4). Fix p = 0, i. e., consider only the ground electronic state, and drop the notational dependence of the potential energy surface on it for simplicity.

# 2.2 Electronic Schrödinger equations and potential energy surfaces

The solution of the time-independent Schrödinger equation for the nuclei (2.5) requires first the construction of the linear operator  $H_{\text{nuc}}$ . This, in turn, requires finding the lowest eigenvalue of the parametric Schrödinger equation (2.4). Since the nuclear geometries x take values in an uncountably infinite set, and an analytic solution is often not achievable, this task is prohibitive. It is, therefore, replaced by computing the smallest eigenvalue of (2.4) for some selected finite set of the nuclear coordinates, and inferring the ground-state energy for other nuclear coordinates. In particular, given a set of nuclear geometries  $\hat{x} = \{x^k\}_k$ , and the corresponding electronic energies  $\hat{E} = \{E_{el}^k\}_k^1$  the goal is to be able to infer the electronic energy for a nuclear geometry  $x \notin \hat{x}$ . While this task can be performed *via* an interpolation procedure of the dataset  $\hat{z} = \{(x^k, E_{el}^k)\}_k$ , this is not recommended in practice, since solving (2.4) for the set  $\hat{x}$  is not an exact process and often contains errors. Instead, this task is often performed *via* a linear or nonlinear regression where both the input x and output  $E_{el}$  are interpreted as random variables and their causal relation is assumed to be governed by a probability distribution  $\mathcal{P}(x, E_{el})$ . Given a class of hypothesis functions  $\mathfrak{H}$  one would then choose a function  $h^*$  that

<sup>&</sup>lt;sup>1</sup>Superscript refers to the observation, i. e.,  $E^k$  refers to the *k*th empirical observation of the random variable *E*.

reproduces the dataset and generalizes well beyond it. In practice, the following problem is solved

$$\frac{1}{|D|} \sum_{(x^k, E^k_{\text{el}}) \in \hat{z}} l(E^k_{\text{el}}, h(x^k)) + \gamma \beta(h(x^k)) \longrightarrow \min_{h \in \mathfrak{H}},$$
(2.6)

where *l* is the loss function quantifying the discrepancy between the true energies and the predicted ones, and  $\beta(h)$  is a regularization term that constraints the complexity of the hypothesis *h* with  $\lambda \in \mathbb{R}^+$ . While this task is clearly easier than solving (2.4) for any value of *x*, it is still complex since generating the dataset is a dimension-dependent problem. In particular, one needs exponentially more computational resources to generate the target values  $\mathcal{E}$  for an increasing dimensionality of the molecular system. Moreover, the quality of the minimizer *h*<sup>\*</sup> of (2.6) highly depends on the choice of *D*. In Chapter 3 I discuss an active learning framework to extend the regression problem defined in (2.6) into an optimization over the choice of the training dataset *D* as well.

### 2.3 Schrödinger equations for the nuclei

The nuclear Schrödinger equation (2.5) is the second step to characterize quantum properties of molecules in the Born-Oppenheimer picture. Similar to the electronic Schrödinger equation (2.4), the numerical treatment of this equation for polyatomic molecules is not straightforward due to its high dimensionality, since the nuclear geometries lie in  $\mathbb{R}^{3N_n}$ . Moreover, several applications require the computations of many excited states that are often highly-oscillatory. These reasons render numerical simulations challenging.

A first step of common methods to solve (2.5) is to write its weak formulation. Multiplying by a test function v with enough regularity and taking the inner product in  $L^2$ , one has

$$\langle v, (T_{\text{nuc}} + E_{\text{el},p})\phi_{\text{nuc},j} \rangle = E_{\text{nuc},j} \langle v, \phi_{\text{nuc},j} \rangle.$$

Choose a trial solution  $\gamma_j \in H^2$  and drop the dependence on *j* for simplicity of notation. In the special case  $v = \gamma$  where  $\|\gamma\| = 1$ , the weak formulation of (2.5) reads

$$\langle G\gamma, G\gamma \rangle + \langle \gamma, E_{\rm el}\gamma \rangle = E_{\rm nuc},$$
 (2.7)

where  $G = \sum_{j=1}^{N_n} \frac{\hbar}{\sqrt{2M_j}} \nabla$ . Define

$$\epsilon(\gamma) := \langle G\gamma, G\gamma \rangle + \langle \gamma, E_{\rm el}\gamma \rangle \tag{2.8}$$

to be the linear form corresponding to the weak formulation (2.7). It turns out that, under proper assumptions on the potential  $E_{el}$ , this linear form has a minimum and this minimum corresponds to the smallest eigenvalue of (2.5).

#### 2.3.1 A variational formulation

Several weak formulations of differential equations correspond to a variational principle, i.e., the weak solution can be obtained by minimizing some energy functional, such as, e.g., the Poisson equation with homogeneous Dirichlet boundary conditions. One can, indeed, construct similar results for TISEs. The correct assumptions on the potential function to obtain these results read as follows [63].

Assumption 2.2 (Conditions on the potential function).

$$\begin{cases} \text{for all } d \ge 3: & E_{el} \in L^{\infty}(\mathbb{R}^d) + L^{d/2}(\mathbb{R}^d) \\ d = 2, \text{ for all } \varepsilon > 0: & E_{el} \in L^{\infty}(\mathbb{R}^d) + L^{1+\varepsilon}(\mathbb{R}^d) \\ d = 1: & E_{el} \in L^{\infty}(\mathbb{R}^d) + L^{1}(\mathbb{R}^d) \end{cases}$$

The variational result is based on the following two lemmas. The first lemma tells us that  $\epsilon(\gamma)$  is bounded from below.

**Lemma 2.1** ([63]). Under Assumption 2.2 there exist C, D > 0 with

$$\epsilon(\gamma) \geq C \int_{\Omega_{nuc}} |\nabla \gamma|^2 \, d\mu - D \|\gamma\|^2.$$

In particular,

$$E_0 := \inf\{\epsilon(\gamma) \mid \|\gamma\| = 1\} > -\infty.$$
(2.9)

The following result establishes a correct setting that guarantees the existence of a minimum of  $\epsilon(\gamma)$ .

Lemma 2.2 ([63]). Under Assumption 2.2 the potential energy

$$P(\gamma) = \int_{\Omega_{nuc}} |\gamma|^2 E_{el} \, d\mu$$

is weakly continuous in  $H^1$ . In other words, if  $\gamma_j \to \gamma$  weakly in  $H^1$ , then  $P(\gamma_j) \to P(\gamma)$ as  $j \to \infty$ .

The stage is now ready to state the main theorem.

**Theorem 2.2** (Existence of a minimum of the energy functional [63]). *Under Assumption 2.2 and assuming that* 

$$E_0 = \inf\{\epsilon(\gamma) \mid \gamma \in H^1, \|\gamma\| = 1\} < 0$$

there exists  $\gamma_0 \in H^1$ , with  $\|\gamma_0\| = 1$  and  $\epsilon(\gamma(0)) = E_0$ . Moreover,  $(\gamma_0, E_0)$  solves the Schrödinger equation (2.5) in the weak sense.

*Proof.* A proof is provided in the supplementary material section at the end of this chapter.  $\Box$ 

**Remark 2.1.** It follows that  $E_0 = \inf_{\psi \in D(H), \|\psi\|=1} \epsilon(\gamma)$  since D(H) is dense in  $H^1(D(H)$  is dense in  $H^2$  which is dense in  $H^1$ ).

Theorem 2.2 guarantees that the smallest eigenvalue of the Schrödinger equation can be computed using a variational method. The same can be done for higher eigenvalues. However, here the minimization is performed over spaces that are orthogonal to the eigenfunctions corresponding for smaller eigenvalues. Define

$$E_1 = \inf \{ \epsilon(\gamma) \mid \gamma \in H^1, \|\gamma\| = 1, \langle \gamma, \gamma_0 \rangle = 0 \}.$$

Continuing recursively, define

$$E_k = \inf\{\epsilon(\gamma) \mid \gamma \in H^1, \|\gamma\|_2 = 1, \langle \gamma, \gamma_k \rangle = 0 \text{ for all } k = 0, \dots, k-1\}.$$
(2.10)

The following theorem provide an equivalent of Theorem 2.2 for larger eigenvalues.

**Theorem 2.3** ([63]). Under Assumption 2.2 and assuming  $E_k < 0$ , then, the infimum in (2.10) is attained and the minimizer  $\gamma_k$  is such that  $H\gamma_k = E_k\gamma_k$ .

**Remark 2.2.** While the analysis in this section was developed for the specific nuclear Schrödinger equation, it is, in general, valid for any time-independent Schrödinger equation where the potential satisfies Assumption 2.2 [63].

### 2.3.2 Spectral discretization and the curse of dimensionality

The established variational formulation of the TISE allows for practical algorithms for computing its eigenpairs. In particular, given a certain approximation space, one would then choose the function that minimizes the established energy functional. This is often referred to as the Rayleigh-Ritz method [30]. Discretizing the problem and constructing approximation spaces can be done *via* standard methods, such as finite volume or finite differences methods [24–27]. However, due to the often oscillatory nature of the solutions of TISEs and to the need to model a vast number of eigenfunctions, a more common approach to discretize the equations is based on spectral methods [28, 29] where one approximates the eigenfunctions in a finite linear span of globally supported sequence of functions  $(\varphi_n)_n$ 

$$\phi_{\operatorname{nuc},j}(x) \approx \sum_{n=1}^{N} c_{n,j} \varphi_n(x) \quad \text{for all } j = 0, 1, \dots$$
(2.11)

To have convergence guarantees as  $N \rightarrow \infty$  and quantify the convergence order one would have to choose sequences with some density properties, i. e., the linear span of the sequence should be dense in some target functional spaces. In

TABLE 2.1: The size of a truncated Hermite basis N that is required to compute the ground-state energy of a perturbed Harmonic oscillator problem in 1,2 and 3 dimensions to a relative absolute error  $< 10^{-1}$ .

d	1	2	3
Ν	3	45	286

molecular physics, the functional space is often  $L^2$ . Indeed, with such density properties, the method is reliable and well-posed for solving a variety of differential equations [28], such as, e. g., determining the spectra of unbounded operators.

In spite of these positive features, the method suffers from the *curse of dimensionality*, as the size of the basis N needed to converge a certain amount of eigenpairs scales exponentially with the number of nuclei  $N_n$  in the system. To exemplify this problem I solved the Schrödinger equation for a perturbed quantum harmonic oscillator problem, i. e., the TISE with a Hamiltonian H = T + V where

$$V = \frac{1}{4}|x|^4 + \frac{1}{2}|x|^2.$$

Table 2.1 shows the size of truncated Hermite basis that is needed in order to converge the ground-state energy in several dimensions. Clearly, N increases rapidly as a function of the dimension of the problem.

Another problem is the scaling with respect to the number of eigenpairs that one would want to approximate. In practice, eigenfunctions to (2.5) corresponding to high eigenvalues are highly oscillatory and, hence, the accuracy of spectral methods decrease nonlinearly with the number of required states. Table 2.2 shows the size of the Hermite basis required in order to converge the energies of the first three states for the perturbed Harmonic oscillator problem in 2 dimensions. Clearly, increasingly more linear terms are needed to converge higher states<sup>2</sup>.

One way to mitigate the curse of dimensionality with respect to the linear parameter *N* is to allow for the use of an adaptive sequence of functions, i. e., a

<sup>&</sup>lt;sup>2</sup>Technical details to reproduce this example are provided in the supplementary materials at the end of the chapter.

TABLE 2.2: The size of truncated Hermite basis N that is required to compute the energies of the first three states of a perturbed Harmonic oscillator problem in 3 dimensions to a relative absolute error  $< 10^{-1}$ .

State	1	2	3
Ν	45	300	555

sequence  $(\varphi_n^{\theta})_n$  that depends on free parameters  $\theta$ 

$$\phi_{\operatorname{nuc},j}(x) \approx \sum_{n=1}^{N} c_{n,j} \varphi_n^{\theta}(x) \quad \text{for all } j = 0, 1, \dots$$
(2.12)

One would then minimize the established upper bound with respect to both the linear parameters,  $c_{n,j}$  and the nonlinear ones,  $\theta$ . Allowing for a set of functions that is more suitable for the problem lessens the need for an exponentially big set of fixed functions. However, the introduction of arbitrarily adaptive functions without proven density properties make it difficult, if not impossible, to obtain convergence guarantees. Neural networks have been lately under intensive investigations to this end <sup>3</sup>. However, their use is not straightforward and often requires a lot of engineering efforts. In Chapter 4 I propose and investigate another alternative that is based on carefully deforming fixed bases into adaptive ones.

**Remark 2.3.** Solving (2.5) can be simplified by differentiating between three kinds of motions in polyatomic molecules: a translational motion where the Cartesian coordinates of all atoms are shifted by the same quantity and in the same direction, a rotational motion of all the atoms, and a vibrational motion where distances between atoms change. Under the potential energy term consider in (2.5), the eigenvalues remain invariant with respect to a translational motion of the nuclei. Thus, one can consider only the rotational and vibrational (rovibrational) motions considered when solving (2.5). The

<sup>&</sup>lt;sup>3</sup>Generally one uses one neural network, i. e., N = 1.

equation describing the rovibrational motion can also be solved in three steps, where, first, the (2.5) is solved for the vibrational degrees of freedom, then for the rotational degrees, and eigenfunctions of both motions are then used to solve the overall rovibrational equation [64]. In this work I only consider the vibrational Schrödinger equation and hence reduce the space of the coordinates to a set  $\Omega_{nuc} \subset \mathbb{R}^{3N_n-6}$ . The set  $\Omega_{nuc}$  contains often two kinds of coordinates, radial coordinates describing the distance between pairs of nuclei having values in  $\mathbb{R}_{>0}$  and angular coordinates having values in  $(0, \pi)$  describing angles between pairs of radial vectors. This set of coordinate is often called an internal set of coordinates. Moving from a Cartesian coordinate systems to the internal one is described in Appendix B.

## Supplementary material

The simulation ran to produce Table 2.1 and Table 2.2 was performed as follows. A Bubnov-Galerkin numerical scheme was used to discretize the Schrödinger equation, where both the test and trial functions were modeled by Hermite functions. For d > 1, the basis was generated from Hermite functions using the truncated direct product (4.16) with  $w_i = 1$ , for all i = 1, ..., d and  $n_i$  were varied from 1 to 30 for all i = 1, ..., d. N = 1, ..., 50 was used for d = 1 and d = 2 while I used N = 1, ..., 22 for d = 3. Hermite quadrature points were used to compute the integrals with 50 quadrature points per dimension. Quadrature points corresponding to weights  $< 10^{-34}$  were removed. I considered for true solutions those converged with the largest basis.

*Proof of Theorem* 2.2. Set  $\Omega = \Omega_{\text{nuc}}$  for simplicity of notation. Let  $\gamma_j$  be a sequence in  $H^1$  with  $\|\gamma_j\| = 1$  and  $\epsilon(\gamma_j) \to E_0$  as  $j \to \infty$ . Lemma 2.1 implies that

$$\epsilon(\gamma_j) \geq \frac{1}{2} \|\nabla \gamma_j\|^2 - C,$$
which implies the boundedness of  $\|\nabla \gamma_j\|_{L^2}^2$ , which, in turns, implies the boundedness of  $\gamma_j$  in  $H^1$  for all j. By the Banach-Alaoglu theorem there exists a subsequence  $\gamma_{n_j}$  and  $\gamma_0 \in H^1$  such that  $\gamma_{n_j} \to \gamma_0$  weakly in  $H^1$ . The weak limit in the norm can only get smaller, thus

$$\|\gamma_0\| \leq 1, \|\nabla\gamma_0\| \leq \liminf_{j \to \infty} \|\nabla\gamma_{n_j}\|.$$

Since  $\|\gamma_0\| \leq 1$  and  $E_0 = \epsilon(\gamma_0)$  it holds

$$E_{0} \|\gamma_{0}\|^{2} \leq \epsilon(\gamma_{0})$$

$$= \|\nabla\gamma_{0}\| + P(\gamma_{0})$$

$$= \|\nabla\gamma_{0}\| + \lim_{j \to \infty} P(\gamma_{n_{j}})$$

$$\leq \liminf_{j \to \infty} (\|\nabla\gamma_{n_{j}}\| + P(\gamma_{n_{j}}))$$

$$\leq \liminf_{j \to \infty} (\|\nabla\gamma_{n_{j}}\| + P(\gamma_{n_{j}}))$$

$$= \leq \liminf_{j \to \infty} \epsilon(\gamma_{n_{j}})$$

$$= E_{0}$$

where in the second equality I used Equation 2.9. Since  $E_0 \le 0$  by assumption, one deduces that  $\|\gamma_0\| \ge 1$ . But  $\|\gamma_0\| \le 1$ . Hence,  $\|\gamma_0\| = 1$ .

To show that  $\gamma_0$  solves the Schrödinger equation take a perturbation  $\gamma_{\delta} = \gamma_0 + \delta f$  of  $\gamma_0$  with  $\delta \in \mathbb{R}$ . Define  $R(\delta) = \frac{\epsilon(\gamma_{\delta})}{\|\gamma_{\delta}\|^2}$ .

Since  $E_0$  is a minimizer of  $\epsilon$  it holds that  $R(\delta)$  attains a minimum at  $\delta = 0$ . Thus,

$$0 = \frac{dR(\delta)}{\delta}|_{0}$$
  
=  $\frac{1}{\|\gamma_{\delta}\|^{4}} \left(\frac{d\epsilon(\gamma_{\delta})}{d\delta} - \frac{d\|\gamma_{\delta}\|^{2}_{L^{2}(\mathbb{R}^{d})}}{d\delta} \frac{\epsilon(\gamma_{\delta})}{\|\gamma_{\delta}\|^{2}_{L^{2}(\mathbb{R}^{d})}}\right)|_{\delta=0}$ 

Therefore,

$$0 = \left(\frac{d\epsilon(\gamma_{\delta})}{d\delta} - \frac{d\|\gamma_{\delta}\|^2}{d\delta}E_0\right)|_{\delta=0},$$

where

$$\begin{aligned} \frac{d\epsilon(\gamma_{\delta})}{d\delta} &= \frac{d}{d\delta} (\int_{\Omega} |\nabla\gamma_0|^2 \, d\mu + \delta^2 \int_{\Omega} |\nabla\gamma_0|^2 \, d\mu + 2\delta \int_{\Omega} \nabla f . \nabla\gamma_0 \, d\mu + \\ &+ \int_{\Omega} (\gamma_0 + \delta f)^2 V \, d\mu), \end{aligned}$$

hence,

$$\frac{d\epsilon(\gamma_{\delta})}{d\delta}|_{\delta=0} = 2\int_{\Omega} \nabla f \cdot \nabla \gamma_0 \, d\mu + 2\int_{\Omega} f \, \gamma_0 \, V \, d\mu.$$

Similarly,

$$E_0 \frac{d \|\gamma_{\delta}\|^2}{d\delta}|_{\delta=0} = 2E_0 \int_{\Omega} \gamma_0 f \, d\mu.$$

Whence,

$$0 = \int_{\Omega} (-\Delta + V - E_0) \gamma_0 f \, d\mu \quad \text{for all } f \in C_c^{\infty},$$

i. e.,  $(E_0, \gamma_0)$  solves the Schrödinger equation in the weak sense.

## Chapter 3

# Active learning for constructing potential energy surfaces<sup>1</sup>

I outlined in Section 2.2 the problem of constructing potential energy surfaces, i.e., of inferring solutions of the parametric Schrödinger equation (2.4) in a statistical manner. Here, the relationship between the nuclear geometries x and the electronic energies *E* is assumed to be governed by a probability distribution  $\mathcal{P}_{z_{\ell}}$  called the generating distribution, with z := (x, E). Given access to an initial set of independent observations  $\hat{z} = \{z^k := (x^k, E^k)\}_k, \hat{z} \sim \mathcal{P}_z$  of nuclear geometries  $\hat{x} = \{x^k\}_k$  and their corresponding electronic energies  $\hat{E} = \{E^k\}_k^2$  one aims at approximating  $\mathcal{P}_z$ by solving the regression problem defined in (2.6).

In recent years, many machine learning (ML) models [65–69] have been used to model the hypothesis class  $\mathfrak{H}$ . The most extensively used models include permutationally invariant polynomials [70–74], neural network (NN)s [32, 75–80], Gaussian process (GP)s [81–87], and other kernel methods [16, 88–90].

The quality of the optimizer  $h^*$  that solves (2.6) highly depends on the quality and size of the dataset  $\hat{z}$ . The more data there is to learn from, the more accurate the model is. However, computing the labels  $\hat{E}$  of the input nuclear geometries  $\hat{x}$  is computationally expensive. It requires solving the electronic Schrödinger

<sup>&</sup>lt;sup>1</sup>This chapter is, in parts, based on this publication: Y. Saleh, V. Sanjay, A. Iske, A.Yachmenev, J. Küpper, J. Chem. Phys. 155, 144109 (2021). My contribution to this publication was the development and implementation of several active learning algorithms and writing the manuscript. <sup>2</sup>In what follows denote the target in  $\hat{z}$  by  $E^k$  instead of  $E^k_{el}$  for notational simplicity.

equation (2.4) for a system of  $N_l$  electrons. High-accuracy solvers scale as  $O(n^7)$  with the size of the system n [23]. Hence, one would want to minimize the size of the dataset  $\hat{z}$  that is needed to obtain an optimized model with a certain accuracy. In other words, one would want to solve

$$\frac{1}{|\hat{z}|} \sum_{(x^k, E^k) \in \hat{z}} l(E^k, h(x^k)) + \gamma \beta(h(x^k)) \longrightarrow \min_{\hat{z}, h \in \mathfrak{H}}.$$
(3.1)

One learning paradigm that tackles such problems is active learning (AL) [91], in contrast to the common passive learning (PL) paradigm, where the training dataset is given *a priori*. Trying to minimize the size of the training dataset for constructing potential energy surfaces (PESs), and reducing the amount of humans' intervention in an AL paradigm became increasingly popular during the last few years [92–103].

Formally, in AL, the training dataset  $\hat{y} := \{y^i := (q^i, E^i)\}_i$  is sampled from a probability distribution  $\mathcal{P}_y$  that is not necessarily equal to  $\mathcal{P}_z$ . However, both distributions are assumed to follow the same conditional probability. Denoting by  $p_z$ ,  $p_y$  the probability densities of  $\mathcal{P}_z$ ,  $\mathcal{P}_y$ , respectively, one has  $p_z = p_x p_{E|x}$ ,  $p_y =$  $p_q p_{E|x}$  [104, 105], where  $p_x$ ,  $p_q$  denote probability densities over the nuclear geometries x and  $p_{E|x}$  denotes the conditional probability, i. e., the probability of energies given a certain nuclear geometry x. The aim of AL is to learn densities  $p_q$ as to solve (3.1). Algorithms that aim at constructing such  $p_q$  are called *query/policy strategies/algorithms*. A deep theoretical understanding of optimal query algorithms is generally absent. Such algorithms are often constructed based on heuristic arguments such as uncertainty sampling [91], where datapoints corresponding to high uncertainties in their predictions are added to the training dataset.

A guiding concept in defining statistical learning algorithms is empirical risk minimization principles [106]. Such principles are based on the observation that the objective function in standard learning algorithms, such as (2.6), is an upper bound of the true risk. In this chapter, a similar theoretical insight into constructing query algorithms is provided. It is shown that upper bounds on the generalization error in PL can be extended naturally to AL (Theorem 3.2). This allows one to propose the following empirical risk minimization principle to define

query algorithms. The probability distribution  $\mathcal{P}_q$  defining the query algorithm should maintain a small integral probability metric from the generating probability distribution while still assigning a significant measure to scarce regions of the generating distribution. This result can be seen as a general formulation of [105] and [104], where specific integral probability metrics were used. Directly optimizing the derived upper bound is rather impractical, whence I review some more practical algorithms, referred to as pool-based algorithms [91]. I survey the state-of-the-art AL algorithms employed for construction of PESs. I propose a novel AL learning algorithm (Algorithm 3) that complies with the derived upper bound. It is a regression version of random query by forest [107].

I validate the proposed algorithm for modeling the PESs of weakly-bound molecules<sup>3</sup>. Such a task is complex [108–111], since higher levels of theory need to be employed to produce correct asymptotic behavior of the training dataset [112]. Furthermore, the landscape of these PESs is complex because of the loosely bound character of intermolecular interactions. Thus, a larger number of grid points is generally required to sample the complete configuration space. Moreover, due to the importance of dynamical electron correlation (dispersion) and its slow basis-set convergence [113], calculations for the noncovalent long-range parts of the PES are generally more costly than the ones at short-range. In particular, I model the PES of pyrrole(H<sub>2</sub>O). I show that the proposed algorithm reduces the computational costs of constructing the PES of pyrrole(H<sub>2</sub>O). It leads to a roughly two times faster convergence with respect to the size of training dataset than other commonly used AL algorithms.

#### 3.1 Formal setting and notation

Throughout this chapter consider a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . For an open  $\mathbb{X} \subseteq \mathbb{R}^m$  (given a molecular system with  $N_n$  nuclei, one often has  $m = 3N_n$ ) I model the nuclear geometries x as a random vector  $x : \Omega \to \mathbb{X}$ . Endow  $\mathbb{X}$  with the Borel  $\sigma$ -algebra and the Lebesgue measure  $\mu^m$ . Set  $\mathcal{P}_x = x_{\#}\mu^m$ . Further,

<sup>&</sup>lt;sup>3</sup>Weakly-bound molecules are complexes characterized by relatively low interaction energies between the constituent molecules

assume that  $\mathcal{P}_x$  is absolutely continuous with respect to  $\mu^m$  and denote by  $p_x$ the corresponding density. The nuclear geometries chosen by a query strategy are modeled by a random vector  $q : \Omega \to \mathbb{X}$  with a distribution  $\mathcal{P}_q = q_{\#}\mu^m$ and a density  $p_q$ . Similarly, for an open  $\mathbb{E} \subseteq \mathbb{R}$  consider the electronic energies to be a random variable  $E : \Omega \to \mathbb{E}$ . Endowing  $\mathbb{E}$  with the Borel  $\sigma$ -algebra and the Lebesgue measure  $\mu$  set  $\mathcal{P}_E = E_{\#}\mu$  and denote by  $p_E$  the corresponding probability density. Given the random vector x and the random variable E define  $z := (x, E), z : \Omega \to \mathbb{X} \times \mathbb{E} =: \mathbb{Z}$  to be the joint random vector. Similarly, set,  $y := (q, E), y : \Omega \to \mathbb{Z}$ . Endowing  $\mathbb{Z}$  with the product  $\sigma$ -algebra and the Lebesgue-measure  $\mu^{m+1}$  set  $\mathcal{P}_z = z_{\#}\mu^{m+1}, \mathcal{P}_y = y_{\#}\mu^{m+1}$  and denote by  $p_z, p_y$  the corresponding probability densities.

# 3.2 Empirical risk minimization principles and generalization errors

The goal in a supervised learning task is to infer the probability distribution  $P_z$  governing the relationship between two random variables *x*, *E*. Given a hypothesis class  $\mathfrak{H}$ , this translates into solving

$$R_{\mathcal{P}_{z}}(h) = \int_{\mathbb{Z}} l \, d\mathcal{P}_{z} \longrightarrow \min_{h \in \mathfrak{H}}, \tag{3.2}$$

where  $R_{\mathcal{P}_z}(h)$  is called the *true risk* or *generalization error* of a hypothesis *h* and  $l : \mathbb{E} \times \mathbb{X} \to \mathbb{R}_{>0}$  is a *loss function* quantifying the discrepancy between true values *E* and predicted ones h(x). In other words, one tries to find the function *h* that would minimize the discrepancy between the predicted values and the true values along the joint probability of the problem [68, 106].

However, in practice, one has access only to a dataset  $\hat{z}$  of a finite size. It is reasonable to try to solve (2.6) where

$$\hat{R}_{\hat{z}\sim\mathcal{P}_{z}}(h) := \frac{1}{|\hat{z}|} \sum_{(x^{k},E^{k})\in\hat{z}} l(E^{k},h(x^{k}))$$

is called the *empirical risk* of a hypothesis h with respect to a dataset  $\hat{z}$ . Let  $h^*$ ,  $\hat{h}^*$  denote the solutions of (3.2) and (2.6), respectively. Approximating  $h^*$  by  $\hat{h}^*$  is called the *inductive empirical risk minimization principle* [106]. Characterizing under what conditions such an approximation is valid is one of the main elements of a learning theory [68, 106].

I start this section by introducing an upper bound to the generalization error in a PL setting, i. e., where a dataset  $\hat{z}$  is given *a priori*, and optimization is performed only over the hypothesis class  $\mathfrak{H}$ . To this end one needs the following concepts.

**Definition 3.1** (Representativeness of a dataset [68]). Given a loss function *l* and a hypothesis class  $\mathfrak{H}$ , define the representativeness of a dataset  $\hat{z}$ 

$$\operatorname{Rep}(\hat{z}) := \sup_{h \in \mathfrak{H}} (R_{\mathcal{P}_z}(h) - \hat{R}_{\hat{z} \sim \mathcal{P}_z}(h)),$$

i.e., the representativeness of a dataset is the biggest generalization error achievable over a certain hypothesis class. Now consider the practical problem of having to compute the representativeness of a dataset  $\hat{z}$ . This is not doable since the computation of  $R_{\mathcal{P}_z}(h)$  needs access to the true distribution of data which is not available. However, an estimate of the representativeness can be obtained by diving the dataset  $\hat{z}$  into two sets  $\hat{z}_1, \hat{z}_2$  and computing the empirical estimate

$$\hat{\operatorname{Rep}}(\hat{z}) = \sup_{h \in \mathfrak{H}} (\hat{R}_{\hat{z}_1 \sim \mathcal{P}_z}(h) - \hat{R}_{\hat{z}_2 \sim \mathcal{P}_z}(h)).$$

To compactify this notation, assume that  $|\hat{z}_1| = |\hat{z}_2| = \frac{m}{2}$  and let  $\sigma = (\sigma_1, ..., \sigma_m)$  be such that  $\sigma_i = 1$  if  $(x^i, E^i) \in \hat{z}_1$  and  $\sigma_i = -1$  if  $(x^i, E^i) \in \hat{z}_2$ . Then, the empirical representativeness can be simplified to

$$\hat{\operatorname{Rep}}(\hat{z}) = \frac{2}{m} \sup_{h \in \mathfrak{H}} \sum_{i}^{m} \sigma_{i} l(E^{i}, h(x^{i})).$$

*Rademacher complexity* generalizes this idea by considering the average empirical representativeness for a random choice of  $\sigma$  with  $\text{Prob}[\sigma_i = 1] = \text{Prob}[\sigma_i = -1] = 0.5$ . This can be understood as taking the average empirical representativeness

over all possible choices of  $\hat{z}_1, \hat{z}_2$  of equal sizes.

**Definition 3.2** (Rademacher complexity [68]). Given a loss function l and a hypothesis class  $\mathfrak{H}$  define the Rademacher complexity

$$\begin{aligned} \operatorname{Rad}(\hat{z}) &:= \frac{1}{2m} \mathbb{E}_{\sigma \sim \{\pm 1\}^m} [\operatorname{Rep}(\hat{z})] \\ &= \frac{1}{m} \mathbb{E}_{\sigma \sim \{\pm 1\}^m} [\sup_{h \in \mathfrak{H}} \sum_{i}^m \sigma_i l(E^i, h(x^i))]. \end{aligned}$$

#### 3.2.1 Upper bound to the true risk in passive learning

**Theorem 3.1** (upper bound on the true risk [68]). Assume that  $|l(E^i, h(x^i))| \le c < \infty$  for all  $(x^i, E^i) \in \hat{z}, h \in \mathfrak{H}$ . Then, with probability of at least  $1 - \delta$ , for all  $h \in \mathfrak{H}$ 

$$R_{\mathcal{P}_z}(h) \le \hat{R}_{\hat{z} \sim \mathcal{P}_z}(h) + 2Rad(\hat{z}) + 4c\sqrt{\frac{2\ln(\frac{4}{\delta})}{|\hat{z}|}}.$$
(3.3)

In particular, this holds for  $\hat{h}^*$  that solves (2.6).

**Remark 3.1.** Note that, in practice, the upper bound derived on the generalization error justifies approximating solutions of (3.2) with solutions of (2.6). The upper bound in (3.3) contains the empirical risk and the Rademacher complexity. The latter is connected to the complexity of the optimizer  $h^*$  [114]. Hence, solving (2.6) is actually equivalent to minimizing an upper bound to the generalization error.

Similarly, an upper bound on the generalization error in AL can provide some insight into a good choice of the probability distribution  $\mathcal{P}_q$  of a query strategy.

#### 3.2.2 Upper bound to the true risk in active learning

Here I derive an upper bound to the generalization error in AL. As in PL case, the upper bound would depend on the empirical risk, i.e., the training error, and on the complexity of the hypothesis class. However, the distribution of the training data  $\mathcal{P}_y$  in AL is not equal to the generating distribution  $\mathcal{P}_z$ . Therefore, it is reasonable to expect the upper bound on the generalization error in AL to depend on some notion of a distance between  $\mathcal{P}_y$  and  $\mathcal{P}_z$ . While there are several possibilities of defining distances between measures, integral probability metrics follow naturally in the settings under consideration.

Consider the measure space  $(X, \mathcal{B}(X))$  (see Section 3.1).

**Definition 3.3** (Integral probability metric [115, 116]). Given a class of real-valued bounded measurable functions  $\mathcal{F}$  on  $(\mathbb{X}, \mathcal{B}(\mathbb{X}))$ , the *integral probability metric* between two measures P, Q on  $(\mathbb{X}, \mathcal{B}(\mathbb{X}))$  is defined as

$$d_{\mathcal{F}}(P,Q) := \sup_{f \in \mathcal{F}} |\int_{\mathbb{X}} f \, dP - \int_{\mathbb{X}} f \, dQ|.$$
(3.4)

**Remark 3.2.** Note that (3.4) is generally not a metric but a pseudometric, since  $d_{\mathcal{F}}(\mathcal{P}, \mathcal{Q}) = 0$  does not imply  $\mathcal{P} = \mathcal{Q}$ . However, for the two special choices of  $\mathcal{F}$  that I later discuss (3.4) is a metric. Moreover, note that, in general, the distance between two probability measures is allowed to be infinite [116, 117].

The stage is now ready to state an upper bound on the generalization error in an AL setting in terms of the integral probability metric. Later, I specify some useful function classes  $\mathcal{F}$ . Recall that, in AL, the distribution of the training data and the generating distribution have the same conditional probability, i. e.,  $p_z = p_x p_{E|x}$  and  $p_y = p_q p_{E|x}$ .

Let  $\hat{x} \sim \mathcal{P}_x$ ,  $\hat{q} \sim \mathcal{P}_q$  be two finite datasets. And let  $\hat{z}$ ,  $\hat{y}$  denote the same datasets with the corresponding labels, i. e., the electronic energies.

**Theorem 3.2.** Define  $\alpha := \int_{\mathbb{E}} l \, d\mathcal{P}_{E|x}$ . Given a function class  $\mathcal{F}$  and assuming  $\alpha \in \mathcal{F}$  the following holds with probability at least  $1 - \delta$ 

$$R_{\mathcal{P}_{z}}(h) \leq \hat{R}_{\hat{y} \sim \mathcal{P}_{y}}(h) + d_{\mathcal{F}}(\mathcal{P}_{x}, \mathcal{P}_{q}) + 2Rad(\hat{y}) + 4c\sqrt{\frac{2\ln(\frac{4}{\delta})}{|\hat{y}|}}.$$
 (3.5)

*Proof.* for all  $h \in \mathfrak{H}$  assume that the loss function is bounded for all elements of  $\hat{q}$ . Using Theorem 3.1 the following holds with probability  $1 - \delta$ 

$$R_{\mathcal{P}_y}(h) \leq \hat{R}_{\hat{y} \sim \mathcal{P}_y}(h) + 2\mathrm{Rad}(\hat{y}) + 4c\sqrt{\frac{2\ln(rac{4}{\delta})}{|\hat{y}|}}.$$

Adding  $R_{\mathcal{P}_z}(h)$  to both sides and rearranging

$$R_{\mathcal{P}_z}(h) \le (R_{\mathcal{P}_z}(h) - R_{\mathcal{P}_y}(h)) + \hat{R}_{\hat{y} \sim \mathcal{P}_y}(h) + 2\operatorname{Rad}(\hat{y}) + 4c\sqrt{\frac{2\ln(\frac{4}{\delta})}{|\hat{y}|}}$$

Consider the first term and note that

$$\begin{aligned} R_{\mathcal{P}_z}(h) - R_{\mathcal{P}_y}(h) &= \int_{\mathbb{Z}} l \, d\mathcal{P}_z - \int_{\mathbb{Z}} l \, d\mathcal{P}_y \\ &= \int_{\mathbb{X}} \int_{\mathbb{E}} l \, p_{E|x} \, d\mu p_x \, d\mu^m - \int_{\mathbb{X}} \int_{\mathbb{E}} l \, p_{E|x} \, d\mu \, p_q \, d\mu^m \\ &= \int_{\mathbb{X}} \alpha \, p_x \, d\mu^m - \int_{\mathbb{X}} \alpha \, p_q \, d\mu^m. \end{aligned}$$

Since  $\alpha \in \mathcal{F}$ , it holds

$$\begin{aligned} R_{\mathcal{P}_z}(h) - R_{\mathcal{P}_y}(h) &\leq \sup_{g \in \mathcal{F}} \left| \int_{\mathbb{X}} g \, p_x \, d\mu^m - \int_{\mathbb{X}} g \, p_q d\mu^m \right| \\ &= d_{\mathcal{F}}(\mathcal{P}_x, \mathcal{P}_q). \end{aligned}$$

Theorem 3.2 establishes an upper bound to the true risk in terms of any generic integral probability metric. By imposing some conditions on the loss function and the conditional probability one can derive from the upper bound in (3.5) various upper bounds. For example, under some conditions on on the loss function  $\mathcal{F} = C_b$ , i. e., the space of continuous and bounded functions. This is, indeed, a good choice since it allows for a unique identification of identical probability measures. In other words, given two probability measures  $\mathcal{P}$ ,  $\mathcal{Q}$ , one can show [116] that

 $\mathcal{P} = \mathcal{Q}$  if and only if

$$\int_{\mathbb{X}} f \, d\mathcal{P} = \int_{\mathbb{X}} f \, d\mathcal{Q} \quad \text{for all } f \in C_b.$$

Another possibility would be to impose conditions on the loss function l such that  $\alpha$  is a Lipschitz continuous function with a Lipschitz constant less than one. In such a case, and assuming that  $\mathcal{P}_x$ ,  $\mathcal{P}_q$  have bounded supports, the Kantorovich metric is recovered, which is the dual of the Wasserstein 1-distance between two probability measures [118]. The following proposition specifies conditions on the loss function l in order recover these two special cases.

**Proposition 3.1.** (*i*) Let  $l(E_0, \cdot) : \mathbb{R}^m \to \mathbb{R}$  be continuous for all  $E_0 \in \mathbb{R}$  and  $l(\cdot, x_0) : \mathbb{R} \to \mathbb{R}$  be bounded for all  $x_0 \in \mathbb{R}^m$ . It follows that  $\alpha \in C_b(\mathbb{X})$ .

(*ii*) Let  $l(E_0, \cdot) : \mathbb{R}^m \to \mathbb{R}$  be Lipschitz-continuous with Lipschitz constant L for all  $E_0 \in \mathbb{R}$  and  $l(\cdot, x_0) : \mathbb{R} \to \mathbb{R}$  be bounded for all  $x_0 \in \mathbb{R}^m$  by some B > 0. Furthermore, assume that  $\max(L, 2B) \leq 1$ . Then  $\alpha$  is Lipschitz-continuous with Lipschitz constant  $\leq 1$ .

*Proof.* Proof is provided in the supplementary material at the end of this section.  $\Box$ 

In light of Proposition 3.1, Theorem 3.2 can be regarded as a more generic form of [105] and [104], where an upper bound on the generalization error in AL in terms of the Wasserstein distance and using a reproducing kernel Hilbert space<sup>4</sup>, respectively, were derived.

The result established in Theorem 3.2 can be used to derive an empirical risk minimization principle for AL. Denote by  $\hat{z}^{(0)} \sim \mathcal{P}_z$  a set of already labeled data. The upper bound suggests a query strategy that chooses data  $\hat{y}$  and a hypothesis *h* such that the empirical risk  $\hat{R}_{\hat{y}\cup\hat{z}^{(0)}}(h)$  is small. For a fixed hypothesis *h*, examples that correspond to a high loss function will then be added to the training set in order to have an overall low empirical loss. Assume that the fixed hypothesis *h* is obtained through a standard training procedure using the already

<sup>&</sup>lt;sup>4</sup>While the authors here assume  $\mathcal{F}$  to be a reproducing kernel Hilbert space, the conditions they impose on the loss function *l* do not guarantee that.

labelled data  $\hat{z}^{(0)}$ . Since the initially labeled data  $\hat{z}^{(0)}$  follows  $\mathcal{P}_z$ , the regions in *X*, on which one expects the empirical loss to be high are actually outliers, i. e., they have a small measure under  $\mathcal{P}_x$ . On the other hand, the upper bound suggests sampling a set  $\hat{q}$  that does not have a big integral probability metric from a dataset  $\sim \mathcal{P}_x$ . Hence, the upper bound suggests sampling points that are both *representative* of the underlying distribution of the unlabeled data  $\mathcal{P}_x$  and points that are outliers, and in a sense, *informative* for any hypothesis *h*. This argumentation is in perfect accordance to a vast literature on the need for querying informative and representative samples in an AL strategy [104, 105, 119–121].

An algorithm to directly minimize the upper bound in (3.3) can be formulated, although it is not clear what statistical distance is best to employ. One point to take into account here is the computational costs of estimating the chosen integral probability metric. For example, methods to estimate the Wasserstein distance are often rather expensive. For constructing PESs, I found out that directly minimizing such upper bounds is extremely impractical and leads to poor results. I observed query strategies that indirectly minimize the upper bound to work better. I outline some of them in the next subsection.

# 3.3 Practical pool-based active learning

One of the most practical and successful AL frameworks is called pool-based active learning [91]. Here, AL is performed in an iterative manner where the query algorithm is given access to a pool of unlabeled data  $\hat{s}^{(0)} = \{(x^i, )\}_{i=1}^l$  and is then asked to query a set of samples *B* from this pool and add it to the already labeled data  $\hat{z}^{(0)} = \{(x^i, E^i)\}_{i=1}^m$ , where  $m \ll l$ . Then, a test is run to judge whether the currently labeled data are enough for a reliable prediction. If not, another set *B* of datapoints are sampled from the pool, labeled, and added to the labeled data. This process continues until one has enough datapoints. Algorithm 1 summarizes this procedure.

The procedure performed to judge whether a certain amount of data is enough for a reliable modeling is often based on training a ML model on the labeled data and evaluating its accuracy on a test set. Note here that, due to the iterative nature Fix batch size  $N_B$ , t = 1; **Input:** Pool of unlabeled data  $\hat{s}^{(t-1)}$ , an initial labeled data.  $\hat{z}^{(0)}$  **while** *performance is unsatisfactory* **do a**) Select a batch *B* of size  $N_B$  from  $\hat{s}^{(0)}$ . **b**) Label this set to obtain  $\hat{z}^{(t)}$ . **c**) Set  $\hat{s}^{(t)} = \hat{s}^{(t-1)} \setminus B$ ,  $\hat{z}^{(t)} = \hat{z}^{(t-1)} \cup \hat{z}^{(t)}$ . **d**) t = t + 1. **end** 

**Algorithm 1:** Basic steps of a generic pool-based AL strategy. In each active learning iteration,  $N_B$  datapoints are chosen from the pool, labeled, and added to the training data.

of pool-based AL, the execution speed and the scaling with the amount of data is a main concern in designing query algorithms.

A simple example of a query algorithm is uniform random sampling (RS) from the pool. Note that such a strategy is representative by construction and results in a small statistical distance in the upper bound (3.5). However, uniform random points from the pool are not informative and would hence lead to a high empirical error.

Another criterion for defining a query algorithm is prediction uncertainty, where an ML model predicts the targets of unlabeled datapoints from the pool, and those corresponding to the highest uncertainties in their predictions are queried. For probabilistic models like GPs, the uncertainties can be directly calculated [95, 97, 98, 102]. For the ML models that do not offer a direct way to compute uncertainties, these can be inferred by training a *diverse* ensemble of models on the currently labeled training set and selecting the points about which the models disagree the most. This algorithm is called query by committee (QBC) [122]. This procedure is formalized in Algorithm 2.

Note that diversity of the ML models is crucial in this algorithm. If the models are not diverse, their predictions for a certain unlabeled datapoint would be almost the same and hence one would not be able to infer the uncertainty. Practically, the diversity of models is introduced through random perturbations to the learning process. For example, when the ML models consisting the ensemble are NNs,

Fix the number of models *n* in the ensemble. **Result:** A batch *B* from  $\hat{s}^{(t)}$ . **Input:**  $\hat{s}^{(t)}, \hat{z}^{(t)}, N_B$ . a) Train an ensemble  $\{T_i\}_{i=1}^n$  of models on data  $\hat{z}^{(t)}$ . b) Compute predictions  $\hat{E}_i = T_i(x^k)$  for all  $x^k \in \hat{s}^{(t)}$ , for all *i*. c) Compute the community disagreement  $q(x^k) = \text{std}(\hat{E}_i)$  for all  $x^k \in \hat{s}^{(t)}$ . d) Take  $N_B$  elements from the unlabeled data that have the highest *q*.

**Algorithm 2:** Basic steps of a query by committee algorithm for regression problems. n models are trained on the labeled datasets and asked to make predictions on the whole unlabeled dataset. The dataset B chosen to be labeled are those that maximize the standard deviation (std) of the prediction among the n models.

diversity can be achieved by randomly initializing their weights, and choosing different architectures and regularization parameters.

Uncertainty-based algorithms aim to minimize the empirical risk in (3.5) by querying points corresponding to high uncertainties, which correspond to underpopulated/sparse areas of  $\mathcal{P}_x$ , i. e., outliers [91, 123], which is a clear downside. Interestingly, the vast majority of AL applications to PESs used uncertainty-based sampling [92, 94–100]. When all sparse regions of the pool can be clearly identified, e. g., as points with high energy, or if prior knowledge about the minima and saddle points exists, the downside of uncertainty-based algorithms can be solved by introducing a weighting function [93, 102]. In a more general setting, one can combine the uncertainty-based query algorithm with a molecular-dynamics sampler starting from various known critical points of the PES [70, 94, 99].

I choose, however, to correct this behavior at a more fundamental level by constructing a probability density function from the QBC-estimated uncertainties. Then, querying grid points is performed through random sampling according to this density function. Algorithm 3 formalizes this idea. In contrast to QBC, points with small uncertainties may still be queried if they fall in high-density regions. In other words, Algorithm 3 respects the statistical information in the pool that is defined *a priori* by the expert. Note that accounting for the statistical

Fix the number of models *n* in the ensemble. **Result:** *B* elements from  $\hat{s}^{(t)}$ . **Input:**  $\hat{s}^{(t)}, \hat{z}^{(t)}, N_B$ . a) Train an ensemble  $\{T_i\}_{i=1}^n$  of *n* models on data  $\hat{z}^{(t)}$ . b) Compute predictions  $\hat{E}_i = T_i(x^k)$  for all  $x^k \in \hat{s}^{(t)}$ , for all *i*. c) Compute the community disagreement  $q(x^k) = \text{std}(\hat{E}_i)$  for all  $x^k \in \hat{s}^{(t)}$ . d) Compute the weights:  $L(x) = \frac{q(x) - q_{\min}}{q_{\max} - q_{\min}}$ , and the sampling probability  $p(x) = \frac{L(x)}{\Sigma_x L(x)}$  where  $q_{\min} = \min_{x \in \hat{s}^{(t)}} q(x)$  and  $q_{\max} = \max_{x \in \hat{s}^{(t)}} q(x)$ . e) sample  $N_B$  elements from the unlabeled data with probabilities p(x).

Algorithm 3: Stochastic query by committee algorithm: Data to query are chosen by sampling according to a probability distribution that gives more weights to datapoints whose predictions are uncertain. Uncertainty is inferred by a standard query by committee algorithm.

information in the pool using QBC can also be performed by considering only a few unlabeled datapoints sampled independently of the input distribution as candidates to query [120], which is very similar in spirit to Algorithm 3. However, I empirically observed Algorithm 3 to work better than this approach.

# **3.4** Simulations on pyrrole(H<sub>2</sub>O)

In what follows I apply the acRS algorithm, Algorithm 2 and Algorithm 3 for building a PES for pyrrole $(H_2O)$  molecules with a reduced number of datapoints.

Due to the highly fluxional nature of the hydrogen bond in  $pyrrole(H_2O)$ , the intermolecular motions are highly delocalized, rendering the calculation and representation of the PES very challenging. The intramolecular vibrations in the pyrrole and water moieties can be described with a relatively simple, though multidimensional, single-minimum PES and thus, for simplicity of calculations, were not considered here. The structures of pyrrole and water monomers were fixed to the experimentally determined values [124, 125], see supplementary material at the end of this chapter, and varied the six intermolecular coordinates, shown



FIGURE 3.1: Internal intermolecular coordinates *R*,  $\theta$ ,  $\phi$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$  of pyrrole(H<sub>2</sub>O).



FIGURE 3.2: The probability density distribution of the energies corresponding to all molecular geometries in the pool. The histogram was calculated for a bin width of 34.5 cm<sup>-1</sup> and has a peak at 1600 cm<sup>-1</sup>, corresponding to the dissociation limit of pyrrole(H<sub>2</sub>O).

in Figure 3.1. These are defined as follows: the relative position of water with respect to pyrrole is described by the three spherical coordinates R = [0.2, 1] nm,  $\theta = [0, \pi], \phi = [0, \pi]$  and the relative orientation of water is defined by the three Euler angles  $\alpha = [0, \pi], \beta = [0, \pi], \gamma = [0, \pi]$ . The angles  $\phi, \alpha$ , and  $\gamma$  were restricted to the ranges  $[0, \pi]$  exploiting the  $C_{2v}(M)$  symmetry of the complex.

The pool of molecular configurations was generated *a priori* as the direct product of one-dimensional grids for every degree of freedom and contained 57500 different molecular geometries covering the potential energy up to 5000 cm<sup>-1</sup> above dissociation. All coordinates were sampled more densely in the vicinity of the equilibrium geometry. Also, the angular coordinates were sampled more densely for small radial distances  $R \leq 500$  pm with a sparser grid for 500 <  $R \leq 1000$  pm. This led to a nonuniform distribution of energies in the pool, shown in Figure 3.2. Note that a direct-product grid is not essential for the accumulation of the pool of unlabeled geometries and the test dataset. Here, it was used mainly because it allows the coverage of the whole configuration space that is relevant for the subsequent quantum dynamics' simulations, and hence prevents biases and holes in the pool and test data. While this method is not arbitrarily extendable to systems with more degrees of freedom, other pool accumulation methods [94] could be used without modifications to the stochastic query by forest (SQBF) approach. The electronic structure calculations employed the density-fitting explicitly-correlated DF-MP2-F12 level of theory [126–128] in the frozen-core approximation using aug-cc-pVDZ-F12 [129] atomic orbital, cc-pVDZ-F12+/OPTRI [130] resolution of the identity, and aug-cc-pVDZ/MP2FIT [131] density fitting bases. The geminal exponent was fixed at 1.0. The electronic structure calculations, i. e., solving the electronic Schrödinger equation (2.4) were carried out using Molpro [132–134]. A subset of 10 % of the total number of points in the pool was randomly selected as a test set and taken out of the pool (OOP). 5 % of the remaining data was randomly selected as a validation set. I employed two different machine learning models, RFR, and NN, to fit the data. Exponential functions of interatomic distances were used, with all distances considered, as molecular descriptors, see supplementary materials at the end of the chapter.

Furthermore, I tested the SQBF algorithm and an NN model on the PES of the  $N_4$  molecule using previously reported electronic structure data [135].

In Algorithm 2 and Algorithm 3 it remains to specify the ensemble of models used to estimate uncertainty. While any ensemble of ML models can be used, I propose to use the trees of a random forest regressor (RFR) as members of this ensemble. I argue that choosing regression trees for inferring uncertainty is advantageous because of relatively low-training complexity and a straightforward diversification-ability. The reader is referred to Appendix D for more information on RFR models. The RFR combined with Algorithm 3 gives rise to a regression version of the stochastic query by forest algorithm (SQBF) [107], employed in this study.

Note that, in Algorithm 3, the balance between sampling points from the sparse and high-density regions is controlled by the function L, which is linear with respect to the community disagreement. The probability of a point being sampled decreases linearly with the decrease of the point's uncertainty. One can have more freedom on this balance by considering powers of this function, i. e.,  $L^{\alpha}$ , where  $\alpha \in \mathbb{R}^+$ . For  $\alpha \in (0, 1)$ , the algorithm will sample more points with low uncertainty and conversely less for  $\alpha \in (1, \infty)$ . We performed a heuristic study of the effect of different powers  $\alpha$ . At each AL iteration, I ran SQBF algorithm for different values of  $\alpha \in \{0.5, 0.75, 1, 1.25, 1.75\}$ . The  $\alpha$  that led to the largest improvement in generalization error was picked the corresponding queried points were collected. I proceeded to AL using this batch as part of the pool. The whole procedure was repeated at every AL iteration. I found only minor improvements of the accuracy when using multiple, optimized, values of  $\alpha$ . I explored a few other heuristics of similar nature, but none of them yielded significantly better results. Hence, throughout the paper I report results obtained with a single value of  $\alpha = 1$ .

#### 3.4.1 Performance

For  $pyrrole(H_2O)$  I compared the performance of the RS, QBC, and SQBF AL query algorithms considering the convergence rate and the fitting accuracy. All

TABLE 3.1: Out-of-the-pool RMS errors (in  $cm^{-1}$ ) of the random forest regressor and neural network models, listed as RFR/NN, computed for various fractions of the total pool data collected by the different AL query.

AL query	20 %	40 %	60 %	80 %	100 %
RS	183/57	117/31	81/20	52/15	39/11
QBC	141/43	74/21	49/13	41/11	38/11
SQBF	88/27	37/14	36/12	38/11	39/11

query algorithms started from the same fixed amount of m = 2458 labeled samples and queried the same equal number of m samples at every AL iteration. For every iteration and query algorithm, I used the RFR and NN models to fit the data. The fitting error is defined as the RMS error of the ML models in predicting the energies on the OOP dataset. This dataset was the same for all query algorithms and followed the joint distribution of the problem  $\mathcal{P}$ . The accuracy of a model on this dataset is an estimate of the generalization error.

The fitting errors of the RFR and NN models for different query algorithms are plotted in Figure 3.3 as functions of the AL iteration, i. e., size of labeled data. The SQBF strategy with RFR model leads to the fastest convergence of the error. QBC strategy outperforms RS. Similar convergence behavior of different query algorithms can be observed for the NN model. For our dataset, the fitting error of NN was smaller than that of RFR for all AL iterations and for all strategies by an average factor of 3.3. Table 3.1 summarizes these results. The better performance of NNs is partially due to the fact that NNs are easier to train to higher accuracy and can approximate complex functions with a better control on the bias-variance trade-off, which was enabled by using an early stopping criterion on the validation set, see the supplementary material at the end of the chapter. The AL iterations were terminated when the pool became empty. In practice, the iterations are to be terminated when the derivative of the fitting error with respect to the amount of labeled data is less than a predefined value [97] or simply when the fitting error of the model is small enough.



FIGURE 3.3: RMS error of out-of-the-pool datasets using (up) the random forest regressor and (down) the neural network models for the RS (triangles), QBC (circles), and SQBF (squares) query. The SQBF has the fastest convergence. A neural network model, trained on 30 % of the total amount of datapoints in the pool achieves an RMS error of 16 cm<sup>-1</sup>. The RMS error on the full dataset is 11 cm<sup>-1</sup>. The neural networks trained on data collected by the QBC or RS algorithms show worse performance. The same convergence patterns hold when using a random forest regressor to train on the data instead of a neural network, albeit at overall somewhat slower convergence.

TABLE 3.2: RMS mean errors and standard deviations using the available data (in cm<sup>-1</sup>) of NN using the AL SQBF algorithm with a NN to fit the data (present work) and latin hypercube sampling with GPs to fit the data [136], applied to PES data of the N<sub>4</sub> molecule [135].

No. training data	NN ( $cm^{-1}$ )	$GP(cm^{-1})$	
240	$36518\pm2697$	$13300\pm2770$	
480	$26207\pm1871$	$10027\pm1371$	
720	$11192\pm1200$	$8401 \pm 1102$	
960	$8111\pm 668$	$7544 \pm 972$	
1200	$6201\pm462$	$6806\pm962$	
1680	$4704\pm612$	±	
1800	$4494 \pm 633$	$5551 \pm 951$	
1920	$4284\pm658$	±	
2400	$3557\pm675$	$5012\pm832$	

Similarly, for the N<sub>4</sub> molecule the SQBF algorithm was used to query geometries from the pool of 16421 molecular geometries reported [135]. The OOP and validation datasets were each generated using 10 % of the uniform-randomly sampled pool data. An initial batch of 240 geometries was uniform-randomly sampled from the pool and the SQBF algorithm queried 240 geometries at each AL iteration. The same molecular descriptor as described above was used to transform the data and an NN model was used for fitting; details on the NN design are provided in the supplementary materials. This procedure was repeated 100 times and the mean and standard deviation of the resulting NN errors on the entire dataset as a function of the number of training examples is reported in Table 3.2. The SQBF results are compared with the ensemble of 100 GPs used to fit the data collected by the Latin hypercube sampling algorithm [136]. The GP method shows a better performance for the first few AL iterations. I attribute this to the fact that it is hard to prevent overfitting with a neural network with a very small set of randomly selected training data. However, already at 1200 training points the two models result in comparable accuracy. With 1680 training points, our SQBF/NN approach achieves the same accuracy as GP with 2400 points, which corresponds to a 30 % reduction in the size of the training dataset. All the following further



investigations are performed for pyrrole(H<sub>2</sub>O).

FIGURE 3.4: Normalized probability density distributions of the number of data points  $N/N_{tot}$  across the potential energies plotted for the data collected by the RS, QBC, and SQBF query at different AL iterations corresponding to 20 %, 40 %, and 60 % of the total pool. The bin width of the histograms is 34.5 cm<sup>-1</sup>.

#### 3.4.2 Distribution of queried data

In Figure 3.4 I plotted the normalized distributions of the samples' electronic energies of pyrrole( $H_2O$ ) collected by different AL query algorithms at three different iterations corresponding to 20 %, 40 %, and 60 % of the total pool. Compare these with the distribution of energies in the total pool Figure 3.2, which has a peak around 1600 cm<sup>-1</sup>, corresponding to the dissociation limit of pyrrole(H<sub>2</sub>O). The densities were computed using 200 equally-sized bins covering the energy range from 0 to 6874 cm<sup>-1</sup> and normalized to the bin width of 34.5 cm<sup>-1</sup>. Evidently, the probability density of data sampled by the acRS query most closely resembles the pool distribution. On the other hand, it is clear that the QBC algorithm samples more data with higher energies, whereas SQBF keeps a balance between both the RS and QBC tendencies. As the number of the labeled data increases, all probability density distributions become more similar to the distribution in the pool.

It is reasonable to expect that a model built on a dataset sampled by QBC algorithm will tend to have a better performance for the high-energy regions. This is demonstrated in Figure 3.5 showing the 2D histograms of OOP energies and the absolute errors of the RFR and NN models in predicting these energies, plotted for different query algorithms. The histograms were computed using 20 and 50 equally-sized bins for the energy and absolute errors, respectively. The size of the training dataset here corresponds to 40 % of the pool's size. We clearly see that RS achieves good accuracy for the points with low energies, QBC works best for the points with high energies, and the SQBF maintains a more regular accuracy across the whole energy spectrum.

#### 3.4.3 Batch size and size of initially labeled dataset

I repeated the above calculations with a smaller batch size of 122 points instead of the initially used 2458, starting from the same initially labeled dataset. The convergence of the RFR fitting error with the number of training data is plotted in Figure 3.6 for different query algorithms. Here, note that both QBC and SQBF strategies benefit slightly from using a smaller batch size. This is in accordance with previous studies that showed a decreasing performance of QBC with increasing batch size, which is due to collecting many similar samples [137].

I also studied the effect of changing the size of initially labeled dataset. Figure 3.7 shows the RFR fitting errors for different query algorithms obtained from initial datasets of 100 and 2458 samples with the batch size of 2458. Observe that



FIGURE 3.5: 2D histograms of discrepancies between the predictions of random forest regressor and neural network models (trained on 40 % of the pool) and the potential energy of the out-of-the-pool data for different query algorithms; 20 and 50 bins were used for energy and absolute error, respectively. Models trained on data collected by QBC tend to perform better on high-energy regions than on low-energy regions. The opposite is true for RS. In contrast, models trained on data collected by SQBF have a more uniform accuracy across the whole energy spectrum.

RS query algorithm outperforms QBC, and that the accuracy of QBC declines significantly. This suggests that with a fewer number of initially labeled data, an AL strategy should focus on collecting grid points from dense regions of the configuration space rather than sampling points with high uncertainties in their



FIGURE 3.6: Effect of the size of initially labeled data on the outof-the-pool error of an RFR model trained using data collected by the RS (blue, triangles), QBC (red, circles), and SQBF (orange, squares) query. Solid (points) and dashed lines correspond to 100 and 2458 initially labeled data, respectively.

predictions. Notably, the SQBF performance is not affected by the size change.

### 3.5 Summary and Conclusion

The first principles calculations of molecular PESs, especially for molecules with many fluxional degrees of freedom, are computationally expensive. One of the major bottlenecks originates from the need to solve the high-dimensional electronic Schrödinger equation (2.4) for tens and hundreds of thousands of different molecular geometries. In particular, standard methods for such calculations suffer from the curse of dimensionality which render them prohibitive. Algorithms that allow to reduce the number of necessary single-point calculations with controlled accuracy of the resulting PES are thus highly demanded. For small molecules, grid reduction algorithms were found beneficial in calculations employing highlevel electron correlation, bases, and relativistic corrections, which are usually computationally affordable only for a relatively small number of points [138–141].



FIGURE 3.7: Effect of the size of initially labeled data on the outof-the-pool error of an RFR model trained using data collected by the RS (blue, triangles), QBC (red, circles), and SQBF (orange, squares) query. Solid (points) and dashed lines correspond to 100 and 2458 initially labeled data, respectively.

I presented in this chapter a theoretical insight into AL, a learning paradigm that allows one to perform statistical inference while minimizing the number of required training datasets. In particular, I presented an upper bound on the generalization error in AL. It suggests that AL algorithms should sample datasets corresponding to high uncertainties in their predictions while not deviating much from the true distribution of the data. I then surveyed practical query algorithms and their applications for constructing PESs. I proposed a regression version of SQBF, a pool-based AL algorithm to generate a compact grid of molecular geometries and the RFR and NN ML-models to construct the six-dimensional intermolecular PES of the weakly-bound pyrrole(H<sub>2</sub>O) complex. I argued that this algorithm is in accordance to the empirical risk minimization principle (3.5). The proposed algorithm led to a roughly two times faster convergence with respect to the number of grid points than the commonly used QBC algorithm to represent the PES to an accuracy of about 16 cm<sup>-1</sup>.

Furthermore, the PES fitted on the data sampled by SQBF exhibited a more uniform accuracy across the whole energy spectrum in comparison to QBC. I empirically showed that the SQBF method is not very sensitive to a variation of parameters such as the size of initially labeled data and size of the batch.

In addition, the proposed method is computationally cheap and scales well with the size of the labeled data N, i. e., as  $\Theta(M \cdot K \cdot \tilde{N} \log_2^2 \tilde{N})$ , where K, M denote the number of random features sampled at each splitting and the number of trees, respectively,  $\tilde{N} = 0.632N$  (see Appendix D). This makes the method attractive for developing universal ML-potentials where large datasets are needed [142–144]. In the case when the accuracy of the RFR is not sufficient for the application of the PES, I showed that the data can be used equally-well by other ML models like NNs. An alternative would be to employ Algorithm 3 with any other ensemble of models or even with a model that offers a direct computation of uncertainty.

Overall, the presented procedure is general and can be applied to the PESs of any polyatomic molecule. It can also be used to model other physical properties like dipole-moment or polarizability surfaces. The major advantage of the proposed method over more popular QBC approach is the heuristic sampling procedure that preserves the distribution of data in the pool while keeping the uncertainty as the primal selection criterion. I believe that in the future the general approach can be improved even further by a better tuned balance between uncertainty and representativeness.

#### Supplementary Material

The regression trees used to implement the QBC and SQBF algorithms were both built using the scikit-learn (sklearn) Python package [145]. All AL algorithms used here were written based on the Libact Python package [146]. For both the QBC and SQBF algorithms I used an ensemble of 100 trees. The training during all AL iterations used an exponential function of the intermolecular distances as a molecular descriptor:  $1 - \exp(-(r - r_0))$  where  $r, r_0$  denote the actual distance and equilibrium distance between two nuclei, respectively. The perturbation of the learning process is controlled through two parameters: (i) a bootstrapping parameter  $\gamma$  that determines the fraction of data sampled by each tree and (ii) the number of features  $\beta$  sampled randomly by each tree. For the batch sizes used in the simulations we experimented with several combinations of these parameters and obtained the best convergence for  $\gamma = 0$ ,  $\beta = 12$  for simulations on pyrrole(H<sub>2</sub>O) and  $\gamma = 0$ ,  $\beta = 4$  for simulations on N<sub>4</sub>. The same parameters were used for uncertainty estimation in both the QBC and SQBF algorithms. Minimal cost complexity pruning was used to reduce the overfitting of RFR with complexity parameter c = 0.01. Since I require that Algorithm 3 queries exactly |B| geometries, one may run into a situation where the number of entries with non-zero probabilities of the distribution p is less than |B|. In such a case I chose to query the elements with the highest uncertainty. The effect of this choice on the simulations conducted in the manuscript is negligible since this case was only encountered once in one of the last AL iterations.

The NN used is a multilayer perceptron and training was implemented using the Python Tensorflow package [147]. The NN has three hidden layers with 256, 512, and 256 neurons, respectively, and a single neuron output layer. The second and third layers were  $l_2$ -regularized with a regularization parameter of  $10^{-5}$ . All hidden layers used "ReLU" as the activation function. The ReLU activation function could be substituted with a smooth approximation such as Softplus to yield a smooth PES. The same aforementioned molecular descriptor was used. The networks were trained for 250 epochs using the Adam optimization algorithm [148], with an initial learning rate of 0.0025 and a decaying learning rate schedule ( $lr_{current} = 0.9825 \times lr_{previous}$ ). An early stopping callback was employed on the validation set that was taken out-of-the-pool with patience of 25. The NN hyperparameters were set to obtain a sufficiently accurate NN, with test error of around 10 cm<sup>-1</sup> when using all the training data.

*Proof of Proposition 3.1.* To prove (i) let  $\epsilon > 0$ . Since *l* is continuous with respect to its second argument there exists  $\delta > 0$  such that  $|x - x_0| \le \delta$  implies  $|l(E_0, x) - \delta| \le \delta$ 

 $l(E_0, x_0) \le \epsilon$  for all  $x_0 \in \mathbb{X}, E_0 \in \mathbb{E}$ . Note that

$$\begin{aligned} |\alpha(x) - \alpha(x_0)| &= \int_{\mathbb{E}} l(.,x) p_{E|x} \, d\mu - \int_{\mathbb{E}} l(.,x_0) p_{E|x_0} \, d\mu \\ &= \int_{\mathbb{E}} \left( l(.,x) - l(.,x_0) \right) p_{E|x} \, d\mu + \int_{\mathbb{E}} l(.,x_0) \left( p_{E|x} - p_{E|x_0} \right) \, d\mu \\ &\leq \epsilon + \int_{\mathbb{E}} l(.,x_0) \left( p_{E|x} - p_{E|x_0} \right) \, d\mu \\ &\leq \underbrace{\epsilon + 2B}_{:=\gamma}, \end{aligned}$$

i. e., for any  $\gamma > 0, x_0 \in \mathbb{X}$  there exists  $\delta > 0$  such that  $|x - x_0| \leq \delta \implies |\alpha(x) - \alpha(x_0)| \leq \gamma$  and hence,  $\alpha$  is continuous.

Since *l* is bounded with respect to its second argument and  $p_{E|x}$  is a finite measure for any  $x \in X$  it follows that  $\alpha$  is bounded.

To prove (*ii*) take  $x_1, x_2 \in X$  and note that

$$\begin{aligned} |\alpha(x_2) - \alpha(x_1)| &\leq \int_{\mathbb{E}} |l(\cdot, x_2) - l(\cdot, x_2)| |p_{E|x_2}| \, d\mu + \int_{\mathbb{E}} |l(\cdot, x_1) \left( p_{E|x_2} - p_{E|x_1} \right)| \, d\mu \\ &\leq L|x_2 - x_1| + 2B \\ &\leq \max(L, 2B) |x_2 - x_1|. \end{aligned}$$

Thus,  $\alpha$  is Lipschitz-continuous with Lipschitz constant  $\leq 1$  if max $(L, 2B) \leq 1$ .  $\Box$ 

## Chapter 4

# Spectral learning: augmenting bases with normalizing flows<sup>1</sup>

The discretization scheme (2.11) proposed to model eigenpairs of (2.5) belongs to the popular class of *spectral methods*, where the target function is expanded by a linear span from a globally defined basis of some function space. Spectral methods were studied quite extensively [28–30], and enjoy an increasing popularity in various applications from computational and engineering sciences. This is mainly due to some favorable approximation to properties of spectral methods, such as their relatively high accuracy and fast convergence for smooth solutions [29, 30]. Moreover, spectral methods are generally a popular and effective choice for modeling highly-oscillatory functions [30]. In fact, spectral methods are the basic building block for a variety of advanced variational techniques to solve Schrödinger equations for the nuclear motion, such as TROVE [149–152], MULTI-MODE [153], TheoRets [154], GENIUSH [155, 156], and others [157, 158]. Despite such favorable approximation properties, spectral methods have high memory requirements and their convergence rate degrades exponentially at increasing

<sup>&</sup>lt;sup>1</sup>This chapter is, in parts, based on the publications: Y. Saleh, A. Iske, A.Yachmenev, J. Küpper, *Proc. Appl. Math. Mech.* **23** (1), *e*202200239 (2023), Y. Saleh, A. Iske, A.Yachmenev, J. Küpper, in preparation (2023), and Y. Saleh, A. F. Corral, A. Iske, A.Yachmenev, J. Küpper, in preparation (2023). Notation was modified when necessary. My contribution to these publications was the development, analysis and implementation of the underlying methods and writing the manuscripts.

54

problem dimension<sup>2</sup>. This phenomenon, referred to as *the curse of dimensionality*, leads to severe limitations, e.g., in applications of quantum mechanics and dynamics, where the systems of interest are inherently high-dimensional.

An alternative linear expansion in an adaptive sequence of functions (2.12) may lessen the need for exponentially many functions at increasing problem dimension. Recently, adaptive nonlinear models, such as neural networks, have been under intensive investigations for approximating solutions to (partial) differential equations [35–37, 41–49, 51]. Their efficiency in approximating high-dimensional functions for challenging applications [159], ranging from image recognition to natural language processing, hints at a great potential for solving high-dimensional differential equations, in particular equations that lend themselves to variational formulations. One important problem class is that of infinite dimensional eigenvalue problems, e.g., static Schrödinger equations. Such problems are strongly related to variational simulations of numerous physics phenomena and, moreover, they often demand solutions for many eigenvalues which correspond to highlyoscillatory functions. Indeed, neural networks were successfully applied to various finite- [41, 43] and infinite-dimensional [42, 46, 47] quantum systems, yielding high accuracies at a lower computational scaling [42, 46], compared to traditional methods. Approximating functions by standard neural network architectures, such as multilayer perceptrons, is, however, rather fragile/not reliable [50] due to the sensitivity of the approximations to the learning parameters, e.g., the network's architecture and training parameters. This results in a need for tedious and elaborate engineering efforts to obtain accurate converged results. This complexity manifests itself clearly in solving static Schrödinger equations [47, 49, 160]. It was shown, for example, that variational schemes that use neural networks to approximate excited states of electronic Schrödinger equations suffer from convergence issues if not initialized to reproduce lower-resolution solutions [47]. I also empirically observed that such standard neural networks are incapable of simultaneously computing many eigenpairs of (2.5). For 3-dimensional molecular systems, I managed to compute only 5 eigenpairs. Trials to compute more

<sup>&</sup>lt;sup>2</sup>See the demonstrative simulations in Chapter 2. For a rigorous numerical analysis, look at, e.g., error bounds for approximating Schwartz functions in the linear span of Hermite functions [15].

eigenpairs often went below the variational minima, a clear hint that the approximation space of the employed neural network falls out of the function space of the variational formulation ( $H^2$  space in case of (2.5), see Theorem 2.2). These empirical difficulties in utilizing neural networks for solving differential equations is often accompanied by a lack of practical convergence guarantees. Formally, neural networks can be used to uniformly approximate any continuous function in  $L^{\infty}(K)$ , where K is compact [161]. Neural networks can be used to approximate any Sobolev function as well. Here, convergence guarantees and rates can be derived with errors measured in  $L^p$  spaces assuming that certain Sobolev embedding conditions hold [162]. However, convergence rates here suffer from the curse of dimensionality. Various universal approximation theorems with dimensionindependent convergence properties were derived mainly in Barron space<sup>3</sup> [50, 51]. However, solutions to differential equations often lie in Sobolev spaces. This renders the use of such approximation theorems less straightforward. Indeed, several results on analyzing neural networks for solving differential equations assume that the solutions and the data of the equation lie in Barron spaces [53–55].

In this chapter I propose and study a special construction of reliable nonlinear approximators. It is, conceptually, based on carefully deforming standard bases *via* special neural networks. The following example illustrates the basic idea.

**Example 4.1.** [Approximating a Gaussian - a motivating example of spectral learning] Denote by  $(\gamma_n)_{n \in \mathbb{N}_{>0}}{}^4$ ,  $\gamma_n : \mathbb{R} \to \mathbb{R}$  the sequence of Hermite functions. This is defined by

$$\gamma_n(x) := a_n \mathfrak{h}_n(x) \exp(-x^2/2), \tag{4.1}$$

where  $\mathfrak{h}_n$  denotes the *n*th Hermite polynomial, and  $a_n$  is a normalizing coefficient. Figure 4.1 shows these functions for n = 1, ..., 5. Let *f* be a normalized Gaussian function centered around a point *a*, i.e.,  $f(x) = \frac{\sqrt{2}}{\sqrt{\sqrt{\pi}}} \exp(-(x-a)^2/2)$ ,  $x \in \mathbb{R}$ , and consider approximating *f* in the linear

<sup>&</sup>lt;sup>3</sup>Function spaces that are tailored to neural networks.

56

span of  $(\gamma_n)_{n \leq N}$  for some  $N \in \mathbb{N}_{>0}$ . Hermite functions  $(\gamma_n)_n$  is an orthonormal basis of  $L^2(\mathbb{R})$ . Thus, such an approximation problem is well-posed in the sense that convergence guarantees can be obtained as N goes to infinity. It appears, however, that an exaggerate approximation method is being used to approximate f since  $\gamma_1$  is actually a Gaussian function. In other words, it does not seem that one needs a very large N to obtain a good approximation. However, it turns out that the approximation error in  $L^2$  becomes small from  $N > \frac{e}{2}a^2$  onward<sup>5</sup> [15]. In other words, the number of functions one needs to have a good approximation depends nonlinearly on the center of the Gaussian a.

Consider now the slightly modified basis  $(\gamma_n^h)_n$  where  $\gamma_n^h = \gamma_n \circ h$  and h(x) = x - a. Note that  $f = \gamma_1^h$ , i.e., one needs only one function of the sequence  $(\gamma_n^h)_n$  to reproduce the target function exactly.



FIGURE 4.1: The figure shows Hermite functions (4.1) for n = 1, ..., 5.

<sup>&</sup>lt;sup>4</sup>From here onward I write  $(\gamma_n)_n$  for notational simplicity.

<sup>&</sup>lt;sup>5</sup>To see this observe Lemma 4.2. Compute the quantity *Af* and use Stirling's formula for *N*!

The previous example shows that the composition operation with an appropriate function can dramatically decrease the computational costs of linear approximation methods. The main thesis of the second contribution of my work is to learn such appropriate functions.

While bases of generic functional spaces can be considered, I restrict the choice to the relevant case of  $L^2(\mu)$  where  $\mu$  denotes the Lebesgue measure. In particular, given a basis  $(\gamma_n)_n$  of  $L^2(\mu)$  I seek approximate solutions of generic approximation problems and (2.5) in particular *via* the ansatz

$$\phi_{N,h} = \sum_{n \le N} c_n \gamma_n \circ h \mid \det Dh \mid^{1/2},$$
(4.2)

where *h* is a function that belongs to a hypothesis class of bijections  $\mathfrak{H}$ , and multiplying by the determinant of the Jacobian of *h* serves to conserve the possible orthonormality of  $(\gamma_n)_n$ . For a fixed  $N \in \mathbb{N}_{>0}$  consider the family of sequences induced by  $\mathfrak{H}$ 

$$\left\{ (\gamma_n \circ h \mid \det Dh \mid^{1/2})_n \mid h \in \mathfrak{H} \right\}.$$
(4.3)

Formally, applying the non-linear ansatz (4.2) is equivalent to identifying a specific  $h^*$  from  $\mathfrak{H}$  that yields the optimal approximation of the target function within the linear span of the truncated sequence  $(\gamma_n \circ h^* | \det Dh^* |^{1/2})_{n < N}$ .

While this scheme has already been proposed in [163] with  $\mathfrak{H}$  being a class of normalizing flows [164, 165], it was only applied to toy quantum models, and no convergence analysis was performed<sup>4</sup>. This chapter aims at providing a rigorous theoretical analysis of approximation schemes based on (4.3) through answers to three main questions.

First, under what conditions on  $\mathfrak{H}$ , if any, are members of the family (4.3) bases of  $L^2(\mu)$ ? If such conditions exist and can be practically imposed, approximations in the linear span of any truncated sequence  $(\gamma_n \circ h | \det Dh|^{1/2})_{n \leq N}$  would be well-posed. I provide an answer based on studying the push-forward measures

<sup>&</sup>lt;sup>4</sup>In the original paper [163] such approximating schemes were proposed to solve Schrödinger equations, where they were referred to as *quantum flows*. I recognize the applicability of such models in approximation problems unrelated to differential equations or quantum mechanics, and therefore refrain from using this terminology.

 $h_{\#}\mu$ . It is shown that essential boundedness conditions on their Radon-Nikodym derivatives are sufficient to this end. Furthermore, it is shown that common implementations of normalizing flows, such as invertible residual networks, satisfy these conditions. I, therefore, proceed to implement a Bubnov-Galerkin numerical scheme for solving differential equations using (4.3). This scheme is called *spectral learning*.

Second, what conditions, if any, can one impose on  $\mathfrak{H}$  such that approximations in the linear span of any truncated sequence  $(\gamma_n \circ h |\det Dh|^{1/2})_{n \leq N}$  converge, in some sense, as the truncation parameter N goes to infinity? I call such a convergence a linear convergence. I answer this question for approximating Schwartz functions and static Schrödinger equations, in particular (2.5), under some assumptions on the Hamiltonian. I provide convergence guarantees in  $L^2$  and characterize the convergence order. I show that standard numerical analysis can be recovered from these results by setting h = id.

Third, does spectral learning achieve faster linear convergence than standard spectral methods? More specifically, does there exist  $h^* \in \mathfrak{H}$  such that an approximation in the linear space of  $(\gamma_n \circ h^* |\det Dh^*|^{1/2})_{n \leq N}$  converges faster than an approximation in the linear space of  $(\gamma_n)_{n \leq N}$  as N goes to infinity? I answer this question positively by studying the push-forward measures  $h_{\#}\mu$  induced by all  $h \in \mathfrak{H}$ .

In addition to the theoretical results, I report simulations I performed to solve (2.5) for polyatomic molecules demonstrating the theoretical findings on the advantages of using spectral learning. In particular, numerical simulations show a two-order increase in accuracy upon the use of approximation schemes based on (4.3).

While the dimension *d* does not play a crucial role in the discussion throughout this chapter, the theoretical results are, nevertheless, presented for an arbitrary dimension. The aim here is to allow for future studies on the effect of dimension on convergence rates for solving approximation problems using augmented bases. This chapter assumes familiarity with neural networks and gradient descent optimization algorithms. It also assumes familiarity with basic measure theory. Fundamental results are, however, provided in Appendix D.
# 4.1 Augmenting expressivity via normalizing flows

Normalizing flows [164, 165] are a powerful tool for generative and discriminative modeling of probability distributions, i. e., in data-based modeling one can use normalizing flows to infer the labels of unseen data and to generate new data examples that have, approximately, the same probability distribution of the training data. Normalizing flows are based on augmenting the expressivity of base distribution that are often easy to evaluate and sample from. I start by pushing this idea to augmenting the expressivity of bases of  $L^2(\mu)$ .

## 4.1.1 Augmenting the expressivity of base distributions

Consider the standard example of supervised learning, where one aims at inferring the probability distribution  $\mathcal{P}_{x,E}$  governing the relation between two random variables or vectors x, E from a finite dataset that is sampled from this distribution. One can approximate  $\mathcal{P}_{x,E}$  via a trial probability distribution, e. g., a Gaussian  $\mathcal{P}_0$ . The approximation process consists of optimizing the parameters of  $\mathcal{P}_0$ , i. e., its width and center, in order to minimize some loss function<sup>5</sup>. However, in cases of very simple base distribution, such as Gaussians, there is not much one can do to approximate potentially complex distributions  $\mathcal{P}_{x,E}$ . One way forward is to increase the expressivity of such base distributions by composing them with a function h, i. e.,

$$\mathcal{P}_h := \left| (\mathcal{P}_0 \circ h) \right| \det Dh |,$$

where the multiplication with the determinant of the Jacobian guarantees that  $\mathcal{P}_h$  integrates to 1, i. e., that it is a valid probability distribution. The task, then, is to find an *h* in a hypothesis class  $\mathfrak{H}$  that would minimize some discrepancy between  $\mathcal{P}_h$  and  $\mathcal{P}_{x,E}$ .

In practice, the class  $\mathfrak{H}$  is modeled by smooth bijections [164, 165] to allow for generative and discriminative modeling of  $\mathcal{P}_{x,E}$ . Such classes are referred to as *normalizing flows*. A common approach to model  $\mathfrak{H}$  is *via* neural networks. There are several ways to produce invertible neural networks [164, 165]. The following

<sup>&</sup>lt;sup>5</sup>See also the introduction to supervised learning in Chapter 3.

example illustrates one way to construct a normalizing flow using standard neural networks.

**Example 4.2.** [Invertible residual neural networks] Consider, residual neural network (ResNet), i. e., a neural network composed of concatenated blocks of the form

$$h(x) = x - K(x),$$
 (4.4)

where *K* is a standard concatenation of linear layers and nonlinear activation functions. It can be shown that such a model is invertible if *K* in each block is Lipschitz continuous with a Lipschitz constant < 1. To guarantee that all linear transformations  $W_i$  in a neural network satisfy this condition one can divide by the biggest singular value  $\sigma_i$  [166], i. e.,

$$ilde{W}_i = \left\{ egin{array}{c} c rac{W_i}{\sigma_i} \, : rac{c}{\sigma_i} < 1 \ W_i ext{ otherwise } \end{array} 
ight.$$

where c < 1 is a hyperparameter. The overall neural network is, hence, invertible upon using, e.g., a *Lipswish* nonlinearity

$$\sigma(x) = \frac{1}{1.1} \cdot \frac{x}{1 + \exp(-x)} \,. \tag{4.5}$$

Next, augmenting the expressivity of bases for solving differential equations is proposed and discussed.

### **4.1.2** Augmenting the expressivity of bases of $L^2(\mu)$

Given a basis  $(\gamma_n)_n$  of square integrable real-valued functions  $L^2(\mu)$  and a hypothesis class  $\mathfrak{H}$  I study the family of sequences (4.3) induced by  $\mathfrak{H}$ . An answer to the first question posed in the introduction of this chapter is provided.

While most of the following results hold for an abstract measure space under the sole constraint of the measure being  $\sigma$ -finite, the analysis is restricted to the

relevant case of  $(\Omega \subseteq \mathbb{R}^d, \mathcal{B}(\Omega), \mu)$  (see Remark D.1) where  $\mathcal{B}$  denotes the Borel  $\sigma$ -Algebra generated by  $\Omega$ . Denote by  $\langle ., . \rangle$  its inner-product.

An important property of mappings for the discussion of this chapter is that of non-singularity.

**Definition 4.1.** A measurable mapping  $h : \Omega \to \Omega$  is said to be non-singular if  $\mu(h^{-1}(B)) = 0$  whenever  $\mu(B) = 0$  for all  $B \in \mathcal{B}$ .

Throughout this chapter let

 $\mathfrak{H} = \{h : \Omega \longrightarrow \Omega \mid h \text{ is a non-singular bijective measurable mapping}\}$ 

be a hypothesis class. Note that  $\mathfrak{H}$  induces a class of linear operators  $\{C_h \mid h \in \mathfrak{H}\}$ on  $L^2(\mu)$  that send any function f into the linear space of all measurable functions on  $\Omega$  defined by

$$f \circ h$$
 for all  $h \in \mathfrak{H}$ .

Since  $\mu$  is  $\sigma$ -finite the non-singularity assumption of h guarantees the existence of the Radon-Nikodym derivative  $\frac{dh_{\#}\mu}{d\mu}$  of the push-forward measures  $h_{\#}\mu$  for all  $h \in \mathfrak{H}$ , and for all  $B \in \mathcal{B}(\Omega)$  (see Lemma D.1 and Theorem D.1). To guarantee wellposedness of the following definition assume that  $\frac{dh_{\#}\mu}{d\mu} \neq 0 \mu$ - almost everywhere and note that this implies  $\mu \ll h_{\#}\mu$ . Since also  $h_{\#}\mu \ll \mu$  it follows that  $h_{\#}\mu$  is equivalent to  $\mu$ .

Next, I restate (4.3) in an abstract measure-theoretic formalism.

**Definition 4.2** (augmented sequence of functions). For  $h \in \mathfrak{H}$ , define the *augmented* sequence of functions  $(\gamma_n^h)_n$  with

$$\gamma_n^h := (\underbrace{\gamma_n \circ h}_{:=\tilde{\gamma}_n^h}) |\frac{dh_{\#}\mu}{d\mu}|^{-1/2}.$$
(4.6)

The family of sequences  $\{(\gamma_n^h)_n \mid h \in \mathfrak{H}\}$  is called *a family/class of augmented* sequences (induced by \mathfrak{H}).

**Remark 4.1.** [On notation] In what follows I will use the hypothesis class  $\mathfrak{H}$  to define several operations and families. To this end, the following terminology is used. I say that a certain property holds for  $\mathfrak{H}$  whenever this property holds for all  $h \in \mathfrak{H}$ . For example,  $\mathfrak{H}$  is Lipschitz or smooth means that any  $h \in \mathfrak{H}$  is Lipschitz, smooth, respectively.

Fix  $h \in \mathfrak{H}$ . To develop a framework, where all sequences in an  $\mathfrak{H}$ -augmented family (Definition 4.2) form bases of  $L^2(\mu)$  one needs well-defined projections  $\langle f, \gamma_n^h \rangle$  of any  $f \in L^2(\mu)$ , for all  $n \in \mathbb{N}_{>0}$ . Since  $h : \Omega \to \Omega$  is a measurable bijection on  $(\Omega, \mathcal{B}(\Omega))$  it follows that the inverse is also measurable, i. e.,  $h(A) \in \mathcal{B}(\Omega)$  for any  $A \in \mathcal{B}(\Omega)$  [167]. Therefore, one can write

$$\langle \gamma_n \circ h \frac{1}{(dh_{\#}\mu/d\mu)^{1/2}}, f \rangle = \langle \gamma_n, f \circ h^{-1} (dh_{\#}\mu/d\mu)^{1/2} \rangle$$

Denote by  $\langle ., . \rangle_{h_{\#}^{-1}\mu}$  the  $L^2$  inner-product on the weighted space with respect to the push-forward measure  $h_{\#}^{-1}\mu$ , i. e.,

$$\langle f,g\rangle_{h^{-1}_{\#}\mu} = \int_{\Omega} fg \, dh^{-1}_{\#}\mu.$$

The main result in this section are sufficient conditions on the class  $\mathfrak{H}$  for (4.6) to be a basis of  $L^2(\mu)$  for all  $h \in \mathfrak{H}$ .

**Theorem 4.1** (Augmented basis [168]). Let  $(\gamma_n)_n$  be an orthonormal basis of  $L^2(\mu)$ and  $\mathfrak{H}$  be as above. For all  $h \in \mathfrak{H}$  it holds that

- (i)  $(\tilde{\gamma}_n^h)_n$  (see (4.6)) is an orthonormal basis of  $L^2(h_{\#}^{-1}\mu)$ .
- (ii) If, in addition, the Radon-Nikodym derivative  $\frac{dh_{\#}\mu}{d\mu}$  is bounded  $\mu$ -almost everywhere, and bounded away from zero  $\mu$ -almost everywhere, then  $(\gamma_n^h)_n$  is an orthonormal basis for  $L^2(\mu)$ .

**Remark 4.2.** For all  $h \in \mathfrak{H}$  one can show that (by noting that  $h_{\#}\mu$  is locally finite and applying Theorem D.2)

$$\frac{dh_{\#}\mu}{d\mu} = \frac{1}{|\det Dh|}.$$

By the inverse function theorem

$$\frac{1}{|\det Dh|} = |\det Dh^{-1}|.$$

Assuming  $\mathfrak{H}$  is differentiable and under the assumption (ii) of Theorem 4.1,  $\mathfrak{H}$  is bi-Lipschitz [169]. This remark is important, since it links the assumptions of Theorem 4.1 to invertible ResNets (4.4). For the rest of this chapter set *r*, *R* to be the lower and upper bounds on the determinant, i. e.,

$$r^d \leq |\det Dh| \leq R^d$$
,

and

$$1/R^d \le |\det Dh^{-1}| \le 1/r^d.$$

*Proof of Theorem 4.1.* Fix an arbitrary  $h \in \mathfrak{H}$ . The orthonormality of  $(\tilde{\gamma}_n^h)_n$  can readily be seen by a simple change of variable. To prove that  $(\tilde{\gamma}_n^h)_n$  is a basis for  $L^2(h_{\#}^{-1}\mu)$ , take one  $f \in L^2(h_{\#}^{-1}\mu)$  satisfying  $f \perp \tilde{\gamma}_n^h$  for all n. In this case

$$0 = \langle \tilde{\gamma}_{n}^{h}, f \rangle_{h_{\#}^{-1}\mu}$$
$$= \int_{\Omega} \gamma_{n} f \circ h^{-1} d\mu \quad \text{for all } n$$

Since  $f \in L^2(h_{\#}^{-1}\mu)$ , one has that  $f \circ h^{-1} \in L^2(\mu)$  and since  $(\gamma_n)_n$  is a basis of  $L^2(\mu)$ ,  $C_{h^{-1}}f = f \circ h^{-1} \equiv 0$ . Since  $h^{-1}$  is invertible,  $f \equiv 0$ . Conclusion follows by Proposition A.3.

To prove that  $(\gamma_n^h)_n$  is a basis for  $L^2(\mu)$  take  $f \in L^2(\mu)$  such that

$$0 = \int_{\Omega} f(\gamma_n \circ h) |\det Dh|^{1/2} d\mu$$
$$= \int_{\Omega} (f \circ h^{-1}) \gamma_n |\det Dh^{-1}|^{1/2} d\mu \quad \text{for all } n.$$

Note that

$$\int_{\Omega} f^2 \circ h^{-1} |\det Dh^{-1}| \, d\mu \le \|f \circ h^{-1}\|_{L^2(\mu)}^2 \|\det Dh^{-1}\|_{L^{\infty}(\mu)} < \infty,$$

since the Radon-Nikodym derivative of  $h_{\#}\mu$  is bounded from below and above. Thus,

$$f \circ h^{-1} |\det Dh^{-1}|^{1/2} \in L^2(\mu).$$

Since  $(\gamma_n)_n$  is a basis for  $L^2(\mu)$  one has that  $f \circ h^{-1} |\det Dh^{-1}|^{1/2} = 0$ , and thus,  $f \equiv 0$  and  $(\gamma_n \circ h |\det Dh|^{1/2})_n$  is a basis for  $L^2(\mu)$  by Proposition A.3.

I provide an alternative proof for (ii). Given an arbitrary  $f \in L^2(\mu)$  one saw that  $f \circ h^{-1} |\det Dh^{-1}|^{1/2} \in L^2(\mu)$ . Define

$$f_N = \sum_{n=1}^N \langle f \circ h^{-1} | \det Dh^{-1} |^{1/2}, \gamma_n \rangle \gamma_n.$$

Since  $(\gamma_n)_n$  is a basis for  $L^2(\mu)$ , one has that

$$\lim_{N \to \infty} \int_{\Omega} \sum_{n=1}^{N} |\langle f \circ h^{-1} | \det Dh^{-1} |^{1/2}, \gamma_n \rangle|^2 \gamma_n^2 \, d\mu = \int_{\Omega} f^2 \circ h^{-1} |\det Dh^{-1} | \, d\mu.$$

Hence,

$$\lim_{N \to \infty} \int_{\Omega} \sum_{n=1}^{N} |\langle f, \gamma_n \circ h| \det Dh |^{1/2} \rangle|^2 \gamma_n^2 \circ h| \det Dh| \, d\mu = \int_{\Omega} f^2 \, d\mu.$$

Therefore,  $(\gamma_n^h)_n$  is a basis for  $L^2(\mu)$ .

**Remark 4.3.** [On the restrictivity of the conditions on  $\mathfrak{H}$ ] In order to have augmented sequences of functions with density properties in  $L^2(\mu)$  one had to adopt the restrictive conditions that  $\mathfrak{H}$  is a class of bijections and that

$$\|\det Dh\|_{L^{\infty}(\mu)} < \infty$$
  
 $\|\det Dh^{-1}\|_{L^{\infty}(\mu)} < \infty,$ 

for all  $h \in \mathfrak{H}$ . Assuming  $\mathfrak{H}$  is also differentiable, this means that  $\mathfrak{H}$  is bi-Lipschitz. This is, indeed, a very restrictive condition on the hypothesis class. Enforcing Lipschitz constraints on machine learning models is linked to less expressivity [169, 170]. However, machine learning models with Lipschitz constraints also have some advantages, e.g., they were linked with better generalization capabilities [171]. In addition, models with small Lipschitz constants are more stable during training [172] and less prone to numerical errors.

This constraint is particularly relevant since several of the most popular normalizing flows are bi-Lipschitz functions, e.g., the invertible residual neural network (4.4).

The hypotheses of Theorem 4.1 are connected to the study of composition operators  $C_{\mathfrak{H}}$  on  $L^p(\mu)$  spaces [173, 174]. As these connections might serve for a further development of adaptive bases, I comment on them in the following remark.

**Remark 4.4.** [On connections to composition operators] For an arbitrary  $h \in \mathfrak{H}$ , note that  $\frac{dh_{\#}\mu}{d\mu} < \infty \mu$ -almost everywhere is a necessary and sufficient condition for the induced composition operator  $C_h : L^2(\mu) \to L^2(\mu)$  to be bounded [173, 174]. Hence, the hypothesis of Theorem 4.1 implies that both  $C_h$  and  $C_{h^{-1}}$  are bounded. Moreover, the hypothesis implies that  $C_h$  is an invertible operator and that the inverse  $C_h^{-1} = C_{h^{-1}}$  [173, 174].

The stage is now ready to employ the augmented family  $\{(\gamma_n^h)_n \mid h \in \mathfrak{H}\}$  for approximation problems.

## 4.2 From spectral methods to spectral learning

Spectral methods are a powerful tool for solving approximation problems appearing in partial differential equations. For an open domain  $\Omega \subseteq \mathbb{R}^d$  and a mapping  $u : \Omega \to \mathbb{R}$ , I introduce spectral method for the generic case of solving

$$\begin{cases} \mathcal{L}u &= f \quad x \in \Omega \\ \mathcal{B}u &= 0 \quad x \in \partial \Omega \end{cases}$$
(4.7)

where  $\mathcal{L}, \mathcal{B}$  are some linear operators and  $f : \Omega \to \mathbb{R}$  is given. One can allow for less regular solutions of (4.7) by adopting a weak formulation. This is derived by multiplying by a test function v and integrating

$$\begin{cases} \int_{\Omega} (\mathcal{L}u) v \, d\mu &= \int_{\Omega} f v \, d\mu \\ \int_{d\Omega} (\mathcal{B}u) v \, d\mu &= 0 \end{cases}$$
(4.8)

Note that the strong (4.7) and weak (4.8) formulations are equivalent when the solutions are smooth. However, the weak formulation can allow for less regular solutions by imposing regularity assumptions on the test function v. It can be shown that the weak solutions solve the strong formulation (4.7) in the sense of distributions [29]. Unlike finite difference methods, spectral methods often adopt the weak formulation to construct the solutions [28–30]. Spectral methods are based on approximating solutions of the weak formulation (4.8) by  $\tilde{u}_N$ , a linear combination of elements of a truncated sequence of globally-supported functions

$$\tilde{u}_N(x) = \sum_{n=1}^N c_n \gamma_n(x).$$

 $(\gamma_n)_n$  is a sequence with some density properties, i. e., the linear span of  $(\gamma_n)_n$  is dense in some function space. This allows one to derive convergence guarantees as *N* grows to infinity. One can differentiate between different spectral methods based on the choice of the test function *v* [29]. Choosing  $v = \tilde{u}_N$  with  $\beta\gamma_n(x) = 0$  for all *n* results in a so-called *Bubnov-Galerkin* spectral method. While choosing

 $v = \tilde{u}_N$  where  $\gamma_n$  do not satisfy the boundary conditions is referred to as the *Tau method* [28–30].

**Example 4.3.** [Spectral discretization of the nuclear TISE] Using an orthonormal basis  $(\gamma_n)_n$  of  $L^2(\mu)$  to discretize (2.5) in a Bubnov-Galerkin framework and noting the variational formulations obtained in Theorem 2.2 and Theorem 2.3 one obtains

$$\tilde{H}\tilde{C}_n = \tilde{E_n}\tilde{C_n} \quad \text{for } n = 1, \dots, N,$$
(4.9)

where  $\tilde{H}$  is an  $N \times N$  matrix whose *ij*th entry is  $\tilde{H}[i, j] = \langle \gamma_i, H\gamma_j \rangle$ , and  $\tilde{C}_n$  in (4.9) is a vector of length N. Hence, solving (2.5) boils down to solving the finite dimensional eigenvalue problem (4.9), i. e., to finding all eigenpairs  $(\tilde{E}_n, \tilde{C}_n)$  that satisfy (4.9).

Given a basis  $(\gamma_n)_n$  I showed in the previous section that it is possible to obtain a family  $\{(\gamma_n)_n \mid h \in \mathfrak{H}\}$  where each member of this family is a basis of  $L^2(\mu)$ . Extending spectral methods to an optimization over this family, i. e., allowing for an optimization of the basis, is straightforward. To this end, consider the family of linear spans of truncated basis for some  $N \in \mathbb{N}_{>0}$ 

$$\left\{ \text{span} \left( \gamma_n^h \right)_{n=1}^N \mid h \in \mathfrak{H} \right\}.$$
(4.10)

**Definition 4.3** (Spectral learning). Given a basis  $(\gamma_n)_n$  of  $L^2(\mu)$ , and a hypothesis class  $\mathfrak{H}$  that induces a family of bases of  $L^2(\mu)$ , an approximation paradigm where an approximate solution of a differential equation is looked for in (4.10) is called *spectral learning*.

I now proceed to employ spectral learning to solve concrete approximation problems. I consider approximating eigenvalues of infinite-dimensional operators in general and (2.5) in particular. For this purpose, I also study approximation of Schwartz functions. Here, questions number two and three in the introduction of this chapter guide the analysis. While the analysis can be carried out for a generic basis of  $L^2(\mu)$ , I set  $(\gamma_n)_n$  to be the sequence of Hermite functions (4.1) for simplicity. Note that this sequence of functions is a basis of  $L^2(\mathbb{R})$ . In multiple dimensions I define the direct product basis

$$\gamma_{(n_1,n_2,\ldots,n_d)}(x_1,\ldots,x_d) := \gamma_{n_1}(x_1)\ldots\gamma_{n_d}(x_d).$$

The size of the full direct product increases exponentially with the dimension *d*. I use, therefore, the hyperbolic reduced tensor-product basis [15]

$$\mathcal{N} = \mathcal{N}(d, N) := \{ (n_1, \dots, n_d) : n_j \ge 0, \prod_{j=1}^d (1+n_j) \le N \}.$$
(4.11)

### 4.2.1 Linear convergence analysis

It is shown in this section that the use of any truncated augmented basis  $(\gamma_n^h)_{n \le N}$  for solving approximation problems is well-posed, in the sense that convergence guarantees can be obtained as N grows to infinity.

The reduced direct product (4.11) of Hermite functions is augmented, in the sense of Definition 4.2, by a hypothesis class  $\mathfrak{H}$  that induces a family of bases. The approximation space is of the form (4.10). For a target function f note that  $\mathfrak{H}$  induces a family of projection operators  $\left\{ \tilde{P}_{\mathcal{N}}^{h} \mid h \in \mathfrak{H} \right\}$  where

$$\tilde{P}^{h}_{\mathcal{N}}f(x) := \sum_{n \in \mathcal{N}} \langle f, \tilde{\gamma}^{h}_{n} \rangle_{h^{-1}_{\#} \mu} \tilde{\gamma}^{h}_{n}(x),$$

where the range of  $\tilde{P}^h_{\mathcal{N}}$  is span  $(\tilde{\gamma}^h_n)_{n \in \mathcal{N}}$ . Set  $\tilde{P}^{h, \perp}_{\mathcal{N}} := I - \tilde{P}^h_{\mathcal{N}}$ . Similarly, one has the family  $\left\{ P^h_{\mathcal{N}} \mid h \in \mathfrak{H} \right\}$  where

$$P^h_{\mathcal{N}}f(x) := \sum_{n \in \mathcal{N}} \langle f, \gamma^h_n \rangle \gamma^h_n(x),$$

and  $P_{\mathcal{N}}^{h,\perp} := I - P_{\mathcal{N}}$ . I drop the notational dependence of all these projection operators on *h* for simplicity. In what follows I set  $\|\cdot\| = \|\cdot\|_{L^{2}(\mu)}$ .

I present the following lemma on the relation between the  $L^2(\mu)$  space and weighted spaces induced by  $\mathfrak{H}$ .

**Lemma 4.1.** Let  $\mathfrak{H}$  satisfy the hypothesis (ii) of Theorem 4.1 and let r, R be as in Remark 4.2. It holds that  $L^2(\mu) = L^2(h_{\#}^{-1}\mu) = L^2(h_{\#}\mu)$  for all  $h \in \mathfrak{H}$ .

*Proof.* Provided in the supplementary material at the end of this chapter.  $\Box$ 

The following quantity is required for the linear convergence analysis.

**Definition 4.4** (Dirac ladder operator). The Dirac ladder operator  $A_i$  is given by

$$A_j = \frac{1}{\sqrt{2}} (q_j + d/dx_j), \tag{4.12}$$

where  $x_j$  is the *j*th component of  $x \in \mathbb{R}^d$ ,  $q_j$  is the momentum operator defined by  $(q_j f)(x) = x_j f(x)$ .

Set  $A_i^*$  to be the adjoint of  $A_i$  on the Schwartz space S, i. e.,

$$\langle A_i^*\psi, \gamma \rangle = \langle \psi, A_i\gamma \rangle$$
 for all  $\gamma, \psi \in \mathcal{S}$ .

For a multi-index  $\sigma = (\sigma_1, ..., \sigma_d)$ , set  $A^{\sigma} = A_1^{\sigma_1} ... A_d^{\sigma_d}$ , i.e., apply the Diracladder operator a different number of times in all dimensions.

Note that Hermite functions satisfy the following useful recurrence relations [15] for any  $n \in \mathbb{N}_{>0}$ 

$$\begin{cases} \gamma_{n+1} &= \frac{1}{\sqrt{n+1}} A^* \gamma_n \\ \gamma_{n-1} &= \frac{1}{\sqrt{n}} A \gamma_n. \end{cases}$$
(4.13)

The eigenfunctions that solve (2.5) belong to  $L^2(\mu)$ . However, it is more convenient to work with Schwartz functions. This is not a problem due to the density of Schwartz functions in  $L^2$  with respect to the  $L^2$  norm. I start the analysis by providing convergence guarantees for approximating Schwartz functions *via* augmented Hermite functions.

To characterize the conditions on  $\mathfrak{H}$  for such results one needs the following definition.

**Definition 4.5** (Symbols of the Schwartz space [175]). Let  $h \in \mathfrak{H}$  be smooth. *h* is called a *symbol for Schwartz space* if  $f \circ h \in S$  for all  $f \in S$ .

**Lemma 4.2.** Suppose that  $\mathfrak{H}$  satisfy the hypothesis of Theorem 4.1. Let r be as in Remark 4.2. Furthermore, assume that  $h^{-1}$  is smooth and that it is a symbol for S for all  $h \in \mathfrak{H}$ . For every fixed non-negative integer  $s \leq N$ , for all  $f \in S$  and for all  $h \in \mathfrak{H}$ , it holds that

$$||f - \tilde{P}_{\mathcal{N}}f|| \le (1+s)^{sd/2} \frac{1}{r^d} N^{-s/2} \max_{|\sigma|_{\infty} \le s} ||A^{\sigma}C_{h^{-1}}f||,$$

where the maximum is taken over all  $\sigma = (\sigma_1, ..., \sigma_d)$  with  $0 \le \sigma_j \le s$  for each j.

*Proof.* The proof extends Lubich's result on approximating Schwartz functions by Hermite functions [15]. Fix an arbitrary  $h \in \mathfrak{H}$ . Note that  $f \in L^2(h_{\#}^{-1}\mu)$  for any  $f \in L^2(\mu)$ . For every multi-index  $n = (n_1, \ldots, n_d)$  define the multi-index  $\sigma(n)$ by the condition  $\sigma(n)_j = n_j - (n_j - s)_+$  (with  $a_+ = \max\{a, 0\}$ ) fo all  $j = 1, \ldots, d$ . Note that

$$\begin{split} \tilde{P}_{\mathcal{N}}^{\perp}f &= \sum_{n \notin \mathcal{N}} \langle f, \tilde{\gamma}_{n}^{h} \rangle_{h_{\pi}^{-1} \mu} \tilde{\gamma}_{n}^{h} \\ &= \sum_{n \notin \mathcal{N}} \langle f \circ h^{-1}, \gamma_{n} \rangle \tilde{\gamma}_{n}^{h} \\ &= \sum_{n \notin \mathcal{N}} a_{n,s} \langle f \circ h^{-1}, (A^{*})^{\sigma(n)} \gamma_{n-\sigma(n)} \rangle \tilde{\gamma}_{n}^{h} \\ &= \sum_{n \notin \mathcal{N}} a_{n,s} \langle A^{\sigma(n)} f \circ h^{-1}, \gamma_{n-\sigma(n)} \rangle \tilde{\gamma}_{n}^{h}, \end{split}$$

where in the first equality I used result (i) from Theorem 4.1, in the third equality I used the recurrence relations (4.13), and in the fourth equality I used the assumption that  $h^{-1}$  is a symbol for S. One has

$$a_{n,s} = \prod_{j=1}^{d} \frac{1}{\sqrt{(1 + (n_j - 1)_+) \dots (1 + (n_j - s)_+)}}}$$

Note that, for  $n \notin \mathcal{N}$ , and  $b = 1, \ldots, s$ 

$$\prod_{j=1}^{d} (1 + (n_j - b)_+) = \prod_{j=1}^{d} (1 + (n_j - b)_+) \frac{1 + n_j}{1 + n_j}$$

$$> N \prod_{j=1}^{d} \frac{1 + (n_j - b)_{-j}}{1 + n_j}$$
  
$$> N \prod_{j=1}^{d} \frac{1}{1 + n_j}$$
  
$$> N(1 + b)^{-d}$$
  
$$> N(1 + s)^{-d}.$$

Hence,

$$|a_{n,s}|^2 \le N^{-s}(1+s)^{sd}$$
.

Taking inner product in  $L^2(h_{\#}^{-1}\mu)$ 

$$\begin{split} \|\tilde{P}_{\mathcal{N}}^{\perp}f\|_{L^{2}(h_{\#}^{-1}\mu)}^{2} &= \sum_{n \notin \mathcal{N}} |a_{n,s}|^{2} |\langle A^{\sigma(n)}f \circ h^{-1}, \gamma_{n-\sigma(n)} \rangle|^{2} \\ &\leq N^{-s}(1+s)^{sd} \sum_{n \notin \mathcal{N}} |\langle A^{\sigma(n)}f \circ h^{-1}, \gamma_{n-\sigma(n)} \rangle|^{2} \\ &\leq N^{-s}(1+s)^{sd} \max_{|\sigma|_{\infty} \leq s} \|A^{\sigma}f \circ h^{-1}\|^{2}. \end{split}$$

Since  $|\det Dh| \ge r^d$  almost everywhere one has that  $r^d \|\tilde{P}_N^{\perp} f\|_{L^2(\mu)}^2 \le \|\tilde{P}_N^{\perp} f\|_{L^2(h_{\#}^{-1}\mu)}^2$  (see Lemma 4.1).

Note that  $||P_{\mathcal{N}}^{\perp}f|| \leq c(r, R) ||\tilde{P}_{\mathcal{N}}^{\perp}f||$  assuming an upper bound on the  $L^2$  norm of  $D^2h^{-1}$ . Similar result can be proved for approximating eigenvalues of Schrödinger operators. One needs the following technical lemma.

**Lemma 4.3.** Let  $\mathfrak{H}$  satisfy hypothesis of Theorem 4.1 with r, R as in Remark 4.2. Then  $\psi(\det Dh^{-1})^{1/2} \in S$  for all  $\psi \in S$  and  $h \in \mathfrak{H}$ .

*Proof.* Provided in the supplementary materials at the end of this chapter.

I impose the following constraints on the Hamiltonian.

**Assumption 4.1.** Assume that the Hamiltonian is compact and that its spectrum is of cardinality  $K < \infty$ . Further, assume that the true eigenvalues  $\{E_k\}_{k=1}^K$  and the approximate ones  $\{\tilde{E}_k\}_{k=1}^K$  are ordered, i. e.,  $E_1 < E_2 < \ldots$  and  $\tilde{E}_1 < \tilde{E}_2 < \ldots$ 

Considering the augmented basis defined above, the following theorem provides convergence guarantees for approximating eigenvalues of linear Hamiltonians in quantum mechanics using Hermite functions that are augmented by a class of normalizing flows  $\mathfrak{H}$ .

**Theorem 4.2.** Let  $\mathfrak{H}$  satisfy the hypotheses of Lemma 4.2 and assume  $\mathfrak{H}$  is a symbol for S. Under Assumption 4.1 let

$$f^* = \operatorname{argmax}_{\|f\|=1, f \in D(H)} \|P_{\mathcal{N}}^{\perp} H f\|.$$

For an arbitrarily small  $\epsilon > 0$  let

$$\tilde{f} \in S$$
 be such that  $\|Hf^* - \tilde{f}\| \leq \epsilon$ .

*Then, for any non-negative integer*  $s \leq N$  *and all*  $h \in \mathfrak{H}$  *one has* 

$$|\tilde{E}_k - E_k| \le \sqrt{2}N^{-s} \frac{1}{r^d} (1+s)^{sd} \max_{|\sigma|_\infty \le s} ||A^{\sigma}C_{h^{-1}}\tilde{f}|| + \mathcal{O}(\epsilon)$$

$$(4.14)$$

for all k = 1, ..., K, where the right-hand side goes to zero as N goes to infinity, and  $\sigma$  is as in Lemma 4.2.

*Proof.* The proof extends the analysis of spectral methods for eigenvalue problems provided in [176] to augmented spectral methods. Fix an arbitrary  $h \in \mathfrak{H}$ . Assume that the operator  $P_{\mathcal{N}}HP_{\mathcal{N}}$  is self-adjoint. Under Assumption 4.1 and by Theorem A.4 one has

$$|\hat{E}_k - E_k| \le ||H - P_{\mathcal{N}}HP_{\mathcal{N}}||$$
  
=  $\sup_{||u||=1, u \in D(H)} ||(H - P_{\mathcal{N}}HP_{\mathcal{N}})u||$ .

Define  $v := (H - P_N H P_N)u$ . One has  $||v|| = ||P_N v + P_N^{\perp}v||$  and  $||v||^2 \le ||P_N v||^2 + ||P_N^{\perp}v||^2$ . Plugging in the definition of v

$$\begin{split} \|P_{\mathcal{N}}v\| &= \|(P_{\mathcal{N}}H - P_{\mathcal{N}}HP_{\mathcal{N}})u\| \\ &= \|P_{\mathcal{N}}HP_{\mathcal{N}}^{\perp}u\| \\ &\leq \|P_{\mathcal{N}}\|\|HP_{\mathcal{N}}^{\perp}\|\|u\| \\ &\leq \|HP_{\mathcal{N}}^{\perp}\|\|u\|, \end{split}$$

since  $||P_{\mathcal{N}}|| = \sup_{u \in D(H)} \frac{\langle P_{\mathcal{N}}u, P_{\mathcal{N}}u \rangle}{\langle u, u \rangle} = \sup_{u \in D(H)} \frac{\sum_{i < N} |c_i|^2}{\sum_i |c_i|^2} \le 1$ . For the second term one has

$$\begin{split} \|P_{\mathcal{N}}^{\perp}v\| &= \|(P_{\mathcal{N}}^{\perp}H - \underbrace{P_{\mathcal{N}}^{\perp}P_{\mathcal{N}}HP_{\mathcal{N}}}_{=0})u\| \\ &= \|P_{\mathcal{N}}^{\perp}Hu\| \\ &\leq \|P_{\mathcal{N}}^{\perp}H\|\|u\| \\ &\leq \|P_{\mathcal{N}}^{\perp}H\|\|u\| \\ &= \|HP_{\mathcal{N}}^{\perp}\|\|u\|. \end{split}$$

Noting that  $P_N^{\perp}H = HP_N^{\perp*}$  and that adjoint operators have the same operator norm one has

$$\begin{split} |\tilde{E}_k - E_k| &\leq \sqrt{2} \|P_{\mathcal{N}}^{\perp} H\| \\ &= \sqrt{2} \sup_{\|f\| = 1, f \in D(H)} \|P_{\mathcal{N}}^{\perp} H f\| \\ &= \sqrt{2} \|P_{\mathcal{N}}^{\perp} H f^*\|. \end{split}$$

By the density of S in  $L^2$  it holds that there exists  $\tilde{f} \in S$  such that  $||Hf^* - \tilde{f}|| \le \epsilon$ . Therefore,

$$\begin{split} |\tilde{E}_k - E_k| &\leq \sqrt{2} \Big( \|P_{\mathcal{N}}^{\perp} (Hf^* - \tilde{f})\| + \|P_{\mathcal{N}}^{\perp} \tilde{f}\| \Big) \\ &\leq \sqrt{2} \Big( \|P_{\mathcal{N}}^{\perp}\| \epsilon + \|P_{\mathcal{N}}^{\perp} \tilde{f}\| \Big) \end{split}$$

 $\leq \mathcal{O}(\epsilon) + \sqrt{2} \| P_{\mathcal{N}}^{\perp} \tilde{f} \|.$ 

74

Using Lemma 4.2 for approximating  $\tilde{f}$  (4.14) is attained. It remains to check that  $P_{\mathcal{N}}HP_{\mathcal{N}}$  is self-adjoint where H = T + V. Clearly,  $P_{\mathcal{N}}vP_{\mathcal{N}}$  is self-adjoint and  $P_{\mathcal{N}}TP_{\mathcal{N}}$  is self-adjoint if T is self-adjoint, which is correct on Schwartz-functions. Since  $\gamma_n \in S$  for all  $n \in \mathbb{N}_{>0}$  and h is a symbol for S then  $\tilde{\gamma}_n^h \in S$  and  $\tilde{\gamma}_n^h(\det Dh)^{1/2} \in S$  by Lemma 4.3.

These results are, indeed, generalizations of spectral methods.

**Remark 4.5.** Using a normalizing flow h = id the augmented basis  $(\gamma_n^h)_n$  is exactly Hermite functions. In this case, Lemma 4.2 and Theorem 4.2 are convergence guarantees of using spectral methods for approximating Schwartz functions, and eigenfunctions to Schrödinger equations, respectively. Hence, Lemma 4.2 and Theorem 4.2 can be regarded as an extension of the numerical analysis of spectral methods to augmented spectral methods, i. e., spectral learning.

The upper bounds in Lemma 4.2 and Theorem 4.2 explain some limitations of spectral methods and spectral learning.

**Remark 4.6.** [On difficulty of approximating high-dimensional and highlyoscillatory functions] The upper bounds in Lemma 4.2 and Theorem 4.2 depend nonlinearly on the dimension of the problem *d*. This is a major drawback of linear approximation concepts. Moreover, it can be shown that the quantity  $||A^{\sigma}C_{h^{-1}}f||$  increases for an increasing oscillatory behavior of the target function *f*. This means that the convergence rate degrades at increasing oscillation behavior.

However, the derived upper bounds depend on a hypothesis class *via*  $||A^{\sigma}C_{h^{-1}}f||$ . In the next subsection I discuss the possibility of minimizing the upper bound and its scaling with *d* by optimization over  $\mathfrak{H}$ .

I conclude the linear analysis by showing that invertible ResNets satisfy the assumptions of Lemma 4.2, and hence, maintain the convergence property as *N* 

grows to infinity. I have already shown that invertible residual networks satisfy hypothesis of Theorem 4.1. It remains to check that invertible residual networks are symbols for S. The following theorem characterizes symbols. *d* is set to one for simplicity.

**Theorem 4.3** (Charachterization of symbols [175]). A function  $\phi \in C^{\infty}(\mathbb{R})$  is a symbol for  $S(\mathbb{R})$  if and only if the following conditions are satisfied

• For all  $j \in \mathbb{N}_{>0}$  there exist C, p > 0 such that, for every  $x \in \mathbb{R}$ 

$$|\phi^{(j)}(x)| \le C(1+\phi(x)^2)^p.$$

• There exists k > 0 such that  $|\phi(x)| \ge |x|^{1/k}$  for all  $|x| \ge k$ .

**Assumption 4.2.** Set  $\mathfrak{H}$  to be the set of invertible ResNets of the form (4.4) and fix the number of their layers to one. Fix an arbitrary  $h \in \mathfrak{H}$ . Write  $h^{-1} = x - k$  and assume that k is a decreasing function of x.

Without loss of generality I set  $h^{-1}(0) = 0$ .

**Theorem 4.4.**  $h^{-1}$  is a symbol for S.

*Proof.* Provided in the supplementary material at the end of the chapter.  $\Box$ 

The next section deals with the third and final question posed in the introduction of this chapter.

#### 4.2.2 Faster convergence rates *via* normalizing flows

So far I proposed to compose standard bases  $(\gamma_n)_n$  of  $L^2(\mu)$  with a class of mappings  $\mathfrak{H}$  to augment their expressivity. I showed that, under some conditions on  $\mathfrak{H}$ , sequences in the induced family  $\{(\gamma_n \circ h \mid \det Dh \mid 1/2)_n \mid h \in \mathfrak{H}\}$  form bases of  $L^2(\mu)$ . Later, I proposed to use such bases in the framework of spectral methods to solve approximation problems. The density property of this family allowed one to derive convergence guarantees as the size of the truncated basis *N* goes to infinity. However, such convergence guarantees are already achievable for the

76

special normalizing flow h = id. The question that arises now is whether there exist special functions h such that a faster linear convergence can be achieved using augmented bases than a standard basis. Example 4.1 indeed hints at such a possibility.

Next, I show that this is achievable as well for approximating Schwartz functions and eigenvalues of Schrödinger operators. In particular, I show that the upper bounds in Lemma 4.2 and Theorem 4.2 admit minima over the hypothesis class  $\mathfrak{H}$ , and that a minimizer  $h^* \neq I$ . In other words, a faster linear convergence is achievable *via* the use of a truncated basis  $(\gamma_n^{h^*})_{n \leq N}$  compared to the use of  $(\gamma_n)_{n \leq N}$ .

Instead of considering an optimization problem over the derived upper bound that depends on h via  $||A^{\sigma}C_{h^{-1}f}||$ , I consider a simpler and more informative upper bound. Without loss of generality set  $s = 1, \sigma = (1, 0, ..., 0)$  and consider  $f \in C_c^{\infty 6}$ . By considering approximations to compactly supported functions instead of Schwartz functions one needs no longer assume that  $\mathfrak{H}^{-1}$  is a symbol for S to prove Lemma 4.2. Indeed, for all  $f \in C_c^{\infty}$  one has that  $C_{h^{-1}}f \in C_c^{\infty}$  for all  $h \in \mathfrak{H}$  due to the non-singularity of these mappings. Without loss of generality, let  $\operatorname{supp}(f) = U, U \in \mathcal{B}(\Omega)$  with  $\mu(U) = 1$ . Assume that for all  $x \in U$  there exists c > 0 such that  $||Dh^{-1}(x)||_{\mathcal{F}} < c$ . Under such settings, for every  $h \in \mathfrak{H}$  it holds

$$\begin{split} \|AC_{h^{-1}}f\|_{L^{2}(\mu)}^{2} &= \int_{\Omega} ((|x|^{2}+1)C_{h^{-1}}f)^{2} d\mu + \int_{\Omega} |\frac{d}{dx_{1}}C_{h^{-1}}f|^{2} d\mu \\ &\leq \int_{\Omega} ((|x|^{2}+1)C_{h^{-1}}f)^{2} d\mu + \int_{\Omega} |DC_{h^{-1}}f|^{2} d\mu \\ &\leq C_{1} \int_{U} |DC_{h^{-1}}f|^{2} d\mu \\ &= C_{1} \int_{U} |C_{h^{-1}}Df \cdot Dh^{-1}|^{2} d\mu \\ &= C_{1} \int_{U} |Df|^{2} \|Dh^{-1}\|_{\mathcal{F}} dh_{\#}^{-1}\mu \\ &\leq C_{2} \int_{U} |Df|^{2} dh_{\#}^{-1}\mu \end{split}$$

<sup>6</sup>Since  $C_c^{\infty}$  functions are dense in S with respect to the norm  $||f||_{a,b} = \sup_{x \in \mathbb{R}^d} |x^a D^b f|$ , for all  $a, b \in \mathbb{N}_{\geq 0}^d$ .

$$\leq C_2 \|Df\|_{L^{\infty}(\mu)} \int_{U} |Df| dh_{\#}^{-1}\mu,$$

where I used Poincare's inequality.

I consider an optimization problem over this upper bound, i. e.,

$$\inf\left\{\underbrace{\int_{U}|Df|\,dh_{\#}^{-1}\mu}_{l(h):=}\mid h\in\mathfrak{H}\right\},\tag{4.15}$$

where  $\mathfrak{H}$  satisfies the hypothesis of Theorem 4.1 with *r*, *R* as in Remark 4.2, and *l* stands for a loss function.

**Remark 4.7.** [On relationship of the optimization problem to total variation] The main motivation of considering an optimization problem over  $\mathfrak{H}$  for a looser upper bound than the ones derived in Lemma 4.2 and Theorem 4.2 is, indeed, simplicity. Proving the existence of minimizers is easier in this setting since the optimization problem is over induced measures  $h_{\#}^{-1}\mu$ . Under such setting there is a variety of similar useful results, especially from the field of optimal transportation [177]. This upper bound comes also with the advantage of more interpretability. Note that the loss function l(h) is actually the total variation of a function f with respect to the induced measure  $h_{\#}^{-1}\mu$ . The total variation of a function serves as an appropriate measure of oscillatory behavior. Hence, this looser upper bound explains the limited capabilities of approximating functions of large l(h), i. e., of large total variation. Specifically, a larger oscillatory behavior means a larger total variation, which translates into a slower convergence of the underlying truncated basis.

The optimization problem in (4.15) can be understood as looking for an induced measure, with respect to which f has a small total variation.

Consider  $(U, \mathcal{B}(U), \mu)$  as an abstract measure space, where  $\mu$  is the restriction of the Lebesgue measure on U, and denote by  $\mathcal{M}(U)$  the space of Radon measures on U. Endow  $\mathcal{M}(U)$  with the total variation norm  $\|\lambda\|_{\mathcal{M}} = |\mu|(U)$ . Note that, under the assumption  $\mu(U) = 1$ , the induced measures  $h_{\#}^{-1}\mu$  for all  $h \in \mathfrak{H}$  are

actually probability measures. Denote by  $\mathbb{P}(U)$  the set of probability measures on U.

The main result in this section requires the introduction of some lemmas.

**Lemma 4.4** (Pushforward measures are Radon measures). *The induced measure*  $h_{\#}^{-1}\mu$  for all  $h \in \mathfrak{H}$  where  $\mathfrak{H}$  is the hypothesis class in the optimization problem (4.15) is a *Radon measure*.

*Proof.* Note that  $h_{\#}^{-1}\mu(U) = \mu(h(U))$ . Since *h* is measurable, one has that  $h(U) \in \mathcal{B}$ . Since  $\mu$  is inner-regular it follows that  $h_{\#}^{-1}\mu$  is inner-regular as well.

Now take  $x \in U$  and a compact neighborhood of it,  $N_x$  and note that the image of  $N_x$  under *h* is compact due to the Lipschitz continuity of *h*. One has

$$h_{\#}^{-1}\mu(N_x) = \mu(h(N_x)) < \infty,$$

since  $\mu$  is locally finite. Hence,  $h_{\#}^{-1}\mu$  is Radon (see Definition D.7).

The following result shows the existence on invertible mappings that induce certain measures on  $(U, \mathcal{B}(U))$ . It follows as a direct corollary from a more abstract result (Theorem D.4).

**Lemma 4.5** (Existence of mappings that induce certain measures<sup>7</sup>). Let v be a measure on  $(U, \mathcal{B}(U))$  such that  $v(U) = \mu(U), v(\{x\}) = 0$  for all  $x \in U$ . There exists a bi-measurable mapping  $T : U \to U$  such that  $T_{\#}\mu = v$ .

*Proof.* There exists a bi-measurable map  $T_1 : U \to [0,1]$  such that  $T_{1\#}\mu = \mu|_{[0,1]}$  (see Theorem D.4). Similarly, under the hypothesis on  $\nu$  there exists a bi-measurable  $T_2 : U \to [0,1]$  such that  $T_{2\#}\nu = \mu_{[0,1]}$ . Set  $T : U \to U$ ,  $T = T_2^{-1} \circ T_1$ . One has  $T_{\#}\mu = \nu$ .

Note that  $\mathcal{M}(U)$  is the topological dual of C(U). The main result depends on the notion of weak-\* convergence of a sequence  $(\lambda_n)_n \in \mathcal{M}(U)$ .

<sup>&</sup>lt;sup>7</sup>I would like to acknowledge the input of Alp Uzman [178] and Oliver Díaz [179] herewith on math stackexchange [180].

**Definition 4.6** (Weak-\* convergence). A sequence  $(\lambda_n)_n \in \mathcal{M}(U)$  converges \*-weakly to  $\lambda \in \mathcal{M}(U)$  if

$$\lim_{n \to \infty} \int_U f \, d\lambda_n = \int_U f \, d\lambda$$

for any  $f \in C(U)$ .

The following two lemmas characterize important results of a \*-weakly converging sequence in  $\mathcal{M}(U)$ .

**Lemma 4.6** ([177]). Let  $f : U \to \mathbb{R} \cup \{\infty\}$  be a lower semi-continuous function. Define  $F : \mathbb{P}(U) \to \mathbb{R} \cup \{\infty\}$ ,  $F(\mu) = \int_U f d\mu$ . Then, F is lower-semi continuous for the \*-convergence in  $\mathcal{M}(U)$ , i. e., for  $(\lambda_n)_n, \lambda \in \mathcal{M}(U)$  such that  $\lambda_n$  converges \*-weakly to  $\lambda$  it holds

$$F(\lambda) \leq \liminf_{n \to \infty} F(\lambda_n).$$

Similarly, if f is upper semi-continuous, it holds

$$F(\lambda) \geq \limsup_{n \to \infty} F(\lambda_n).$$

**Lemma 4.7.** Let  $(\lambda_n)_n \in \mathcal{M}(U)$  be a sequence of finite measures on  $(U, \mathcal{B}(U))$  that converges \*-weakly to a  $\lambda \in \mathcal{M}(U)$ . Assume that  $(\lambda_n)_n$  is absolutely continuous with respect to the Lebesgue measure  $\mu$  with the Radon-Nikodym derivatives satisfying  $\frac{d\lambda_n}{d\mu} \leq L < \infty \mu$ -almost everywhere. Then  $\lambda$  is absolutely continuous with respect to  $\mu$ .

*Proof* [181]. Choose  $k \in C^{\infty}(U)$  and define

$$\gamma(k) := \int_U k \, d\lambda.$$

Note that

$$\begin{aligned} |\gamma| &\leq \limsup_{n \to \infty} \int_{U} k \frac{\partial \lambda_n}{\partial \mu} \ d\mu \\ &\leq M \|k\|_{L^2(\mu)}^2. \end{aligned}$$

Thus,  $\gamma$  can be extended to a linear operator on  $L^2(U)$ . By Riesz representation theorem there exists  $f \in L^2(U)$  such that

$$\gamma(g) = \int_U fg \, d\mu.$$

Therefore,  $\lambda$  is absolutely continuous with respect to  $\mu$  with Radon-Nikodym derivative *f*.

The stage is ready for the main result.

**Theorem 4.5** (Minimizer of (4.15)). *The optimization problem defined in* (4.15) *admits a minimizer.* 

*Proof.* Note that l(h) is bounded from below. Let  $(h_n)_n \in \mathfrak{H}$  be a minimizing sequence, i.e.,

$$\lim_{n\to\infty} l(h_n) = E_0 \text{ with } E_0 = \inf \Big\{ l(h) \mid h \in \mathfrak{H} \Big\}.$$

By Lemma 4.4,  $h_{\#}^{-1}\mu$  is a Radon measure. One has that

$$\begin{split} \|h_n^{-1} \mu\|_{\mathcal{M}(U)} &= |\int_U dh_{\#}^{-1} \mu/d\mu \, d\mu| \\ &= \|dh_{\#}^{-1} \mu/d\mu\|_{L^{\infty}(\mu)} |\mu|(U) \\ &< \infty, \end{split}$$

i. e.,  $(h_n^{-1} \mu)_n$  is bounded in  $\mathcal{M}$ . Since  $\mathcal{M}$  is the topological dual of C(U) which is a separable space, by the Banach-Alaoglu theorem (Theorem E.2) there exists a subsequence of induced measures  $(h_{nj}^{-1} \mu)_j$  that converges \*-weakly to  $\alpha^* \in \mathcal{M}(U)$ . By Lemma 4.7,  $\alpha^*$  is absolutely continuous and by Lemma 4.6 it holds

$$\int_{U} |Df| \, d\alpha^* \leq \liminf_{j \to \infty} \int_{U} |Df| \, dh_{n_j}^{-1} \# \mu$$
$$\leq R^d \int_{U} |Df| \, d\mu.$$

Thus,  $\frac{\partial \alpha^*}{\partial \mu} \leq R^d$ . Similarly,  $\frac{\partial \alpha^*}{\partial \mu} \geq r^d$ . It remains to show that there exists  $h^* \in \mathfrak{H}$  such that  $h_{\#}^{-1^*} \mu = \alpha^*$ . Since  $\alpha^*$  is absolutely continuous with respect to  $\mu$  one has that  $\alpha^*(\{x\}) = 0$  for all  $x \in U$ . Note that, by Lemma 4.6

$$\begin{aligned} \alpha^*(U) &= \int_U d\alpha^* \\ &\leq \int_U \liminf_{j \to \infty} dh_{nj_{\#}}^{-1} \mu \\ &= 1. \end{aligned}$$

Similarly,  $\alpha^*(U) \ge 1$ . Hence,  $\alpha^*(U) = 1$ . It follows by Lemma 4.5 that there exists a bi-measurable  $h^* : U \to U$  such that  $h_{\#}^{-1^*} \mu = \alpha^*$ . Moreover,  $h^* \in \mathfrak{H}$ .

Choosing r > 1 in Lemma 4.2 one can conclude that a minimizer is not equal to identity.

This completes the analysis defined by the three questions in the introduction of this chapter. Next, I report simulations on computing nuclear spectra of polyatomic molecules, where results confirm the theoretical analysis.

## 4.3 Computing the spectra of polyatomic molecules

Here, I report simulations to approximate solutions of the vibrational Schrödinger equation (2.5) for hydrogen sulfide  $H_2S$  in a Bubnov-Galerkin framework (4.9). The aim is to study convergence patterns of spectral learning (Definition 4.3) and compare it to convergence of standard spectral methods.

The problem under consideration is three-dimensional. The coordinates are defined by the internal vibrational degrees of freedom  $(r_1, r_2, \theta)$  (see Appendix D) of H<sub>2</sub>S, where  $r_1, r_2 \in (0, \infty)$  denote the stretching coordinates, and  $\theta \in (0, \pi)$  denotes the bending coordinate.

For standard spectral methods Hermite functions (primitive basis) are used, and for spectral learning Hermite functions are augmented with an invertible ResNet (basis!augmented). Note here that invertible ResNets satisfy the assumptions of Theorem 4.1. However, the minimization problem (4.15) was considered

N	1	3	5	7	9	11	13	15	17	19	21	23
$ \mathcal{N} $	2	8	20	40	70	112	168	240	330	440	571	723

TABLE 4.1: Size of the index set  $\mathcal{N}$  for different values of the polyad number.

in a class of more general mappings  $\mathfrak{H}$ , i. e., mappings that are not necessarily differentiable. The 3-dimensional approximating functions were constructed from the one-dimensional Hermite functions (4.1) *via* the reduced tensor product

$$\mathcal{N}(N) = \{ (n_1, n_2, n_3) : n_j \ge 0 \text{ for } j = 1, 2, 3, w_1 n_1 + w_2 n_2 + w_3 n_3 \le N \}, \quad (4.16)$$

where N = 1, 2, ..., 18,  $w_1, w_2$  correspond to the two stretching modes and  $w_3$  corresponds to the bending mode. The aim here is to study the approximate energies for an increasing N, i. e., an increasing size of linear expansion. I chose  $w_1 = w_2 = 2, w_3 = 1$ . This choice is suitable since the density of the bending states is two times larger than that of the stretching states. I refer to N by *polyad number*. Table 4.1 shows the size  $|\mathcal{N}(N)|$  of the set  $\mathcal{N}(N)$  for different values of N. Note that the number of approximate eigenvalues for each N is equal to  $|\mathcal{N}(N)|$ . Define the set of approximate eigenvalues

$$\mathcal{E}(N) = \{\tilde{E}_n\}_{n \in \mathcal{N}(N)}, N = 1, 2, \dots, 18$$

To find the coefficients  $\tilde{C}_n$  and the approximate eigenvalues  $\tilde{E}_n$  of (4.9) I used a direct eigensolver. For augmented Hermite functions, the coefficients and the approximate eigenvalues depend on the parameters of the neural network. Following the variational formulation derived for Schrödinger equation in Theorem 2.2 one has that

$$\sum_{n\in\mathcal{N}(N)}\tilde{E}_n\geq\sum_{n\in\mathcal{N}(N)}E_n,$$

i. e., the sum of the approximate eigenvalues is always bigger than that of the true eigenvalues. I defined the sum of the approximate eigenvalues as a loss function

and used a first order optimizer to optimize the parameters of the normalizing flow.

To compute the matrix elements of the projected Hamiltonian  $\hat{H}$ , I used Gauss-Hermite quadrature.

The eigenvalues were calculated in the units of inverse centimeters. For spectroscopic applications, a high accuracy correspond to  $< 1 \text{ cm}^{-1}$  error.

More extensive information on the architecture of the neural network, training procedure and numerical integrations are provided in the supplementary material at the end of this chapter.

#### 4.3.1 Convergence of the numerical schemes

Here, the convergence of the approximate eigenvalues, i. e., approximate vibrational energies, that solve (4.9) as a function of the truncation parameter N for the two discretization schemes is reported. Since I consider many eigenvalues, it is more convenient to consider the convergence of energy bands. While there are many ways to group the energies, grouping them according to their polyad number is the one often used in physics literature.

To formally introduce this grouping define subsets of the set of indices  $\mathcal{N}(N)$  as follows

$$\mathcal{N}_i(N) := \{ n \in \mathcal{N} : n = (n_1, n_2, n_3), w_1 n_1 + w_2 n_2 + w_3 n_3 = i \}$$

for  $i = 1, ..., N_i$ . Note that  $\mathcal{N}(N) = \bigcup_{i=1,...,N} \mathcal{N}_i(N)$ . Similarly, define the energy subsets

$$\mathcal{E}_i(N) := \{\tilde{E}_n\}_{n \in \mathcal{N}_i(N)}.$$

I denote by  $\tilde{\mathcal{E}}_i(N)$  the average energy over the band  $\mathcal{E}_i(N)$ , i. e.,

$$ilde{\mathcal{E}}_i(N) = rac{\sum_{E \in \mathcal{E}_i(N)} E}{|\mathcal{E}_i(N)|}.$$

To characterize the convergence of both discretization schemes I studied the quantities

$$\Delta \tilde{\mathcal{E}}_i(N) := \tilde{\mathcal{E}}_i(N) - \tilde{\mathcal{E}}_i(N+1), \ i = 1, \dots, M \le N, N = 1, \dots, 17.$$

A decreasing  $\Delta \tilde{\mathcal{E}}_i(N)$  means that the eigenvalues are converging. I studied this quantity for M = 9. In other simpler words, I studied the convergence of the average (over group) energy of the first nine groups of energies as a function of the truncation parameter *N*.

Figure 4.2 shows  $\Delta \tilde{\mathcal{E}}_i$  for the primitive and augmented basis. For both schemes  $\Delta \tilde{\mathcal{E}}_i$  decreases as a function of the polyad number *N*. However, the this quantity converges faster as a function of *N* for the augmented basis and increasing the polyad number results in an improvement of less than 1 cm<sup>-1</sup>, the usually required accuracy in my field of application.

I concluded that the 40 eigenvalues corresponding to polyad numbers up to 7 are converged for the augmented basis. I note that the absolute values for these eigenvalues are lower than that of the primitive basis for the same polyad number. Hence, I consider the approximate eigenvalues obtained using the augmented basis to be the true reference values.

### 4.3.2 Error quantification

I compared the lowest 40 converged eigenvalues with reference calculations using symmetry adapted vibrational basis functions (TROVE) [149]. I note that the Hamiltonian in TROVE calculations includes a pseudopotential that I excluded from the calculations due to the high expenses of its computation and its relatively small contribution to the overall Hamiltonian.

Table 4.2 shows the absolute error between the approximate eigenvalues<sup>8</sup> and the reference ones. All energies differ from the reference ones by less than 1 cm<sup>-1</sup>. This indicates a very good agreement with reference simulations.

<sup>&</sup>lt;sup>8</sup>To be accurate I compared relative eigenvalues, i. e., I subtracted the smallest eigenvalue from all eigenvalues before comparing to reference simulations. This is a common procedure in physics literature.



FIGURE 4.2: Convergence of the first 9 energy bands for two discretization schemes, Hermite functions (primitive basis) and augmented Hermite functions (augmented basis).

Next, I compared the absolute error in the approximate eigenvalue bands for both discretization schemes as a function of the polyad number. The correct energies here are considered to be those of augmented Hermite functions for a polyad number 18. Figure 4.3 shows these results. One observes that the approximate eigenvalues of the augmented basis converge fast as a function of N, while only the first band in the primitive basis reaches the desired error tolerance of  $< 1 \text{ cm}^{-1}$ . Moreover, approximation by Hermite functions is particularly poor for larger-eigenvalue bands, which correspond to more oscillatory functions. For augmented Hermite functions, the convergence of these energies is slower as well, but the desired accuracy can, nevertheless, be achieved. Moreover, the deterioration in approximation capabilities for higher energy bands is smaller than that encountered when using Hermite functions. This indicates that augmentation with invertible neural networks can improve approximations of highly-oscillatory functions.



FIGURE 4.3: The absolute error (in cm<sup>-1</sup>) in the computation of the the first 7 vibrational polyad bands (in different colors) of H<sub>2</sub>S using augmented Hermite functions (augmented basis) and Hermite functions (primitive basis).

### 4.3.3 Loss function

Figure 4.4 shows the convergence of the total loss function for polyad number N = 7, i. e., the sum of all the 40 eigenvalues at each nonlinear training iteration t, plotted for the two discretization schemes. The fast convergence of the loss function in the augmented scheme demonstrates the high quality of the inductive bias provided by Hermite functions. The smoothness of the convergence can be partially attributed to that fact that the sequence of augmented functions used for the approximation of energies is always a basis in the limit  $N \rightarrow \infty$  for any values of the parameters of the ResNet (Theorem 4.1).

## 4.3.4 Accessing of the approximate eigenfunctions

To assess the quality of the approximate eigenfunctions I computed another observable. A physical quantity of interest is the electric dipole moment  $\mu$  which takes



FIGURE 4.4: The training loss, i. e., the sum of the approximate eigenvalues as a function of training iteration t for the two discretization schemes.

values in  $\mathbb{R}^3$ . It can be calculated from electronic structure theory [23]. Denote by  $h^*$  the optimized neural network for the polyad number 18. The projection of the dipole moment on the approximate eigenfunctions is a matrix *D* whose *nm*th element is given by

$$D[n,m] = \tilde{C}_n \tilde{C}_m \langle \gamma_n^{h^*}, |\mu| \gamma_m^{h^*} \rangle, \qquad (4.17)$$

where  $\tilde{C}_n$ ,  $\tilde{C}_m$  are the coefficients of the approximate eigenfunctions. I computed this quantity for n = m for the first 35 converged eigenfunctions and compared against results from TROVE [149]. Table 4.3 summarizes this comparison. Note again that TROVE calculations include a pseudopotential that I ignore. Nevertheless, results show an agreement up to the second decimal number.

## 4.4 Summary and outlook

Motivated by the curse of dimensionality and the difficulty of approximating highly-oscillatory functions encountered in spectral methods, I discussed the

augmentation of standard bases of  $L^2(\mu)$  with a hypothesis class  $\mathfrak{H}$ . I derived sufficient conditions on  $\mathfrak{H}$  to guarantee that the resulting augmented sequences form bases of  $L^2(\mu)$  (Theorem 4.1). An interesting question here is whether one can derive weaker conditions to this end. Furthermore, I showed that invertible ResNets satisfy these assumptions.

These results allowed one to provide convergence guarantees for approximating Schwartz functions (Lemma 4.2) and eigenvalues of bounded operators (Theorem 4.2) in the linear span of truncated augmented bases as the truncation parameter grows to infinity. While the analysis here was limited to augmented Hermite functions, extension to other bases is straightforwardly achievable by noting the underlying recurrence relations of the bases. For example, Legendre polynomials in one dimension  $(P_n)_n$  satisfy the recurrence relation  $P_{n-1}(x) = BP_n(x)$  where the operator *B* reads  $B(f) = (q - \frac{q^2 - 1}{2} \frac{d}{dx})(f)$ . Similar to the recurrence relations of Hermite functions (4.13), this can be used to derive convergence guarantees for augmented Legendre functions upon noting the space on which  $B^*$  is defined and imposing correct conditions on  $\mathfrak{H}$  to be a symbol of this space. It is also possible to analyze augmented bases for a generic underlying basis assuming it satisfies some recurrence relations. Another restriction of the current analysis is that Theorem 4.2 is limited to bounded quantum Hamiltonian operators. While a convergence behavior is empirically observed also for unbounded Hamiltonians [168, 182], the proof technique used in Theorem 4.2 is not suitable to show that.

I derived a looser upper bound to the approximation errors in Lemma 4.2 and Theorem 4.2 in terms of the total variation of the target function f with respect to push-forward measures induced by  $\mathfrak{H}$ . This upper bound shows that both Hermite functions and augmented Hermite functions suffer from a slower convergence at an increasing oscillation of f. However, I showed that this upper bound admits a minimum over  $\mathfrak{H}$  (Theorem 4.5). This means that faster convergence can be achieved for augmented Hermite functions. One interesting research problem is to quantify the improvement in accuracy possible upon the use of normalizing flows and how this relates to their parameters. This can be possibly accomplished *via* constructive approximation theories of normalizing flows [183]. Moreover, further research is needed to see whether normalizing flows can directly decrease the scaling of the computational costs with the dimensionality of the problem.

I reported simulations to compute the vibrational spectrum of  $H_2S$ . The results agreed qualitatively with the theoretical insights. In particular, I demonstrated numerically the convergence of a numerical scheme based on augmented Hermite functions as a function of the truncation parameter. Moreover, I showed that this convergence is indeed faster than what is achievable *via* standard spectral methods. In particular, one gained a 2-order of magnitude improvement in accuracy. The increased accuracy is particularly important for approximating highly-oscillatory functions. Furthermore, the results agreed well with reference values computed by the TROVE variational method [149–152]. I note that high-dimensional quadrature rules should be further discussed in order to successfully apply spectral learning for bigger systems. In the performed calculations, I used Gauss quadratures. These are not suitable for higher dimensions because their size grows exponentially with the number of dimensions. A possible remedy may be to use sparse grid approaches, such as Smolyak grids [158]. Stochastic estimations of integrals, such as Monte-Carlo methods, may provide a dimension-independent scaling at the expense of lower accuracy. Another approach may be the use of collocation methods, which are equivalent to solving the Schrödinger equation by demanding that it is satisfied at a set of points, i.e., no integration is necessary [184].

Overall, the proposed spectral learning framework is theoretically sound, in the sense that it can be well-analyzed for solving differential equations. It improves the approximation capabilites of standard methods, especially for approximating highly-oscillatory functions. The training procedure in this framework is stable. These advantages make spectral learning a promising tool for investigating large scale computational problems, such as modelling ultrafast molecular dynamics [7–9].

# Supplementary material

The normalizing flow I used for training has the form

$$h^{-1}(x) = \tanh(f_{\gamma}(\tanh^{-1}(x-\beta)/\alpha)) * \alpha + \beta,$$

where *f* is an invertible residual neural network, whose parameters are denoted by  $\gamma$ . The neural network *f* is composed of 2 layers with 64 hidden units, with Lipswish activation functions, i. e., functions of the form

$$\sigma(x) = \frac{1}{1.1} \cdot \frac{x}{1 + \exp(-x)}.$$

Through a fixed scaling procedure, the input to the  $tanh^{-1}$  function is guaranteed to lie within [-1, 1]. To solve the optimization problem I use Adam algorithm [148]. To compute integrals in (4.9) I used Hermite Gauss quadratures. During training, I used different data batches that I constructed by varying the order per dimension of the quadrature in  $\{30, 25, 21, 26, 22, 29\}$ . I ran all calculations using 200 training epochs. The final eigenvalue calculations after the nonlinear training were performed using 70 quadrature points per dimension. The 3–dimensional quadrature grid was generated by taking a direct product on the 1*D* grids while eliminating points corresponding to a quadrature weight of  $< 10^{-34}$ . To ensure that I did not run into an overfitting problem, i. e., the weights of the normalizing flow are not sensitive to the training data, I compared the eigenvalues computed using 29 quadrature points per dimension to 70 quadrature points per dimension. The maximum difference per eigenvalue between the two integration schemes over the first 70 eigenvalues is 0.0196. Hence, I conclude the absence of overfitting.

*Proof of Lemma 4.1.* For all  $f \in L^2(h_{\#}\mu)$  one has

$$\begin{split} \|f\|_{L^{2}(h_{\#}\mu)}^{2} &= \int_{\Omega} |f|^{2} dh_{\#}\mu \\ &= \int_{\Omega} |f|^{2} |\det Dh^{-1}| d\mu \\ &\geq 1/R^{d} \int_{\Omega} |f|^{2} d\mu \\ &= 1/R^{d} \|f\|^{2}, \end{split}$$

i. e.,  $f \in L^2(\mu)$  and thus,  $L^2(h_{\#}\mu) \subseteq L^2(\mu)$ . Similarly, for all  $f \in L^2(\mu)$ 

$$\|f\|_{L^2(h_{\#}\mu)}^2 \le 1/r^d \int_{\Omega} |f|^2 d\mu$$

$$= 1/r^d ||f||^2$$

and thus,  $L^2(\mu) \subseteq L^2(h_{\#}\mu)$ . The same can be done to show that  $L^2(\mu) = L^2(h_{\#}^{-1}\mu)$ .

*Proof of Lemma 4.3.*  $\psi(\det Dh^{-1})^{1/2} \in S$  if there exists  $C_{n,\beta} > 0$  such that

$$\sup_{x\in\Omega}|x^{\beta}D^{\alpha}(\psi(\det Dh^{-1})^{1/2})|\leq C_{n,\beta},$$

where

$$D^{\alpha}(\psi(\det Dh^{-1})^{1/2}) = \sum_{\beta,\beta \leq lpha} {lpha eta eta}(D^{\beta}\psi)(D^{lpha-eta}(\det Dh^{-1})^{1/2}).$$

Let  $f(x) = \sqrt{x}$ , and  $d(x) = \det Dh^{-1}$  and let |v| = n. By Faà di Bruno formula [185] one has

$$D^{v}f \circ d = \sum_{j=1}^{n} f^{(j)} \circ d \sum_{s=1}^{n} \sum_{p_{s}(v,j)} (v!) \prod_{j=1}^{s} \frac{1}{(k_{j}!)[l_{j}!]^{k_{j}}} [D_{x}^{l_{j}}d],$$

where  $p_s(v, j) = \{(k_1, \ldots, k_s; l_1, \ldots, l_s) : k_i > 0, 0 \prec l_1 \prec \cdots \prec l_s. \sum_{i=1}^s k_i = j, \sum_{i=1}^s k_i l_i = v\}$  and where  $k_i$  is a scalar and  $l_i$  is a d-dimensional vector for all *i*. Note that  $f^{(l)} \circ d = (-1)^{(l-1)} \frac{(2l-3)!!}{2l} (d)^{0.5-l}$  (formula for *l*th derivative of square root) which is bounded since  $d \ge r$  everywhere. Thus, the whole derivative is bounded.

*Proof of Theorem* 4.4. Clearly,  $h^{-1}$  is smooth upon the use of smooth activations, e. g., sigmoid functions. Assume that the activation functions are sigmoid and are appropriately scaled to have a Lipschitz constant < 1. For a neural network with one layer one can see

$$|h^{-1}(x)| = |x - k(x)| \ge |x| - |k(x)| \ge h \text{ is bilipschitz} (1 - L)|x| = \frac{1}{2}|x| \ge |x|^{1/3}$$

for all  $|x| \ge 3$ , where I sat L = 1/2. It remains to check the first condition of Theorem 4.3. Note that  $\frac{d^n}{dx^n}h^{-1}(x) \le \frac{d^n}{dx^n}x + |\frac{d^n}{dx^n}\sigma(wx+b)|$  for any  $n \in \mathbb{N}$  and where  $\sigma$  is the sigmoid function. Using Faà di Bruno's formula where  $f = \sigma$  and g = wx + b one can write

$$\frac{d^n}{dx^n}\sigma(wx+b) = \sum_l \frac{n!}{m_1!\dots m_n!}\sigma(\overbrace{m_1+\dots+m_n}^{M_l})(wx+b)w^{m_1},$$

where

$$\sigma^{(M_l)} = \sum_{k=0}^{M_l} \sum_{j=0}^k (-1)^j (j+1)^{(M_l)} \binom{k}{j} \sigma^{k+1}$$
$$\leq \sum_{k=0}^{M_l} \sum_{j=0}^k (j+1)^{(M_l)} \binom{k}{j} 1 = c_l$$

since the sigmoid function is bounded and hence, all derivatives of the normalizing flow are bounded. Thus,

$$\begin{aligned} \frac{d^n}{dx^n} \sigma(wx+b) &\leq \sum_l \frac{n!}{m_1! \dots m_n!} c_l w^{m_1} \\ &\leq L^{m_1} \sum_l \frac{n!}{m_1! \dots m_n!} c_l = C \leq C(1+|h^{-1}(x)|^2). \end{aligned}$$

State	Trove ( $cm^{-1}$ )	Augmented Hermite functions ( $cm^{-1}$ )	Absolute error
1	1182.57	1182.69	0.12
2	2353.91	2354.16	0.25
3	2614.39	2614.27	0.12
4	2628.46	2628.33	0.14
5	3513.7	3514.1	0.4
6	3779.19	3779.22	0.03
7	3789.27	3789.26	0.01
8	4661.61	4662.16	0.56
9	4932.69	4932.84	0.15
10	4939.13	4939.24	0.11
11	5145.03	5144.87	0.17
12	5147.17	5146.94	0.23
13	5243.16	5242.94	0.22
14	5797.21	5797.97	0.76
15	6074.57	6074.99	0.43
16	6077.63	6077.91	0.28
17	6288.14	6288.34	0.2
18	6289.13	6289.18	0.05
19	6385.32	6385.4	0.08
20	6920.08	6921.04	0.96
21	7204.31	7204.74	0.43
22	7204.44	7205.0	0.57
23	7419.85	7420.09	0.24
24	7420.08	7420.21	0.13
25	7516.83	7517.0	0.18
26	7576.42	7576.46	0.04
27	7576.6	7576.6	0.0
28	7752.34	7752.74	0.4
29	7779.35	7779.07	0.28
30	8029.81	8031.1	1.29
31	8318.69	8319.44	0.76
32	8321.87	8322.99	1.13
33	8539.58	8540.43	0.86
34	8539.83	8540.79	0.96
35	8637.16	8638.02	0.85

TABLE 4.3: Absolute error in the dipole moment calculations for  $H_2S$  where the reference calculations are performed using TROVE [149]. Numbers are rounded to the 4th decimal number.

State	Trove ( $cm^{-1}$ )	Augmented Hermite functions ( $cm^{-1}$ )	Absolute error
1	0.9703	0.9745	0.0041
2	0.9751	0.9711	0.0041
3	0.98	0.98	0.0
4	0.9652	1.0034	0.0382
5	0.9674	1.001	0.0336
6	0.985	1.0212	0.0362
7	0.9695	0.9698	0.0004
8	0.9717	0.9801	0.0084
9	0.99	0.9543	0.0357
10	0.9738	0.9216	0.0523
11	0.9762	0.9538	0.0224
12	0.9596	0.937	0.0226
13	0.9602	0.9738	0.0135
14	0.9628	0.9869	0.0241
15	0.9952	0.9265	0.0687
16	0.9782	0.9683	0.01
17	0.9808	0.9366	0.0441
18	0.9633	0.9515	0.0118
19	0.9638	0.9623	0.0015
20	0.9665	0.9626	0.0039
21	1.0005	0.9568	0.0437
22	0.9855	0.9724	0.0132
23	0.9828	1.0537	0.0709
24	0.9671	0.9939	0.0268
25	0.9675	1.0441	0.0766
26	0.9702	0.997	0.0268
27	0.9496	0.9471	0.0025
28	0.9497	0.9319	0.0178
29	0.954	0.9042	0.0498
30	0.9587	0.9348	0.0238
31	1.0058	0.9452	0.0605
32	0.9904	0.976	0.0145
33	0.9875	1.0017	0.0142
34	0.9713	0.9385	0.0328
35	0.9711	0.9418	0.0293
#### Chapter 5

## **Conclusions and outlook**

Numerical modeling of quantum chemical dynamics poses computational challenges at the focus of research efforts in numerical analysis and computer science. For example, constructing potential energy surfaces of weakly-bound complexes, such as pyrrole( $H_2O$ ), requires the use of high-resolution numerical methods to solve the electronic Schrödinger equation. Furthermore, the landscape of these potential energy surfaces is complex due to the loosely bound character of intermolecular interactions. Thus, a large number of points is usually required to sample the complete configuration space. In addition, a correct description of the dissociation dynamics of such molecular systems requires the computation of a vast number of highly-oscillatory eigenfunctions of unbounded linear operators that lie in high dimensions. The present work proposed two novel machine learning algorithms that allow for more accurate quantum simulations of molecules than what is possible using state-of-the-art methods, while enjoying a high-level of robustness.

First, constructing potential energy surfaces of polyatomic molecules in a supervised learning paradigm was considered. Due to the high-computational costs of generating the training dataset, this problem has been under extensive recent investigations in an active learning paradigm that allows for optimizing the choice of the training dataset. However, an understanding of optimal strategies to minimize the size of the training dataset is still lacking, and common approaches are rather heuristic. To this end, I proposed an upper bound to the generalization error

for active learning (Theorem 3.2) which allows for an empirical risk minimization principle. It suggests that optimal strategies should sample points with high uncertainties in their predictions and such that their distribution does not deviate, in the sense of integral probability metrics, much from the true distribution of data. This result can be seen as a general formulation of similar results where distances between probability distributions is measured *via* the Wasserstein metric [105], or the maximum mean discrepancy over reproducing kernel Hilbert spaces [104]. I proposed an algorithm based on this empirical risk minimization principle (Algorithm 3) and used it to reduce the computational costs of constructing the potential energy surface of pyrrole $(H_2O)$ . Simulations [121] (also reported in Section 3.4) showed that accurate potential energy surfaces can be constructed with a roughly two times smaller dataset than what is possible *via* other common active learning algorithms. The proposed algorithm is general, can be applied to any molecular system, and can be combined with quantum chemistry packages to solve the electronic Schrödinger equation. For a thorougher technical discussion into the down-, upsides of the proposed framework and its prospective see Section 3.5.

Second, the present work considered solving static Schrödinger equations that describe nuclear motions in molecules. To this end, it proposed augmenting the expressivity of standard bases of  $L^2$  via composition with invertible neural networks. The work identified sufficient conditions on the neural networks for the resulting sequence of functions to form a basis for  $L^2$  (Theorem 4.1). I put forward a spectral learning framework for solving differential equations (Definition 4.2) where augmented bases are used for spatial discretization of differential equations. As applications of the density property of augmented bases in  $L^2$ , I provided convergence guarantees for approximating Schwartz functions (Lemma 4.2) and eigenvalues of some Schrödinger operators (Theorem 4.2) as the truncation parameters goes to infinity. I showed that the total variation of target functions with respect to the push-forward measures induced by the neural networks that satisfy the assumptions of Theorem 4.1 admit minima. A direct corollary of this result is that spectral learning enjoys a faster convergence for approximating Schwartz functions and eigenvalues of Schrödinger operators than standard methods. I employed spectral learning to compute eigenpairs of the vibrational Schrödinger

equation for polyatomic molecules. Results demonstrate a two-order of magnitude increased accuracy upon the use of neural networks. Spectral learning showed a robust training process with little sensitivity to training parameters. This can be attributed to the fact that the augmented sequence of functions is a basis for any values of the parameters of the neural network. The robust training process renders spectral learning more attractive for approximating many eigenfunctions of operators as opposed to standard neural networks, which are generally fragile. Currently, I am investigating the applicability of spectral learning to higher-dimensional problems, which are inaccessible to standard numerical techniques, e. g., the computation of the vibrational spectra of higher-dimensional weakly-bound systems. For a thorougher technical discussion on the limitations and prospective of the proposed spectral learning see Section 4.4.

Overall, the proposed active and spectral learning algorithms extend the boundaries of accuracy and scalability achievable *via* available state-of-the-art methods. They provide powerful tools for modeling quantum chemical dynamic. Furthremore, they are applicable to a wide range of other domains. In particular, the proposed active learning algorithm (Algorithm 3) can be applied for any standard regression task. Similarly, spectral learning can also be straight-forwardly employed for spatial discretizations of other partial differential equations, such as time-dependent Schrödinger equations<sup>1</sup>, or the Hamilton-Jacobi-Bellman equations from control theory.

<sup>&</sup>lt;sup>1</sup>In fact, this is currently under investigation in the framework of a DASHH project by Álvaro Fernández Corral with the aim of simulating strong field ionization processes.

#### Appendix A

# Hilbert spaces and linear operators thereon

Quantum mechanics is formulated in the language of Hilbert spaces and makes extensive use of (unbounded) linear operators thereon to describe physical measurable quantities. Here, relevant results and definitions are collected with special emphasis on bases of Hilbert spaces, which are extensively used to define the spectral learning paradigm (Definition 4.3). Familiarity with the definition of inner products and Hilbert spaces is assumed. Set  $\mathcal{H}$  to be a complex Hilbert space with inner-product  $\langle ., . \rangle$ , taken anti-linear in its first and linear in its second argument. Set  $\|.\| = \sqrt{\langle ., . \rangle}$ . The reader is referred to [63] for thorougher discussion on mathematical quantum mechanics.

#### A.1 Bases of Hilbert spaces

Bases of functional spaces are a crucial tool in approximation theory, as they allow approximating function of these spaces to an arbitrary accuracy.

**Definition A.1.** Let  $\mathcal{H}$  be a Hilbert space. A sequence  $(\phi_n)_n$  is called orthonormal if  $\langle \phi_n, \phi_m \rangle = \delta_{nm}$ .

**Definition A.2.** An orthonormal sequence  $(\phi_n)_n$  in a Hilbert space  $\mathcal{H}$  is called an orthonormal basis of  $\mathcal{H}$  if for all  $\psi \in \mathcal{H}$  one can write

$$\psi = \sum_n \langle \psi, \phi_n \rangle \phi_n.$$

**Proposition A.1** (Bessel's inequality). Let  $(\phi_n)_n$  be an orthonormal sequence in a Hilbert space  $\mathcal{H}$ . for all  $\psi \in \mathcal{H}$  it holds that

$$\|\psi\|^2 \ge \sum_{n=1}^N |\langle \phi_n, \psi \rangle|^2 \quad \text{for all } N \in \mathbb{N}_{>0}.$$
(A.1)

**Example A.1.** The Hermite functions  $(\gamma_n)_n$  defined by

$$\gamma_n(x) = \mathfrak{h}_n(x) \exp(-x^2/2),$$

where  $\mathfrak{h}_n$  denotes the *n*th Hermite polynomial, is a basis set for the Hilbert space  $L^2(\mathbb{R})$ .

**Example A.2.** The sequence  $(\phi_n)_n, \phi_n(x) = \frac{1}{\sqrt{2\pi}} \exp(inx)$  is a basis set for the Hilbert space  $L^2([0, 2\pi])$ .

**Proposition A.2.** *A Hilbert space is separable if and only if it contains an orthonormal basis.* 

One has the following characterization of orthonormal bases.

**Proposition A.3.** An orthonormal sequence  $(\phi_n)_n$  is an orthonormal basis of  $\mathcal{H}$  if and only if

$$\langle \phi_n, \psi \rangle = 0 \quad \text{for all } n \implies \psi = 0.$$
 (A.2)

*Proof.* Let  $(\phi_n)_n$  be a basis and let  $\langle \phi_n, \psi \rangle = 0$  for all *n* and some  $\psi \in \mathcal{H}$ . By Definition A.2  $\psi = 0$ .

Now take  $\chi \in \mathcal{H}$ . By Bessel's inequality Proposition A.1

$$\sum_{n=1}^N |\langle \chi, \phi_n \rangle|^2 \le \|\chi\|^2.$$

The sequence on the left-hand side is non-decreasing and bounded, thus, its limit exits and

$$\lim_{N\to\infty}\sum_{n=1}^N|\langle\chi,\phi_n\rangle|^2=\sum_n|\langle\chi,\phi_n\rangle|^2.$$

Now define  $\psi = \chi - \sum_n |\langle \chi, \phi_n \rangle| \phi_n$ . Clearly  $\psi \in \mathcal{H}$ . Assuming it satisfies (A.2), it holds that  $\psi = 0$  and thus,  $\chi = \sum_n |\langle \chi, \phi_n \rangle| \phi_n$ , i. e.,  $(\phi_n)_n$  is a basis.

#### A.2 Linear operators on Hilbert spaces

The measurement of physical quantities in quantum mechanics is mathematically formulated as computing the spectrum of unbounded linear operators on Hilbert spaces. Basic definitions on such operators are collected here with a special emphasis on self-adjoint operators.  $\mathcal{L}(\mathcal{H})$  denotes the set of bounded linear operators  $u : \mathcal{H} \to \mathcal{H}$ .

An important consequence of the Riesz representation theorem for Hilbert spaces is the existence and uniqueness of adjoints of bounded linear operators. In particular, for each  $A \in \mathcal{L}(\mathcal{H})$  there exists a unique bounded operator  $A^*$  such that

$$\langle A\psi, \phi \rangle = \langle \psi, A^* \phi \rangle$$
 for all  $\psi, \phi \in \mathcal{H}$ .

 $A^*$  is called the adjoint of A.

**Definition A.3** (Symmetric and self-adjoint operators). A linear operator H :  $D(H) \subseteq \mathcal{H} \rightarrow \mathcal{H}$  is called *symmetric* if

$$\langle H\psi, \phi \rangle = \langle \psi, H\phi \rangle$$
 for all  $\psi, \phi \in D(H)$ .

It is called *self-adjoint* if for any  $\phi$ ,  $\eta \in \mathcal{H}$  the relation

$$\langle H\psi, \phi \rangle = \langle \psi, \eta \rangle$$
 for all  $\psi \in D(H)$ 

implies  $\phi \in D(H)$  and  $\eta = H\phi$ .

Every self-adjoint operator is symmetric, but the converse is not true for unbounded operators. Every self-adjoint operator is *closed*, i. e., for any sequence  $(\phi_n)_n$  in D(H), the convergence  $\phi_n \rightarrow \phi$ ,  $H\phi_n \rightarrow \eta$  implies  $\phi \in D(H)$  and  $\eta = H\psi$ .

**Definition A.4** (Unitary operators). An operator U on  $\mathcal{H}$  is *unitary* if it preserves the inner product

$$\langle U\psi, U\phi \rangle = \langle \psi, \phi \rangle$$
 for all  $\psi, \phi \in \mathcal{H}$ ,

or equivalently if  $||U\psi|| = ||\psi||$  for all  $\psi \in \mathcal{H}$ .

**Definition A.5** (Direct sum of Hilbert spaces). Let  $\mathcal{H}_1$ ,  $\mathcal{H}_2$  be two Hilbert spaces. Then, their *direct sum* is defined as

$$\mathcal{H}_1 \oplus \mathcal{H}_2 := \mathcal{H}_1 \times \mathcal{H}_2$$
,

equipped with the scalar product

$$\langle \phi, \psi 
angle_{\mathcal{H}_1 \oplus \mathcal{H}_2} := \langle \phi_1, \psi_1 
angle_{\mathcal{H}_1} + \langle \phi_1, \psi_1 
angle_{\mathcal{H}_2}.$$

 $(\mathcal{H}_1 \oplus \mathcal{H}_2, \langle, \rangle_{\mathcal{H}_1 \oplus \mathcal{H}_2})$  is a Hilbert space.

**Definition A.6** (Graph of an operator, closed operator, closure). • The graph of a linear operator  $T : D(T) \subseteq \mathcal{H} \rightarrow \mathcal{H}$  is the space

$$G(T) = \{ (\phi, T\phi) \in \mathcal{H} \oplus \mathcal{H} \mid \phi \in D(T) \} \subset \mathcal{H} \oplus \mathcal{H}.$$

- An operator *T* is called *closed* if *G* is a closed subset of  $\mathcal{H} \oplus \mathcal{H}$ .
- An operator *T* is called *closable* if it admits a closed extension. In such a case, the smallest closed extension *T* is called the *closure* of *T*.

**Remark A.1.** An operator  $T : D(T) \subseteq \mathcal{H} \to \mathcal{H}$  is said to be densely defined if D(T) is dense in  $\mathcal{H}$ .

One has the following criterion for self-adjoint operators.

**Theorem A.1.** Let  $T : D(T) \subseteq \mathcal{H} \to \mathcal{H}$  be a densely defined and symmetric linear operator. Then, the following are equivalent.

- T is self-adjoint.
- *T* is closed and  $ker(T^* \pm i) = \{0\}$ .

In what follows results concerning spectra of linear operators are introduced.

**Definition A.7** (Resolvent, resolvent set and spectrum). Let  $T : D(T) \subseteq \mathcal{H} \to \mathcal{H}$  be a linear operator on  $\mathcal{H}$ . The *resolvent set* of *T* is defined as

 $\rho(T) := \{z \in \mathbb{C} \mid (T - z) : D(T) \to \mathcal{H} \text{ is a bijection with continuous inverse} \}.$ 

For  $z \in \rho(T)'$  define the *resolvent* of *T* at *z* as

$$R_z(T) := (T-z)^{-1} \in \mathcal{L}(\mathcal{H}).$$

The *spectrum* of *T* is defined as the compliment of the resolvent set

$$\sigma(T) := \mathbb{C} \setminus \rho(T).$$

**Remark A.2.** For closed operators, the continuity requirement in the definition of the resolvent set can be dropped as a consequence of the closed graph theorem, stating that a linear map  $T : X \to Y$  between two Banach spaces X, Y is continuous if and only if it is closed.

**Proposition A.4.** *If T is not closed, then*  $\rho(T) = \emptyset$ *.* 

**Definition A.8** (Partition of the spectrum of a closed operator). Let (T, D(T)) be a closed, linear operator. One can partition its spectrum  $\sigma(T)$  according to the following criteria.

- $\sigma_p(T) = \{z \in \mathbb{C} \mid T z \text{ is not injective}\}$  is called the *point spectrum*, and it coincides with the set of eigenvalues of the operator.
- σ<sub>c</sub>(T) = {z ∈ C | T − z is injective, not surjective, with dense range} is called the *continuous spectrum*.
- $\sigma_r(T) = \{z \in \mathbb{C} \mid T z \text{ is injective, not surjective, with no dense range}\}$  is called the *residual spectrum*.

**Example A.3.** Consider the position operator *q* with the domain

$$D(q) = \{ \psi \in L^2(\mathbb{R}) \mid x\psi \in L^2(\mathbb{R}) \}.$$

Define by

 $q:\psi\to x\psi.$ 

The operator  $(q - z)^{-1}$  is equivalent to the multiplication by  $(x - z)^{-1}$ , which is bounded for all  $z \in \mathbb{C} \setminus \mathbb{R}$ . Thus,  $\sigma(q) = \mathbb{R}$ .

The map  $(q - \lambda)$  has a dense range for all  $\lambda \in \mathbb{R}$ . To see this, define for all  $\psi \in L^2$ 

$$\phi_n := \chi_{\mathbb{R} \setminus [\lambda - rac{1}{n}, \lambda + rac{1}{n}]} rac{\psi}{1 - \lambda}.$$

It holds that  $(x - \lambda)\phi_n \to \psi$  in  $L^2$  and hence, the range of  $x - \lambda$  is dense. Thus,  $\sigma(q) = \sigma_c(q) = \mathbb{R}$ .

One has the following results on the spectrum of self-adjoint operators.

**Theorem A.2.** Let  $T : D(T) \subseteq \mathcal{H} \to \mathcal{H}$  be a symmetric operator. *T* is self-adjoint if and only if  $\sigma(T) \in \mathbb{R}$ .

**Remark A.3.** Quantum mechanics postulates that the outcomes of physical measurements can be represented by the spectrum of linear operators. Since the outcomes of physical measurements are always real quantities, one

would want to describe measurements with operators having real spectrums, whence the importance of self-adjoint operators in quantum mechanics.

One has the following result on the spectrum of compact self-adjoint operators.

**Theorem A.3** (Spectral theorem for compact, self-adjoint operators). Let  $A : \mathcal{H} \to \mathcal{H}$  be a compact, self-adjoint operator. There is an orthonormal basis of  $\mathcal{H}$  consisting of eigenvectors of A. The nonzero eigenvalues of A form a finite or countably infinite set  $\{\lambda_k\}_k$  of real numbers and

$$A = \sum_{k} \lambda_k P_k, \tag{A.3}$$

where  $P_K$  is the orthogonal projection onto the finite-dimensional eigenspace of eigenvectors with eigenvalues  $\lambda_k$ . If the number of nonzero eigenvalues is countably infinite, then the series in (A.3) converges to A in the operator norm.

I finish this appendix by reporting a result that characterizes the effect of perturbations to self-adjoint operators on their spectrum.

**Theorem A.4** (Weyl's inequality [186]). Let  $V_1$ ,  $V_2$  be self-adjoint compact operators and denote by  $E_n(V_1)$ ,  $E_n(V_2)$  their respective eigenvalues which are assumed to be positive. Set  $V = V_1 + V_2$  and denote by  $E_n(V)$  its eigenvalues. It holds that

 $|E_n(V) - E_n(V_1)| \le ||V_2||.$ 

#### Appendix **B**

## From Cartesian to internal coordinates of molecules

Consider a molecule composed of N nuclei. Its state can be described by the position p and velocity v of each atom, thus requiring 6N variables. However, the number of independent variables can be effectively reduced. Assume that the bonds between nuclei are always fixed, i. e., the distances between nuclei do not change. In such a case, one is left with two kinds of motion for the molecule, rotational and translational. Unless an external field is present, properties of molecules, e.g., spectra, are invariant with respect to the translational motion. A common way to effectively describe motions in molecules is using Euler angles, a simple physical extension of spherical polar coordinates. One starts by fixing a Cartesian coordinate system to the center of mass of the molecule. This coordinate system is called a *body-fixed frame*. One aligns the *body-fixed frame* with a *space fixed frame*, also called *laboratory axis*, so that the body-fixed *x*, *y*, and *z* axes coincide with the space-fixed X, Y, and Z axis. Secondly, the body and its frame are rotated actively over a positive angle  $\alpha$ , around the z-axis. Thirdly, one rotates the body and its frame over a positive angle  $\beta$  around the *y*-axis. The z-axis of the body-fixed frame has after these two rotations the longitudinal angle  $\alpha$  and the colatitude angle  $\beta$ , both with respect to the space-fixed frame. A last rotation around its z-axis is necessary to specify its orientation completely. The last rotation angle is

called  $\gamma$ . The total matrix of the three consecutive rotations is the product

$$D = \begin{pmatrix} \cos(\alpha) & -\sin(\alpha) & 0\\ \sin(\alpha) & \cos(\alpha) & 0\\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos(\beta) & 0 & \sin(\beta)\\ 0 & 1 & 0\\ -\sin(\beta) & 0 & \cos(\beta) \end{pmatrix} \begin{pmatrix} \cos(\gamma) & -\sin(\gamma) & 0\\ \sin(\gamma) & \cos(\gamma) & 0\\ 0 & 0 & 1 \end{pmatrix}.$$

Thus, the positions of nuclei of the molecule in the laboratory frame are given by

$$q = R + D \cdot r, \tag{B.1}$$

where  $R \in \mathbb{R}^3$  denotes the position of the center of mass, while  $r \in \mathbb{R}^{3N-6}$  describes the positions of the nuclei with respect to the body-fixed frame. In the rigid model, this vector assumes a fixed value. So, one can see that the total system has 6 degrees of freedom, 3 rotational and 3 translational. In a non-rigid model, one has 3N - 6 more degrees of freedom, known as *vibrational degrees of freedom*. As an example let us consider a water molecule. It has three nuclei and 2 bonds. Thus, one has have 3 vibrational degrees of freedom. Two corresponding to the change of the bonds-length and one corresponding to the change of the angle between the bonds. For such a model equation B.1 is still valid but *r* is not any more constant. Hence, to describe the whole system one needs 3N degrees of freedom.

#### Appendix C

## **Random forest regressors**

Regression trees are a non-parametric way of solving a regression problem. They are based on the intuition that the output value can be inferred by partitioning the input space. In particular, for solving a regression problem with datapairs  $\{(x_i, E_i)\}_{i=1}^l$ , a *Tree-Regressor* aims at finding *J* distinct and non-overlapping regions  $R_1 \dots R_J$  in the feature space that minimize

$$\sum_{j=1}^{J} \sum_{i; x_i \in R_j} (E_i - \bar{E}_{R_j})^2,$$

where

$$\bar{E}_{R_j} = \sum_{i; X_i \in R_j} E_i / m_j$$

is the average target value of  $m_j$  examples in region  $R_j$  [187]. This problem is NP-complete [188]. Therefore, only near-optimal solutions are considered by restricting ourselves to hyper-rectangular regions and using recursive binary splitting, a greedy algorithm to obtain a near-optimal segmentation. A prediction for a new input x is done by assigning the input to one of the regions. The prediction for this input is then the average value of all examples in the training dataset that fall in this region. A major drawback of tree regressors is their large variance [187]. A powerful approach to mitigate this problem is to consider an ensemble of trees. The key idea is that averaging a set of independent random

variables, which have comparable variances, reduces their overall variance [187]. In an ensemble method, a random perturbation is introduced to the learning process in order to produce different learners from the same training set. Thus, taking the average of the predictions of the ensemble would result in a reduction of variance. Such a random perturbation can be introduced by bootstraping, which gives rise to bootstrap aggregation (bagging) methods. Here, *B* different bootstrapped datasets of size  $m_b$  are generated. A tree is built on each model. For a new data point, a prediction is made by taking the average of the predictions of all trees

$$T_{\text{bag}}(x) = \frac{1}{B} \sum_{i=1}^{B} T_i(x).$$

A further random perturbation in tree models can be introduced by considering, at each split, only a randomly drawn subset of all possible features. This gives rise to the *random forest regressor* (RFR) [189]. Thus, one can see that RFR employs a 2-fold randomization procedure. The ensemble can be made even more diverse by introducing further randomization in the learning process, e. g., extremely randomized trees [190]. RFR is an inherent ensemble method encompassing diverse learners. This makes the model a very attractive option for a query by committee-based algorithm like Algorithm 2 or Algorithm 3.

Another advantage of using the trees of an RFR in Algorithm 2 and Algorithm 3 is its relatively low training complexity. An AL paradigm is a dynamic paradigm that needs to be performed iteratively until one is satisfied with the performance. One wishes to be able to perform these iterations quickly. Otherwise, the time saved from performing redundant electronic structure calculations would be wasted in performing AL iterations. Building an RFR is relatively cheap. It has an average time complexity of  $\Theta(M \cdot K \cdot \tilde{N} \log_2^2 \tilde{N})$  [191] where *K*, *M* denote the number of random features sampled at each splitting and the number of trees, respectively, and  $\tilde{N} \approx 0.632N$  with the number *N* of training examples.<sup>1</sup> This should be compared to the computational cost of training a Gaussian process, which scales

<sup>&</sup>lt;sup>1</sup>The probability of not selecting a point in *n* draws of *n* samples with replacement is  $(1 - 1/n)^n$ , which converges in the limit of  $n \to \infty$  to  $e^{-1}$ . Hence, bootstrap samples draw, on average,  $1 - 1/e \approx 63.2$  % of unique samples [191, 192].

as  $\mathcal{O}(N^3)$  [84]. With extremely randomized trees the average time complexity for training is  $\Theta(M \cdot K \cdot N \log_2 N)$  [190]. The average inference complexity of RFR is  $\Theta(M \log N)$  [190]. Thus, one AL iteration scales as  $\mathcal{O}(M * K * \tilde{N} \log_2^2 \tilde{N})$  with the number of so-far-labeled data N. The complexity of RFR is asymptotically inferior to that of a neural network (NN), which has a training time complexity<sup>2</sup> of  $\mathcal{O}(N_e N(\sum_i^{l-1} N^i N^{i+1}))$  with the number of epochs  $N_e$  needed for the NN to converge and the number of neurons  $N^i$  in layer i. However, in the data size regime of our application, the computational costs of an RFR are smaller than that of the NN.

<sup>&</sup>lt;sup>2</sup>This bound can be straightforwardly obtained by noting that matrix multiplications are the most expensive computations in the forward and backward passes of the NN training. We assume here that matrix multiplication scales as  $\Theta(N^3)$ .

#### Appendix D

## **Measure theory**

The proposed spectral learning paradigm (Definition 4.3) and the discussion of active learning are formulated using the language of measure theory. Some relevant definitions and results are collected here. The reader is referred to [193] for a deeper discussion on standard measure theory and to [194] for extensive results on probability measure spaces.

For what follows set  $(X, \mathcal{A}, \lambda)$  to be an abstract measure space  $\mathcal{M}$  to be the families of real-valued measurable functions  $u : (X, \mathcal{A}) \to (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , where  $\mathcal{B}$  denotes the Borel  $\sigma$ -algebra.

**Remark D.1.** Given a measure space  $(X, A, \lambda)$  and  $A \in A$  one would sometimes want to define a new measure space  $(A, A_A, \lambda_A)$  on A. This appears, e.g., when  $X = \mathbb{R}$ . A formal way to do that is to set

$$\mathcal{A}_A = \{ M \in \mathcal{A}, M \subseteq A \},\$$

and

 $\lambda_A = \lambda|_A.$ 

**Definition D.1** (Bi-measurable functions). A bijection  $f : X \to X$  is said to be a *bi-measurable* or an *isomorphism* if both f and  $f^{-1}$  are measurable.

#### D.1 Push-forward measures

In Chapter 4 normalizing flows are required to be transformations that are wellbehaved in the following sense.

**Definition D.2** (Non-singular measurable transformations). A measurable transformation or mapping  $h : X \to X$  is said to be *non-singular* if  $\lambda(T^{-1}(S)) = 0$  whenever  $\lambda(S) = 0$ .

Measurable transformations induce measures in the following way.

**Definition D.3** (Push-forward measure). Given a measure space  $(X, A, \lambda)$  and a measurable mapping  $h : X \to X$  the *push-forward measure of*  $\lambda$  is

$$h_{\#}\lambda(A) = \lambda(h^{-1}(A))$$
 for all  $A \in \mathcal{A}$ .

The push-forward measure is sometimes denoted by  $\lambda h^{-1}$ .

Push-forward measures relate to composition operators in the following sense.

**Example D.1.** [Integrating a function with respect to a push-forward measure] Given a measurable mapping  $g : X \to X$  it holds

$$\int_X g \, dh_* \lambda = \int_X g \, d\lambda h^{-1}$$
$$= \int_X g \circ h \, d\lambda.$$

#### D.2 The Radon-Nikodym theorem

Radon-Nikodym theorem is an important result that establishes relationships between measures on the same measurable space. It has important applications, e.g., in probability theory.

The following is an important property of measures.

**Definition D.4.** Let  $\lambda$ ,  $\nu$  be two measure on the measurable space (X, A).  $\nu$  is said to be *absolutely continuous* with respect to  $\lambda$  if

for all 
$$A \in \mathcal{A}$$
 with  $\lambda(A) = 0 \implies \nu(A) = 0$ .

If this holds write  $\nu \ll \lambda$ .

The push-forward of a measure  $\lambda$  is absolutely continuous assuming some regularity on the underlying transformation.

**Lemma D.1** ([173]). *If a measurable transformation*  $h : X \to X$  *is non-singular, then*  $h_{\#}\lambda \ll \lambda$ .

Radon-Nikodym theorem characterizes absolute continuity of measures.

**Theorem D.1** (Radon-Nikodym [193]). Let  $\lambda$ ,  $\nu$  be two measures on a measurable space  $(X, \mathcal{A})$ . If  $\lambda$  is  $\sigma$ -finite the following is equivalent:

- $\nu \ll \lambda$ .
- $\nu(A) = \int_A f d\lambda$  for all  $A \in A$  for some almost everywhere unique  $f \in \mathcal{M}, f \ge 0$ .

f is called the Radon-Nikodym derivative and is denoted by  $\frac{dv}{d\lambda}$ .

As an example, one can see that the integration against a push-forward measure induced by a non-singular transformation can be written using the Radon-Nikodym derivative

$$\int_X g \, dh_{\#} \lambda = \int_X g \frac{dh_{\#} \lambda}{d\lambda} \, d\lambda.$$

Finally, I collect some fundamental results for the case  $X = \mathbb{R}$ , A = B, and  $\lambda$  is the Lebesgue measure.

**Definition D.5.** A measure  $\lambda$  on  $(\mathbb{R}, \mathcal{B})$  is said to be *locally finite* if for all  $x \in X$  there exists a neighborhood  $N_x$  of x such that  $\lambda(N_x) < \infty$ .

The following theorem is a corollary of Lebesgue's differentiation theorem which links Radon-Nikodym derivatives with definitions of derivatives from standard calculus.

**Theorem D.2** (Lebesgue differentiation [193]). Let v be a locally finite measure on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  that is absolutely continuous with respect to  $\lambda$ . Set

$$D\mu := \lim_{r \to 0} \frac{\nu(B_r(x))}{\lambda(B_r(x))}.$$

*Then*  $D\mu$  *exists Lebesgue almost everywhere and coincides almost everywhere with the Radon-Nikodym derivative*  $\frac{d\nu}{d\lambda}$ .

### D.3 Radon measures and a Riesz representation theorem

A very useful measure theoretic tool is the *Radon measure*, since it characterizes the topological duals of important functional spaces. Such results are important, e. g., in the field of optimal transportation, where they allow for proving that the Kantorovitch formulation of optimal transport admit minima [177]. For what follows, set  $X = \mathbb{R}$ , and  $\mathcal{A} = \mathcal{B}(\mathbb{R}) = \mathcal{B}$ .

**Definition D.6.** A measure  $\lambda$  on  $(\mathbb{R}, \mathcal{B})$  is said to be *inner-regular* if for any  $U \in \mathcal{B}$  one has

$$\lambda(U) = \sup\{K : K \subseteq U, K \text{ is compact}\}.$$

**Definition D.7** (Radon measure). A *Radon measure* on  $(\mathbb{R}, \mathcal{B})$  is a measure that is both inner-regular and locally finite. Denote by  $\mathcal{M}(\mathbb{R}) = \mathcal{M}$  the set of Radon measures on  $(\mathbb{R}, \mathcal{B})$ .

**Remark D.2.** Note that Radon measure are usually introduced for locally compact Hausdorff spaces but analysis here is restricted to the only relevant case of  $X = \mathbb{R}$ .

Example D.2. [The Lebesgue measure is Radon] note that the Lebesgue

measure restricted to  $(X, \mathcal{B})$  is locally-finite and inner-regular. Hence, it is a Radon measure.

The importance of Radon measures stems from their relations to the topological duals of spaces of continuous functions. Equip  $\mathcal{M}$  with the total variation norm

$$\|\gamma\|_{\mathcal{M}} := |\gamma|(X).$$

Equip the space  $C(\mathbb{R})$  with the sup-norm.

**Theorem D.3** ([193]). The topological dual of C can be identified with M with the duality pairing

$$\gamma(f) = \langle \gamma, f \rangle_{\mathcal{M},\mathcal{C}} = \int_{\mathbb{R}} f d\gamma, \quad \text{for } \gamma \in \mathcal{M}, f \in \mathcal{C}.$$

#### **D.4 Probability theory**

Finally, some relevant results from probability measure theory are collect. Consider  $(X, \mathcal{A})$  to be a measure space. Denote by  $\mathbb{P}$  an arbitrary probability measure and by  $\mathbb{P}(X)$  the set of probability measures over X. It is common to refer to any  $A \in \mathcal{A}$  by an *event*.

**Definition D.8** (Independent events). Two events *A*, *B* are said to be *independent* if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B).$$

Given an index set *I*, a family  $(A_i)_{i \in I}$  of events is said to be *independent* if

$$\mathbb{P}(\cap_{j\in J}A_j) = \prod_{j\in J}\mathbb{P}(A_j) \quad \text{for all } J\subset I.$$

One can extend the notion of independence of family of events to independence of family of classes of events.

**Definition D.9.** Let *I* be an index set and consider  $\mathcal{E}_i \in \mathcal{A}$  for all  $i \in I$ . The family  $(\mathcal{E}_i)_{i \in I}$  is called *independent* if, for any finite subset  $J \in I$  and any choice of

 $E_j \in \mathcal{E}_j, j \in J$ , one has

$$\mathbb{P}(\cap_{j\in J}E_j)=\prod_{j\in J}\mathbb{P}(E_j).$$

**Definition D.10** (Random variables/distributions). (i) A measurable mapping  $V : X \to \mathbb{R}$  is called a *real random variable*.

(ii) The push-forward measure  $P_V := V_{\#}\mathbb{P}$  induced by a real random variable *V* is called the *distribution of V*.

The following is an important family of random variables that one often encounters in machine learning and statistics.

**Definition D.11** (Independent and identically distributed random variables). Let *I* be an index set and  $(V_i)_{i \in I}$  be a family of real random variables. Endow  $\mathbb{R}$  by the Borel- $\sigma$  algebra  $\mathcal{B}$ .

(i) The family  $(V_i)_{i \in I}$  is said to be *identically distributed* if

$$P_{V_i} = P_{V_i}$$
 for all  $i, j \in I$ .

(ii) The family  $(V_i)_{i \in I}$  is said to be *independent* if the family of generated sigma algebras  $(\sigma(V_i))_{i \in I}$ , where  $\sigma(V_i) = V_i^{-1}(\mathcal{B})$  is independent.

A family of real random variables satisfying both conditions is said to be *i.i.d.* In such a case set  $P = P_{V_i}$ .

Given  $\nu \in \mathbb{P}(X)$  an interesting question is whether one can find a measurable mapping  $T : X \to [0, 1]$  such that the induced push-forward measure  $T_{\#\nu}$  coincides with  $\mu|_{[0,1]}$ , i. e., the Lebesgue measure on [0,1]. This question turns out to be essential in the nonlinear analysis of spectral learning carried out in Chapter 3. The following theorem provides an answer under some particular settings.

**Theorem D.4** ([167] (Theorem 17.41)). *Let X* be a standard Borel space and assume that  $\nu \in \mathbb{P}(X)$  satisfies

$$\nu(\{x\}) = 0$$
 for all  $x \in X$ .

Then, there exists a bi-measurable mapping  $T: X \rightarrow [0, 1]$  such that

 $T_{\#}\nu = \mu|_{[0,1]}.$ 

#### Appendix E

## **Functional analysis**

Some important results to prove the existence of minimizers to the variational formulation of Schrödinger equation Theorem 2.2 and optimization problem defined in (4.15) are collected here.

Denote by *X* a normed vector space, and by *X*<sup>\*</sup> its topological dual.

**Theorem E.1** (Banach-Alaoglu). Let X be separable, then the closed unit ball in  $X^*$  is compact in the weak-\* topology.

The following is a very important corollary.

**Theorem E.2.** *Any bounded sequence in* X\* *has a weak*\* *converging subsequence.* 

## Bibliography

- [1] A. H. Zewail, "Femtochemistry: Atomic-Scale Dynamics of the Chemical Bond", *J. Phys. Chem. A* **104**, 5660–5694 (2000).
- [2] A. A. Ischenko, P. M. Weber, and R. J. D. Miller, "Capturing Chemistry in Action with Electrons: Realization of Atomically Resolved Reaction Dynamics", *Chem. Rev.* **117**, 11066–11124 (2017).
- [3] E. T. Karamatskos, S. Raabe, T. Mullins, A. Trabattoni, P. Stammer, G. Goldsztejn, R. R. Johansen, K. Długołęcki, H. Stapelfeldt, M. J. J. Vrakking, S. Trippel, A. Rouzée, and J. Küpper, "Molecular movie of ultrafast coherent rotational dynamics of OCS", *Nat. Commun.* **10**, 3364 (2019), arXiv: 1807.01034 [physics].
- [4] J. C. Polanyi and A. H. Zewail, "Direct Observation of the Transition State", Acc. Chem. Res. 28, 119–132 (1995).
- [5] C. I. Blaga, J. Xu, A. D. DiChiara, E. Sistrunk, K. Zhang, P. Agostini, T. A. Miller, L. F. DiMauro, and C. D. Lin, "Imaging ultrafast molecular dynamics with laser-induced electron diffraction", *Nature* 483, 194–197 (2012).
- [6] J. L. McHale, *Molecular spectroscopy*, CRC Press, 2017.
- [7] A. Owens, A. Yachmenev, S. N. Yurchenko, and J. Küpper, "Climbing the Rotational Ladder to Chirality", *Phys. Rev. Lett.* **121**, 193201 (2018), arXiv: 1802.07803 [physics].

- [8] T. Endo, S. P. Neville, V. Wanie, S. Beaulieu, C. Qu, J. Deschamps, P. Lassonde, B. E. Schmidt, H. Fujise, M. Fushitani, A. Hishikawa, P. L. Houston, J. M. Bowman, M. S. Schuurman, F. Légaré, and H. Ibrahim, "Capturing roaming molecular fragments in real time", *Science* **370**, 1072–1077 (2020).
- [9] T. Mullins, E. T. Karamatskos, J. Wiese, J. Onvlee, A. Rouzée, A. Yachmenev, S. Trippel, and J. Küpper, "Picosecond pulse-shaping for strong threedimensional field-free alignment of generic asymmetric-top molecules", *Nat. Commun.* **13**, 1431 (2022), arXiv: 2009.08157 [physics].
- [10] R. Tóbiás, T. Furtenbacher, I. Simkó, A. G. Császár, M. L. Diouf, F. M. J. Cozijn, J. M. A. Staa, E. J. Salumbides, and W. Ubachs, "Spectroscopic-network-assisted precision spectroscopy and its application to water", *Nat. Commun.* 11, 1708 (2020).
- [11] A. Campargue, S. Kassi, A. Yachmenev, A. A. Kyuberis, J. Küpper, and S. N. Yurchenko, "Observation of electric-quadrupole infrared transitions in water vapor", *Phys. Rev. Research* 2, 023091 (2020), arXiv: 2001.02922 [physics].
- [12] A. G. Császár, I. Simkó, T. Szidarovszky, G. C. Groenenboom, T. Karman, and A. van der Avoird, "Rotational–vibrational resonance states", *Phys. Chem. Chem. Phys.* 22, 15081–15104 (2020).
- [13] T. de Jongh, M. Besemer, Q. Shuai, T. Karman, A. van der Avoird, G. C. Groenenboom, and S. Y. T. van de Meerakker, "Imaging the onset of the resonance regime in low-energy NO-He collisions", *Science* 368, 626–630 (2020).
- [14] S. N. Yurchenko, J. Tennyson, J. Bailey, M. D. J. Hollis, and G. Tinetti, "Spectrum of hot methane in astronomical objects using a comprehensive computed line list", *PNAS* **111**, 9379–9383 (2014).
- [15] C. Lubich, *From quantum to classical molecular dynamics: reduced models and numerical analysis*, European Mathematical Society, 2008.
- [16] A. Iske, Approximation Theory and Algorithms for Data Analysis, Springer International Publishing, 2018.

- [17] S. Teufel, Adiabatic perturbation theory in quantum dynamics, Springer Science & Business Media, 2003.
- [18] M. Johny, C. A. Schouder, A. Al-Refaie, L. He, J. Wiese, H. Stapelfeldt, S. Trippel, and J. Küpper, *Molecular sunscreen: water protects pyrrole from radiation damage*, submitted, 2020, arXiv: 2010.00453 [physics].
- [19] J. Onvlee, S. Trippel, and J. Küpper, "Ultrafast light-induced dynamics in solvated biomolecules: The indole chromophore with water", *Nat. Commun.* 13, 7462 (2022), arXiv: 2103.07171 [physics].
- [20] J. Tennyson and S. N. Yurchenko, "ExoMol: molecular line lists for exoplanet and other atmospheres", *Mon. Not. R. Astron. Soc.* 425, 21–33 (2012).
- [21] B. Engquist, A. Fokas, E. Hairer, and A. Iserles, *Highly oscillatory problems*, 366, Cambridge University Press, 2009.
- [22] M. Born and R. Oppenheimer, "Zur Quantentheorie der Molekeln", Ann. Physik 84, 457–484 (1927).
- [23] T. Helgaker, P. Jorgensen, and J. Olsen, *Molecular electronic-structure theory*, John Wiley & Sons, 2014.
- [24] R. J. LeVeque, *Finite volume methods for hyperbolic problems*, vol. 31, Cambridge university press, 2002.
- [25] E. Godlewski and P.-A. Raviart, *Numerical approximation of hyperbolic systems of conservation laws*, vol. 118, Springer Science & Business Media, 2013.
- [26] W. Gautschi, *Numerical analysis*, Springer Science & Business Media, 2011.
- [27] B. Wendroff, "Difference Methods for Initial-Value Problems (Robert D. Richtmyer and K. W. Morton)", SIAM Rev. Soc. Ind. Appl. Math. 10, 381–383 (1968).
- [28] D. Gottlieb and S. A. Orszag, *Numerical analysis of spectral methods: theory and applications*, SIAM, 1977.
- [29] C. Canuto, M. Y. Hussaini, A. Quarteroni, and T. A. Zang, Spectral methods: fundamentals in single domains, Springer Science & Business Media, 2007.

- [30] J. P. Boyd, *Chebyshev and Fourier Spectral Methods*, Applied Mathematical Sciences, Springer, 2000.
- [31] S. Manzhos and T. Carrington, "Neural Network Potential Energy Surfaces for Small Molecules and Reactions", *Chem. Rev.* **121**, 10187 (2021).
- [32] J. Behler, "Constructing high-dimensional neural network potentials: A tutorial review", Int. J. Quantum Chem. 115, 1032–1050 (2015).
- [33] T. B. Blank, S. D. Brown, A. W. Calhoun, and D. J. Doren, "Neural network models of potential energy surfaces", J. Chem. Phys. 103, 4129–4137 (1995).
- [34] K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko, and K.-R. Müller, "Schnet–a deep learning architecture for molecules and materials", *J. Chem. Phys.* 148, 241722 (2018).
- [35] W. E and B. Yu, "The deep Ritz method: a deep learning-based numerical algorithm for solving variational problems", *Commun. Math. Stat.* 6, 1–12 (2018).
- [36] E. Kharazmi, Z. Zhang, and G. E. Karniadakis, "Variational physics-informed neural networks for solving partial differential equations", arXiv (2019), arXiv: 1912.00873 [cs].
- [37] S. Cuomo, V. S. Di Cola, F. Giampaolo, G. Rozza, M. Raissi, and F. Piccialli, "Scientific Machine Learning through Physics-Informed Neural Networks: Where we are and What's next", *arXiv* (2022), arXiv: 2201.05624 [cs].
- [38] M. Hutzenthaler, A. Jentzen, T. Kruse, T. Anh Nguyen, and P. von Wurstemberger, "Overcoming the curse of dimensionality in the numerical approximation of semilinear parabolic partial differential equations", *Proc. R. Soc. A.* 476, 20190630 (2020).
- [39] P. Grohs, F. Hornung, A. Jentzen, and P. Von Wurstemberger, "A proof that artificial neural networks overcome the curse of dimensionality in the numerical approximation of Black-Scholes partial differential equations", *arXiv* (2018), arXiv: 1809.02362 [physics].

- [40] J. Han, A. Jentzen, et al., "Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations", *Comm. Math. Stat.* 5, 349–380 (2017).
- [41] G. Carleo and M. Troyer, "Solving the quantum many-body problem with artificial neural networks", *Science* **355**, 602–606 (2017).
- [42] J. Hermann, Z. Schätzle, and F. Noé, "Deep-neural-network solution of the electronic Schrödinger equation", *Nat. Chem.* **12**, 891–897 (2020).
- [43] K. Choo, G. Carleo, N. Regnault, and T. Neupert, "Symmetries and manybody excitations with neural-network quantum states", *Phys. Rev. Lett.* 121, 167204 (2018).
- [44] J. Han, J. Lu, and M. Zhou, "Solving high-dimensional eigenvalue problems using deep neural networks: A diffusion Monte Carlo like approach", J. *Comput. Phys.* 423, 109792 (2020).
- [45] J. Kessler, F. Calcavecchia, and T. D. Kühne, "Artificial neural networks as trial wave functions for quantum monte carlo", *Adv. Theory Simul.* 4, 2000269 (2021).
- [46] D. Pfau, J. S. Spencer, A. G. Matthews, and W. M. C. Foulkes, "Ab initio solution of the many-electron Schrödinger equation with deep neural networks", *Phys. Rev. Research* 2, 033429 (2020).
- [47] M. Entwistle, Z. Schätzle, P. A. Erdman, J. Hermann, and F. Noé, "Electronic excited states in deep variational Monte Carlo", *Nat. Commun.* 14, 274 (2023), arXiv: 2203.09472 [physics].
- [48] J. S. Spencer, D. Pfau, A. Botev, and W. M. C. Foulkes, *Better, faster fermionic neural networks*, 2020, arXiv: 2011.07125 [physics].
- [49] J. Hermann, J. Spencer, K. Choo, A. Mezzacapo, W. Foulkes, D. Pfau, G. Carleo, and F. Noé, *Ab-initio quantum chemistry with neural-network wave-functions*, 2022, arXiv: 2208.12590 [physics].
- [50] W. E, C. Ma, S. Wojtowytsch, L. Wu, et al., "Towards a Mathematical Understanding of Neural Network-Based Machine Learning: what we know and what we don't", (2020), arXiv: 2009.10713 [math].

- [51] W. E, C. Ma, and L. Wu, "The Barron Space and the Flow-Induced Function Spaces for Neural Network Models", *Constr. Approx.* 55, 369–406 (2022).
- [52] R. Parhi and R. D. Nowak, "Banach space representer theorems for neural networks and ridge splines", *J. Mach. Learn. Res.* 22, 1960–1999 (2021).
- [53] Z. Chen, J. Lu, Y. Lu, and S. Zhou, "A Regularity Theory for Static Schrödinger Equations on  $\mathbb{R}^d$  in Spectral Barron Spaces", (2022), arXiv: 2201.10072 [cs].
- [54] Y. Lu, J. Lu, and M. Wang, "A priori generalization analysis of the deep ritz method for solving high dimensional elliptic partial differential equations", *Conference on learning theory*, PMLR, 2021, 3196–3241.
- [55] Z. Chen, J. Lu, and Y. Lu, "On the Representation of Solutions to Elliptic PDEs in Barron Spaces", *Advances in Neural Information Processing Systems*, ed. by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, vol. 34, Curran Associates, Inc., 2021, 6454–6465.
- [56] G. Ongie, R. Willett, D. Soudry, and N. Srebro, "A function space view of bounded norm infinite width relu nets: The multivariate case", (2019), arXiv: 1910.01635 [cs].
- [57] J. W. Siegel and J. Xu, "Sharp bounds on the approximation rates, metric entropy, and n-widths of shallow neural networks", *Found. Comput. Math.* 1–57 (2022).
- [58] W. E and S. Wojtowytsch, "Representation formulas and pointwise properties for Barron functions", *Calc. Var.* 61, 1–37 (2022).
- [59] D. Elbrächter, D. Perekrestenko, P. Grohs, and H. Bölcskei, "Deep neural network approximation theory", *IEEE Trans. Inf. Theory* 67, 2581–2623 (2021).
- [60] Z.-Q. J. Xu, Y. Zhang, T. Luo, Y. Xiao, and Z. Ma, "Frequency principle: Fourier analysis sheds light on deep neural networks", (2019), arXiv: 1901. 06523 [cs].
- [61] P. W. Atkins and R. S. Friedman, *Molecular Quantum Mechanics*, 3rd ed., Oxford University Press, 1997.

- [62] T. Kato, "On the convergence of the perturbation method. I", Prog. Theor. Phys. 4, 514–523 (1949).
- [63] M. Porta, *Mathematical Quantum Theory*, Lecture notes (accessed 2022-03-15), 2019.
- [64] B. T. Sutcliffe and J. Tennyson, "A general treatment of vibration-rotation coordinates for triatomic molecules", *Int. J. Quantum Chem.* **39**, 183–196 (1991).
- [65] A. Blum, J. Hopcroft, and R. Kannan, *Foundations of Data Science*, Cambridge University Press, 2020.
- [66] M. Lotz, *Mathematics of Machine Learning*, Lecture notes (accessed 2021-03-30), 2018.
- [67] Y. Yao, *A mathematical Introduction to Data Science*, Lecture notes (accessed 30-03-2021), 2019.
- [68] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*, Cambridge University Press, 2014.
- [69] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*, vol. 1, MIT Press, 2016.
- [70] X. Wang, P. L. Houston, and J. M. Bowman, "A new (multi-reference configuration interaction) potential energy surface for H<sub>2</sub>CO and preliminary studies of roaming", *Phil. Trans. R. Soc. A* 375, 20160194 (2017).
- [71] B. J. Braams and J. M. Bowman, "Permutationally invariant potential energy surfaces in high dimensionality", *Int. Rev. Phys. Chem.* 28, 577–606 (2009).
- [72] Z. Xie and J. M. Bowman, "Permutationally invariant polynomial basis for molecular energy surface fitting via monomial symmetrization", J. Chem. Theory Comput. 6, 26–34 (2010).
- [73] C. Qu, Q. Yu, and J. M. Bowman, "Permutationally Invariant Potential Energy Surfaces", Annu. Rev. Phys. Chem. 69, 151–175 (2018).

- [74] R. Conte, C. Qu, P. L. Houston, and J. M. Bowman, "Efficient Generation of Permutationally Invariant Potential Energy Surfaces for Large Molecules", *J. Chem. Theory Comput.* 16, 3264–3272 (2020).
- [75] T. Morawietz, V. Sharma, and J. Behler, "A neural network potential-energy surface for the water dimer based on environment-dependent atomic energies and charges", J. Chem. Phys. 136, 064103 (2012).
- [76] J. Behler and M. Parrinello, "Generalized neural-network representation of high-dimensional potential-energy surfaces", *Phys. Rev. Lett.* 98, 146401 (2007).
- [77] O. T. Unke and M. Meuwly, "PhysNet: A Neural Network for Predicting Energies, Forces, Dipole Moments, and Partial Charges", J. Chem. Theory Comput. 15, 3678–3693 (2019).
- [78] S. Manzhos, R. Dawes, and T. Carrington, "Neural network-based approaches for building high dimensional and quantum dynamics-friendly potential energy surfaces", *Int. J. Quantum Chem.* **115**, 1012–1020 (2014).
- [79] B. Jiang, J. Li, and H. Guo, "Potential energy surfaces from high fidelity fitting of *ab initio* points: The permutation invariant polynomial - neural network approach", *Int. Rev. Phys. Chem.* 35, 479–506 (2016).
- [80] C. Schran, J. Behler, and D. Marx, "Automated Fitting of Neural Network Potentials at Coupled Cluster Accuracy: Protonated Water Clusters as Testing Ground", J. Chem. Theory Comput. 16, 88–99 (2019).
- [81] A. Kamath, R. A. Vargas-Hernández, R. V. Krems, T. Carrington, and S. Manzhos, "Neural networks vs Gaussian process regression for representing potential energy surfaces: A comparative study of fit quality and vibrational spectrum accuracy", J. Chem. Phys. 148, 241702 (2018).
- [82] A. P. Bartók, M. C. Payne, R. Kondor, and G. Csányi, "Gaussian Approximation Potentials: The Accuracy of Quantum Mechanics, without the Electrons", *Phys. Rev. Lett.* **104**, 136403 (2010).
- [83] C. Qu, Q. Yu, B. L. V. Hoozen, J. M. Bowman, and R. A. Vargas-Hernández, "Assessing Gaussian Process Regression and Permutationally Invariant Polynomial Approaches To Represent High-Dimensional Potential Energy Surfaces", J. Chem. Theory Comput. 14, 3381–3396 (2018).
- [84] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*, ed. by T. Dietterich, Cambridge, MS, USA: the MIT Press, 2006.
- [85] H. Sugisawa, T. Ida, and R. V. Krems, "Gaussian process model of 51dimensional potential energy surface for protonated imidazole dimer", *J. Chem. Phys.* 153, 114101 (2020).
- [86] J. Dai and R. V. Krems, "Interpolation and extrapolation of global potential energy surfaces for polyatomic systems by Gaussian processes with composite kernels", J. Chem. Theory Comput. 16, 1386–1395 (2020).
- [87] R. Vargas-Hernández, Y Guan, D. Zhang, and R. Krems, "Bayesian optimization for the inverse scattering problem in quantum reaction dynamics", *New J. Phys.* 21, 022001 (2019).
- [88] O. T. Unke and M. Meuwly, "Toolkit for the Construction of Reproducing Kernel-Based Representations of Data: Application to Multidimensional Potential Energy Surfaces", J. Chem. Inf. Model. 57, 1923–1931 (2017).
- [89] P. O. Dral, A. Owens, S. N. Yurchenko, and W. Thiel, "Structure-based sampling and self-correcting machine learning for accurate calculations of potential energy surfaces and vibrational levels", J. Chem. Phys. 146, 244108 (2017).
- [90] D. Koner and M. Meuwly, "Permutationally Invariant, Reproducing Kernel-Based Potential Energy Surfaces for Polyatomic Molecules: From Formaldehyde to Acetone", J. Chem. Theory Comput. 16, 5474–5484 (2020).
- [91] B. Settles, Active learning literature survey, tech. rep. 1648, University of Wisconsin-Madison Department of Computer Sciences, 2009.
- [92] A. A. Peterson, R. Christensen, and A. Khorshidi, "Addressing uncertainty in atomistic machine learning", *Phys. Chem. Chem. Phys.* 19, 10978–10985 (2017).

- [93] Q. Lin, Y. Zhang, B. Zhao, and B. Jiang, "Automatically growing global reactive neural network potential energy surfaces: A trajectory-free active learning strategy", J. Chem. Phys. 152, 154104 (2020).
- [94] L. Zhang, D.-Y. Lin, H. Wang, R. Car, and W. E, "Active learning of uniformly accurate interatomic potentials for materials simulation", *Phys. Rev. Mater.* 3, 023804 (2019).
- [95] E. Uteva, R. S. Graham, R. D. Wilkinson, and R. J. Wheatley, "Active learning in Gaussian process interpolation of potential energy surfaces", J. *Chem. Phys.* 149, 174114 (2018).
- [96] T. D. Loeffler, T. K. Patra, H. Chan, M. Cherukara, and S. K. Sankaranarayanan, "Active learning the potential energy landscape for water clusters from sparse training data", *J. Phys. Chem. C* **124**, 4907–4916 (2020).
- [97] Y. Zhai, A. Caruso, S. Gao, and F. Paesani, "Active learning of many-body configuration space: Application to the Cs<sup>+</sup>–water MB-nrg potential energy function as a case study", J. Chem. Phys. 152, 144103 (2020).
- [98] J. Vandermause, S. B. Torrisi, S. Batzner, Y. Xie, L. Sun, A. M. Kolpak, and B. Kozinsky, "On-the-fly active learning of interpretable Bayesian force fields for atomistic rare events", *npj Comput Mater* 6, 20 (2020).
- [99] M. Gastegger, J. Behler, and P. Marquetand, "Machine learning molecular dynamics for the simulation of infrared spectra", *Chem. Sci.* 8, 6924–6935 (2017).
- [100] J. S. Smith, B. Nebgen, N. Lubbers, O. Isayev, and A. E. Roitberg, "Less is more: Sampling chemical space with active learning", *J. Chem. Phys.* 148, 241733 (2018).
- [101] G. Sivaraman, A. N. Krishnamoorthy, M. Baur, C. Holm, M. Stan, G. Csányi, C. Benmore, and Á. Vázquez-Mayagoitia, "Machine-learned interatomic potentials by active learning: amorphous and liquid hafnium dioxide", *npj Comput. Mater.* 6, 1–8 (2020).

- [102] Y. Guan, S. Yang, and D. H. Zhang, "Construction of reactive potential energy surfaces with Gaussian process regression: Active data selection", *Mol. Phys.* **116**, 823–834 (2018).
- [103] Q. Lin, L. Zhang, Y. Zhang, and B. Jiang, "Searching Configurations in Uncertainty Space: Active Learning of High-Dimensional Neural Network Reactive Potentials", J. Chem. Theory Comput. 17, 2691–2701 (2021).
- [104] Z. Wang and J. Ye, "Querying Discriminative and Representative Samples for Batch Mode Active Learning", ACM Trans. Knowl. Discov. Data 9, 1–23 (2015).
- [105] C. Shui, F. Zhou, C. Gagné, and B. Wang, "Deep Active Learning: Unified and Principled Method for Query and Training", *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, ed. by S. Chiappa and R. Calandra, vol. 108, Proceedings of Machine Learning Research, PMLR, 2020, 1308–1318.
- [106] V. Vapnik, *The Nature of Statistical Learning Theory*, Information Science and Statistics, Springer New York, 1995.
- [107] A. Borisov, E. Tuv, and G. Runger, "Active batch learning with stochastic query-by-forest (SQBF)", Active Learning and Experimental Design workshop In conjunction with AISTATS 2010, PMLR, 2011, 59–69.
- [108] O. Akin-Ojo and K. Szalewicz, "Potential energy surface and second virial coefficient of methane-water from *ab initio* calculations", *J. Chem. Phys.* **123**, 134311 (2005).
- [109] M. P. Metz, K. Szalewicz, J. Sarka, R. Tóbiás, A. G. Császár, and E. Mátyus, "Molecular dimers of methane clathrates: *ab initio* potential energy surfaces and variational vibrational states", *Phys. Chem. Chem. Phys.* 21, 13504–13525 (2019).
- [110] Q. Ma and H.-J. Werner, "Accurate Intermolecular Interaction Energies Using Explicitly Correlated Local Coupled Cluster Methods [PNO-LCCSD(T)-F12]", J. Chem. Theory Comput. 15, 1044–1052 (2019).

- [111] M. P. Metz and K. Szalewicz, "Automatic Generation of Flexible-Monomer Intermolecular Potential Energy Surfaces", J. Chem. Theory Comput. 16, 2317–2339 (2020).
- [112] S. N. Vogels, T. Karman, J. Kłos, M. Besemer, J. Onvlee, A. van der Avoird, G. C. Groenenboom, and S. Y. T. van de Meerakker, "Scattering resonances in bimolecular collisions between NO radicals and H<sub>2</sub> challenge the theoretical gold standard", *Nat. Chem.* **10**, 435–440 (2018).
- [113] B. Brauer, M. K. Kesharwani, S. Kozuch, and J. M. L. Martin, "The S66×8 benchmark for noncovalent interactions revisited: explicitly correlated *ab initio* methods and density functional theory", *Phys. Chem. Chem. Phys.* 18, 20905–20925 (2016).
- [114] P. L. Bartlett and S. Mendelson, "Rademacher and Gaussian complexities: Risk bounds and structural results", J. Mach. Learn. Res. 3, 463–482 (2002).
- [115] B. K. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf, and G. R. Lanckriet, "On integral probability metrics, φ-divergences and binary classification", (2009), arXiv: 0901.2698 [cs].
- [116] A. Müller, "Integral probability metrics and their generating classes of functions", Adv. Appl. Probab. 29, 429–443 (1997).
- [117] S. T. Rachev, Probability metrics and the stability of stochastic models, vol. 269, Wiley, 1991.
- [118] R. M. Dudley, *Real Analysis and Probability*, 2nd ed., Cambridge Studies in Advanced Mathematics, Cambridge University Press, 2002.
- [119] S. Dasgupta, "Two faces of active learning", *Theor. Comput. Sci.* 412, 1767– 1781 (2011).
- [120] Y. Freund, H. S. Seung, E. Shamir, and N. Tishby, "Information, prediction, and query by committee", *Advances in neural information processing systems*, Morgan Kaufmann publishers, 1993, 483–490.

- [121] Y. Saleh, V. Sanjay, A. Iske, A. Yachmenev, and J. Küpper, "Active learning of potential-energy surfaces of weakly bound complexes with regressiontree ensembles", J. Chem. Phys. 155, 144109 (2021), arXiv: 2104.00708 [physics].
- [122] H. S. Seung, M. Opper, and H. Sompolinsky, "Query by committee", Proceedings of the fifth annual workshop on Computational learning theory, ACM, 1992, 287–294.
- [123] S. Kee, E. del Castillo, and G. Runger, "Query-by-committee improvement with diversity and density in batch active learning", *Inf. Sci.* 454-455, 401– 418 (2018).
- [124] M. J. Tubergen, A. M. Andrews, and R. L. Kuczkowski, "Microwave spectrum and structure of a hydrogen-bonded pyrrole-water complex", *J. Phys. Chem.* 97, 7451–7457 (1993).
- [125] U. Nygaard, J. Nielsen, J. Kirchheiner, G. Maltesen, C. S, J. Rastrup-Andersen, and G. Sørensen, "Microwave Spectra of isotopic pyrroles: Molecular structure, dipole moment and <sup>14</sup>N quadrupole coupling constants of pyrrole", J. Mol. Struct. 3, 491–506 (1969).
- [126] H.-J. Werner and F. R. Manby, "Explicitly correlated second-order perturbation theory using density fitting and local approximations", *J. Chem. Phys.* 124, 054114 (2006).
- [127] F. R. Manby, H.-J. Werner, T. B. Adler, and A. J. May, "Explicitly correlated local second-order perturbation theory with a frozen geminal correlation factor", J. Chem. Phys. 124, 094103 (2006).
- [128] H.-J. Werner, T. B. Adler, and F. R. Manby, "General orbital invariant MP2-F12 theory", J. Chem. Phys. 126, 164102 (2007).
- [129] K. A. Peterson, T. B. Adler, and H.-J. Werner, "Systematically convergent basis sets for explicitly correlated wavefunctions: The atoms H, He, B–Ne, and Al–Ar", J. Chem. Phys. 128, 084102 (2008).

- [130] R. A. Shaw and J. G. Hill, "Approaching the Hartree-Fock Limit through the Complementary Auxiliary Basis Set Singles Correction and Auxiliary Basis Sets", J. Chem. Theory Comput. 13, 1691–1698 (2017).
- [131] F. Weigend, A. Köhn, and C. Hättig, "Efficient use of the correlation consistent basis sets in resolution of the identity MP2 calculations", *J. Chem. Phys.* 116, 3175–3183 (2002).
- [132] H.-J. Werner, P. J. Knowles, F. R. Manby, J. A. Black, K. Doll, A. Heßelmann, D. Kats, A. Köhn, T. Korona, D. A. Kreplin, Q. Ma, T. F. Miller, A. Mitrushchenkov, K. A. Peterson, I. Polyak, G. Rauhut, and M. Sibaev, "The Molpro quantum chemistry package", J. Chem. Phys. 152, 144107 (2020).
- [133] H.-J. Werner, P. J. Knowles, G. Knizia, F. R. Manby, and M. Schütz, "Molpro: a general-purpose quantum chemistry program package", WIREs Comput. Mol. Sci. 2, 242–253 (2012).
- [134] H.-J. Werner, P. J. Knowles, G. Knizia, F. R. Manby, M. Schütz, P. Celani, W. Györffy, D. Kats, T. Korona, R. Lindh, A. Mitrushenkov, G. Rauhut, K. R. Shamasundar, T. B. Adler, R. D. Amos, S. J. Bennie, A. Bernhardsson, A. Berning, D. L. Cooper, M. J. O. Deegan, A. J. Dobbyn, F. Eckert, E. Goll, C. Hampel, A. Hesselmann, G. Hetzer, T. Hrenar, G. Jansen, C. Köppl, S. J. R. Lee, Y. Liu, A. W. Lloyd, Q. Ma, R. A. Mata, A. J. May, S. J. McNicholas, W. Meyer, T. F. Miller III, M. E. Mura, A. Nicklass, D. P. O'Neill, P. Palmieri, D. Peng, K. Pflüger, R. Pitzer, M. Reiher, T. Shiozaki, H. Stoll, A. J. Stone, R. Tarroni, T. Thorsteinsson, M. Wang, and M. Welborn, *MOLPRO, version*, *a package of ab initio programs*, see https://www.molpro.net, Stuttgart, Germany.
- [135] Y. Paukku, K. R. Yang, Z. Varga, and D. G. Truhlar, "Global ab initio groundstate potential energy surface of N<sub>4</sub>", *J. Chem. Phys.* **139**, 044309 (2013).
- [136] J. Cui and R. V. Krems, "Efficient non-parametric fitting of potential energy surfaces for polyatomic molecules with Gaussian processes", J. Phys. B: At. Mol. Opt. Phys. 49, 224001 (2016).

- [137] J. T. Ash, C. Zhang, A. Krishnamurthy, J. Langford, and A. Agarwal, "Deep Batch Active Learning by Diverse, Uncertain Gradient Lower Bounds", *International Conference on Learning Representations (ICLR)*, 2020.
- [138] K. A. Peterson, D. Feller, and D. A. Dixon, "Chemical accuracy in ab initio thermochemistry and spectroscopy: current strategies and future challenges", *Theor. Chem. Acc.* 131, 1079 (2012).
- [139] A. Owens, S. N. Yurchenko, A. Yachmenev, J. Tennyson, and W. Thiel, "Accurate *ab initio* vibrational energies of methyl chloride", *J. Chem. Phys.* 142, 244306 (2015).
- [140] A. Yachmenev, S. N. Yurchenko, T. Ribeyre, and W. Thiel, "High-level ab initio potential energy surfaces and vibrational energies of H<sub>2</sub>CS", J. Chem. Phys. 135, 074302 (2011).
- [141] P. O. Dral, A. Owens, A. Dral, and G. Csányi, "Hierarchical machine learning of potential energy surfaces", J. Chem. Phys. 152, 204110 (2020).
- [142] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. Von Lilienfeld, "Quantum chemistry structures and properties of 134 kilo molecules", *Sci. Data* 1, 140022 (2014).
- [143] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. Von Lilienfeld, "Fast and accurate modeling of molecular atomization energies with machine learning", *Phys. Rev. Lett.* **108**, 058301 (2012).
- [144] J. S. Smith, O. Isayev, and A. E. Roitberg, "ANI-1, A data set of 20 million calculated off-equilibrium conformations for organic molecules", *Sci. Data* 4, 170193 (2017).
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel,
   M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,
   D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn:
   Machine Learning in Python", J. Mach. Learn. Res. 12, 2825–2830 (2011).
- [146] Y.-Y. Yang, S.-C. Lee, Y.-A. Chung, T.-E. Wu, S.-A. Chen, and H.-T. Lin, *libact: Pool-based Active Learning in Python*, tech. rep., National Taiwan University, 2017, arXiv: 1710.00379 [cs].

- [147] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*, Software available from tensorflow.org, 2015.
- [148] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization", 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, ed. by Y. Bengio and Y. LeCun, 2015.
- S. N. Yurchenko, A. Yachmenev, and R. I. Ovsyannikov, "Symmetry adapted ro-vibrational basis functions for variational nuclear motion calculations: TROVE approach", *J. Chem. Theory Comput.* 13, 4368 (2017), arXiv: 1708. 07185 [physics].
- [150] A. Yachmenev and J. Küpper, "Communication: General variational approach to nuclear-quadrupole coupling in rovibrational spectra of polyatomic molecules", J. Chem. Phys. 147, 141101 (2017), arXiv: 1709.08558 [physics].
- [151] A. Yachmenev and S. N. Yurchenko, "Automatic differentiation method for numerical construction of the rotational-vibrational Hamiltonian as a power series in the curvilinear internal coordinates using the Eckart frame", *J. Chem. Phys.* **143**, 014105 (2015).
- [152] S. N. Yurchenko, W. Thiel, and P. Jensen, "Theoretical ROVibrational Energies (TROVE): A robust numerical approach to the calculation of rovibrational energies for polyatomic molecules", J. Mol. Spectrosc. 245, 126–140 (2007).
- [153] J. M. Bowman, S. Carter, and X. Huang, "MULTIMODE: A code to calculate rovibrational energies of polyatomic molecules", *Int. Rev. Phys. Chem.* 22, 533–549 (2003).

- [154] M. Rey, A. V. Nikitin, and V. G. Tyuterev, "Complete nuclear motion Hamiltonian in the irreducible normal mode tensor operator formalism for the methane molecule", J. Chem. Phys. 136, 244106 (2012).
- [155] E. Mátyus, G. Czakó, and A. G. Császár, "Toward black-box-type full- and reduced-dimensional variational (ro)vibrational computations", J. Chem. Phys. 130, 134112 (2009).
- [156] G. Avila and E. Mátyus, "Toward breaking the curse of dimensionality in (ro)vibrational computations of molecular systems with multiple largeamplitude motions", J. Chem. Phys. 150, 174107 (2019).
- [157] X.-G. Wang and T. Carrington, "Computing rovibrational levels of methane with curvilinear internal vibrational coordinates and an Eckart frame", J. *Chem. Phys.* 138, 104106 (2013).
- [158] G. Avila and T. Carrington, "Solving the Schrödinger equation using Smolyak interpolants", J. Chem. Phys. 139, 134114 (2013).
- [159] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning", *Nature* **521**, 436–444 (2015).
- [160] A. Cuzzocrea, A. Scemama, W. J. Briels, S. Moroni, and C. Filippi, "Variational principles in quantum Monte Carlo: The troubled story of variance minimization", J. Chem. Theory Comput. 16, 4203–4212 (2020).
- [161] G. Cybenko, "Approximation by superpositions of a sigmoidal function", Math. Control Signals Syst. 2, 303–314 (1989).
- [162] J. W. Siegel, "Optimal Approximation Rates for Deep ReLU Neural Networks on Sobolev Spaces", (2022), arXiv: 2211.14400 [cs].
- [163] K. Cranmer, S. Golkar, and D. Pappadopulo, *Inferring the quantum density matrix with machine learning*, 2019, arXiv: 1904.05903 [physics].
- [164] I. Kobyzev, S. J. Prince, and M. A. Brubaker, "Normalizing flows: An introduction and review of current methods", *IEEE Trans. Pattern Anal. Mach. Intell.* 43, 3964–3979 (2020).

- [165] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan, "Normalizing Flows for Probabilistic Modeling and Inference", *J. Mach. Learn. Res.* 22, (2022).
- [166] J. Behrmann, W. Grathwohl, R. T. Q. Chen, D. Duvenaud, and J.-H. Jacobsen, "Invertible Residual Networks", *Proceedings of the 36th International Conference on Machine Learning*, ed. by K. Chaudhuri and R. Salakhutdinov, vol. 97, Proceedings of Machine Learning Research, PMLR, 2019, 573–582.
- [167] A. Kechris, Classical descriptive set theory, vol. 156, Springer Science & Business Media, 2012.
- [168] Y. Saleh, A. Iske, A. Yachmenev, and J. Küpper, "Augmenting basis sets by normalizing flows", *Proc. Appl. Math. Mech.* 23, e202200239 (2023), arXiv: 2212.01383 [math].
- [169] A. Verine, B. Negrevergne, F. Rossi, and Y. Chevaleyre, "On the expressivity of bi-Lipschitz normalizing flows", (2021), arXiv: 2107.07232 [stat].
- [170] R. Cornish, A. Caterini, G. Deligiannidis, and A. Doucet, "Relaxing bijectivity constraints with continuously indexed normalising flows", *International conference on machine learning*, PMLR, 2020, 2133–2143.
- [171] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky, "Spectrally-normalized margin bounds for neural networks", *Advances in Neural Information Processing Systems*, ed. by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, vol. 30, Curran Associates, Inc., 2017.
- [172] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks", (2018), arXiv: 1802.05957 [stat].
- [173] R. K. Singh and J. S. Manhas, Composition operators on function spaces, Elsevier, 1993.
- [174] R. K. Singh, "Invertible composition operators on  $L^2(\lambda)$ ", *Proc. Amer. Math. Soc.* **56**, 127–129 (1976).
- [175] A. Galbis and E. Jordá, "Composition operators on the Schwartz space", *Rev. Mat. Iberoam.* 34, 397–412 (2018).

- [176] G. Vainikko, "Evaluation of the error of the Bubnov Galerkin method in an eigenvalue problem", *USSR Comput. Math. Math. Phys.* **5**, 1–31 (1965).
- [177] C. Villani, *Topics in Optimal Transportation*, Providence, RI: American Mathematical Society, 2003.
- [178] A. Uzman, Existence of a mapping that induces a certain push-forward measure, Mathematics Stack Exchange, (accessed February 17, 2023).
- [179] O. Diaz, Existence of a mapping that induces a certain push-forward measure, Mathematics Stack Exchange, (accessed: February 17, 2023).
- [180] Y. Saleh, Existence of a mapping that induces a certain push-forward measure, Mathematics Stack Exchange, (accessed February 17, 2023).
- [181] Shalop, Does weak convergence with uniformly bounded densities imply absolute continuity of the limit?, Mathematics Stack Exchange, (accessed February 15, 2023).
- [182] Y. Saleh, Alvaro Fernàndez-Corral, A. Iske, A. Yachmenev, and J. Küpper, "Normalizing flows for modeling excited states of molecular Schrödinger equations", (2023).
- [183] Z. Kong and K. Chaudhuri, "Universal Approximation of Residual Flows in Maximum Mean Discrepancy", (2021), arXiv: 2103.05793 [stat, cs].
- [184] W. Yang and A. C. Peet, "The collocation method for bound solutions of the Schrödinger equation", *Chem. Phys. Lett.* **153**, 98–104 (1988).
- [185] G Constantine and T Savits, "A multivariate Faa di Bruno formula with applications", *Trans. Amer. Math. Soc.* **348**, 503–520 (1996).
- [186] M. S. Birman and M. Z. Solomjak, Spectral theory of self-adjoint operators in Hilbert space, vol. 5, Springer Science & Business Media, 2012.
- [187] G. James, D. Witten, T. Hastie, and R. Tibshirani, An introduction to statistical learning, vol. 112, Springer, 2013.
- [188] H. Laurent and R. L. Rivest, "Constructing optimal binary decision trees is NP-complete", *Inf. Process. Lett.* 5, 15–17 (1976).
- [189] L. Breiman, "Random forests", Mach. Learn. 45, 5–32 (2001).

- [190] P. Geurts, D. Ernst, and L. Wehenkel, "Extremely randomized trees", *Mach. Learn.* **63**, 3–42 (2006).
- [191] G. Louppe, "Understanding random forests", PhD thesis, University of Liège, 2014, arXiv: 1407.7502 [stat.ML].
- [192] I. H. Witten, E. Frank, and M. A. Hall, *Data Mining: Practical Machine Learning Tools and Techniques*, Morgan Kaufmann Series in Data Management Systems, Morgan Kaufmann, 2017.
- [193] R. L. Schilling, *Measures, integrals and martingales*, Cambridge University Press, 2017.
- [194] A. Klenke, *Probability theory: a comprehensive course*, Springer Science & Business Media, 2013.

## List of illustrations

### Algorithms

1	Basic steps of a generic pool-based active learning strategy.	35
2	Query by committee algorithm.	36
3	Stochastic query by committee algorithm	37

## Figures

3.1	Intermolecular coordinates	38
3.2	Probability density distribution of the energies of $\ensuremath{\text{pyrrole}}(H_2O)$	38
3.3	Root-mean-square error for random forest regressor and neural	
	network models for different active learning algorithms	42
3.4	Normalized probability density distributions for data collected by	
	different active learning algorithms	44
3.5	2D histograms of discrepancies between the predictions of random	
	forest regressor and neural network models for different active	
	learning algorithms	46
3.6	Root-mean-square error for random forest regressor and neural	
	network models for different active learning algorithms - effect of	
	the size of the batch	47

3.7	Root-mean-square error for random forest regressor and neural	
	network models for different active learning algorithms - effect of	
	the size of the initially-labelled data	48
4.1	Hermite functions	56
4.2	Convergence of bands of vibrational energies for H <sub>2</sub> S molecule as a	
	function of the polyad number $N$	85
4.3	Error in the approximate vibrational energies of $H_2S$ as a function	
	of the polyad number $N. \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$	86
4.4	Training loss for computing vibrational spectra of $H_2S$	87

### Tables

2.1	Size of basis to ensure convergence for an increasing problem	
	dimensionality	20
2.2	Size of basis necessary to converge excited states	21
3.1	Out-of-the-pool root-mean-square errors for $pyrrole(H_2O)$ for a	
	random forest regressor and neural models for different active	
	learning algorithms	41
3.2	Root-mean-square mean errors and standard deviations for $\mathrm{N}_4$	
	molecule of a neural network model	43
4.1	Size of the index set ${\mathcal N}$ for different values of the polyad number.	82
4.2	Discrepancy between vibrational calculations of H <sub>2</sub> S via an	
	augmented basis and TROVE calculations	93
4.3	Absolute error in the dipole moment calculations for $H_2S$	94

## Acronyms

- AL active learning (p. 26)
  BO Born-Oppenheimer (p. 13)
  GP Gaussian process (p. 25)
  ML machine learning (p. 25)
  NN neural network (p. 25)
  OOP out of the pool (p. 25)
  PES potential energy surface (p. 4)
  PL passive learning (p. 26)
- **QBC** query by committee (p. 35)
- ResNet residual neural network (p. 60)
- **RFR** random forest regressor (p. 40)
- **RS** random sampling (p. 35)
- SQBF stochastic query by forest (p. 39)
- **TISE** time-independent Schrödinger equation (p. 10)

## Notation and terminology

#### Sets

- (i)  $\mathbb{N}_{\geq 0}$  = set of non-negative integers.  $\mathbb{N}_{\geq 0}^d$  = the *d*-fold Cartesian product of  $\mathbb{N}_{>0}$ . Similarly, denote by  $\mathbb{N}_{>0}$  the set of positive integers.
- (ii) Given a set D, denote by |D| its cardinality.
- (iii) Given two sets *A*, *B*, denote by  $A \setminus B$ , and  $A \cup B$  the set difference between, and the union of *A*, *B*, respectively.  $A \times B$  denotes the Cartesian product.

#### Enumeration

- (i) I use enumeration over  $\mathbb{N}_{>0}$  for sequences, i. e., I write,  $(a_n)_{n \in \mathbb{N}_{n>0}}$  for arbitrary sequences. The indexing is, however, left off for notational simplicity, i. e., I write  $(a_n)_n$ .
- (ii) I use  $\sum_n$  to denote  $\sum_{n \in \mathbb{N}_{>0}}$ .
- (iii) Sometimes I define sequences indexed by an index set *N*. I write  $(a_n)_{n \in N}$ .
- (iv) I sometimes call a sequence of elements a family of elements.

#### Matrices

- (i) tr(A) = trace of a matrix A.
- (ii)  $\det A = \det A$  a matrix A.

- (iii)  $A^T$  = transpose of a matrix A.
- (iv)  $||A||_{\mathcal{F}} =$  Forbenius norm of a matrix, i. e.,

$$\|A\|_{\mathcal{F}} = \sqrt{\sum_{i \le N} \sum_{j \le M} |a_{ij}|^2}$$

, where *N*, *M* denote the number of rows, columns in *A*, respectively.

#### Geometry

- (i)  $\mathbb{R}^d = d \text{dimensional real Euclidean space. } \mathbb{R}_{\geq 0} = \text{set of non-negative real numbers. } \mathbb{R}^d_{\geq 0} = \text{the set of positive real numbers.}$
- (ii)  $e_i = (0, ..., 0, 1, ..., 0) = ith$  standard coordinate vector.
- (iii) A typical point in  $\mathbb{R}^d$  is  $x = (x_1, \dots, x_d)$ .
- (iv) For an open  $\Omega \in \mathbb{R}^d$  denote by  $\partial \Omega$  its boundary.  $\overline{\Omega} = \Omega \cup \partial \Omega =$ closure of  $\Omega$ .
- (v)  $B_r(x) = \{y \in \mathbb{R}^d \mid |x y| < r\}$  = open ball in  $\mathbb{R}^d$  with center x and radius r > 0.
- (vi) Given  $a = (a_1, \ldots, a_d)$  and  $b = (b_1, \ldots, b_d)$  in  $\mathbb{R}^d$  set

$$a.b = \sum_{i=1}^{d} a_i b_i, \ |a| = \left(\sum_{i=1}^{d} a_i^2\right)^{1/2}$$

#### **Differential operators**

Assume  $f : \Omega \to \mathbb{R}^d$ .

- (i)  $\frac{\partial f}{\partial x_i}(x) = \lim_{h \to 0} \frac{f(x+he_i) f(x)}{h}$ , provided the limit exists.
- (ii) Multi-index notation:

• A vector  $\alpha = (\alpha_1, \dots, \alpha_d), \alpha \in \mathbb{N}^d$  is called a *multi-index* of order

$$|\alpha| = \sum_{i=1}^d \alpha_i.$$

Set

$$\alpha! = \prod_{n=1}^d (\alpha_n!).$$

Given  $l = (l_1, \ldots, l_d) \in \mathbb{R}^d$  set

$$l^{\alpha} = \prod_{n=1}^{d} l_n^{\alpha_n}$$

• given a multi-index  $\alpha$ 

$$D^{\alpha}f(x) = \frac{\partial^{|\alpha|}f(x)}{\partial x_1^{\alpha_1}\dots \partial x_d^{\alpha_d}} = \frac{\partial^{\alpha_1}}{\partial x_1^{\alpha_1}}\dots \frac{\partial^{\alpha_d}}{\partial x_d^{\alpha_d}}f(x),$$

• if  $k \in \mathbb{N}$ 

$$D^k f(x) := \{ D^\alpha f(x) \mid |\alpha| = k \},\$$

the set of all partial derivatives of order *K*.

- $|D^k f| = \left( \sum_{|\alpha|=k} |D^{\alpha} f|^2 | \right)^{1/2}.$
- If k = 1 consider the elements of Df to be arranged in a vector, i. e.,

$$Df := (\frac{df}{dx_1}, \dots, \frac{df}{dx_d}) =$$
 gradient vector.

Similarly, if k = 2 regard the elements of  $D^2 f$  as elements of a matrix, i. e., the *Hessian matrix*.

•  $\Delta f = \operatorname{tr}(D^2 f) = \operatorname{Laplacian} \operatorname{of} f$ .

#### Functions and functional spaces

- (i) For an open  $\Omega \subseteq \mathbb{R}^d$ , and two functions  $f, g : \Omega \to \mathbb{R}$  I write  $f \equiv g$  to mean that f is identically equal to v, i. e., the two functions agree for all values of their arguments.
- (ii) For an open  $\Omega \subseteq \mathbb{R}^d$ , and two functions  $f : \Omega \to \mathbb{R}$ , and  $g : \Omega \to \Omega$ ,  $f \circ g$  denotes the composition of f with g.  $f \equiv g$  means that f is identically equal to v, i. e., the two functions agree for all values of their arguments.
- (iii) A function *f* is said to be *smooth* provided it is infinitely differentiable.
- (iv) Denote by supp(f) the support of a function f.
- (v) For an open  $\Omega \in \mathbb{R}^d$  function  $f : \Omega \times \Omega \to \mathbb{R}$ ,  $(x, y) \mapsto f(x, y)$  and a fixed  $y \in \Omega$  I write  $f(\cdot, y)$  to denote the mapping  $\Omega \ni x \mapsto f(x, y) \in \mathbb{R}$ .
- (vi) Let  $\Omega \subseteq \mathbb{R}^d$  be open. I use the following notation for functional spaces.
  - $C^k(\Omega) = \{f : \Omega \to \mathbb{R}^d \mid f \text{ is } k \text{times continuously differentiable} \}.$
  - $C_b(\Omega) = \{f : \Omega \to \mathbb{R}^d \mid f \text{ is continuous and bounded}\}.$
  - $C^{\infty}(\Omega) = \{f : \Omega \to \mathbb{R}^d \mid f \text{ is infinitely differentiable}\}.$
  - $C_c^{\infty}(\Omega) = \{ f : \Omega \to \mathbb{R}^d \mid f \in C^{\infty} \text{and } f \text{ is compactly supported} \}.$
  - Consider a measure space (Ω, A, µ). Define

$$L^{p}(\mu) := \{f : \Omega \to \mathbb{R} \mid f \text{ is } \mu - \text{measurable, } \|f\|_{L^{p}(\Omega)}\} < \infty,$$

where

$$||f||_{L^{p}(\mu)} := \left(\int_{\Omega} |f|^{p} d\mu\right)^{1/p} \ (1 \le p < \infty).$$

 $L^{\infty}(\mu):=\{f:\Omega\to\mathbb{R}\mid f\text{ is Lebesgue measurable, }\|f\|_{L^{\infty}(\mu)}<\infty\},$  where

$$||f||_{L^{\infty}(\mu)} := \inf\{c > 0 \mid \mu(\{|u| \ge c\}) = 0\}.$$

When the choice of either the measure  $\mu$  is clear the explicit dependence of  $L^p$  spaces thereon is omitted, i. e., sometimes I write,  $||f||_{L^p}$ . Similarly, I sometimes write  $L^p(\mathcal{A})$  or  $L^p(\Omega)$  if I intend to stress the underlying space or  $\sigma$ -algebra.  $H^k(\Omega)$ : the Sobolev space of functions that are  $L^2$  together with their weak derivatives up to the *k*th order.

• 
$$S = \{f : \Omega \to \mathbb{R} \mid f \in C^{\infty}, \|f\|_{a,b} < \infty \text{ for all } a, b \in \mathcal{N}^d\}, \text{ where}$$
$$\|f\|_{a,b} = \sup_{x} |x^a(D^b f)(x)|.$$

 $\mathcal{S}$  is called the space of Schwartz functions.

(vii) Given a measure space  $(X, \Sigma, \mu)$  and a measurable function *f* denote by

$$\int_X f \, d\mu$$

the integral of *f* over X with respect to the measure  $\mu$ .

(viii) For a set *E* denote by  $\chi_E$  the *indicator function* of *E*, i. e.,

$$\chi_E(x) = \begin{cases} 1 \text{ if } x \in E \\ 0 \text{ if } x \notin E \end{cases}$$

#### **Operators**

Let *T* be a linear operator between two vector spaces *X*, *Y*.

- (i) Denote by D(T) its domain.
- (ii) ker(T) = kernel of T, i. e.,

$$\ker(T) = \{ v \in X \mid Tx = 0 \}$$

#### Measure and probability theory

I generally use calligraphic letters to denote  $\sigma$ -algebras. By  $\mathcal{B}(X)$  denote the Borel  $\sigma$ -algebra over a set X. Unless otherwise specified denote by  $\mu$  the Lebesgue measure.

Let (*X*, A,  $\mu$ ) denote a  $\sigma$ -finite measure space.

- (i) Given a measurable mapping  $T : X \to X$  denote by  $T_{\#}X$  the push-forward measure.
- (ii) Given another *σ*−finite measure space (*Y*, *B*, *µ*) denote by *µ*<sup>2</sup> the product measure defined on the measurable space (*X* × *Y*, *A* ⊗ *B*), where *A* ⊗ *B* denotes the product *σ*−algebra. Similarly, *µ<sup>m</sup>* for *m* ∈ ℕ<sub>>0</sub> is the product measure on the *m*− product measurable spaces. When the dimensionality of the underlying space is clear, or not relevant for the discussion, I simply write *µ* for *µ<sup>m</sup>*.

Let  $(X, \mathcal{A}, \mathcal{P})$  be a probability measure space. Generally, I use small letters to denote random variables (Definition D.10). Given a random variable x, I use  $\mathcal{P}_x$  to denote the probability distribution of x. Write  $\hat{x}$  to denote a finite dataset of evaluations that have distributions  $\mathcal{P}_x$ . I refer to elements of  $\hat{x}$  by *observations* or (*training*) *examples*, in accordance with machine learning literature. I write  $\hat{x} \sim \mathcal{P}_x$  to say that the observations have distribution  $\mathcal{P}_x$ .

Denote the expected value of *x* by  $\mathfrak{E}_{\sim \mathcal{P}}[x]$ , i. e.,

$$\mathfrak{E}_{\sim \mathcal{P}}[x] = \int_X x \, d\mathcal{P}.$$

#### Machine learning

In Chapter 3 and Chapter 4 I denote by  $\mathcal{H}$  the hypothesis class. By this I mean the set of all admissible functions that is used for learning, e.g., for a neural network of a fixed architecture,  $\mathcal{H}$  denotes all possible functions that one obtains for different values of the parameters.

## Quantum physics

Following quantum physics' terminology I sometimes refer to eigenvalues and eigenfunctions of Hamiltonians (2.2) by *energies* and *eigenfunctions*, respectively.

# A statement on data availability and reproducibility

The data and codes for reproducing results on actively learning potential energy surfaces, originally published in [121] and reported in Chapter 3, are available at https://github.com/CFEL-CMI/Active-Learning-of-PES.

The codes for reproducing results on augmenting bases for solving quantum problems, originally published in [168], are available at https://github.com/CFEL-CMI/FlowBasis.

# Publications and conference contributions

The following is a complete and chronologically ordered list of publications and conference contributions related to the present thesis.

#### **Publications**

- Y. Saleh, A. Iske, A. Yachmenev, and J. Küpper, Spectral learning for solving differential equations, in preparation (2023).
- Y. Saleh, A. F. Corral, A. Iske, A. Yachmenev, and J. Küpper, Normalizing flows for modelling excited states of molecular Schrödinger equations, in preparation (2023).
- Y. Saleh, A. Iske, A. Yachmenev, and J. Küpper, Augmenting basis sets by normalizing flows, *Proc. Appl. Math. Mech.* 23 (1), e202200239 (2023).
- Y. Saleh, V. Sanjay, A. Iske, A. Yachmenev, and J. Küpper, Active learning of potential-energy surfaces of weakly bound complexes with regression-tree ensembles, *J. Chem. Phys.* **155**, 144109 (2021).

#### **Conference contributions**

- Deutsche Physikalische Gesellschaft (DPG) SAMOP Frühjahrstagung, *invited talk* (2023).
- Machine learning in engineering summer school (MLE), TU Hamburg, *invited lecture* (2022).

- Conference on computational methods in applied mathematics, TU Wien, *contributed talk* (2022).
- Gesellschaft für angewandte Mathematik und Mechanik (GAMM) Jahrestagung, TU Aachen, *contributed talk* (2022).
- Hausdorff center for mathematics workshop: synergies between data science and PDE analysis, Universität Bonn, *contributed talk* (2022).
- Hausdorff school: foundational methods in machine learning (2022).
- Opening symposium of the Center for Data and Computing in Natural Science (CDCS), *poster* (2022).
- Deutsche Physikalische Gesellschaft (DPG) SAMOP Frühjahrstagung, contributed talk (2022).
- European CFEL and DESY photon science users' meeting, *poster* (2022).
- Helmholtz H3 hackathon (2021).
- Warsaw summer school for quantum physics and chemistry, *poster*, University of Warsaw, *poster* (2021).
- Machine learning for quantum X, online, *invited talk* (2021).
- Bunsentagung, *contributed talk*, online (2021).

## Index

approximation problem, 29, 56 well-posed, 56 augmented sequence, 61 Barron space, 5, 55 basis, 6, 60 augmented, 62 direct product, 68 orthonormal, 56 primitive, 81 truncated, 66 Born-Oppenheimer approximation, 3, 10 chemical reaction, 1 composition operator, 61 configuration space, 27 curse of dimensionality, 4, 19, 54 Dirac ladder operator, 69 dissociation dynamics, 3 eigenpair, 9 electric dipole moment, 86 empirical risk, 29

empirical risk minimization principle, 6, 26 active learning, 31 passive learning, 30 energy/eigenvalue,9 generalization error, 28 generating distribution, 25 ground state, 15 Hamiltonian, 9 bounded, 80 compact, 72 electronic, 13 kinetic energy operator, 12 linear, 12 molecular, 10 potential energy operator, 12 self-adjoint, 11 symmetric, 11 unbounded, 12 harmonic oscillator, 20 Hermite functions, 55 highly-oscillatory functions, 74 hot exoplanet, 3

hypothesis class, 15, 25 imaging experiment, 2 infinite dimensional eigenvalue problem, 3, 9 informative, 34 integral probability metric, 31 Kantorovich metric, 33 learning algorithm, 5 active, 6 fragile, 54 passive, 26 spectral, 6, 67 supervised, 6 linear convergence, 58 Lipswish nonlinearity, 60 loss function, 28 machine learning model, 25 Gaussian process, 25 kernel method, 25 neural network, 25 polynomial, 25 random forest regressor, 40 residual neural network, 60 mapping, 61 bi-Lipschitz, 63 bi-measurable, 78 bijective, 61 differentiable, 63 invertible, 78 non-sigular, 61

measure, 27, 57  $\sigma$ -finite. 60 absolutely continuous, 79 finite, 51 induced, 78 Lebesgue, 27, 57 probability, 31, 78 push-forward, 61 Radon, 77 normalizing flow, 6, 57 invertible residual neural network, 60 nuclear/molecular geometry, 15 polyad number, 82 polyatomic molecule, 16  $pyrrole(H_2O), 37$ H<sub>2</sub>S, 81  $N_4, 40$ pool, 34 potential energy surface, 4, 15 pump-probe, 1 quantum chemistry, 13 quantum state/eigenfunction, 9 quantum-molecular movie, 1 query algorithm, 26 query by committee, 35 random sampling, 35 stochastic query by forest, 40 Rademacher complexity, 29

Radon-Nikodym derivative, 61 random variable, 15 random vector, 27 Rayleigh-Ritz principle, 19 regression, 15 representative, 34 representativeness, 29 reproducing kernel Hilbert space, 33 rotational motion, 21 spectral method, 6, 66 Bubnov-Galerkin, 66 Tau, 67 strong formulation, 66 symbol for the Schwartz space, 70 test function, 66 time-dependent Schrödinger equation, g time-independent Schrödinger equation, 9 electronic, 4, 13 nuclear, 16 total variation, 77 transition state, 1 translational motion, 21 true risk, 28 ultrafast, 1 uncertainty, 35 variational principle, 17 vibrational motion, 21

Wasserstein distance, 33 weak formulation, 16, 66 weak—\* convergence, 79 weakly-bound molecule, 27

## **Eidesstattliche Versicherung**

Hiermit versichere ich an Eides statt, dass ich die vorliegende Dissertationsschrift selbst verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere, dass die eingereichte schriftliche Fassung mit der elektronischen Fassung der Dissertation übereinstimmt. Die Dissertation wurde in der vorgelegten oder einer ähnlichen Form nicht schon einmal in einem früheren Promotionsverfahren angenommen oder als ungenügen beurteilt.

Hamburg, den 20. März, 2023

Yahya SALEH