

SOLVING THE QUANTUM MANY-BODY PROBLEM
WITH NEURAL-NETWORK QUANTUM STATES

By

Jane Mee Kim

A DISSERTATION

Submitted to
Michigan State University
in partial fulfillment of the requirements
for the degree of

Physics — Doctor of Philosophy
Computational Math, Science and Engineering — Dual Major

2023

ABSTRACT

SOLVING THE QUANTUM MANY-BODY PROBLEM WITH NEURAL-NETWORK QUANTUM STATES

By

Jane Mee Kim

Quantum many-body systems provide a rich framework for exploring and understanding the fundamental laws of physics. By studying the collective behavior and emergent phenomena that arise from the intricate microscopic interactions among particles, we can deepen our understanding of quantum mechanics and gain insight into its broader implications for macroscopic observations. However, the exponential scaling of the Hilbert space dimension coupled with the presence of strong interparticle correlations make the accurate description of these systems difficult, if not impossible, for traditional analytical or perturbative techniques.

In this interdisciplinary approach, we aim to solve the quantum many-body problem by utilizing neural-network quantum states. These states, as their name suggests, are rooted in artificial neural networks and serve as a novel alternative to conventional parameterizations of trial wave functions within a variational Monte Carlo calculation. In addition to reviewing key concepts in quantum many-body theory and machine learning, we investigate a diverse set of continuous-space systems with varying levels of complexity. Starting from a pedagogical overview of an exactly solvable system of bosons in one dimension, we work our way up to strongly-interacting fermionic systems in three dimensions, including dilute neutron matter and ultra-cold Fermi gases. The highly non-perturbative interactions featured in these systems motivate the development of innovative neural-network quantum states, capable of discovering strong correlations while maintaining required symmetries and boundary

conditions. We accompany this study with the description of two distinct implementations of neural-network quantum states, each with their unique goals and strategies. Our findings indicate that neural-network quantum states provide a powerful and flexible strategy for investigating a wide range of quantum systems without relying on prior knowledge.

This thesis is dedicated to my parents, who have been a constant source of love and support despite their distance, and to Elliot, my life partner and best friend, who fills my days with joy and laughter.

ACKNOWLEDGMENTS

Meeting Morten Hjorth-Jensen during my senior year at MSU was an incredibly fortunate event that profoundly influenced my academic journey. Knowing he would make an excellent mentor, he was a large part of the reason I decided to stay at MSU to pursue my PhD. His expertise, patience, and dedication have played a pivotal role in fostering my passion and shaping my scientific identity. Morten, I am forever grateful for your kindness, as well as your ability to strike a delicate balance between providing steadfast support and encouraging independent thinking. Please know that you have had an immense impact on my growth as a researcher and a human being.

Scott Pratt, I distinctly remember approaching your office one day after encountering the statement " $x = x + 1$ " and thinking to myself, "Well, that's simply not true!" It seems trivial now, but this moment marked my very first introduction to coding. What's more, I knew nothing about research at the time, but you convinced me to give it a shot. Thank you for taking a chance on me, and thank you for your many words of wisdom.

To my esteemed committee members, Filomena Nunes, Artemis Spyrou, Huey-Wen Lin, Michael Murillo, and Stuart Tessmer (and a special mention to Matthew Hirn, who is no longer at MSU), I am grateful for your support and guidance throughout the years. You have not only fulfilled your roles as committee members but also served as invaluable role models for me. I would also like to express my gratitude to the nuclear theory group, particularly Dean Lee, Scott Bogner, and Heiko Hergert, for the countless stimulating discussions we have had over the years. I hope that our paths will cross often in the future.

Alessandro Lovato and Bryce Fore, it has been an absolute delight collaborating with both of you. I have learned so much during this past year and a half, and I eagerly look

forward to our continued collaborations in the future. Gabriel Pescia, Jannes Nys, and Giuseppe Carleo, I am deeply appreciative of the opportunity to work alongside each of you and to have had the unforgettable experience of meeting in person. The time spent together in Lausanne holds a special place in my heart, and I hope to have the chance to visit again in the future.

Throughout my nearly decade-long journey at MSU, my interests and career aspirations have undergone complete transformations. In the face of these changes, I am deeply grateful for the presence of Kim Crosslan in the P&A department. Kim, speaking not only for myself but for all graduate students, you are an incredibly important source of stability, and your impact is immeasurable. Thank you for the many times you have saved me.

My dearest Elliot, I can not imagine what graduate school would have been like if you had not been there for every second of it. Thank you for being you, and for helping me through stressful times. Umma and Appa, I am excited for you to finally see what I've been working on for the past few years. Thank you for always encouraging me to work through tough obstacles to strive towards my goals, hearing 'hwaiting!' is often exactly what I needed. Therese (Erin), thank you for your supportive words over the years. I miss you all deeply.

TABLE OF CONTENTS

LIST OF TABLES	x
LIST OF FIGURES	xi
Chapter 1 Introduction	1
1.1 Challenges of the Quantum Many-Body Problem	1
1.2 Neural-Network Quantum States	2
1.3 Algorithm Development Through the Lens of Nuclear Theory	2
1.4 Structure of Dissertation	3
Chapter 2 The Quantum Many-Body Problem	5
2.1 Formalism	6
2.2 Indistinguishable Particles	12
2.3 The Schrödinger Equation	17
2.4 The Variational Principle	20
2.5 Ab Initio Methods	23
2.5.1 Configuration Interaction	24
2.5.2 Hartree-Fock Theory	25
2.5.3 Many-Body Perturbation Theory	30
Chapter 3 Quantum Monte Carlo	34
3.1 Monte Carlo Methods	35
3.1.1 Markov Chain Monte Carlo Sampling	35
3.1.1.1 Metropolis-Hastings	37
3.1.1.2 Importance Sampling	41
3.1.2 Integration	45
3.2 Variational Monte Carlo	46
3.2.1 Trial Wave Functions	49
3.2.1.1 Kato's Cusp Condition	49
3.2.1.2 Jastrow Factor	50
3.2.1.3 Slater Determinant	54
3.2.1.4 Backflow Transformations	56
3.2.1.5 Bardeen-Cooper-Schreifer Wave Function	58
3.2.1.6 Pfaffian Wave Function	60
3.2.2 Optimization	62
3.2.2.1 Stochastic Gradient Descent	63
3.2.2.2 Momentum	64
3.2.2.3 Root Mean Squared Propagation	65
3.2.2.4 Adaptive Moment Estimation	65
3.2.2.5 Stochastic Reconfiguration	66
3.3 Diffusion Monte Carlo	68

Chapter 4 Machine Learning	72
4.1 The Curse of Dimensionality	72
4.2 Cost Functions	74
4.3 Supervised Learning	75
4.4 Unsupervised Learning	77
4.5 Reinforcement Learning	79
4.6 Transfer Learning	81
4.7 Artificial Neural Networks	82
4.7.1 Boltzmann Machines	82
4.7.1.1 Gaussian-Binary Restricted Boltzmann Machines	84
4.7.1.2 Multivariate Gaussian-Binary Restricted Boltzmann Machines	90
4.7.1.3 Multivariate Gaussian-Uniform Restricted Boltzmann Ma- chines	97
4.7.2 Feedforward Neural Networks	101
4.7.2.1 Backpropagation	102
4.7.2.2 Universal Approximation Theorem	104
4.7.3 Deep Sets	105
4.7.4 Graph Neural Networks	109
4.7.5 Attention Mechanisms	111
Chapter 5 Implementation	114
5.1 NeuralAnsatz: C++ Software for Localized Systems	114
5.1.1 Trainer	114
5.1.2 Systems	117
5.1.3 Wave Functions	120
5.1.3.1 Inputs	121
5.1.3.2 Continuous Restricted Boltzmann Machines	129
5.1.3.3 Feedforward Neural Networks	132
5.1.3.4 Deep Sets	136
5.1.3.5 Slater Determinants	137
5.1.3.6 Products	140
5.1.4 Hamiltonian	140
5.1.5 Sampler	142
5.1.6 Optimizer	142
5.1.7 Cost Functions	143
5.1.8 Density	145
5.1.9 Running NeuralAnsatz	145
5.2 Python Software for Periodic Systems	147
5.2.1 JAX: High-Performance Array Computing	147
5.2.1.1 Building Neural Networks	148
5.2.2 Symmetries and Boundary Conditions	149
5.2.3 Computing the Pfaffian	151
5.2.4 Stochastic Reconfiguration with RMSprop Regularization	152
Chapter 6 The Calogero-Sutherland Model	153

6.1	Hamiltonian	153
6.2	Neural-Network Quantum States	155
6.2.1	Feedforward Neural Networks	156
6.2.2	Continuous Restricted Boltzmann Machines	156
6.3	Results	157
6.4	Conclusions	164
Chapter 7	Dilute Neutron Matter	165
7.1	Introduction	166
7.2	Method	168
7.3	Results and Discussion	171
7.4	Conclusion	173
7.5	Acknowledgements	174
Chapter 8	Homogeneous Electron Gas	179
8.1	Introduction	180
8.2	Exact Backflow Transformations	182
8.3	Message-Passing Neural-Network Quantum States	184
8.4	Hamiltonian	186
8.5	Results	188
8.5.1	Energy Benchmarks for Small Systems	189
8.5.2	Energy Benchmarks for Large Systems	189
8.6	Conclusions	193
8.7	Acknowledgements	194
Chapter 9	Ultra-cold Fermi Gases	195
9.1	Introduction	196
9.2	Hamiltonian	200
9.3	Neural-Network Quantum States	201
9.3.1	Pfaffian-Jastrow-Backflow	202
9.3.2	Fixed-Node Slater-Jastrow	205
9.3.3	Slater-Jastrow-Backflow	206
9.3.4	Message-Passing Neural Network	210
9.4	Variational Monte Carlo and Training	214
9.5	Diffusion Monte Carlo	215
9.6	Results	217
9.6.1	Energy	217
9.6.2	Pair Distribution Functions	220
9.6.3	Pairing Gap	225
9.7	Conclusions and perspectives	226
9.8	Acknowledgments	228

LIST OF TABLES

Table 8.1: Total energy per particle in Hartree for unpolarized system of $N = 14$ particles. WAPNet and FermiNet are alternative NQS architectures optimized via VMC. We include FCI and DCD results as benchmarks from quantum chemistry.	192
Table 8.2: Total energy per particle in Hartree for the unpolarized system of $N = 54$ particles. We include FCI and DCD results as benchmarks from quantum chemistry.	192
Table 9.1: Energy per particle for various values of μ and the corresponding values of r_e . The values with asterisks (*) are extrapolations from the linear fits shown in Fig. 9.4. The parameter $v_0 = 1$ is fixed.	220
Table 9.2: Energies per particle and parameters for the two-body potential in Eq. (9.2) giving different scattering lengths with the same effective range $k_F r_e = 0.2$	223
Table 9.3: Energies per particle for different numbers of particles, with $k_F r_e = 0.2$. Values with asterisks (*) indicate translation-invariant calculations.	225
Table 9.4: Energies per particle for different numbers of particles, with $k_F r_e = 0.2$. Values with asterisks (*) indicate translation-invariant calculations.	226

LIST OF FIGURES

Figure 3.1: The cyclic workflow of the variational Monte Carlo algorithm.	48
Figure 3.2: Cartoon of filling single-particle states for bosons and fermions.	56
Figure 3.3: The effects of a symmetric Jastrow factor and a backflow transformation on the nodes of a fermionic wave function.	57
Figure 3.4: Cartoon of the simple stochastic gradient descent algorithm (blue) and the effect of including momentum (orange) near a narrow minimum.	64
Figure 3.5: Visualization of how the Stochastic Reconfiguration smoothens the energy landscape.	66
Figure 4.1: Cartoon of the curse of dimensionality; as the dimension increases, the density of a fixed number of data points decreases exponentially.	73
Figure 4.2: Comparison of a random data set (left) with no information, a data set with clear clusters (middle), and a data set with a clear lower-dimensional manifold (right). The latter two are examples of real-life data sets that contain information, i.e. structure.	74
Figure 4.3: The two main types of supervised learning: regression (left) and classification (right).	76
Figure 4.4: A depiction of principal component analysis (left) and clustering (right), two types of unsupervised learning.	77
Figure 4.5: A depiction of the variational Monte Carlo algorithm in the context of agent-based modeling, a commonly employed reinforcement learning framework.	80
Figure 4.6: A general Boltzmann machine in which all visible and hidden nodes are fully-connected.	83
Figure 4.7: A standard restricted Boltzmann machine, where the hidden nodes are binary and the visible nodes can be either binary or Gaussian. There are no connections between nodes within the same layer.	85
Figure 4.8: A multivariate Gaussian-binary restricted Boltzmann machine, a step between a standard Gaussian-binary restricted Boltzmann machine and a general Boltzmann machine.	92

Figure 4.9: The conditional probability of a single binary hidden node h_j activating, given the state of the visible nodes \mathbf{v} , for a multivariate Gaussian-binary restricted Boltzmann machine. The blue line represents the probability if the hidden nodes are allowed to take values of 0 and 1 (Eq. 4.32), while the orange line represents the probability if the hidden nodes are allowed to take values of -1 and 1 (Eq. 4.39). The x-axis is a transformation of the visible nodes given by Eq. 4.28.	96
Figure 4.10: The conditional probability of a single uniform node h_j activating, given the state of the visible nodes \mathbf{v} , for a multivariate Gaussian-uniform restricted Boltzmann machine. The blue line represents the probability if the uniform nodes are allowed to take values between 0 and 1 (Eq. 4.43), while the orange line represents the probability if the hidden nodes are allowed to take values of -1 and 1 (Eq. 4.46). The x-axis is a transformation of the visible nodes given by Eq. 4.28. We also plot the softplus function (Eq. 4.17) for comparison, as it exhibits similar character. Probability values greater than 1 are assumed to mean the hidden node h_j is always activated.	99
Figure 4.11: Comparison of the four functions we called $f(x)$ during our derivations of the multivariate Gaussian-binary and Gaussian-uniform restricted Boltzmann machines. These functions appear in the log-likelihoods $\log P(\mathbf{v})$.	100
Figure 4.12: A deep feedforward neural network with six inputs (blue) and three outputs (orange). There are two hidden layers with eight nodes each (gray).	103
Figure 4.13: An example of a permutation-invariant Jastrow wave function constructed from the single output of a feedforward neural network. To enforce permutation-invariance, the input vectors can be sorted according to some rule before concatenating them into a single vector.	105
Figure 4.14: A depiction of a permutation-invariant Jastrow wave function constructed from the single output of a Deep Set. As Deep Sets enforce permutation invariance at the architecture level, no prior sorting of the inputs is required.	107
Figure 4.15: One layer of a graph neural network that produces an output graph that preserves the structure of the original input graph. Three types of transformations can be considered: global (top), edge (middle), and vertex (bottom) transformations.	109
Figure 4.16: A visualization of the attention scores between a single element of a query vector and all the elements of a key vector. Larger attention scores are depicted with darker, bolder lines.	111

Figure 5.1: The highest-level structure of the <code>NeuralAnsatz</code> code, an object-oriented C++ software for neural-network quantum states.	115
Figure 5.2: A specific example of the hierarchical structure provided by abstract base classes and derived classes. The classes that appear lower on this graph inherit properties from the classes above it.	119
Figure 6.1: The regularized two-body potential for the Calogero-Sutherland model, where x_0 is the regularization parameter (Eq. 6.4). The long-range behavior of the original potential (black, dashed line) is preserved, but the height of the potential decreases with larger values of x_0 . Here, we have used an interaction parameter value of $\beta = 2$	155
Figure 6.2: Training curves for feedforward neural networks with different activation functions. In this initial pretraining stage, $x_0 = 0.5$	158
Figure 6.3: Training curves for the different variants of Gaussian-binary and Gaussian-uniform restricted Boltzmann machines, with $x_0 = 0.5$	160
Figure 6.4: The converged energy as a function of the regularization parameter x_0 for feedforward neural networks with different activation functions.	161
Figure 6.5: The converged energy as a function of the regularization parameter x_0 for the various Gaussian-binary and Gaussian-uniform restricted Boltzmann machines.	162
Figure 6.6: One-body densities during three different points during training: the initial stage, after pretraining with $x_0 = 0.5$, and after training with $x_0 = 0.01$. The black dashed line represents the one-body distribution from the exact ground state wave function.	163
Figure 7.1: NQS training data in neutron matter at $\rho = 0.04 \text{ fm}^{-3}$ (data points) compared with Hartree-Fock (dotted line), conventional VMC (dashed line), constrained-path ADMC (dash-dotted line) and unconstrained-path ADMC results (solid line).	176
Figure 7.2: Low-density neutron-matter $\not\pi$ EFT equation of state as obtained with the hidden-nucleon NQS for $\not\pi$ EFT potential "o" (blue circles) and $\not\pi$ EFT potential "a" (orange squares) compared with interactions which retain pion-exchange terms (green band). We see that the "o" potential is in excellent agreement with the π -full interactions while the "a" potential has a slightly stiffer equations of state due primarily to a more repulsive three-body force.177	

Figure 7.3: Spin-singlet and triplet two-body distribution functions at $\rho = 0.01 \text{ fm}^{-3}$ (panel a) and $\rho = 0.04 \text{ fm}^{-3}$ (panel b). The NQS calculations (solid symbols) are compared with non-interacting Fermi Gas results (solid lines). ¹⁷⁸	198
Figure 8.1: Energy differences between ground-state energies, obtained with other methods and with the MP-NQS, in units of the thermodynamic Hartree Fock energy, for various densities, polarizations, and system sizes. (Top) $N = 14, 19$, (Bottom): $N = 54$ particles. Error bars are too small to be visible for most densities. The corresponding numerical data can be found in the Supplemental Material.	190
Figure 8.2: Spin-averaged radial distribution function for the HEG system with $N = 128$ electrons at low densities ($r_s = 110, 200$). Error bars are smaller than the symbols. The unpolarized fluid is obtained from [110] for $r_s = 110$	191
Figure 9.1: A cartoon of the BCS-BEC crossover. Moving from left to right, the attractive interaction between opposite-spin fermions increases. However, in the BEC regime, the attraction binds the pairs so tightly that they behave as weakly repulsive bosons. The region between the weakly attractive Cooper pairs and the weakly repulsive dimers is known as the unitary limit.	197
Figure 9.2: Schematic representation of a message-passing neural network with T iterations. Dashed lines represent the concatenation operations, while solid lines represent the parameterized transformations (linear transformations and nonlinear feedforward neural networks). Messages, highlighted in pink, mediate the exchange of information between the one- and two-body streams, in blue. A yellow box indicates a single iteration of the network.	205
Figure 9.3: Ground-state energies per particle as a function of the MPNN depth T for the SJ-PW (blue squares), SJ-BF (orange circles), and PJ-BF (green triangles) ansätze. The interaction parameters are set to $v_0 = 1$ and $\mu = 5$, corresponding to an effective range of $r_e k_F = 0.4$. The DMC benchmark energies with and without pairing are displayed as solid and dashed lines, respectively.	218
Figure 9.4: Ground-state energies per particle as a function of the effective range. The DMC-BCS benchmark energies (blue circles) and the Pfaffian-Jastrow with backflow (PJ-BF) energies (orange triangles) are extrapolated to zero effective range using linear fits (dashed lines). The shaded regions are the error bands for the DMC-BCS and PJ-BF energies.	219

Figure 9.5: Opposite-spin pair densities as a function of small k_{Fr} at unitarity ($v_0 = 1$) and $\mu = 5$ (blue squares), $\mu = 10$ (orange circles), and $\mu = 20$ (green triangles).	221
Figure 9.6: Opposite-spin pair densities in the crossover region for the BCS phase $1/ak_F = -0.5$ (blue squares), unitarity $1/ak_F = 0$ (orange circles), and BEC phase $1/ak_F = 0.5$ (green triangles). The effective range of all cases are fixed $k_{Fre} = 0.2$. See Table 9.3 for the corresponding values of v_0 and μ	222
Figure 9.7: Upper panel: Energy per particle in the BCS-BEC crossover region as a function of the scattering length a for a fixed effective range $k_{Fre} = 0.2$. Lower panel: Difference between Pfaffian-Jastrow with backflow (PJ-BF) and DMC-BCS benchmark energies. See Table 9.3 for the corresponding values of v_0 and μ	224

Chapter 1

Introduction

1.1 Challenges of the Quantum Many-Body Problem

One of the fundamental challenges in addressing the quantum many-body problem lies in the extensive amount of information required to fully characterize a quantum state. In contrast to classical mechanics, where the focus is to determine the trajectories of the particles in the system, in quantum mechanics, our objective is to determine the relative probabilities of *all* potential states occurring for all particles.

This distinction arises due to the probabilistic nature of quantum mechanics, where particles can exist in multiple states simultaneously and their behaviors are described by the evolution of wave functions. As the number of particles increases, the dimension of the Hilbert space grows exponentially, resulting in a staggering number of possible states. This exponential growth poses a formidable computational challenge, as the explicit representation and manipulation of the wave function becomes intractable even for systems of modest size.

Consequently, quantum many-body theory is driven by the use of efficient and physically motivated approximations. Specific strategies vary depending on the type of problem at hand, but ideally, we want to encode as much relevant information about the system as possible in the fewest number of free parameters or degrees of freedom. In continuous-space, the exponential scaling problem is magnified by an infinite basis. Nonetheless, resolving

the spatial distribution of a quantum system remains of utmost interest, as it is often more intuitive to interpret physical phenomena in terms of position and spatial relationships.

1.2 Neural-Network Quantum States

In this work, we choose to employ the variational Monte carlo algorithm combined with neural-network quantum states to solve continuous-space quantum many body problems. Since their initial application by Carleo and Troyer in 2017 [1], neural-network quantum states have been applied to problems ranging from condensed matter physics [2] to quantum chemistry [3] and nuclear physics [4]. The field has quickly gained traction, attesting to the immense potential demonstrated by neural-network quantum states. As we continue to develop and refine these techniques, exciting opportunities arise for further breakthroughs.

1.3 Algorithm Development Through the Lens of Nuclear Theory

The advancements in neural-network quantum states often draw inspiration from well-established techniques in machine learning. While the possibility of adapting unexplored machine learning methods piques curiosity, it is also interesting to consider the advantages that a physics-oriented background can offer in tackling these challenges.

The unique form of nuclear interactions, in particular, serves as motivation to design neural-network quantum states that possess a greater degree of generality compared to states designed for interactions commonly encountered in condensed matter systems. Characterized by their explicit dependence on spin and isospin, nuclear Hamiltonians provide illuminating

insights into the intricate mechanisms that govern the universe, encompassing the full range of observable matter—from the microscopic domain of atomic nuclei to the vast reaches of neutron stars.

With the intention to eventually apply our neural-network quantum states to nuclear systems, we carefully consider designs that are compatible with the exchange of spin and isospin. For instance, in our investigation of ultra-cold Fermi gases in Sec. 9, we adopt a completely generalized approach to the Pfaffian wave function, which is commonly decomposed into spin-singlet and spin-triplet contributions for electronic structure problems. By embracing a more inclusive approach, we discover a remarkably versatile and scalable ansatz capable of capturing strong pairing correlations not only in condensed matter systems, but also in nuclear systems.

1.4 Structure of Dissertation

This dissertation begins with a comprehensive overview of the quantum many-body problem in Chapter 2, including the establishment of notation and vocabulary used in later chapters. In Chapter 3, we discuss Quantum Monte Carlo methods, serving as a foundational framework to facilitate the incorporation of machine learning techniques for enhanced flexibility. Then in Chapter 4, we provide a broad background on machine learning, focusing on the mathematics of artificial neural networks. Chapter 5 presents two distinct approaches for the implementation of neural-network quantum states: one in C++ and the other in Python. The former allows for a deeper understanding of how the neural-network quantum states are trained, as gradients are computed analytically in terms of the trainable parameters, while the latter affords greater freedom in design, as gradients are computed numerically. Chap-

ters 6-9 provide illustrative examples of neural-network quantum states employed in different quantum systems, with the latter three chapters being a collection of papers. Each of these examples will provide conclusions and perspectives on neural-network quantum states as they are applied in the specific application and as a whole.

Chapter 2

The Quantum Many-Body Problem

The collective behavior of a quantum many-body system often displays remarkable deviations from the simple summation of the individual constituents. Strong correlations and intricate interactions between particles lead to emergent phenomena that cannot be adequately described by classical mechanics alone. These phenomena, such as phase transitions and the formation of superfluid states, underscore the importance of quantum mechanics in unraveling the complex behaviors and novel properties of quantum many-body systems.

In this study, our objective is to solve the time-independent Schrödinger equation for a nonrelativistic system of identical particles in continuous space. We will focus on determining the ground state, which plays a pivotal role in understanding the system's stability, symmetries, and equilibrium properties. Furthermore, the ground state serves as a reference point for all excited states, enabling meaningful comparisons to experimental observations.

The purpose of this chapter is to establish the foundation for our investigation. We will cover the essential aspects, significance, and challenges associated with quantum many-body problems, as well as examples of commonly employed ab initio methods. This discussion will help provide a comprehensive and contextualized understanding of the field and motivate our subsequent efforts to enhance quantum Monte Carlo methods with neural networks.

2.1 Formalism

Dirac notation, also known as bra-ket notation, offers a concise and abstract representation of quantum states, operators, and measurements. It is intimately linked to the concept of Hilbert space, a complex inner product space denoted as \mathcal{H} . A general quantum state is represented as a ket vector that lives in Hilbert space $|\Psi\rangle \in \mathcal{H}$. For each ket $|\Psi\rangle$, there exists a corresponding linear map from Hilbert space to the field of complex numbers $\langle\Psi| : \mathcal{H} \rightarrow \mathbb{C}$, which we call a bra vector. One can interpret a bra vector as a measurement device or probing tool that extracts useful information regarding a particular state. By applying a bra $\langle\Psi|$ to the left-hand side of another ket $|\Phi\rangle$, we define the inner product $\langle\Psi|\Phi\rangle$, a complex scalar that characterizes the overlap between the states $|\Psi\rangle$ and $|\Phi\rangle$. This inner product is essential for computing probabilities and expectation values of observables, as it enables us to properly define the notions of distance and orthogonality between states in Hilbert space.

Operators are linear transformations $\hat{A} : \mathcal{H} \rightarrow \mathcal{H}$ that act on a state to produce another state. While some operators, like time-evolution or rotation operators, do not correspond directly to measurable quantities, they still have mathematical and physical significance for describing the dynamics and symmetries of quantum states. Observables, such as position, momentum, and spin, are represented by Hermitian operators. Formally, they satisfy $\langle\hat{A}\Psi|\Phi\rangle = \langle\Psi|\hat{A}\Phi\rangle$ for any $|\Psi\rangle, |\Phi\rangle \in \mathcal{H}$. Hermitian operators have two important features:

1. *Real eigenvalues.* All possible measurements are real, implying that all expectation values $\langle\hat{A}\rangle \equiv \langle\Psi|\hat{A}|\Psi\rangle/\langle\Psi|\Psi\rangle$ are also real. The denominator ensures that the average is taken over a normalized probability distribution. These real expectation values play a crucial role in making comparisons and predictions with experimental results.
2. *Orthogonal eigenvectors.* The eigenvectors corresponding to distinct eigenvalues form

a set of orthonormal basis vectors $\{|\alpha\rangle\}$ that span the Hilbert space, yielding the completeness relation

$$\hat{1} = \sum_{\alpha} |\alpha\rangle\langle\alpha|. \quad (2.1)$$

Any state can be expanded as a linear combination of the eigenvectors

$$|\Psi\rangle = \left(\sum_{\alpha} |\alpha\rangle\langle\alpha| \right) |\Psi\rangle = \sum_{\alpha} \langle\alpha|\Psi\rangle |\alpha\rangle. \quad (2.2)$$

If the states are continuous, we convert the sum into an integral

$$\hat{1} = \int d\alpha |\alpha\rangle\langle\alpha|, \quad |\Psi\rangle = \int d\alpha \langle\alpha|\Psi\rangle |\alpha\rangle. \quad (2.3)$$

Since we will work with both continuous and discrete states, this distinction is of minor significance. It will be left to the reader to convert sums into integrals or vice versa, whenever they are not explicitly written.

Therefore, one should identify the observables which provide the most convenient set of basis states $\{|\alpha\rangle\}$ to represent the many-body state $|\Psi\rangle$. Determining the coefficients (or functions) $\langle\alpha|\Psi\rangle$ is equivalent to finding a representation of the abstract state $|\Psi\rangle$.

The best choice of basis for a many-body system depends on the specific nature of the particles and the chosen method for solving the problem. In preparation for our later treatment of N particles in d spatial dimensions, let us define our single-particle basis states as

$$|\mathbf{x}_i\rangle = \begin{cases} |\mathbf{r}_i\rangle, & \text{for bosons,} \\ |\mathbf{r}_i\rangle \otimes |\mathbf{s}_i\rangle \equiv |\mathbf{r}_i, \mathbf{s}_i\rangle, & \text{for fermions,} \end{cases} \quad (2.4)$$

where $\mathbf{r}_i \in \mathbb{R}^d$ are the coordinates of the i th particle, $\mathbf{s}_i \in \mathbb{R}^s$ contains the spin-like degrees of freedom (if there are any), and $i = 1, 2, \dots, N$. The spatial parts are eigenstates of the single-particle position operator¹ $\hat{\mathbf{r}}_i$,

$$\hat{\mathbf{r}}_i |\mathbf{r}_i\rangle = \mathbf{r}_i |\mathbf{r}_i\rangle, \quad (2.5)$$

where the inner product between any two eigenstates is a d -dimensional Dirac delta function

$$\langle \mathbf{r}_i | \mathbf{r}'_i \rangle = \delta(\mathbf{r}_i - \mathbf{r}'_i). \quad (2.6)$$

The spin parts are written as

$$|\mathbf{s}_i\rangle = \begin{cases} |s_i^z\rangle, & \text{for spin-1/2 fermions,} \\ |s_i^z, t_i^z\rangle = |s_i^z\rangle \otimes |t_i^z\rangle, & \text{for nucleons,} \end{cases} \quad (2.7)$$

where $|s_i^z\rangle$ are the eigenstates of the single-particle spin- z operator \hat{s}_i^z ,

$$\hat{s}_i^z |s_i^z\rangle = s_i^z |s_i^z\rangle, \quad \langle s_i^z | s_i^{z'} \rangle = \delta_{s_i^z s_i^{z'}}, \quad (2.8)$$

and $|t_i^z\rangle$ are the eigenstates of the single-particle isospin- z operator \hat{t}_i^z ,

$$\hat{t}_i^z |t_i^z\rangle = t_i^z |t_i^z\rangle, \quad \langle t_i^z | t_i^{z'} \rangle = \delta_{t_i^z t_i^{z'}}, \quad (2.9)$$

¹The position operator for the i th particle is commonly denoted as $\hat{\mathbf{X}}_i$ or $\hat{\mathbf{x}}_i$, but we chose $\hat{\mathbf{r}}_i$ to stay consistent with our definition of \mathbf{x}_i and \mathbf{r}_i .

and δ_{ij} denotes the Kronecker delta. The resulting single-particle Hilbert spaces are

$$\mathcal{H}_i = \begin{cases} L^2(\mathbb{R}^d), & \text{for spin-0 bosons,} \\ L^2(\mathbb{R}^d) \otimes \mathbb{C}^2, & \text{for spin-1/2 fermions,} \\ L^2(\mathbb{R}^d) \otimes \mathbb{C}^2 \otimes \mathbb{C}^2, & \text{for nucleons,} \end{cases} \quad (2.10)$$

and a many-body configuration can be written as the tensor product of the single-particle basis states

$$|\mathbf{X}\rangle = |\mathbf{x}_1\rangle \otimes |\mathbf{x}_2\rangle \otimes \cdots \otimes |\mathbf{x}_N\rangle \equiv |\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\rangle \in \mathcal{H}, \quad (2.11)$$

living in a tensor product of single-particle Hilbert spaces

$$\mathcal{H} = \mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \cdots \otimes \mathcal{H}_N. \quad (2.12)$$

Assuming the dimensions of all the single-particle Hilbert spaces are the same, it is now clear that the dimension of the many-body Hilbert space scales exponentially with the number of particles

$$\dim(\mathcal{H}) = \prod_{i=1}^N \dim(\mathcal{H}_i) = \dim(\mathcal{H}_1)^N. \quad (2.13)$$

This is the primary source of difficulty in quantum many-body problems. The exponential scaling of the Hilbert space dimension applies for even the simplest spin systems in which $\dim(\mathcal{H}_1) = 2$. Our challenge is exacerbated by the fact that the particles in our systems of interest have continuous spatial degrees of freedom, resulting in an infinite $\dim(\mathcal{H}_1)$. Accordingly, quantum many-body methods are often built around ways different ways of

mitigating this problem, such as truncating Hilbert space to a finite subspace or exploring only relevant areas of Hilbert space stochastically.

Our many-body state $|\Psi\rangle$ can be expanded in terms of the many-body configurations defined in Eq. (2.11),

$$|\Psi\rangle = \int d\mathbf{X} \langle \mathbf{X} | \Psi \rangle |\mathbf{X}\rangle \quad (2.14)$$

where our goal is to determine the wave function $\Psi(\mathbf{X}) \equiv \langle \mathbf{X} | \Psi \rangle$, the probability amplitude of finding the system in a certain configuration $|\mathbf{X}\rangle$. Since we may have a mixture of continuous and discrete degrees of freedom, it is sometimes convenient to decompose $|\mathbf{X}\rangle$ as

$$|\mathbf{X}\rangle = |\mathbf{R}\rangle \otimes |\mathbf{S}\rangle, \quad (2.15)$$

$$|\mathbf{R}\rangle = |\mathbf{r}_1\rangle \otimes |\mathbf{r}_2\rangle \otimes \cdots \otimes |\mathbf{r}_N\rangle \equiv |\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N\rangle, \quad (2.16)$$

$$|\mathbf{S}\rangle = |\mathbf{s}_1\rangle \otimes |\mathbf{s}_2\rangle \otimes \cdots \otimes |\mathbf{s}_N\rangle \equiv |\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_N\rangle, \quad (2.17)$$

instead of how it was originally expressed in Eq. (2.11). Then we can easily convert the integral over X into a discrete sum over the spin degrees of freedom and an integral over the continuous spatial degrees of freedom

$$\int d\mathbf{X} \longrightarrow \sum_{\mathbf{S}} \int d\mathbf{R} = \sum_{\mathbf{s}_1} \sum_{\mathbf{s}_2} \cdots \sum_{\mathbf{s}_N} \int d^d \mathbf{r}_1 d^d \mathbf{r}_2 \cdots d^d \mathbf{r}_N, \quad (2.18)$$

with the latter omitted in the bosonic case.

The single-particle momentum operator

$$\hat{\mathbf{p}}_i = -i\hbar \nabla_i, \quad (2.19)$$

where $\nabla_i \equiv \frac{\partial}{\partial r_i}$ is the gradient with respect to position the i th particle, acts on the wave function as

$$\hat{p}_i |\Psi\rangle = \int dX \langle X | \hat{p}_i | \Psi \rangle |X\rangle = \int dX \left(-i\hbar \nabla_i \langle X | \Psi \rangle \right) |X\rangle \quad (2.20)$$

The single-particle eigenstates are

$$\hat{p}_i |\mathbf{p}_i\rangle = \mathbf{p}_i |\mathbf{p}_i\rangle, \quad \langle \mathbf{p}_i | \mathbf{p}'_i \rangle = \delta(\mathbf{p}_i - \mathbf{p}'_i), \quad (2.21)$$

with the continuous momenta $\mathbf{p}_i \in \mathbb{R}^d$ becoming discrete in the presence of periodicity. Expressed in position space, the momentum eigenfunctions are plane waves

$$\varphi_{\mathbf{p}_i}(\mathbf{x}_i) \equiv \langle \mathbf{x}_i | \mathbf{p}_i \rangle = \frac{1}{(2\pi\hbar)^{d/2}} e^{i\mathbf{p}_i \cdot \mathbf{r}_i / \hbar}, \quad (2.22)$$

which will become relevant in our studies of infinite matter.

Just as we can expand states in terms of our working basis, we can expand operators as well,

$$\hat{A} = \left(\int d\mathbf{X} |\mathbf{X}\rangle \langle \mathbf{X}| \right) \hat{A} \left(\int d\mathbf{X}' |\mathbf{X}'\rangle \langle \mathbf{X}'| \right) = \int d\mathbf{X} \int d\mathbf{X}' \langle \mathbf{X} | \hat{A} | \mathbf{X}' \rangle |\mathbf{X}\rangle \langle \mathbf{X}'|, \quad (2.23)$$

where the coefficients $\langle \mathbf{X} | \hat{A} | \mathbf{X}' \rangle$ are called matrix elements of the operator \hat{A} , even though our basis contains continuous components. If the operator is diagonal in our basis, we can

write

$$\begin{aligned}
\hat{A} &= \int d\mathbf{X} \int d\mathbf{X}' \langle \mathbf{X} | \hat{A} | \mathbf{X}' \rangle | \mathbf{X} \rangle \langle \mathbf{X}' | \\
&= \int d\mathbf{X} \int d\mathbf{X}' \langle \mathbf{X} | \hat{A} | \mathbf{X}' \rangle \delta(\mathbf{X} - \mathbf{X}') | \mathbf{X} \rangle \langle \mathbf{X}' | \\
&= \int d\mathbf{X} \langle \mathbf{X} | \hat{A} | \mathbf{X} \rangle | \mathbf{X} \rangle \langle \mathbf{X} | \\
&= \int d\mathbf{X} \hat{A}(\mathbf{X}) | \mathbf{X} \rangle \langle \mathbf{X} |
\end{aligned} \tag{2.24}$$

with $\delta(\mathbf{X} - \mathbf{X}') \equiv \delta(\mathbf{R} - \mathbf{R}') \delta_{\mathbf{S}, \mathbf{S}'}$ representing the product of a Dirac delta function over the spatial degrees of freedom and a Kronecker delta over the spin degrees of freedom, if it applies. It is common to leave off the explicit dependence on \mathbf{X} in the function $\hat{A}(\mathbf{X}) \equiv \langle \mathbf{X} | \hat{A} | \mathbf{X} \rangle$, since this representation of \hat{A} is particularly simple.

2.2 Indistinguishable Particles

The behavior of quantum systems exhibit fundamental differences from classical systems, primarily due to the concept of particle indistinguishability. In classical mechanics, particles are always distinguishable; even when they possess identical masses, charges, or other properties, they can be assigned different labels to unambiguously describe and predict their motion. As a result, classical systems feature well-defined trajectories for each particle, albeit often involving highly nonlinear coupled differential equations.

When we established our mathematical framework in the previous section, we implicitly assumed that the particles were distinguishable. This is evident in Eq. (2.11) and Eq. (2.12), for example, where each particle and single-particle Hilbert space was given an index corresponding to their ordering. Now, let us consider what happens if we assume the particles are

indistinguishable. More specifically, we assume that the expectation value of any observable \hat{A} is invariant under any permutation $\hat{P} \in S_N$ of the N particle indices

$$\frac{\langle \Psi | \hat{A} | \Psi \rangle}{\langle \Psi | \Psi \rangle} = \frac{\langle \hat{P} \Psi | \hat{A} | \hat{P} \Psi \rangle}{\langle \hat{P} \Psi | \hat{P} \Psi \rangle} = \frac{\langle \Psi | \hat{P}^{-1} \hat{A} \hat{P} | \Psi \rangle}{\langle \Psi | \Psi \rangle}, \quad (2.25)$$

where we have used that $\hat{P}^\dagger = \hat{P}^{-1}$ in the last equality. The above equation implies that $[\hat{A}, \hat{P}] = 0$.

Since any permutation in the symmetric group S_N can be written as a product of transpositions, it is sufficient to investigate the action of the transposition operator \hat{P}_{ij} , also known as the particle exchange operator, which trades the labeling of the i th and j th particle.

$$\hat{P}_{ij} |\mathbf{X}\rangle = \hat{P}_{ij} |\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N\rangle = |\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N\rangle. \quad (2.26)$$

Applying the same particle exchange operator twice restores the original ordering

$$\hat{P}_{ij}^2 |\mathbf{X}\rangle = \hat{P}_{ij} |\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N\rangle = |\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N\rangle = |\mathbf{X}\rangle, \quad (2.27)$$

implying the eigenvalues of \hat{P}_{ij} are ± 1 . Thus, the possible effects of exchanging the particle ordering

$$\hat{P}_{ij} |\Psi_{\mathcal{S}}\rangle = |\Psi_{\mathcal{S}}\rangle, \quad \text{for bosons,} \quad (2.28)$$

$$\hat{P}_{ij} |\Psi_{\mathcal{A}}\rangle = -|\Psi_{\mathcal{A}}\rangle, \quad \text{for fermions,} \quad (2.29)$$

give rise to two fundamentally different classes of identical particles: bosons, which have purely symmetric states and obey Bose-Einstein statistics, and fermions, which have purely

antisymmetric states and obey Fermi-Dirac statistics. The subscripts \mathcal{S} and \mathcal{A} stand for symmetric and antisymmetric, respectively. By the spin-statistics theorem, the difference between the two categories is attributed to an intrinsic property called spin. Bosons, such as photons and gluons, have integer values of spin, while fermions, such as electrons, protons, and neutrons, have half-integer values. We will only consider spin-0 and spin-1/2 particles in this work, as shown already in Eq. (2.10). In Sec. 9, we will discover that the line distinguishing bosons and fermions is sometimes less distinct than it initially appears.

It is straightforward to see how the action of the particle exchange operator can be generalized to an arbitrary permutation $\hat{P} \in S_N$. The sign, or parity, of the permutation is $\sigma(P) = (-1)^p$, where p is the number of transpositions required to decompose the permutation into a product of transpositions². Then the permutation acts as

$$\hat{P}|\Psi_{\mathcal{S}}\rangle = |\Psi_{\mathcal{S}}\rangle, \quad \text{for bosons,} \quad (2.30)$$

$$\hat{P}|\Psi_{\mathcal{A}}\rangle = \sigma(P)|\Psi_{\mathcal{A}}\rangle, \quad \text{for fermions,} \quad (2.31)$$

for the two cases. Our working basis states $|X\rangle$ are not eigenstates themselves, but we can construct symmetric and antisymmetric states from them. In fact, we can define the symmetrization operator

$$\hat{\mathcal{S}} = \frac{1}{N!} \sum_{\hat{P} \in S_N} \hat{P}, \quad (2.32)$$

²There exists an infinite number of ways to decompose a permutation into a product of transpositions. However, an odd permutation will invariably be decomposed into an odd number of transpositions, and an even permutation will be decomposed into an even number of transpositions.

and the antisymmetrization operator

$$\hat{\mathcal{A}} = \frac{1}{N!} \sum_{\hat{P} \in S_N} \sigma(P) \hat{P}, \quad (2.33)$$

that will map any many-body state $|\Psi\rangle \in \mathcal{H}$ to the equivalent symmetric and antisymmetric states

$$\hat{\mathcal{S}}|\Psi\rangle \equiv |\Psi_{\mathcal{S}}\rangle, \quad \hat{\mathcal{A}}|\Psi\rangle \equiv |\Psi_{\mathcal{A}}\rangle. \quad (2.34)$$

Both of these operators are idempotent,

$$\hat{\mathcal{S}}^2 = \hat{\mathcal{S}}, \quad \hat{\mathcal{A}}^2 = \hat{\mathcal{A}}, \quad (2.35)$$

meaning they project generic states in Hilbert space to symmetric and antisymmetric subspaces, respectively,

$$|\Psi_{\mathcal{S}}\rangle \in \mathcal{H}_{\mathcal{S}}, \quad |\Psi_{\mathcal{A}}\rangle \in \mathcal{H}_{\mathcal{A}}. \quad (2.36)$$

Applying the symmetrization operator after already applying the antisymmetrization operator destroys any state, and vice versa

$$\hat{\mathcal{S}}\hat{\mathcal{A}} = \hat{\mathcal{A}}\hat{\mathcal{S}} = 0, \quad (2.37)$$

implying the subspaces $\mathcal{H}_{\mathcal{S}}$ and $\mathcal{H}_{\mathcal{A}}$ are entirely distinct. Furthermore, any two states that are related by a permutation

$$|\Psi'\rangle = \hat{P}'|\Psi\rangle, \quad \hat{P}' \in S_N, \quad (2.38)$$

will map to the same symmetric state

$$|\Psi'_S\rangle = \hat{\mathcal{S}}|\Psi'\rangle = \hat{\mathcal{S}}|\Psi\rangle = |\Psi_S\rangle, \quad (2.39)$$

and to the same antisymmetric state except for a possible sign flip

$$|\Psi'_A\rangle = \hat{\mathcal{A}}|\Psi'\rangle = \sigma(P')\hat{\mathcal{A}}|\Psi\rangle = \sigma(P')|\Psi_A\rangle. \quad (2.40)$$

Therefore, the many-body Hilbert space \mathcal{H} that we established in Eq. (2.12) for distinguishable particles is typically much larger than the direct sum of the subspaces for indistinguishable particles

$$\mathcal{H}_{\mathcal{S}} \oplus \mathcal{H}_{\mathcal{A}} \subseteq \mathcal{H}. \quad (2.41)$$

This is spectacular news! The physically realizable quantum states of our system will either be in $\mathcal{H}_{\mathcal{S}}$ or $\mathcal{H}_{\mathcal{A}}$, which is smaller than the entire Hilbert space \mathcal{H} by a factor that scales as $N!$, the order of the symmetric group S_N . While this does not alleviate the dimensionality problem entirely because our Hilbert space is still infinite dimensional, it does help make the many-body problem more tractable. Enforcing other symmetries, such as parity and time-reversal, can similarly help reduce the effective size of our Hilbert space.

Now that we are able to construct both symmetric and antisymmetric states from a generic one, let us consider what happens if two (or more) indistinguishable particles occupy the same exact single-particle state. Then there exists a particle exchange operator \hat{P}_{ij} for some $i \neq j$ that generates a symmetry of the many-body state,

$$\hat{P}_{ij}|\Psi_{ij}\rangle = |\Psi_{ij}\rangle. \quad (2.42)$$

Here, we have denoted the state as $|\Psi_{ij}\rangle$ to specify that it is invariant under $i \leftrightarrow j$. No issues arise for the corresponding bosonic system according to Eq. (2.39), so bosons are fine with being on top of their identical neighbors. However, for the fermionic system, Eq. (2.40) implies a contradiction

$$\hat{\mathcal{A}}|\Psi_{ij}\rangle = -\hat{\mathcal{A}}|\Psi_{ij}\rangle. \quad (2.43)$$

The above equation can only be possible if $\hat{\mathcal{A}}|\Psi_{ij}\rangle = 0$, leading to the well-known Pauli exclusion principle. This fundamental principle asserts that no two identical fermions can occupy the same quantum state simultaneously. As a result, fermions exhibit distinct behavior from bosons, particularly at low temperatures.

2.3 The Schrödinger Equation

The behavior of a nonrelativistic quantum system is governed by the Schrödinger equation

$$i\hbar \frac{\partial}{\partial t} |\Psi\rangle = \hat{H}|\Psi\rangle, \quad (2.44)$$

where the Hamiltonian \hat{H} is a Hermitian operator that represents the total energy of the system. We will exclusively work with a time-independent Hamiltonian that is diagonal in our working basis $|\mathbf{X}\rangle$, meaning

$$\hat{H} = \int d\mathbf{X} \hat{H}(\mathbf{X}) |\mathbf{X}\rangle \langle \mathbf{X}|, \quad (2.45)$$

where $\hat{H}(\mathbf{X}) \equiv \langle \mathbf{X} | \hat{H} | \mathbf{X} \rangle$ as defined in Eq. (2.24). Our Hamiltonian takes the form

$$\hat{H} = \hat{K} + \hat{V}, \quad (2.46)$$

where the non-relativistic kinetic energy is defined as

$$\hat{K} = \frac{1}{2m} \sum_{i=1}^N \hat{\mathbf{p}}_i^2 = -\frac{\hbar^2}{2m} \sum_{i=1}^N \nabla_i^2 \quad (2.47)$$

and the potential energy $\hat{V} = V(\mathbf{X})$ may include external potentials, two- or three-body interactions, and spin- or isospin-dependence. While the above Hamiltonian is not dependent time, the state $|\Psi\rangle$ still could be. We have suppressed this time-dependence thus far, but we will now make it explicit by writing

$$|\Psi(t)\rangle = \int d\mathbf{X} \langle \mathbf{X} | \Psi(t) \rangle | \mathbf{X} \rangle, \quad (2.48)$$

where the time-dependent wave function can be decomposed as

$$\langle \mathbf{X} | \Psi(t) \rangle \equiv \Psi(\mathbf{X}, t) = \psi(\mathbf{X}) \phi(t), \quad (2.49)$$

following the separation of variables method of solving partial differential equations. Then the Schrödinger equation reads

$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{X}, t) = \hat{H}(\mathbf{X}) \Psi(\mathbf{X}, t) \implies i\hbar \psi(\mathbf{X}) \frac{\partial}{\partial t} \phi(t) = \phi(t) \hat{H}(\mathbf{X}) \psi(\mathbf{X}). \quad (2.50)$$

Rearranging all of the t -dependent components on the left and the X -dependent components on the right,

$$\frac{i\hbar}{\phi(t)} \frac{\partial}{\partial t} \phi(t) = \frac{1}{\psi(\mathbf{X})} \hat{H}(\mathbf{X}) \psi(\mathbf{X}) = E, \quad (2.51)$$

we find that each side of the equation must be equal to some constant E . The left-hand side can be solved immediately

$$i\hbar \frac{\partial}{\partial t} \phi(t) = E\phi(t) \implies \phi(t) = \phi(0)e^{-iEt/\hbar}. \quad (2.52)$$

The right-hand side, however, is an eigenvalue problem

$$\hat{H}(\mathbf{X}) \psi(\mathbf{X}) = E\psi(\mathbf{X}), \quad (2.53)$$

also known as the time-independent Schrödinger equation.

Eigenvalue problems are challenging because they are inherently nonlinear. What's more, the eigenvectors in our situation are actually continuous eigenfunctions rather than finite-dimensional vectors. However, once we identify the solutions of the time-independent Schrödinger equation above, we will automatically obtain the solutions of the full time-dependent Schrödinger equation as well. Let us index the possible solutions with α , which can either label the discrete energy spectrum for a bound system or the continuous spectrum for an unbound system. Then the so-called stationary states are

$$\Psi_\alpha(\mathbf{X}, t) = \psi_\alpha(\mathbf{X}) e^{-iE_\alpha t/\hbar}, \quad (2.54)$$

where we omit writing the constant $\phi(0)$ from Eq. (2.52) because it can be absorbed into

the coefficients of the general solution

$$\Psi(\mathbf{X}, t) = \sum_{\alpha} c_{\alpha} \Psi_{\alpha}(\mathbf{X}, t). \quad (2.55)$$

The stationary states are called "stationary" because their associated probability distributions do not depend on time

$$|\Psi_{\alpha}(\mathbf{X}, t)|^2 = \psi_{\alpha}^{*}(\mathbf{X}) e^{i E_{\alpha} t / \hbar} \psi_{\alpha}(\mathbf{X}) e^{-i E_{\alpha} t / \hbar} = |\psi_{\alpha}(\mathbf{X})|^2. \quad (2.56)$$

Given that

$$|\Psi(\mathbf{X}, t)|^2 = \sum_{\alpha} \sum_{\beta} c_{\alpha}^{*} c_{\beta} \psi_{\alpha}^{*}(\mathbf{X}) \psi_{\beta}(\mathbf{X}) e^{i(E_{\alpha} - E_{\beta})t / \hbar}, \quad (2.57)$$

the same cannot be said about the general solution.

2.4 The Variational Principle

While the solutions of the time-dependent Schrödinger equation are vital for understanding the dynamics of the quantum system, it is a necessary step to first solve the time-independent Schrödinger equation. However, the exact time-independent solutions and their corresponding energies are still difficult, if not impossible, to obtain. One valuable tool for finding approximate solutions is called the variational principle, which states that the expectation value of the Hamiltonian is an upperbound for the true ground state energy of the system E_0 ,

$$E[\Psi] \equiv \langle \hat{H} \rangle = \frac{\langle \Psi | \hat{H} | \Psi \rangle}{\langle \Psi | \Psi \rangle} \geq E_0, \quad (2.58)$$

for any state $|\Psi\rangle$. It is simple to prove. First, we write the eigenstates of the Hamiltonian as $|\Psi_\alpha\rangle = |\alpha\rangle$ and their corresponding energies as E_α . The eigenstates form a complete orthonormal basis with $\langle\alpha|\beta\rangle = \delta_{\alpha\beta}$. If we assume $|\Psi\rangle$ is normalized, $\langle\Psi|\Psi\rangle = 1$, then the expectation value of the energy can be expanded as

$$\begin{aligned} E[\Psi] &= \langle\Psi|\hat{H}|\Psi\rangle \\ &= \sum_{\alpha} \sum_{\beta} \langle\Psi|\alpha\rangle \langle\alpha|\hat{H}|\beta\rangle \langle\beta|\Psi\rangle \\ &= \sum_{\alpha} \sum_{\beta} \langle\Psi|\alpha\rangle \langle\beta|\Psi\rangle E_\alpha \delta_{\alpha\beta} \\ &= \sum_{\alpha} E_\alpha |\langle\alpha|\Psi\rangle|^2. \end{aligned} \tag{2.59}$$

The ground state energy E_0 , by definition, is smaller than all the other energy eigenvalues $\{E_\alpha\}_{\alpha \neq 0}$, leading to the result

$$E[\Psi] = \sum_{\alpha} E_\alpha |\langle\alpha|\Psi\rangle|^2 \geq E_0, \tag{2.60}$$

where equality can only hold for the ground state $|\Psi_0\rangle$.

Similar variational principles can be developed for the excited states. First, let us order our energy eigenstates such that $E_0 \leq E_1 \leq E_2 \dots$, where E_1 is the energy of the first excited state $|\Psi_1\rangle$, E_2 is the energy of the second excited state $|\Psi_2\rangle$, and so on. Then the variational principle for the n th excited state is given by

$$E[\Psi] = \frac{\langle\Psi|\hat{H}|\Psi\rangle}{\langle\Psi|\Psi\rangle} \geq E_n, \tag{2.61}$$

if our state $|\Psi\rangle$ is orthogonal to all the lower-lying states,

$$\langle \Psi_m | \Psi \rangle = 0 \text{ for all } m < n. \quad (2.62)$$

This statement can be proved by the same method as Eq. (2.59). We have used the Latin indices m and n in place of the Greek indices α and β to emphasize that our eigenstates have been ordered according to their corresponding energies.

Another useful fact is that the variance of energy is only zero if and only if our state is an eigenstate of the Hamiltonian

$$\sigma_E^2 = 0 \iff |\Psi\rangle = |\Psi_\alpha\rangle \text{ for some } \alpha. \quad (2.63)$$

The backward direction of the statement above is trivial to prove, so we will just prove the forward direction here. Expanding the variance of the energy in terms of the energy eigenstates, we obtain

$$\begin{aligned} \sigma_E^2 &= \langle (\hat{H} - \langle \hat{H} \rangle)^2 \rangle = \langle \hat{H}^2 \rangle - \hat{H}^2 \\ &= \sum_{\alpha} E_{\alpha}^2 |\langle \alpha | \Psi \rangle|^2 - \left(\sum_{\alpha} E_{\alpha} |\langle \alpha | \Psi \rangle|^2 \right) \left(\sum_{\beta} E_{\beta} |\langle \beta | \Psi \rangle|^2 \right) \\ &= \sum_{\alpha} E_{\alpha}^2 |\langle \alpha | \Psi \rangle|^2 - \sum_{\alpha} \sum_{\beta} E_{\alpha} E_{\beta} |\langle \alpha | \Psi \rangle|^2 |\langle \beta | \Psi \rangle|^2 \\ &= \sum_{\alpha} E_{\alpha} |\langle \alpha | \Psi \rangle|^2 \left(E_{\alpha} - \sum_{\beta} E_{\beta} |\langle \beta | \Psi \rangle|^2 \right). \end{aligned} \quad (2.64)$$

There are only two ways that the above can equal zero *for all* α : either $|\langle \alpha | \Psi \rangle|^2 = 0$ or

$$E_\alpha = \sum_{\beta} E_{\beta} |\langle \beta | \Psi \rangle|^2 = E_\alpha |\langle \alpha | \Psi \rangle|^2 \implies |\langle \alpha | \Psi \rangle|^2 = 1. \quad (2.65)$$

Therefore, the state $|\Psi\rangle$ must be an eigenstate of the Hamiltonian for some α . A nonzero variance of the energy indicates that our system is not in a well-defined state, but rather, a superposition of eigenstates.

We must address an important caveat regarding the variational principle. Specifically, the variational principle only holds true if the state adheres to the correct symmetries of the system. The many-body Hamiltonian often does not provide information about whether the system consists of bosons or fermions. For example, when describing a system of identical charged particles, the Coulomb interaction only depends on the charge and distances between the pairs. If one employs a state intended for bosons when computing the expectation value of the energy for a fermionic system, it can lead to energies much lower than the true ground state energy of the system. Such erroneous results are not violations of the variational principle; instead, they arise due to an incorrect implementation of the principle.

2.5 Ab Initio Methods

A many-body method is called ab initio (or "from first principles") if it only relies on fundamental principles and established laws of nature. It involves solving the Schrödinger equation starting from a microscopic Hamiltonian rather than a Hamiltonian derived from empirical or experimental data. In this section, we will provide examples of common ab initio methods. While some of these methods are not explicitly formulated for real-space systems, they

can often be adapted and applied to such systems with the appropriate modifications.

2.5.1 Configuration Interaction

Full Configuration Interaction is an exact method that aims to solve the time-independent Schrödinger equation through diagonalization of the Hamiltonian in matrix form. We begin by expanding the wave function as a linear combination of orthonormal many-body basis states $\{|\alpha\rangle\}$,

$$|\Psi\rangle = \sum_{\alpha} c_{\alpha} |\alpha\rangle, \quad (2.66)$$

where c_{α} are the coefficients of the expansion and the states $|\alpha\rangle$ are usually taken to be all possible antisymmetrized products of single-particle basis states, also known as Slater determinants, for fermionic systems. By inserting this expansion and the completeness relation into the Schrödinger equation, we obtain

$$\hat{H}|\Psi\rangle = E|\Psi\rangle, \quad (2.67)$$

$$\sum_{\alpha} \hat{H}c_{\alpha} |\alpha\rangle = E \sum_{\alpha} c_{\alpha} |\alpha\rangle \quad (2.68)$$

$$\sum_{\alpha} \sum_{\beta} \langle \beta | \hat{H} | \alpha \rangle c_{\alpha} |\beta\rangle = E \sum_{\alpha} c_{\alpha} |\alpha\rangle \quad (2.69)$$

Since the basis states are orthogonal, multiplying on the left by $\langle \gamma |$ yields

$$\sum_{\alpha} \langle \gamma | \hat{H} | \alpha \rangle c_{\alpha} = E c_{\gamma}, \quad (2.70)$$

which can be written in a more convenient form by organizing the coefficients c_α into a vector and $\langle \gamma | \hat{H} | \alpha \rangle$ into a matrix

$$\begin{bmatrix} \langle 1 | \hat{H} | 1 \rangle & \langle 1 | \hat{H} | 2 \rangle & \cdots & \langle 1 | \hat{H} | \alpha \rangle & \cdots \\ \langle 2 | \hat{H} | 1 \rangle & \langle 2 | \hat{H} | 2 \rangle & \cdots & \langle 2 | \hat{H} | \alpha \rangle & \cdots \\ \vdots & \vdots & \ddots & \vdots & \cdots \\ \langle \alpha | \hat{H} | 1 \rangle & \langle \alpha | \hat{H} | 2 \rangle & \cdots & \langle \alpha | \hat{H} | \alpha \rangle & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_\alpha \\ \vdots \end{bmatrix} = E \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_\alpha \\ \vdots \end{bmatrix} \quad (2.71)$$

Clearly, solving the eigenvalue problem above is only imaginable when the matrix is of finite size. However, even in such cases, the dimensionality of the problem often poses significant challenges and renders it computationally intractable.

One alternative approach is to truncate the many-body basis to a finite, manageable subset $\{|\alpha_1\rangle, |\alpha_2\rangle, \dots, |\alpha_n\rangle\}$, leading to the more commonly-employed method called Configuration Interaction. In truncating our many-body basis, it becomes important to consider which states we include in our subset. Typically, the most relevant states are constructed by identifying a suitable reference state, like the Hartree-Fock solution, and constructing excitations or configurations around it.

2.5.2 Hartree-Fock Theory

Hartree-Fock theory is one of the simplest approximate methods for solving the many-body Schrödinger equation. It provides a mean-field description of the system, where each particle moves in an effective average field created by all other particles. This approach essentially decouples the two-body interaction, resulting in a tractable computational scheme, but it

neglects correlation effects, which can be significant in systems with strong interactions.

Let us assume that the Hamiltonian of a fermionic system contains up to two-body interactions

$$\hat{H}(\mathbf{X}) = \hat{H}_0(\mathbf{X}) + \hat{V}(\mathbf{X}), \quad (2.72)$$

$$\hat{H}_0(\mathbf{X}) = \sum_{i=1}^N \hat{h}_0(\mathbf{x}_i), \quad (2.73)$$

$$\hat{V}(\mathbf{X}) = \sum_{i < j}^N \hat{v}(\mathbf{x}_{ij}) \quad (2.74)$$

where we have divided the terms in the Hamiltonian such that \hat{H}_0 contains only one-body contributions and \hat{V} contains the pair-wise interactions. Accordingly they can be decomposed into individual one-body and two-body operators. The notation \mathbf{x}_{ij} is meant to imply that the two-body interactions are symmetric with respect to i and j .

To achieve the best mean-field description of the interacting problem, our goal is to find an effective one-body Hamiltonian that approximates the original one. Namely,

$$\hat{H}(\mathbf{X}) \approx \sum_{i=1}^N \hat{f}(\mathbf{x}_i), \quad (2.75)$$

where \hat{f} can be further decomposed as $\hat{f}(\mathbf{x}_i) = \hat{h}_0(\mathbf{x}_i) + \hat{v}^{HF}(\mathbf{x}_i)$, for some choice of \hat{v}^{HF} that is to be determined. Then eigenstates of this mysterious one-body Hamiltonian must satisfy

$$\hat{f}(\mathbf{x})\varphi_\alpha(\mathbf{x}) = \varepsilon_\alpha \varphi_\alpha(\mathbf{x}), \quad (2.76)$$

with the corresponding single-particle energies ε_α .

In the Hartree-Fock approximation, the many-body wave function $\Phi^{HF}(\mathbf{X})$ is taken

to be a Slater determinant of these single-particle eigenstates. The Slater determinant is a natural result of applying the antisymmetrization operator (Eq. (2.33)) to a product of single-particle wave functions

$$\begin{aligned}\Phi^{HF}(\mathbf{X}) &= \sqrt{N!} \hat{\mathcal{A}}\{\varphi_1(\mathbf{x}_1)\varphi_2(\mathbf{x}_2) \cdots \varphi_N(\mathbf{x}_N)\} \\ &= \frac{1}{\sqrt{N!}} \det \begin{bmatrix} \varphi_1(\mathbf{x}_1) & \varphi_1(\mathbf{x}_2) & \cdots & \varphi_1(\mathbf{x}_N) \\ \varphi_2(\mathbf{x}_1) & \varphi_2(\mathbf{x}_2) & \cdots & \varphi_2(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_N(\mathbf{x}_1) & \varphi_N(\mathbf{x}_2) & \cdots & \varphi_N(\mathbf{x}_N) \end{bmatrix}. \quad (2.77)\end{aligned}$$

The above assumption only applies for fermions, but an equivalent assumption can be made for bosons by applying the symmetrization operator (Eq. (2.32)) to the product, resulting in a permanent instead of a determinant. To simplify our notation moving forward, let us define $\Phi(\mathbf{X}) \equiv \varphi_1(\mathbf{x}_1)\varphi_2(\mathbf{x}_2) \cdots \varphi_N(\mathbf{x}_N)$ to mean the product of single-particle eigenstates before they have been antisymmetrized. Then $\Phi^{HF}(\mathbf{X}) = \sqrt{N!}\hat{\mathcal{A}}\Phi(\mathbf{X})$.

According to the variational principle, the best set of orthogonal single-particle states $\varphi_\alpha(\mathbf{x})$ would minimize the energy using the Hartree-Fock ansatz $\Phi^{HF}(\mathbf{X})$. The expectation value of the energy is

$$E[\Phi^{HF}] = \langle \Phi^{HF} | \hat{H} | \Phi^{HF} \rangle = \langle \Phi^{HF} | \hat{H}_0 | \Phi^{HF} \rangle + \langle \Phi^{HF} | \hat{V} | \Phi^{HF} \rangle, \quad (2.78)$$

where we have assumed $\Phi^{HF}(\mathbf{X})$ is normalized. Each of these terms can be evaluated by using the known properties of the antisymmetrization operator and the orthogonality of the

single-particle states $\varphi_\alpha(\mathbf{x}) = \langle \mathbf{x} | \alpha \rangle$. The calculations for both are as follows:

$$\begin{aligned}
& \langle \Phi^{HF} | \hat{H}_0 | \Phi^{HF} \rangle \\
&= \int d\mathbf{X} \Phi^{HF*}(\mathbf{X}) \hat{H}_0(\mathbf{X}) \Phi^{HF}(\mathbf{X}) \\
&= N! \int d\mathbf{X} \Phi^*(\mathbf{X}) [\hat{\mathcal{A}} \hat{H}_0(\mathbf{X})] \hat{\mathcal{A}} \Phi(\mathbf{X}) = N! \int d\mathbf{X} \Phi^*(\mathbf{X}) [\hat{H}_0(\mathbf{X}) \hat{\mathcal{A}}] \hat{\mathcal{A}} \Phi(\mathbf{X}) \\
&= N! \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{H}_0(\mathbf{X}) [\hat{\mathcal{A}} \hat{\mathcal{A}}] \Phi(\mathbf{X}) = N! \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{H}_0(\mathbf{X}) \hat{\mathcal{A}} \Phi(\mathbf{X}) \quad (2.79) \\
&= \sum_{i=1}^N \sum_{\hat{P} \in S_N} \sigma(P) \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{h}_0(\mathbf{x}_i) \hat{P} \Phi(\mathbf{X}) = \sum_{i=1}^N \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{h}_0(\mathbf{x}_i) \Phi(\mathbf{X}) \\
&= \sum_{\alpha=1}^N \int d\mathbf{x} \varphi_\alpha^*(\mathbf{x}) \hat{h}_0(\mathbf{x}) \varphi_\alpha(\mathbf{x}) = \sum_{\alpha=1}^N \langle \alpha | \hat{h}_0 | \alpha \rangle,
\end{aligned}$$

$$\begin{aligned}
& \langle \Phi^{HF} | \hat{V} | \Phi^{HF} \rangle \\
&= \int d\mathbf{X} \Phi^{HF*}(\mathbf{X}) \hat{V}(\mathbf{X}) \Phi^{HF}(\mathbf{X}) \\
&= N! \int d\mathbf{X} \Phi^*(\mathbf{X}) [\hat{\mathcal{A}} \hat{V}(\mathbf{X})] \hat{\mathcal{A}} \Phi(\mathbf{X}) = N! \int d\mathbf{X} \Phi^*(\mathbf{X}) [\hat{V}(\mathbf{X}) \hat{\mathcal{A}}] \hat{\mathcal{A}} \Phi(\mathbf{X}) \\
&= N! \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{V}(\mathbf{X}) [\hat{\mathcal{A}} \hat{\mathcal{A}}] \Phi(\mathbf{X}) = N! \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{V}(\mathbf{X}) \hat{\mathcal{A}} \Phi(\mathbf{X}) \\
&= \sum_{i < j}^N \sum_{\hat{P} \in S_N} \sigma(P) \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{v}(\mathbf{x}_{ij}) \hat{P} \Phi(\mathbf{X}) = \sum_{i < j}^N \int d\mathbf{X} \Phi^*(\mathbf{X}) \hat{v}(\mathbf{x}_{ij}) (1 - \hat{P}_{ij}) \Phi(\mathbf{X}) \\
&= \frac{1}{2} \sum_{\alpha\beta}^N \int d\mathbf{x}_1 d\mathbf{x}_2 \left[\varphi_\alpha^*(\mathbf{x}_1) \varphi_\beta^*(\mathbf{x}_2) \hat{v}(\mathbf{x}_{12}) \varphi_\alpha(\mathbf{x}_1) \varphi_\beta(\mathbf{x}_2) - \varphi_\alpha^*(\mathbf{x}_1) \varphi_\beta^*(\mathbf{x}_2) \hat{v}(\mathbf{x}_{12}) \varphi_\alpha(\mathbf{x}_2) \varphi_\beta(\mathbf{x}_1) \right] \\
&= \frac{1}{2} \sum_{\alpha\beta}^N \left[\langle \alpha\beta | \hat{v} | \alpha\beta \rangle - \langle \alpha\beta | \hat{v} | \beta\alpha \rangle \right] \equiv \frac{1}{2} \sum_{\alpha\beta}^N \langle \alpha\beta | \hat{v} | \alpha\beta \rangle_{\mathcal{A}}. \quad (2.80)
\end{aligned}$$

Now that we know

$$E[\Phi^{HF}] = \sum_{\alpha=1}^N \langle \alpha | \hat{h}_0 | \alpha \rangle + \frac{1}{2} \sum_{\alpha\beta}^N \langle \alpha\beta | \hat{v} | \alpha\beta \rangle_{\mathcal{A}}, \quad (2.81)$$

we can now vary the orbitals $\varphi_{\alpha}(\mathbf{x})$ in order to find the minimum $E[\Phi^{HF}]$. Alternatively, we can perform a unitary transformation on our basis, expanding $|\alpha\rangle$ as a linear combination of well-known, orthogonal basis states $|i\rangle$, and vary the coefficients. We will choose the latter method, specifically with the basis states coming from the non-interacting problem,

$$|\alpha\rangle = \sum_i C_{\alpha i} |i\rangle, \quad (2.82)$$

where

$$\hat{f}|\alpha\rangle = \varepsilon_{\alpha}|\alpha\rangle, \quad (2.83)$$

$$\hat{h}_0|i\rangle = \varepsilon_i^0|i\rangle. \quad (2.84)$$

In terms of these expansion coefficients, the expectation value of the energy becomes

$$E[\Phi^{HF}] = \sum_{\alpha=1}^N \sum_{ij} C_{\alpha i}^* C_{\alpha j} \langle i | \hat{h}_0 | j \rangle + \frac{1}{2} \sum_{\alpha\beta}^N \sum_{ijkl} C_{\alpha i}^* C_{\beta j}^* C_{\alpha k} C_{\beta \ell} \langle ij | \hat{v} | kl \rangle_{\mathcal{A}} \quad (2.85)$$

Recognizing that what we really desire is Eq. 2.75, we can then write

$$E[\Phi^{HF}] = \sum_{\alpha=1}^N \langle \alpha | \hat{f} | \alpha \rangle = \sum_{\alpha=1}^N \sum_{ij} C_{\alpha i}^* C_{\alpha j} \langle i | \hat{f} | j \rangle. \quad (2.86)$$

By comparing this equation with the one above it, we find

$$\langle i | \hat{f} | j \rangle = \langle i | \hat{h}_0 | j \rangle + \sum_{\alpha} \sum_{k\ell} C_{\alpha k}^* C_{\beta \ell} \langle ik | \hat{v} | jl \rangle_{\mathcal{A}}, \quad (2.87)$$

which simplifies to

$$\langle i | \hat{f} | j \rangle = \varepsilon_i^0 + \sum_{\alpha} \sum_{k\ell} C_{\alpha k}^* C_{\beta \ell} \langle ik | \hat{v} | jl \rangle_{\mathcal{A}}, \quad (2.88)$$

because the states $|i\rangle$ are eigenstates of \hat{h}_0 . Next, by expanding Eq. 2.83 in terms of $|j\rangle$ and projecting on the left with $\langle i|$, we find that we can write

$$\sum_j C_{\alpha j} \langle i | \hat{f} | j \rangle = \varepsilon_{\alpha} C_{\alpha i}. \quad (2.89)$$

This means that if we organize Eq. 2.86 into a matrix and the coefficients $C_{\alpha i}$ into a vector for each α , the above can be written as a normal eigenvalue problem $FC_{\alpha} = \varepsilon_{\alpha} C_{\alpha}$.

The Hartree-Fock eigenvalue problem can be solved iteratively. First, we provide guesses for the initial $C_{\alpha}^{(0)}$, where the superscript indicates the iteration. Then for each iteration, we construct the matrix $F^{(t)}$, which has elements $F_{ij}^{(t)} = \langle i | \hat{f} | j \rangle$ using the previous coefficients $C_{\alpha}^{(t-1)}$. By solving the eigenvalue problem through diagonalization, we obtain new eigenvectors $C_{\alpha}^{(t)}$ and their corresponding eigenvalues $\varepsilon_{\alpha}^{(t)}$. We continue this process until the eigenvalues converge.

2.5.3 Many-Body Perturbation Theory

Many-body perturbation theory offers a systematic framework for incorporating the effects of interparticle correlations beyond the Hartree-Fock level. The Hamiltonian is assumed to

consist of two parts, an unperturbed Hamiltonian \hat{H}_0 and an interacting Hamiltonian \hat{H}_I ,

$$\hat{H} = \hat{H}_0 + \hat{H}_I, \quad (2.90)$$

where the solutions for the unperturbed case are easy to find

$$\hat{H}_0|\Phi_n\rangle = W_n|\Phi_n\rangle. \quad (2.91)$$

Let us split the completeness relation such that the projector onto the subspace spanned by the ground state is separated from the rest,

$$\hat{1} = \sum_{n=0}^{\infty} |\Phi_n\rangle\langle\Phi_n| = |\Phi_0\rangle\langle\Phi_0| + \sum_{n=1}^{\infty} |\Phi_n\rangle\langle\Phi_n| = \hat{P} + \hat{Q}. \quad (2.92)$$

Then we can expand the ground state $|\Psi_0\rangle$ of the full Hamiltonian \hat{H} as

$$|\Psi_0\rangle = (\hat{P} + \hat{Q})|\Psi_0\rangle = |\Phi_0\rangle + \hat{Q}|\Psi_0\rangle = |\Phi_0\rangle + \sum_{n=1}^{\infty} \langle\Phi_n|\Psi_0\rangle|\Phi_n\rangle, \quad (2.93)$$

where we have made the assumption $\langle\Phi_0|\Psi_0\rangle = 1$. By applying the full Hamiltonian on the above state and projecting with $\langle\Phi_0|$, we obtain

$$\langle\Phi_0|\hat{H}|\Psi_0\rangle = \langle\Phi_0|\hat{H}_0|\Phi_0\rangle + \langle\Phi_0|\hat{H}_I|\Psi_0\rangle = W_0 + \langle\Phi_0|\hat{H}_I|\Psi_0\rangle, \quad (2.94)$$

due to the orthogonality of the eigenstates $|\Phi_n\rangle$. In addition, by applying $\langle\Phi_0|$ to the Schrödinger equation for the full Hamiltonian, we obtain

$$\langle\Phi_0|\hat{H}|\Psi_0\rangle = E_0\langle\Phi_0|\Psi_0\rangle = E_0, \quad (2.95)$$

implying that $E_0 = W_0 + \Delta E$, where $\Delta E \equiv \langle \Phi_0 | \hat{H}_I | \Psi_0 \rangle$. Through simple manipulations of the Schrödinger equation, we find

$$\begin{aligned}\hat{H}|\Psi_0\rangle &= \hat{H}_0|\Psi_0\rangle + \hat{H}_I|\Psi_0\rangle = E_0\hat{H}|\Psi_0\rangle \\ (\hat{H}_0 - W_0)|\Psi_0\rangle &= (E_0 - W_0 - \hat{H}_I)|\Psi_0\rangle \\ (W_0 - \hat{H}_0)|\Psi_0\rangle &= (\hat{H}_I - E_0)|\Psi_0\rangle\end{aligned}\tag{2.96}$$

Since \hat{Q} is idempotent and commutes with the Hamiltonian, applying \hat{Q} to both sides of the above equation yields

$$\begin{aligned}\hat{Q}(\hat{H}_I - \Delta E)|\Psi_0\rangle &= \hat{Q}(W_0 - \hat{H}_0)|\Psi_0\rangle \\ &= \hat{Q}\hat{Q}(W_0 - \hat{H}_0)|\Psi_0\rangle \\ &= \hat{Q}(W_0 - \hat{H}_0)\hat{Q}|\Psi_0\rangle \\ &= \sum_{m,n=1}^{\infty} |\Phi_m\rangle\langle\Phi_m|(W_0 - \hat{H}_0)|\Phi_n\rangle\langle\Phi_n|\Psi_0\rangle \\ &= \sum_{m,n=1}^{\infty} (W_0 - W_m)\delta_{mn}|\Phi_m\rangle\langle\Phi_n|\Psi_0\rangle \\ &= \sum_{n=1}^{\infty} (W_0 - W_n)|\Phi_n\rangle\langle\Phi_n|\Psi_0\rangle \\ &= (W_0 - \hat{H}_0)\hat{Q}|\Psi_0\rangle.\end{aligned}\tag{2.97}$$

Inserting this equation into Eq. 2.93 yields

$$|\Psi_0\rangle = \frac{\hat{Q}}{W_0 - \hat{H}_0}(\hat{H}_I - \Delta E)|\Psi_0\rangle,\tag{2.98}$$

which can be recursively inserting into itself

$$|\Psi_0\rangle = \sum_{n=0}^{\infty} \left[\frac{\hat{Q}}{W_0 - \hat{H}_0} (\hat{H}_I - \Delta E) \right]^n |\Phi_0\rangle. \quad (2.99)$$

Likewise, the implications for the correlation energy is as follows:

$$\Delta E = \sum_{n=0}^{\infty} \langle \Phi_0 | \left[\frac{\hat{Q}}{W_0 - \hat{H}_0} (\hat{H}_I - \Delta E) \right]^n | \Phi_0 \rangle. \quad (2.100)$$

By organizing the terms in the above by order in \hat{H}_I , the correlation energy can be written as

$$\Delta E = \sum_{n=1}^{\infty} \Delta E^{(n)}, \quad (2.101)$$

where the first few terms in the expansion are

$$\Delta E^{(1)} = \langle \Phi_0 | \hat{H}_I | \Phi_0 \rangle \quad (2.102)$$

$$\Delta E^{(2)} = \langle \Phi_0 | \hat{H}_I \frac{\hat{Q}}{W_0 - \hat{H}_0} \hat{H}_I | \Phi_0 \rangle \quad (2.103)$$

$$\Delta E^{(3)} = \langle \Phi_0 | \hat{H}_I \frac{\hat{Q}}{W_0 - \hat{H}_0} \hat{H}_I \frac{\hat{Q}}{W_0 - \hat{H}_0} \hat{H}_I | \Phi_0 \rangle \quad (2.104)$$

$$- \langle \Phi_0 | \hat{H}_I \frac{\hat{Q}}{W_0 - \hat{H}_0} \langle \Phi_0 | \hat{H}_I | \Phi_0 \rangle \frac{\hat{Q}}{W_0 - \hat{H}_0} \hat{H}_I | \Phi_0 \rangle \quad (2.105)$$

Chapter 3

Quantum Monte Carlo

In this chapter, we introduce quantum Monte Carlo (QMC), a diverse family of ab initio many-body methods based on simulating the quantum system and computing its properties stochastically. We will cover two different flavors of QMC methods: variational Monte Carlo (VMC) and diffusion Monte Carlo (DMC). The former will take center stage in this study, as our neural-network quantum states are designed to push the limits of this method. The latter consistently yields state-of-the-art results, so it serves as a benchmark for our neural-network quantum states.

We will begin by establishing the basics of Monte Carlo sampling and integration. Both VMC and DMC use these techniques extensively, as implied by their names. Then, we will delve into the details of the conventional implementation of variational Monte Carlo. This will not only include an extensive list of commonly employed trial wave functions, but also the optimization algorithms used to update the parameters of the wave function. We will use these same optimization methods to train our neural-network quantum states, discussed throughout the remaining chapters of this thesis. We will end with a brief overview of diffusion Monte Carlo. While these methods can be used to solve for excited states, we will focus only on the ground state in our investigation.

3.1 Monte Carlo Methods

3.1.1 Markov Chain Monte Carlo Sampling

Monte Carlo sampling methods are used to generate random samples from a target probability distribution. In our case, the target distribution $P(\mathbf{X})$ is given by the square of the wave function $|\Psi(\mathbf{X})|^2$, which is a function of continuous spatial degrees of freedom, and possibly discrete spin degrees of freedom as well. However, our wave function is not guaranteed to be normalized, so it is more precise to say

$$P(\mathbf{X}) = \frac{1}{Z} |\Psi(\mathbf{X})|^2, \quad (3.1)$$

where $Z = \int d\mathbf{X} |\Psi(\mathbf{X})|^2$ is the normalization constant. It is impractical to compute Z directly since it involves a high-dimensional integral over all the spatial degrees of freedom and a sum over all the spin degrees of freedom. Luckily, there are clever tricks we can use to avoid computing Z , leading to a style of sampling called Markov Chain Monte Carlo (MCMC) sampling.

MCMC algorithms, in general, are designed to address the problem of sampling from high-dimensional and nontrivial distributions. Instead of generating independent samples from $P(\mathbf{X})$, as done with a standard hit-or-miss algorithm for instance, we can produce a discrete sequence of samples $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_t, \dots\}$, where each sample is dependent only on the previous one. Then the transition probability for the Markov process must have the simple form

$$P(\mathbf{X}_t | \mathbf{X}_{t-1}, \mathbf{X}_{t-2}, \dots, \mathbf{X}_0) = P(\mathbf{X}_t | \mathbf{X}_{t-1}), \quad (3.2)$$

for each $t > 0$. Our goal is to accurately approximate the target distribution $P(\mathbf{X})$ by the

distribution of a finite number of samples

$$P(\mathbf{X}) \approx \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{X} - \mathbf{X}_t), \quad (3.3)$$

which is only possible if $P(\mathbf{X})$ is invariant under the action of the transition probability. In other words, $P(\mathbf{X})$ is the stationary distribution for the Markov process characterized by $P(\mathbf{X}'|\mathbf{X})$ if it satisfies

$$P(\mathbf{X}') = \int d\mathbf{X} P(\mathbf{X}'|\mathbf{X})P(\mathbf{X}), \quad (3.4)$$

for all \mathbf{X}' .

In principle, the above stationary condition provides a way to evolve to $P(\mathbf{X})$ starting from any initial distribution $P(\mathbf{X}_0)$. The distribution at a later time t can be obtained by repeatedly applying the transition probability on $P(\mathbf{X}_0)$,

$$P(\mathbf{X}_t) = \int d\mathbf{X}_{t-1} d\mathbf{X}_{t-2} \cdots d\mathbf{X}_0 P(\mathbf{X}_t|\mathbf{X}_{t-1})P(\mathbf{X}_{t-1}|\mathbf{X}_{t-2}) \cdots P(\mathbf{X}_1|\mathbf{X}_0)P(\mathbf{X}_0). \quad (3.5)$$

For large enough t , the Markov chain must converge to the stationary distribution

$$\lim_{t \rightarrow \infty} P(\mathbf{X}_t) = P(\mathbf{X}), \quad (3.6)$$

regardless of the initial state \mathbf{X}_0 . The algorithm for generating the samples is as follows:

1. *Initialize.* Draw the first sample \mathbf{X}_0 from any initial probability distribution $P(\mathbf{X}_0)$.
2. *Iterate.* For each iteration $t = 1, \dots, T$, draw the next sample \mathbf{X}_t from $P(\mathbf{X}_t)$ in Eq. (3.5).

While it is simple to state the algorithm, the second step is somewhat difficult complete

without additional restrictions.

In order to treat the samples in a Markov chain as independent samples from the target distribution, the chain should be fully equilibrated and the autocorrelation between samples should be low. Let us define τ as the number of steps required to sufficiently reach the stationary distribution, also known as the equilibration or burn-in time. In addition, let us define γ as the number of steps required to sufficiently reduce the autocorrelation between samples

$$\langle \mathbf{X}_t \cdot \mathbf{X}_{t+\gamma} \rangle \approx 0. \quad (3.7)$$

Then, our target distribution can be better approximated by

$$P(\mathbf{X}) \approx \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{X} - \mathbf{X}_{\tau+\gamma t}), \quad (3.8)$$

for some large, but finite value of T . This approach allows us to faithfully treat each effective sample $\mathbf{X}_{\tau+\gamma t}$ as independent samples from $P(\mathbf{X})$, at the cost of throwing away $\tau + (\gamma - 1)T$ samples out of the total number of samples $\tau + \gamma T$.

3.1.1.1 Metropolis-Hastings

The Metropolis-Hastings algorithm is a specific type of MCMC sampling algorithm that enforces the detailed balance condition, a stricter restriction than the stationary condition in Eq. (3.4). Namely, it assumes the Markov process is reversible

$$P(\mathbf{X}' | \mathbf{X}) P(\mathbf{X}) = P(\mathbf{X} | \mathbf{X}') P(\mathbf{X}'), \quad (3.9)$$

for every pair of states \mathbf{X}, \mathbf{X}' . By integrating both sides of the detailed balance condition,

$$\int d\mathbf{X} P(\mathbf{X}'|\mathbf{X})P(\mathbf{X}) = \int d\mathbf{X} P(\mathbf{X}|\mathbf{X}')P(\mathbf{X}') = P(\mathbf{X}') \int d\mathbf{X} P(\mathbf{X}|\mathbf{X}') = P(\mathbf{X}'), \quad (3.10)$$

it automatically follows that $P(\mathbf{X})$ is the stationary distribution.

In addition, this approach involves decomposing the transition probability into two components

$$P(\mathbf{X}'|\mathbf{X}) = Q(\mathbf{X}'|\mathbf{X})A(\mathbf{X}'|\mathbf{X}), \quad (3.11)$$

where the proposal probability density $Q(\mathbf{X}'|\mathbf{X})$ suggests a candidate state \mathbf{X}' based on the current state \mathbf{X} , and the acceptance probability $A(\mathbf{X}'|\mathbf{X})$ determines if the proposed state should be accepted or rejected. Plugging this into the detailed balance equation yields

$$\frac{A(\mathbf{X}'|\mathbf{X})}{A(\mathbf{X}|\mathbf{X}')} = \frac{Q(\mathbf{X}|\mathbf{X}')P(\mathbf{X}')}{Q(\mathbf{X}'|\mathbf{X})P(\mathbf{X})}. \quad (3.12)$$

If $A(\mathbf{X}|\mathbf{X}') = 1$, then $A(\mathbf{X}'|\mathbf{X}) = \frac{Q(\mathbf{X}|\mathbf{X}')P(\mathbf{X}')}{Q(\mathbf{X}'|\mathbf{X})P(\mathbf{X})}$. Likewise, if $A(\mathbf{X}'|\mathbf{X}) = 1$ then $A(\mathbf{X}'|\mathbf{X}) = \frac{Q(\mathbf{X}'|\mathbf{X})P(\mathbf{X})}{Q(\mathbf{X}|\mathbf{X}')P(\mathbf{X}')}$. Thus the choice of $A(\mathbf{X}'|\mathbf{X})$ that maximizes the acceptance probability is

$$A(\mathbf{X}'|\mathbf{X}) \equiv \min \left\{ 1, \frac{Q(\mathbf{X}|\mathbf{X}')P(\mathbf{X}')}{Q(\mathbf{X}'|\mathbf{X})P(\mathbf{X})} \right\}. \quad (3.13)$$

In practice, the Metropolis-Hastings algorithm can be realized by separating the transition step into two parts:

1. *Initialize.* Draw the first sample \mathbf{X}_0 from an initial probability distribution $P(\mathbf{X}_0)$.
2. *Iterate.* For each iteration $t = 1, \dots, T$:
 - (a) *Propose.* Draw a candidate state \mathbf{X}' based on the previous state \mathbf{X}_{t-1} according

to the proposal probability density $Q(\mathbf{X}'|\mathbf{X}_{t-1})$.

- (b) *Calculate.* Calculate the acceptance probability of the transition $A(\mathbf{X}'|\mathbf{X}_{t-1})$.
- (c) *Accept-or-reject.* Draw a uniform random number $r \sim \mathcal{U}(0, 1)$ between 0 and 1. If $r < A(\mathbf{X}'|\mathbf{X}_{t-1})$, accept the transition by setting $\mathbf{X}_t = \mathbf{X}'$. Otherwise, reject the transition by setting $\mathbf{X}_t = \mathbf{X}_{t-1}$.

Since our states may involve discrete spin projections \mathbf{S} in addition to the continuous spatial coordinates \mathbf{R} , we will further decompose all the probabilities as

$$P(\mathbf{X}) = P(\mathbf{R})P(\mathbf{S}), \quad (3.14)$$

$$Q(\mathbf{X}'|\mathbf{X}) = Q(\mathbf{R}'|\mathbf{R})Q(\mathbf{S}'|\mathbf{S}), \quad (3.15)$$

$$A(\mathbf{X}'|\mathbf{X}) = A(\mathbf{R}'|\mathbf{R})A(\mathbf{S}'|\mathbf{S}), \quad (3.16)$$

and perform Metropolis steps in \mathbf{R} and \mathbf{S} separately. For a localized system of particles, one can take the initial distribution of the spatial degrees of freedom as

$$P(\mathbf{R}_0) \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I}), \quad (3.17)$$

with a variance σ_0^2 chosen to minimize the equilibration time τ . For an infinite system, one can take the initial distribution as a uniform distribution

$$P(\mathbf{R}_0) \sim \mathcal{U}(-L/2, L/2), \quad (3.18)$$

ranging the entire size of the d -dimensional simulation box of side length L . The proposal probability density for the spatial coordinates is commonly chosen to be a Gaussian distri-

bution

$$Q(\mathbf{R}'|\mathbf{R}) \sim \mathcal{N}(\mathbf{R}, \sigma^2 \mathbf{I}), \quad (3.19)$$

with a variance σ^2 that balances a short autocorrelation time γ with a high acceptance probability $A(\mathbf{R}'|\mathbf{R})$. Since the proposal probability is symmetric, $Q(\mathbf{R}'|\mathbf{R}) = Q(\mathbf{R}|\mathbf{R}')$, the acceptance probability for the move simplifies to

$$A(\mathbf{R}'|\mathbf{R}) = \min \left\{ 1, \frac{P(\mathbf{R}')}{P(\mathbf{R})} \right\} = \min \left\{ 1, \frac{|\Psi(\mathbf{R}', \mathbf{S})|^2}{|\Psi(\mathbf{R}, \mathbf{S})|^2} \right\}. \quad (3.20)$$

For the spin degrees of freedom \mathbf{S} , it is possible to take a similar strategy as above by initializing a random configuration of spins and proposing random spins to flip rather than coordinates to perturb. However, the Hamiltonians we will use always commute with the total spin projection operator on the z -axis

$$\hat{S}^z = \sum_{i=1}^N \hat{s}_i^z, \quad (3.21)$$

meaning the total spin projection $S^z = \sum_{i=1}^N s_i^z$ is conserved. Similarly, the total isospin projection on the z -axis T^z is also preserved for systems of nucleons, as the total isospin projection operator

$$\hat{T}^z = \sum_{i=1}^N \hat{t}_i^z \quad (3.22)$$

commutes with the nuclear Hamiltonian. To avoid sampling unphysical spin configurations, we can restrict our Metropolis walk to preserve the total spin (and isospin) projection as follows:

1. *Initialize.* Generate any spin configuration \mathbf{S}_0 with the desired total spin projection S^z

(and isospin projection T^z). Spin-up particles are assigned a spin value of $+1$, while the spin-down particles are assigned -1 .

2. *Iterate.* For each iteration $t = 1, 2, \dots, T$:

- (a) *Propose.* Choose a random pair of particles (i, j) to exchange spin degrees of freedom. The proposal spin configuration \mathbf{S}' is the same as \mathbf{S}_{t-1} , except with the spins of particle i and j exchanged

$$\begin{aligned} \mathbf{s}'_i &= \mathbf{s}_{j,t-1}, \\ \mathbf{s}'_j &= \mathbf{s}_{i,t-1}. \end{aligned} \tag{3.23}$$

A different pair can be chosen to exchange the isospin degrees of freedom, if desired.

- (b) *Compute.* Compute the acceptance probability of the proposed configuration

$$A(\mathbf{S}'|\mathbf{S}_{t-1}) = \min \left\{ 1, \frac{P(\mathbf{S}')}{P(\mathbf{S}_{t-1})} \right\} = \min \left\{ 1, \frac{|\Psi(\mathbf{R}, \mathbf{S}')|^2}{|\Psi(\mathbf{R}, \mathbf{S}_{t-1})|^2} \right\}. \tag{3.24}$$

- (c) *Accept-or-reject.* Draw a uniform random number $r \sim \mathcal{U}(0, 1)$ between 0 and 1. If $r < A(\mathbf{S}'|\mathbf{S}_{t-1})$, accept the spin exchange by setting $\mathbf{S}_t = \mathbf{S}'$. Otherwise, reject the exchange by setting $\mathbf{S}_t = \mathbf{S}_{t-1}$.

3.1.1.2 Importance Sampling

Importance sampling, a variant of the Metropolis-Hastings algorithm, draws inspiration from the Fokker-Planck equation—a generalization of the diffusion equation. The basic idea of importance sampling is to guide the proposal probability for the spatial degrees of freedom

$Q(\mathbf{R}'|\mathbf{R})$ towards regions with higher probability $P(\mathbf{R})$. For a general, time-dependent probability distribution function, the Fokker-Planck equation of a diffusion process reads

$$\frac{\partial}{\partial t}P(\mathbf{R},t) = D\boldsymbol{\nabla} \cdot (\boldsymbol{\nabla} - \mathbf{F}(\mathbf{R}))P(\mathbf{R},t), \quad (3.25)$$

where $\boldsymbol{\nabla} = (\boldsymbol{\nabla}_1, \boldsymbol{\nabla}_2, \dots, \boldsymbol{\nabla}_N)$ is the combination of all gradient operators for the N particles, $D = \frac{\hbar^2}{2m}$ is a constant diffusion coefficient, and $\mathbf{F}(\mathbf{R})$ is the drift velocity due to an external potential. Here, we have also suppressed all dependence on the spins \mathbf{S} , as they remain constant during these steps. Our goal is find the drift velocity $\mathbf{F}(\mathbf{R})$ such that our probability distribution converges to the stationary distribution determined by our wave function

$$P(\mathbf{R},t) = P(\mathbf{R}) = |\Psi(\mathbf{R})|^2 \implies \frac{\partial}{\partial t}P(\mathbf{R},t) = 0. \quad (3.26)$$

Under this condition, the Fokker-Planck equation simplifies to

$$\boldsymbol{\nabla}^2 P(\mathbf{R}) = \boldsymbol{\nabla} \cdot \mathbf{F}(\mathbf{R})P(\mathbf{R}) = \mathbf{F}(\mathbf{R}) \cdot \boldsymbol{\nabla} P(\mathbf{R}) + P(\mathbf{R})\boldsymbol{\nabla} \cdot \mathbf{F}(\mathbf{R}) \quad (3.27)$$

For the Laplacian to appear on the right-hand side, the drift force should have the form

$$\mathbf{F}(\mathbf{R}) = f(P(\mathbf{R}))\boldsymbol{\nabla} P(\mathbf{R}), \quad (3.28)$$

where f is a scalar function to be determined. Then, Eq. (3.27) can be written as

$$\begin{aligned} \nabla^2 P(\mathbf{R}) &= f(P(\mathbf{R})) \nabla P(\mathbf{R}) \cdot \nabla P(\mathbf{R}) \\ &\quad + P(\mathbf{R}) f(P(\mathbf{R})) \nabla^2 P(\mathbf{R}) \\ &\quad + P(\mathbf{R}) \nabla f(P(\mathbf{R})) \cdot \nabla P(\mathbf{R}). \end{aligned} \quad (3.29)$$

Matching the Laplacian terms, we find that

$$f(P(\mathbf{R})) = \frac{1}{P(\mathbf{R})}, \quad (3.30)$$

which also eliminates the gradient terms. Therefore, the drift velocity must have the form

$$\mathbf{F}(\mathbf{R}) = \frac{1}{P(\mathbf{R})} \nabla P(\mathbf{R}) = 2 \frac{1}{\Psi(\mathbf{R})} \nabla \Psi(\mathbf{R}), \quad (3.31)$$

where we have assumed $\Psi(\mathbf{R})$ is real-valued in the last equality.

The Fokker-Planck equation can be seen as a deterministic description of the stochastic dynamics captured by the corresponding Langevin equation,

$$\frac{\partial}{\partial t} \mathbf{R}(t) = D \mathbf{F}(\mathbf{R}) + \boldsymbol{\xi}, \quad (3.32)$$

a stochastic differential equation that describes the dynamics of a particle under the influence of both deterministic forces and random noise. In the above, $\boldsymbol{\xi} \in \mathbb{R}^{Nd}$ is a random perturbation drawn from a Gaussian distribution $\mathcal{N}(\mathbf{0}, 2D\mathbf{I})$ and $\mathbf{F}(\mathbf{R})$ is the drift force evaluated at the initial configuration, defined in Eq. (3.31). Integrating the Langevin equation over a

short time interval Δt , we obtain

$$\mathbf{R}(t + \Delta t) = \mathbf{R}(t) + D\Delta t \mathbf{F}(\mathbf{R}) + \sqrt{\Delta t} \boldsymbol{\xi}. \quad (3.33)$$

This proposal rule implies the normalized proposal density is given by

$$\begin{aligned} Q(\mathbf{R}' | \mathbf{R}, \Delta t) &= \mathcal{N}(\mathbf{R} + D\Delta t \mathbf{F}(\mathbf{R}), 2D\delta t \mathbf{I}) \\ &= \frac{1}{(4\pi D\Delta t)^{Nd/2}} e^{-(\mathbf{R}' - \mathbf{R} - D\Delta t \mathbf{F}(\mathbf{R}))^2 / 4D\Delta t}, \end{aligned} \quad (3.34)$$

which is no longer symmetric due to the drift term. Then the acceptance probability becomes

$$\begin{aligned} A(\mathbf{R}' | \mathbf{R}, \Delta t) &= \min \left\{ 1, \frac{P(\mathbf{R}')Q(\mathbf{R}|\mathbf{R}', \Delta t)}{P(\mathbf{R})Q(\mathbf{R}'|\mathbf{R}, \Delta t)} \right\} \\ &= \min \left\{ 1, \frac{|\Psi(\mathbf{R}', \mathbf{S})|^2}{|\Psi(\mathbf{R}, \mathbf{S})|^2} \frac{e^{-(\mathbf{R}-\mathbf{R}'-D\Delta t \mathbf{F}(\mathbf{R}'))^2 / 4D\Delta t}}{e^{-(\mathbf{R}'-\mathbf{R}-D\Delta t \mathbf{F}(\mathbf{R}))^2 / 4D\Delta t}} \right\}. \end{aligned} \quad (3.35)$$

The additional weight appearing next to the original acceptance probability can be easily simplified by defining the drift velocity $\mathbf{v}(\mathbf{R}) \equiv \mathbf{F}(\mathbf{R})/2$,

$$\frac{e^{-(\mathbf{R}-\mathbf{R}'-D\Delta t \mathbf{F}(\mathbf{R}'))^2 / 4D\Delta t}}{e^{-(\mathbf{R}'-\mathbf{R}-D\Delta t \mathbf{F}(\mathbf{R}))^2 / 4D\Delta t}} = e^{-\mathbf{v}(\mathbf{R}) - \mathbf{v}(\mathbf{R}') + (\mathbf{R}-\mathbf{R}') + D\Delta t(\mathbf{v}(\mathbf{R}) - \mathbf{v}(\mathbf{R}'))} \quad (3.36)$$

This weight corrects for the sampling bias introduced by our proposal distribution exploring more relevant regions of the probability distribution. An equivalent derivation can be achieved by integrating the Fokker-Planck equation itself for a small time step Δt and using the Green's function method to solve for the time-dependent probability distribution. The resulting Green's function coincides with the proposal distribution in Eq. (3.34).

3.1.2 Integration

Monte Carlo integration is a numerical technique that uses random sampling to efficiently estimate the value of a high-dimensional integral. Recall that we can expand the expectation value of operators \hat{A} in our $|\mathbf{X}\rangle$ basis as

$$\langle \hat{A} \rangle = \frac{\langle \Psi | \hat{A} | \Psi \rangle}{\langle \Psi | \Psi \rangle} = \frac{\int d\mathbf{X} \langle \Psi | \mathbf{X} \rangle \langle \mathbf{X} | \hat{A} | \Psi \rangle}{\int d\mathbf{X} \langle \Psi | \mathbf{X} \rangle \langle \mathbf{X} | \Psi \rangle}, \quad (3.37)$$

Inserting $1 = \langle \mathbf{X} | \Psi \rangle / \langle \mathbf{X} | \Psi \rangle$ into the numerator, we have

$$\langle \hat{A} \rangle = \frac{\int d\mathbf{X} |\Psi(\mathbf{X})|^2 A(\mathbf{X})}{\int d\mathbf{X} |\Psi(\mathbf{X})|^2}, \quad (3.38)$$

where $A(\mathbf{X})$ is the local quantity corresponding to the operator \hat{A}

$$A(\mathbf{X}) \equiv \frac{\langle \mathbf{X} | \hat{A} | \Psi \rangle}{\langle \mathbf{X} | \Psi \rangle}. \quad (3.39)$$

Notice that

$$P(\mathbf{X}) = \frac{|\Psi(\mathbf{X})|^2}{\int d\mathbf{X} |\Psi(\mathbf{X})|^2} \quad (3.40)$$

is none other than the normalized probability distribution given by our wave function from Eq. (3.1). Then the expectation value becomes

$$\langle \hat{A} \rangle = \int d\mathbf{X} P(\mathbf{X}) A(\mathbf{X}), \quad (3.41)$$

where the integral over \mathbf{X} can be appropriately decomposed into a sum over \mathbf{S} and an integral over \mathbf{R} whenever necessary. Any integral with this form can be approximated as a

simple average over local quantities

$$\langle \hat{A} \rangle \approx \frac{1}{T} \sum_{t=1}^T A(\mathbf{X}_t) \equiv \langle A(\mathbf{X}) \rangle, \quad (3.42)$$

where the average is weighted over T samples from the probability distribution $\mathbf{X}_t \sim P(\mathbf{X})$.

By the law of large numbers, the average of the local quantities $\langle A(\mathbf{X}) \rangle$ will converge to the expected value of the operator $\langle \hat{A} \rangle$ as T increases. Notice the notation we have chosen for the expectation value of the operator versus the weighted sum over the local quantities, as it differs from other texts. We express both with angled brackets, but the explicit dependence on \mathbf{X} for the latter is meant to imply \mathbf{X} is pulled from our probability distribution. We choose this notation $\langle A(\mathbf{X}) \rangle$ over the more common $\langle A_L \rangle$ to leave room for other subscripts.

3.2 Variational Monte Carlo

Variational Monte Carlo (VMC) is a direct application of the variational principle from Sec. 2.4 and Monte Carlo integration from Sec. 3.1.2. The core idea is to take some parameterized trial wave function $\Psi_{\boldsymbol{\theta}}(\mathbf{X})$, and minimize the corresponding expectation value of the energy with respect to the variational parameters $\boldsymbol{\theta}$

$$E[\Psi_{\boldsymbol{\theta}}] \equiv \frac{\langle \Psi_{\boldsymbol{\theta}} | \hat{H} | \Psi_{\boldsymbol{\theta}} \rangle}{\langle \Psi_{\boldsymbol{\theta}} | \Psi_{\boldsymbol{\theta}} \rangle}, \quad (3.43)$$

$$\min_{\boldsymbol{\theta}} E[\Psi_{\boldsymbol{\theta}}] \geq E_0. \quad (3.44)$$

The variational principle guarantees that the variational energy $E[\Psi_{\boldsymbol{\theta}}]$ establishes a strict upperbound on the true ground state energy of the system. To find the optimal parameters, we use gradient descent methods, a family of iterative optimization algorithms discussed in

more detail in Sec. 3.2.2. For now, we just compute the gradient of the variational energy with respect to the parameters

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} E[\Psi_{\boldsymbol{\theta}}] &= \frac{1}{\langle \Psi_{\boldsymbol{\theta}} | \Psi_{\boldsymbol{\theta}} \rangle} \left(\langle \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} | \hat{H} | \Psi_{\boldsymbol{\theta}} \rangle + \langle \Psi_{\boldsymbol{\theta}} | \hat{H} | \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} \rangle \right) \\ &\quad - \frac{\langle \Psi_{\boldsymbol{\theta}} | \hat{H} | \Psi_{\boldsymbol{\theta}} \rangle}{\langle \Psi_{\boldsymbol{\theta}} | \Psi_{\boldsymbol{\theta}} \rangle^2} \left(\langle \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} | \Psi_{\boldsymbol{\theta}} \rangle + \langle \Psi_{\boldsymbol{\theta}} | \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} \rangle \right) \\ &= 2 \left(\frac{\langle \Psi_{\boldsymbol{\theta}} | \hat{H} | \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} \rangle}{\langle \Psi_{\boldsymbol{\theta}} | \Psi_{\boldsymbol{\theta}} \rangle} - E[\Psi_{\boldsymbol{\theta}}] \frac{\langle \Psi_{\boldsymbol{\theta}} | \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} \rangle}{\langle \Psi_{\boldsymbol{\theta}} | \Psi_{\boldsymbol{\theta}} \rangle} \right),\end{aligned}\tag{3.45}$$

where we have assumed our trial wave function is real-valued in the last equality. To approximate all the high-dimensional integrals above, we employ Monte Carlo integration. For the variational energy, we simply write

$$E[\Psi_{\boldsymbol{\theta}}] \approx \frac{1}{T} \sum_{t=1}^T E_{\boldsymbol{\theta}}(\mathbf{X}_t) \equiv \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle,\tag{3.46}$$

where the configurations \mathbf{X}_t are sampled from $|\Psi_{\boldsymbol{\theta}}(\mathbf{X})|^2$, T is the number of samples, and the local energy is defined as

$$E_{\boldsymbol{\theta}}(\mathbf{X}) \equiv \frac{\langle \mathbf{X} | \hat{H} | \Psi_{\boldsymbol{\theta}} \rangle}{\langle \mathbf{X} | \Psi_{\boldsymbol{\theta}} \rangle} = \frac{\hat{H} \Psi_{\boldsymbol{\theta}}(\mathbf{X})}{\Psi_{\boldsymbol{\theta}}(\mathbf{X})}.\tag{3.47}$$

For the gradient in Eq. (3.45), we first define the local gradient of the wave function as

$$\mathcal{O}_{\boldsymbol{\theta}}(\mathbf{X}) \equiv \frac{\langle \mathbf{X} | \nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}} \rangle}{\langle \mathbf{X} | \Psi_{\boldsymbol{\theta}} \rangle} = \frac{\nabla_{\boldsymbol{\theta}} \Psi_{\boldsymbol{\theta}}(\mathbf{X})}{\Psi_{\boldsymbol{\theta}}(\mathbf{X})} = \nabla_{\boldsymbol{\theta}} \log \Psi_{\boldsymbol{\theta}}(\mathbf{X}),\tag{3.48}$$

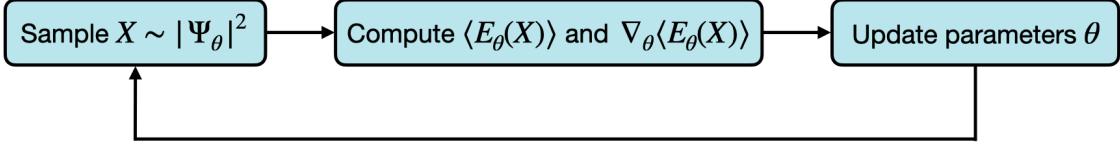


Figure 3.1: The cyclic workflow of the variational Monte Carlo algorithm.

which allows us to write

$$\begin{aligned} \frac{\langle \Psi_\theta | \nabla_\theta \Psi_\theta \rangle}{\langle \Psi_\theta | \Psi_\theta \rangle} &= \frac{\int d\mathbf{X} \langle \Psi_\theta | \mathbf{X} \rangle \langle \mathbf{X} | \nabla_\theta \Psi_\theta \rangle}{\int d\mathbf{X} \langle \Psi_\theta | \mathbf{X} \rangle \langle \mathbf{X} | \Psi_\theta \rangle} = \frac{\int d\mathbf{X} |\Psi_\theta(\mathbf{X})|^2 \mathbf{O}_\theta(\mathbf{X})}{\int d\mathbf{X} |\Psi_\theta(\mathbf{X})|^2} \\ &\approx \frac{1}{T} \sum_{t=1}^T \mathbf{O}_\theta(\mathbf{X}_t) \equiv \langle \mathbf{O}_\theta(\mathbf{X}) \rangle, \end{aligned} \quad (3.49)$$

$$\begin{aligned} \frac{\langle \Psi_\theta | \hat{H} | \nabla_\theta \Psi_\theta \rangle}{\langle \Psi_\theta | \Psi_\theta \rangle} &= \frac{\int d\mathbf{X} \langle \Psi_\theta | \hat{H} | \mathbf{X} \rangle \langle \mathbf{X} | \nabla_\theta \Psi_\theta \rangle}{\int d\mathbf{X} \langle \Psi_\theta | \mathbf{X} \rangle \langle \mathbf{X} | \Psi_\theta \rangle} = \frac{\int d\mathbf{X} |\Psi_\theta(\mathbf{X})|^2 E_\theta(\mathbf{X}) \mathbf{O}_\theta(\mathbf{X})}{\int d\mathbf{X} |\Psi_\theta(\mathbf{X})|^2} \\ &\approx \frac{1}{T} \sum_{t=1}^T E_\theta(\mathbf{X}_t) \mathbf{O}_\theta(\mathbf{X}_t) \equiv \langle E_\theta(\mathbf{X}) \mathbf{O}_\theta(\mathbf{X}) \rangle, \end{aligned} \quad (3.50)$$

where we have once again assumed the wave function is real. Therefore, the gradient can be approximated as

$$\nabla_\theta E[\Psi_\theta] \approx \nabla_\theta \langle E_\theta(\mathbf{X}) \rangle = 2 \left(\langle E_\theta(\mathbf{X}) \mathbf{O}_\theta(\mathbf{X}) \rangle - \langle E_\theta(\mathbf{X}) \rangle \langle \mathbf{O}_\theta(\mathbf{X}) \rangle \right). \quad (3.51)$$

The variational Monte Carlo algorithm is a simple cyclic procedure; we sample configurations from the square of the wave function, estimate the variational energy and its gradient, and update the parameters. There are many ways we can complete the last step of this cycle, which we will cover in Sec. 3.2.2. More information about the implementation will be discussed in Chapter 5.

3.2.1 Trial Wave Functions

The choice of trial wave function, or ansatz, is the most crucial ingredient of a variational Monte Carlo calculation. If chosen incorrectly, the upperbound on the ground state energy provided by the variational principle is effectively meaningless. For this reason, traditional VMC calculations rely on carefully curating a trial wave function to capture the essential features and correlations of the system under study. This involves considering the symmetries, boundary conditions, and known properties of the system, as well as incorporating relevant physical insights and intuition. Consequently, these trial wave functions are intricately tailored to the specific Hamiltonians they were originally devised for, making them ill-suited for application to similar, albeit distinct, Hamiltonians without substantial structural modifications.

In this section, we will showcase a diverse range of conventional trial wave functions and analyze their suitability for different types of systems. Our objective identify the limitations of these methods while appreciating their inherent strengths. Through this exploration, we aim to develop a deeper understanding of the advantages offered by neural networks and their potential to overcome the limitations commonly associated with traditional strategies.

3.2.1.1 Kato's Cusp Condition

In the presence of singularities in the potential $V(\mathbf{X})$, Kato's cusp condition asserts that the local kinetic energy $K_{\theta}(\mathbf{X})$ must precisely counterbalance the diverging potential such that the total local energy remains finite as any two particles approach each other. More

formally, the cusp condition states

$$\lim_{r_{ij} \rightarrow 0} E_{\theta}(\mathbf{X}) = \lim_{r_{ij} \rightarrow 0} (K_{\theta}(\mathbf{X}) + V(\mathbf{X})) < \infty, \quad (3.52)$$

for all pairs i, j , where $r_{ij} = |\mathbf{r}_i - \mathbf{r}_j|$ is defined as the Euclidean distance between particles i and j , and the local kinetic energy is given by

$$K_{\theta}(\mathbf{X}) = \sum_{i=1}^N -\frac{\hbar^2}{2m} \frac{\nabla_i^2 \Psi_{\theta}(\mathbf{X})}{\Psi_{\theta}(\mathbf{X})}. \quad (3.53)$$

For strongly-correlated systems, it is especially important to design a trial wave function that upholds the cusp condition throughout training. Even a finite, but hard interaction potential may warrant a cusp condition to ensure stability throughout the optimization of $\Psi_{\theta}(\mathbf{X})$. A "hard" potential is often used to describe an interaction potential that exhibits abrupt changes in magnitude as the distance between particles varies. These potentials typically have a large scattering length while being short-range in nature.

3.2.1.2 Jastrow Factor

Enforcement of Kato's cusp condition usually comes in the form of a Jastrow factor, also known as a Jastrow wave function. For a system of identical particles, it can be any permutation-invariant function of the particles, as it preserves the symmetry of the overall wave function. The goal of the Jastrow factor is to incorporate the effects of particle-particle interactions into the trial wave function. The form of the Jastrow factor depends on the specific system and the character of the interactions being considered.

For a system of bosons, the ground state wave function must be positive-definite and symmetric with respect to particle exchange. Since we can ignore spin, the generic trial

wave function for bosons can be written as

$$\Psi_{\theta}(\mathbf{R}) = \Phi_0(\mathbf{R}) e^{J(\mathbf{R})}, \quad (3.54)$$

where the first term represents the ground state of the non-interacting problem and the second term is the symmetric Jastrow factor containing the inter-particle correlations. By writing

$$\Phi_0(\mathbf{R}) = \prod_{i=1}^N \xi_0(\mathbf{r}_i), \quad (3.55)$$

the non-interacting problem becomes simple to solve because it can be separated into N identical and independent parts

$$\begin{aligned} \left(\hat{K} + V_1(\mathbf{R}) \right) \Phi_0(\mathbf{R}) &= E_0^{\text{nonint}} \Phi_0(\mathbf{R}), \\ \left(-\frac{\hbar^2}{2m} \nabla_i^2 + v_1(\mathbf{r}_i) \right) \xi_0(\mathbf{r}_i) &= \varepsilon_0 \xi_0(\mathbf{r}_i), \\ E_0^{\text{nonint}} &= \sum_{i=1}^N \varepsilon_0 = N \varepsilon_0, \end{aligned} \quad (3.56)$$

which can often be solved analytically. Each of the N bosons occupy the lowest single-particle state with energy ε_0 , so the non-interacting energy is simply their sum.

Examples of exactly-solvable non-interacting systems include the free system

$$\hat{H}^{\text{free}} = \hat{K}, \quad (3.57)$$

where the single-particle ground state is the trivial plane-wave orbital

$$\psi_0^{\text{free}}(\mathbf{r}_i) = 1, \quad (3.58)$$

with zero energy $\varepsilon_0^{\text{free}} = 0$, in which case, the trial wave function is just the Jastrow factor on its own. Another common example is the non-interacting system of isotropic harmonic oscillators

$$\hat{H}^{\text{ho}} = \hat{K} + \sum_{i=1}^N \frac{1}{2} m \omega^2 \mathbf{r}_i^2, \quad (3.59)$$

where ω is the oscillation frequency. The single-particle ground state wave function is a Gaussian

$$\psi_0^{\text{ho}}(\mathbf{r}_i) = e^{-\frac{m\omega}{2\hbar} \mathbf{r}_i^2} \quad (3.60)$$

with the corresponding energy

$$\varepsilon_0^{\text{ho}} = \frac{d}{2} \hbar \omega. \quad (3.61)$$

Constructing the Jastrow factor is often much more difficult than the non-interacting ground state. Nonetheless, there are some exactly-solvable interactions that serve as great examples for demonstrating the general strategy. Consider the Calogero-Sutherland model of bosons in a one-dimensional harmonic oscillator trap interacting with an inverse-squared potential

$$\hat{H}^{\text{cs}} = \hat{H}^{\text{ho}} + V^{\text{cs}}(\mathbf{R}), \quad V^{\text{cs}}(\mathbf{R}) = \frac{\hbar^2}{m} \sum_{i < j}^N \frac{\beta(\beta - 1)}{x_{ij}^2}, \quad (3.62)$$

where $x_{ij} = |x_i - x_j|$ is the distance between the pair i and j . Say we look for an ansatz with the form

$$\Psi_{\boldsymbol{\theta}}^{\text{cs}}(\mathbf{R}) = \Phi_0^{\text{ho}}(\mathbf{R}) e^{J(\mathbf{R})}, \quad (3.63)$$

where $\Phi_0^{\text{ho}}(\mathbf{R}) \equiv \prod_{i=1}^N \xi_0^{\text{ho}}$ is the ground state of the one-dimensional non-interacting harmonic oscillators. Then according to Kato's cusp condition, we need to find the function

$J(\mathbf{R})$ that satisfies

$$\lim_{x_{ij} \rightarrow 0} E_L(\mathbf{R}) < \infty, \quad (3.64)$$

for all configurations $\mathbf{R} = (x_1, x_2, \dots, x_N)$ and pairs i, j . It is straightforward to show that

$$\frac{1}{\Psi_\theta} \frac{\partial^2}{\partial x_i^2} \Psi_\theta = \frac{1}{\Phi_0^{\text{ho}}} \frac{\partial^2 \Phi_0^{\text{ho}}}{\partial x_i^2} + \frac{2}{\Phi_0^{\text{ho}}} \frac{\partial \Phi_0^{\text{ho}}}{\partial x_i} \frac{\partial J(\mathbf{R})}{\partial x_i} + \frac{\partial^2 J(\mathbf{R})}{\partial x_i^2} + \left(\frac{\partial J(\mathbf{R})}{\partial x_i} \right)^2, \quad (3.65)$$

for any particular particle i . Since the ground state for the non-interacting system is $\Phi_0^{\text{ho}} = \prod_{i=1}^N e^{-\frac{m\omega}{2\hbar} x_i}$ in one dimension, the local kinetic energy then evaluates to

$$K_L(\mathbf{R}) = K_L^{\text{ho}}(\mathbf{R}) + \frac{\hbar^2}{2m} \sum_{i=1}^N \left[\frac{m\omega}{\hbar} x_i \frac{\partial J(\mathbf{R})}{\partial x_i} - \frac{\partial^2 J(\mathbf{R})}{\partial x_i^2} - \left(\frac{\partial J(\mathbf{R})}{\partial x_i} \right)^2 \right], \quad (3.66)$$

where $K_L^{\text{ho}}(\mathbf{R}) = \frac{1}{\Phi_0^{\text{ho}}} \hat{K} \Phi_0^{\text{ho}}$ is the local kinetic energy for the non-interacting harmonic oscillators. In order to cancel the divergences coming from the interaction potential in Eq. (3.62), the right-hand side must scale as $\sim \frac{1}{x_{ij}^2}$ as $x_{ij} \rightarrow \infty$ for all pairs. The simplest Jastrow correlator with the proper scaling behavior is

$$J(\mathbf{R}) = b \sum_{i < j} \log x_{ij}, \quad (3.67)$$

for some constant b . By matching the coefficients in Eq. (3.66), we find that $b = \beta$, and the resulting ansatz is the exact solution for the ground state

$$\Psi_0^{\text{cs}} = \Phi_0^{\text{ho}} \prod_{i < j}^N x_{ij}^\beta, \quad (3.68)$$

with the corresponding energy

$$E_0^{\text{cs}} = E_0^{\text{ho}} + \frac{1}{2}\beta N(N-1)\hbar\omega. \quad (3.69)$$

In this particular example, we were fortunate to uncover an exact solution, making the choice of a parameterized wave function an obvious one. Specifically, we can use the Jastrow factor as defined in Eq. (3.67), with a single variational parameter denoted as b . However, by the time we have found the optimal design for $J(\mathbf{R})$, we essentially solved the entire problem, rendering the application of variational Monte Carlo pointless except for benchmarking purposes.

Although the cusp condition can be solved exactly in the aforementioned Calogero-Sutherland model, this is not the case for the vast majority of interesting systems. Furthermore, determining the appropriate form of the Jastrow factor necessitated a specific choice for the interaction potential. This means that the knowledge we gained throughout this process is not immediately applicable for other problems. We will see that when utilizing neural networks as trial wave functions, it is possible to forgo manually enforcing the cusp condition, as the neural networks can autonomously discover it, allowing a single ansatz to be reused for a large class of problems.

3.2.1.3 Slater Determinant

For systems of fermions, the situation becomes more complicated because the ground state wave function must be antisymmetric with respect to particle exchange, and we have to be mindful of spins. We can still use a symmetric Jastrow factor to capture the correlations in our trial wave function, but we need to replace the symmetric non-interacting ground state

in Eq. 3.54 with an antisymmetric one. Hence, we reintroduce the spins and take our ansatz to be

$$\Psi_{\boldsymbol{\theta}}(\mathbf{X}) = \Phi_0(\mathbf{X}) e^{J(\mathbf{X})}, \quad (3.70)$$

where $\Phi_0(\mathbf{X})$ is a Slater determinant involving the N single-particle spin-orbitals evaluated for each of the N particles

$$\Phi_0(\mathbf{X}) = \det \left[\varphi_{\alpha}(\mathbf{x}_i) \right] = \det \begin{bmatrix} \varphi_1(\mathbf{x}_1) & \varphi_1(\mathbf{x}_2) & \cdots & \varphi_1(\mathbf{x}_N) \\ \varphi_2(\mathbf{x}_1) & \varphi_2(\mathbf{x}_2) & \cdots & \varphi_2(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_N(\mathbf{x}_1) & \varphi_N(\mathbf{x}_2) & \cdots & \varphi_N(\mathbf{x}_N) \end{bmatrix} \quad (3.71)$$

The quantum numbers of the single-particle states are labeled generically by α , while the particles are labeled with the index i . As with bosons, we obtain the single-particle states by solving the non-interacting problem, but because of the Pauli exclusion principle, we must find the spectrum rather than just the ground state. The fermions are filled in a way that minimizes the total energy of the system, starting from the lowest energy state and moving upwards. This filling continues until the highest-occupied state, often referred to as the Fermi level, is reached. The Fermi level represents the boundary between filled and unfilled states in a system of fermions, and its corresponding single-particle energy is denoted as ε_F .

As another simple example, let us consider the circular quantum-dots system,

$$\hat{H}^{\text{qd}} = \hat{H}^{\text{ho}} + V^{\text{qd}}(\mathbf{R}), \quad V^{\text{qd}}(\mathbf{R}) = \frac{e^2}{4\pi\varepsilon_0} \sum_{i < j} \frac{1}{r_{ij}}, \quad (3.72)$$

a two-dimensional system of electrons trapped in a harmonic oscillator well. For convenience,

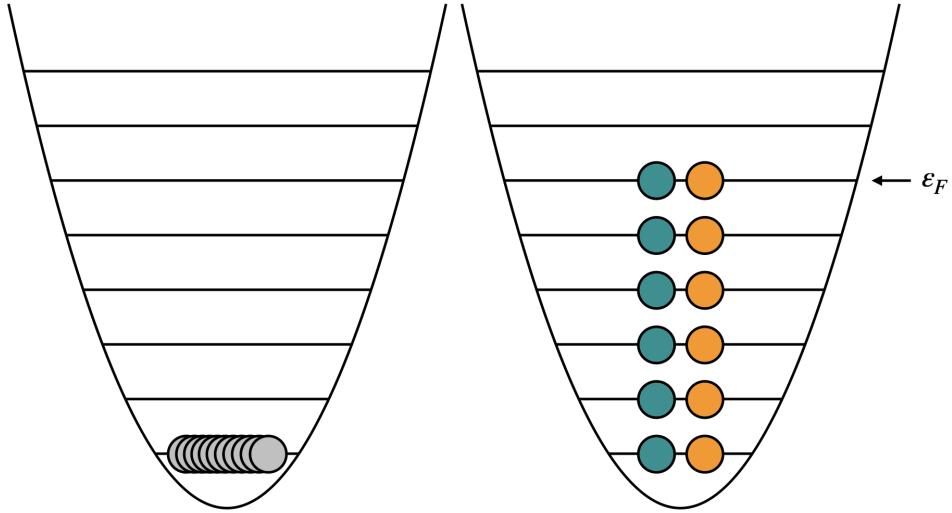


Figure 3.2: Cartoon of filling single-particle states for bosons and fermions.

we set $e = 4\pi\varepsilon_0 = 1$ and assume the system is unpolarized, meaning there is an equal number of spin-up and spin-down electrons, $N_\uparrow = N_\downarrow = N/2$. Then in two-dimensions, the single-particle wave functions for the non-interacting problem are

$$\varphi_\alpha(\mathbf{r}_i) = e^{-\frac{m\omega}{2\hbar} \mathbf{r}_i^2} H_{m\alpha} \left(\sqrt{\frac{m\omega}{\hbar}} x_i \right) H_{n\alpha} \left(\sqrt{\frac{m\omega}{\hbar}} y_i \right) \langle s_i | \sigma_\alpha \rangle, \quad (3.73)$$

where $i = (\mathbf{r}_i, s_i) = (x_i, y_i, s_i)$ labels the single-particle degrees of freedom, $\alpha = (m_\alpha, n_\alpha, \sigma_\alpha)$ labels the single-particle states with $\sigma_\alpha \in \{\uparrow, \downarrow\}$, and $H_n(x)$ are the Hermite polynomials of degree n . By inserting the above expression into Eq. (3.71), we can automatically constrain the antisymmetry of the wave function.

3.2.1.4 Backflow Transformations

The Slater determinant discussed in the preceding section relies on single-particle orbitals derived from the non-interacting problem. Consequently, the nodal structure of the overall wave function remains fixed during the training process, as the Slater determinant contains

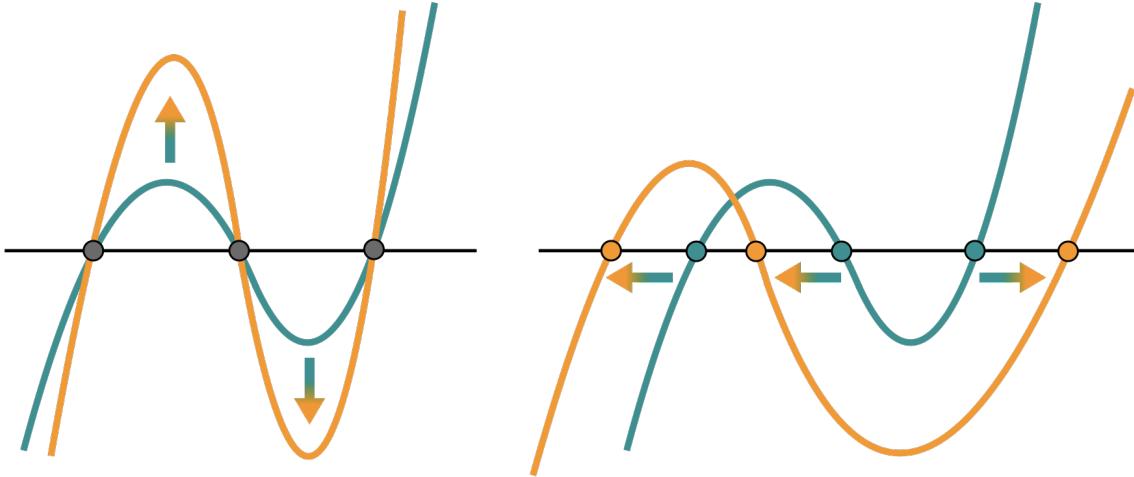


Figure 3.3: The effects of a symmetric Jastrow factor and a backflow transformation on the nodes of a fermionic wave function.

no variational parameters. To improve the nodal structure of the Slater determinant, we can incorporate backflow transformations, which aim to capture the influence of all surrounding particles on the state of an individual particle. While Jastrow factors technically encode some backflow correlations as well, we only utilize positive definite Jastrows that are incapable of changing the nodal structure. For the sake of clarity, we will specifically use the term "backflow transformation" for the technique that facilitates alterations to the nodal structure of the Slater determinant.

A backflow transformation typically alters the single-particle positions as

$$\mathbf{r}_i \mapsto \mathbf{r}_i + \sum_{j \neq i}^N \eta(r_{ij})(\mathbf{r}_i - \mathbf{r}_j), \quad (3.74)$$

where $\eta(r)$ is the backflow correlation function. Then the transformed coordinates, along

with the untransformed spin, serve as input to the non-interacting spin-orbitals $\varphi_\alpha(\mathbf{x}_i)$. Determining the functional form of $\eta(r)$, much like the Jastrow factor, necessitates the imposition of known physical insights about the system. However, accomplishing this task manually can be quite challenging in general.

Alternate forms for the backflow transformation can be chosen, as long as the transformation is permutation equivariant, i.e. for a given particle i , the transformation is invariant under any permutation of all the other particles $j \neq i$. With this constraint, the anti-symmetry of the Slater determinant is preserved. We write this arbitrary transformation as

$$\mathbf{x}_i \mapsto \boldsymbol{\eta}(\mathbf{x}_i, \{\mathbf{x}_{j \neq i}\}), \quad (3.75)$$

which suggests a means to potentially incorporate the spins of all particles into the backflow transformations, in addition to the spatial degrees of freedom.

3.2.1.5 Bardeen-Cooper-Schreifer Wave Function

In fermionic systems with strong pairing correlations, it is often insufficient to construct an antisymmetric ansatz from a Slater determinant of single-particle states. Even with the utilization of multiple Slater determinants or the incorporation of backflow transformations, the antisymmetric component of the wave function remains constrained within the single-particle mean-field framework. Thus, we graduate to the number-projected Bardeen-Cooper-Schreifer (BCS) wave function, an antisymmetrized product of singlet-pair orbitals. This ansatz is suitable for unpolarized systems of fermions with an attractive interaction sufficiently strong enough in the s -wave channel to facilitate the formation of Cooper pairs.

It can be written as the determinant of an $N/2 \times N/2$ matrix,

$$\Phi_0(\mathbf{X}) = \det \left[\xi(\mathbf{r}_i^\uparrow, \mathbf{r}_j^\downarrow) \right] = \det \begin{bmatrix} \xi(\mathbf{r}_1^\uparrow, \mathbf{r}_1^\downarrow) & \xi(\mathbf{r}_1^\uparrow, \mathbf{r}_2^\downarrow) & \cdots & \xi(\mathbf{r}_1^\uparrow, \mathbf{r}_{N/2}^\downarrow) \\ \xi(\mathbf{r}_2^\uparrow, \mathbf{r}_1^\downarrow) & \xi(\mathbf{r}_2^\uparrow, \mathbf{r}_2^\downarrow) & \cdots & \xi(\mathbf{r}_2^\uparrow, \mathbf{r}_{N/2}^\downarrow) \\ \vdots & \vdots & \ddots & \vdots \\ \xi(\mathbf{r}_{N/2}^\uparrow, \mathbf{r}_1^\downarrow) & \xi(\mathbf{r}_{N/2}^\uparrow, \mathbf{r}_2^\downarrow) & \cdots & \xi(\mathbf{r}_{N/2}^\uparrow, \mathbf{r}_{N/2}^\downarrow) \end{bmatrix}, \quad (3.76)$$

where $\xi(\mathbf{r}_i^\uparrow, \mathbf{r}_j^\downarrow)$ is the unique singlet-pairing orbital evaluated for each combination of spin-up and spin-down fermions. It has been shown that the above BCS wave function can be expanded to include single-particle orbitals $\varphi_\alpha(\mathbf{r})$ and remain antisymmetric. Define $P = \min\{N_\uparrow, N_\downarrow\}$ as the number of pairs and $U = |N_\uparrow - N_\downarrow|$ as the number of unpaired particles. Then the augmented BCS wave function becomes the determinant of a $(P + U) \times (P + U)$ matrix,

$$\Phi_0(\mathbf{X}) = \det \begin{bmatrix} \xi(\mathbf{r}_1^\uparrow, \mathbf{r}_1^\downarrow) & \xi(\mathbf{r}_1^\uparrow, \mathbf{r}_2^\downarrow) & \cdots & \xi(\mathbf{r}_1^\uparrow, \mathbf{r}_P^\downarrow) & \varphi_1(\mathbf{r}_1^\uparrow) & \varphi_2(\mathbf{r}_1^\uparrow) & \cdots & \varphi_U(\mathbf{r}_1^\uparrow) \\ \xi(\mathbf{r}_2^\uparrow, \mathbf{r}_1^\downarrow) & \xi(\mathbf{r}_2^\uparrow, \mathbf{r}_2^\downarrow) & \cdots & \xi(\mathbf{r}_2^\uparrow, \mathbf{r}_P^\downarrow) & \varphi_1(\mathbf{r}_2^\uparrow) & \varphi_2(\mathbf{r}_2^\uparrow) & \cdots & \varphi_U(\mathbf{r}_2^\uparrow) \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \xi(\mathbf{r}_P^\uparrow, \mathbf{r}_1^\downarrow) & \xi(\mathbf{r}_P^\uparrow, \mathbf{r}_2^\downarrow) & \cdots & \xi(\mathbf{r}_P^\uparrow, \mathbf{r}_P^\downarrow) & \varphi_1(\mathbf{r}_P^\uparrow) & \varphi_2(\mathbf{r}_P^\uparrow) & \cdots & \varphi_U(\mathbf{r}_P^\uparrow) \\ \varphi_1(\mathbf{r}_1^\downarrow) & \varphi_1(\mathbf{r}_2^\downarrow) & \cdots & \varphi_1(\mathbf{r}_P^\downarrow) & 0 & 0 & \cdots & 0 \\ \varphi_2(\mathbf{r}_1^\downarrow) & \varphi_2(\mathbf{r}_2^\downarrow) & \cdots & \varphi_2(\mathbf{r}_P^\downarrow) & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \varphi_U(\mathbf{r}_1^\downarrow) & \varphi_U(\mathbf{r}_2^\downarrow) & \cdots & \varphi_U(\mathbf{r}_P^\downarrow) & 0 & 0 & \cdots & 0 \end{bmatrix}. \quad (3.77)$$

The BCS wave function provides a more accurate description of the system's collective behavior in the presence of strong pairing correlations, such as in fermionic superfluids.

However, it is limited in its ability to capture triplet correlations and is not applicable to Hamiltonians that involve spin exchange.

3.2.1.6 Pfaffian Wave Function

To overcome some the limitations of the BCS wave function, alternative wave functions based on the Pfaffian have been employed extensively in electronic structure as well as nuclear applications. Similar to the determinant, the Pfaffian of a matrix is a scalar value calculated by taking sums of products of permuted matrix elements. It is defined only for even-dimensional, skew-symmetric matrices with the explicit formula

$$\text{pf}[A] = \frac{1}{2^n n!} \sum_{\hat{P} \in S_{2n}} \sigma(P) \prod_{i=1}^n a_{p(2i-1), p(2i)}, \quad (3.78)$$

where we have assumed A is a $2n \times 2n$ matrix. The Pfaffian allows us to write the most generic antisymmetric wave function constructed from pairing orbitals rather than single-particle ones. Assuming N is even, the Pfaffian wave function takes the form

$$\Phi_0(\mathbf{X}) = \text{pf} [\phi(\mathbf{x}_i, \mathbf{x}_j)] = \text{pf} \begin{bmatrix} 0 & \phi(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1, \mathbf{x}_N) \\ \phi(\mathbf{x}_2, \mathbf{x}_1) & 0 & \cdots & \phi(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N, \mathbf{x}_1) & \phi(\mathbf{x}_N, \mathbf{x}_2) & \cdots & 0 \end{bmatrix}, \quad (3.79)$$

where $\phi(\mathbf{x}_i, \mathbf{x}_j)$ is an *antisymmetric* pairing spin-orbital, $\phi(\mathbf{x}_i, \mathbf{x}_j) = -\phi(\mathbf{x}_j, \mathbf{x}_i)$. If we expand the matrix above to include unpaired single-particle spin-orbitals, mirroring the approach we used for the BCS wave function, the wave function becomes the Pfaffian of an

$(N + U) \times (N + U)$ matrix

$$\Phi_0(\mathbf{X}) = \text{pf} \begin{bmatrix} \boldsymbol{\phi} & \boldsymbol{\varphi} \\ -\boldsymbol{\varphi}^T & \mathbf{0} \end{bmatrix}, \quad (3.80)$$

where $\boldsymbol{\phi}$ is the same $N \times N$ matrix displayed in Eq. (3.79), $\boldsymbol{\varphi}$ is an $N \times U$ matrix constructed from the occupied single-particle orbitals

$$\boldsymbol{\varphi} \equiv \begin{bmatrix} \varphi_1(\mathbf{x}_1) & \varphi_2(\mathbf{x}_1) & \cdots & \varphi_U(\mathbf{x}_1) \\ \varphi_1(\mathbf{x}_2) & \varphi_2(\mathbf{x}_2) & \cdots & \varphi_U(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_1(\mathbf{x}_N) & \varphi_2(\mathbf{x}_N) & \cdots & \varphi_U(\mathbf{x}_N) \end{bmatrix}, \quad (3.81)$$

and U is the number of unpaired particles. The pairing spin-orbital is commonly decomposed into explicit singlet and triplet contributions in order to write down the spatial dependence more easily

$$\begin{aligned} \phi(\mathbf{x}_i, \mathbf{x}_j) = & \xi(\mathbf{r}_i, \mathbf{r}_j) \langle s_i, s_j | (|\uparrow\downarrow\rangle - |\downarrow\uparrow\rangle) / \sqrt{2} \\ & + \chi^{\uparrow\uparrow}(\mathbf{r}_i, \mathbf{r}_j) \langle s_i, s_j | \uparrow\uparrow\rangle \\ & + \chi^{\uparrow\downarrow}(\mathbf{r}_i, \mathbf{r}_j) \langle s_i, s_j | (|\uparrow\downarrow\rangle + |\downarrow\uparrow\rangle) / \sqrt{2} \\ & + \chi^{\downarrow\downarrow}(\mathbf{r}_i, \mathbf{r}_j) \langle s_i, s_j | \downarrow\downarrow\rangle. \end{aligned} \quad (3.82)$$

The singlet pairing orbital $\xi(\mathbf{r}_i, \mathbf{r}_j)$ is the same as $\xi(\mathbf{r}_i^\uparrow, \mathbf{r}_j^\downarrow)$ in the BCS wave function, and must be even with respect to the spatial coordinates. The newly introduced triplet pairing orbitals $\chi^{\uparrow\uparrow}(\mathbf{r}_i, \mathbf{r}_j)$, $\chi^{\uparrow\downarrow}(\mathbf{r}_i, \mathbf{r}_j)$, $\chi^{\downarrow\downarrow}(\mathbf{r}_i, \mathbf{r}_j)$ must all be odd.

When the Hamiltonian does not exchange spin, in condensed matter problems for instance, we can further decompose the $\boldsymbol{\phi}$ and $\boldsymbol{\varphi}$ matrices into separate blocks corresponding

to the singlet and triplet contributions, and the spin-up and spin-down unpaired contributions, respectively. This singlet-triplet-unpaired (STU) Pfaffian ansatz, developed in Ref. [5], is written as

$$\Phi_0(\mathbf{X}) = \text{pf} \begin{bmatrix} \chi^{\uparrow\uparrow} & \xi^{\uparrow\downarrow} & \varphi^\uparrow \\ -\xi^{\uparrow\downarrow T} & \chi^{\downarrow\downarrow} & \varphi^\downarrow \\ -\varphi^{\uparrow T} & -\varphi^{\downarrow T} & \mathbf{0} \end{bmatrix}, \quad (3.83)$$

where we have omitted the $\chi^{\uparrow\downarrow}(\mathbf{r}_i, \mathbf{r}_j)$ term in Eq. (3.82). The singlet block $\xi^{\uparrow\downarrow}$ is an $N_\uparrow \times N_\downarrow$ matrix constructed from applying the singlet pairing orbital $\xi(\mathbf{r}_i^\uparrow, \mathbf{r}_j^\downarrow)$ to all combinations of the spin-up and spin-down particles. The triplet blocks $\chi^{\uparrow\uparrow}$ and $\chi^{\downarrow\downarrow}$ have sizes $N_\uparrow \times N_\uparrow$ and $N_\downarrow \times N_\downarrow$, respectively. They are both skew-symmetric and constructed by applying the triplet pairing orbitals $\chi^{\uparrow\uparrow}(\mathbf{r}_i^\uparrow, \mathbf{r}_j^\uparrow)$ and $\chi^{\downarrow\downarrow}(\mathbf{r}_i^\downarrow, \mathbf{r}_j^\downarrow)$ to the appropriate same-spin pairs. Finally, the unpaired blocks φ^\uparrow and φ^\downarrow can be viewed as a reorganization of the rows in φ such that the spin-up particles and spin-down particles are placed into separate blocks.

The STU wave function is a convenient approach if the spins of the particles can be fixed during the Monte Carlo simulation. Otherwise, the more general Pfaffian wave function can be used to handle the exchange of spins. Either way, the pairing spin-orbital $\phi(\mathbf{x}_i, \mathbf{x}_j)$ is decomposed according to Eq. (3.82) so that we can leverage our intuition in its design. In our investigation of neural-network quantum states, we will find that such a decomposition is not necessary. Instead, we can construct the most general antisymmetric pairing orbital that allows the spatial and spin degrees of freedom to influence one another arbitrarily.

3.2.2 Optimization

After selecting a specific parameterization of the trial wave function, we require methods by which the variational parameters can be changed in order to minimize the energy. These

optimization methods operate iteratively, gradually refining the solution over time. The first-order optimization methods we will discuss, known as gradient descent methods, coincide with the very same techniques commonly employed in numerous machine learning problems. However, we will also introduce one additional second-order optimization method, known as Stochastic Reconfiguration, which is specifically tailored for the variational Monte Carlo method. It is akin to the Natural Gradient method in the context of machine learning, but as our application of neural networks will always involve minimizing the energy, as opposed to a general objective function, we opt to explore all optimization methods within the context of variational Monte Carlo rather than machine learning.

3.2.2.1 Stochastic Gradient Descent

Stochastic Gradient Descent (SGD) is a simple method that utilizes a single estimate of the gradient $\nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle$ to update the parameters $\boldsymbol{\theta}$. At any given iteration of the optimization procedure, the parameters are transformed following the opposite direction of the gradient

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle, \quad (3.84)$$

where η is a small learning rate and the gradient is defined in Eq. (3.51). The learning rate can be scheduled to decrease with the number of optimization steps, with typical values between 10^{-5} and 10^{-2} .

Since our measures of the gradient are approximate, our trajectory through the energy landscape can be noisy. On one hand, this noise can sometimes prevent us from getting trapped in a local minimum, as there is always a chance of a perturbation that propels us out of it. On the other hand, the presence of this noise can make the overall training

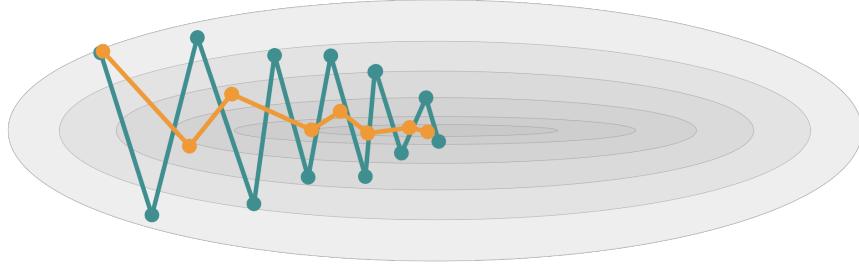


Figure 3.4: Cartoon of the simple stochastic gradient descent algorithm (blue) and the effect of including momentum (orange) near a narrow minimum.

procedure slow and inaccurate, as we may find ourselves frequently oscillating back and forth in a particularly narrow or shallow valley.

3.2.2.2 Momentum

The simple SGD algorithm can be improved by including momentum, a moving average of the gradients. Momentum is implemented by first defining a vector \mathbf{m} that stores the moving average of gradients, and initializing it to zero. Then parameter updates are performed as

$$\mathbf{m} \leftarrow \beta \mathbf{m} + (1 - \beta) \nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle, \quad (3.85)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{m}, \quad (3.86)$$

where η is the learning rate and β is a constant hyperparameter that controls the influence of the previous gradients on the current update, usually set around $\beta = 0.9$. By averaging over a history of gradients, the optimization trajectory becomes smoother as oscillations in opposite directions tend to cancel out. This smoothing effect helps to stabilize the parameter updates and prevent the algorithm from getting trapped in erratic behavior. Moreover, in situations where the gradients become small, such as in plateaus or saddle points, momentum helps the training process to persist in the "right" direction.

3.2.2.3 Root Mean Squared Propagation

Root Mean Squared Propagation (RMSprop) adaptively adjusts the learning rate for each parameter based on the magnitudes of the recent gradients. We begin by initializing a new vector $\mathbf{v} = \mathbf{0}$, and updating the parameters as

$$\mathbf{v} \leftarrow \beta \mathbf{v} + (1 - \beta) (\nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle)^2, \quad (3.87)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\eta}{\sqrt{\mathbf{v} + \epsilon}} \nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle, \quad (3.88)$$

where the square in Eq. (3.87) and the square root in Eq. (3.88) are both element-wise operations. The hyperparameter β is typically set to 0.9. The hyperparameter ϵ provides numerical stability and is set around 10^{-8} . Unlike the previous methods, RMSprop is less sensitive to choices of the learning rate η . By normalizing the gradients in this way, RMSprop mitigates the problem of exploding gradients, as extreme fluctuations are alleviated, making the optimization process more stable. At the opposite extreme, RMSprop also prevents gradients from becoming too small, allowing training to continue effectively.

3.2.2.4 Adaptive Moment Estimation

Adaptive Moment Estimation (ADAM) combines the concepts of RMSprop and momentum, and it remains one of the most popular first-order optimization methods since its development in 2017Ref?. In addition to computing the moving averages of gradients and the squares of gradients, the ADAM algorithm incorporates a learning rate correction to compensate for the bias caused by initializing these averages to zero. This correction is dependent on the optimization iteration t , so we begin by setting $t = 0$, $\mathbf{m} = \mathbf{0}$, and $\mathbf{v} = \mathbf{0}$. Then the

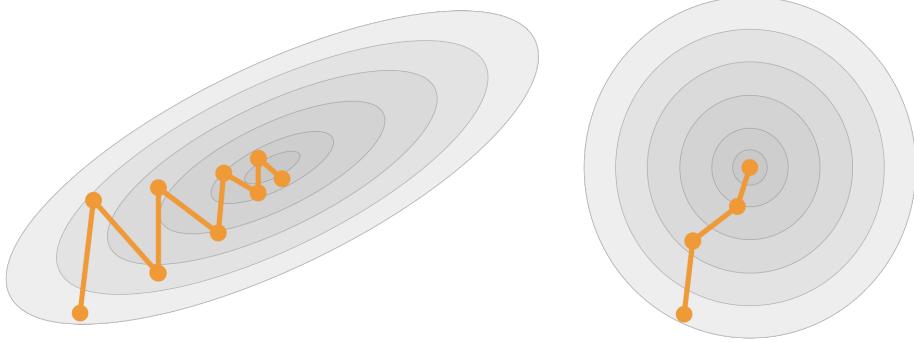


Figure 3.5: Visualization of how the Stochastic Reconfiguration smoothens the energy landscape.

parameters are updated as

$$\mathbf{m} \leftarrow \alpha \mathbf{m} + (1 - \alpha) \nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle, \quad (3.89)$$

$$\mathbf{v} \leftarrow \beta \mathbf{v} + (1 - \beta) (\nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle)^2, \quad (3.90)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \left(\frac{1 - \beta^t}{1 - \alpha^t} \right) \frac{\mathbf{m}}{\sqrt{\mathbf{v}} + \epsilon}, \quad (3.91)$$

$$t \leftarrow t + 1. \quad (3.92)$$

The hyperparameters are usually set to $\alpha = 0.9$, $\beta = 0.999$, and $\epsilon = 10^{-8}$. ADAM has gained popularity across diverse optimization problems, thanks to its efficient and robust performance.

3.2.2.5 Stochastic Reconfiguration

Stochastic Reconfiguration [6] (SR) is a second-order optimization method specific to the variational Monte Carlo (VMC) method. Instead of manipulating the gradients according to their history, the SR algorithm manipulates the gradients according to the *curvature* of the energy landscape. It can alternatively be viewed as stretching and squeezing the landscape itself, making it smoother or more isotropic in certain areas. SR provides a more favorable

terrain for finding the global minimum and improves the exploration of the parameter space.

The SR update is derived by determining the parameter change, denoted as $\delta\boldsymbol{\theta}$, that best reproduces a small step $\delta\tau$ in an imaginary-time evolution. The concept of imaginary-time propagation, which will be further elaborated on in Sec. 9.5, involves applying the operator $e^{-\delta\tau\hat{H}}$ to states non-orthogonal to the ground state, thereby driving them closer to the ground state. We can approximate the imaginary-time propagator to first order in $\delta\tau$ and apply it to our original state $|\Psi_{\boldsymbol{\theta}}\rangle$, resulting in a new state

$$|\Phi\rangle = e^{-\delta\tau\hat{H}}|\Psi_{\boldsymbol{\theta}}\rangle \approx (1 - \delta\tau\hat{H})|\Psi_{\boldsymbol{\theta}}\rangle = |\Psi_{\boldsymbol{\theta}}\rangle - \delta\tau\hat{H}|\Psi_{\boldsymbol{\theta}}\rangle. \quad (3.93)$$

Our goal is to match this state to one obtained by changing the parameters $\boldsymbol{\theta}$. To do this, we expand around our original state to first order in $\delta\boldsymbol{\theta}$,

$$|\Psi_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}\rangle \approx |\Psi_{\boldsymbol{\theta}}\rangle + (\delta\boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle. \quad (3.94)$$

For $|\Phi\rangle$ and $|\Psi_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}\rangle$ to coincide, we need their projections onto the original state be the same, keeping in mind it may not be normalized,

$$\frac{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}+\delta\boldsymbol{\theta}}\rangle}{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle} = \frac{\langle\Psi_{\boldsymbol{\theta}}|\Phi\rangle}{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle}, \quad (3.95)$$

$$\frac{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle}{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle} + (\delta\boldsymbol{\theta}) \cdot \frac{\langle\Psi_{\boldsymbol{\theta}}|\nabla_{\boldsymbol{\theta}}\Psi_{\boldsymbol{\theta}}\rangle}{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle} = \frac{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle}{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle} - \delta\tau \frac{\langle\Psi_{\boldsymbol{\theta}}|\hat{H}|\Psi_{\boldsymbol{\theta}}\rangle}{\langle\Psi_{\boldsymbol{\theta}}|\Psi_{\boldsymbol{\theta}}\rangle}, \quad (3.96)$$

Similarly, the gradient with respect to the parameters must also coincide,

$$\frac{\langle \Psi_{\theta} | \nabla_{\theta} | \Psi_{\theta+\delta\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle} = \frac{\langle \Psi_{\theta} | \nabla_{\theta} | \Phi \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle}, \quad (3.97)$$

$$\frac{\langle \Psi_{\theta} | \nabla_{\theta} \Psi_{\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle} + (\delta\theta) \cdot \frac{\langle \nabla_{\theta} \Psi_{\theta} | \nabla_{\theta} \Psi_{\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle} = \frac{\langle \Psi_{\theta} | \nabla_{\theta} \Psi_{\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle} - \delta\tau \frac{\langle \nabla_{\theta} \Psi_{\theta} | \hat{H} | \Psi_{\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle}, \quad (3.98)$$

$$(3.99)$$

Combining Eqs. (3.96) and (3.99) reveals the following update rule for the parameters:

$$\theta \leftarrow \eta S^{-1} \nabla_{\theta} \langle E_{\theta}(\mathbf{X}) \rangle, \quad (3.100)$$

where S is the quantum geometric tensor [7] with the following matrix elements

$$S_{ij} = \frac{\langle \partial_i \Psi_{\theta} | \partial_j \Psi_{\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle} - \frac{\langle \partial_i \Psi_{\theta} | \Psi_{\theta} \rangle \langle \Psi_{\theta} | \partial_j \Psi_{\theta} \rangle}{\langle \Psi_{\theta} | \Psi_{\theta} \rangle} \quad (3.101)$$

and ∂_i denotes a derivative with respect to the i th variational parameter θ_i .

3.3 Diffusion Monte Carlo

Diffusion Monte Carlo (DMC) is a virtually exact stochastic projector method based on solving the imaginary-time many-body Schrödinger equation. By transforming $t \mapsto -i\tau\hbar$, the Schrödinger equation becomes

$$-\frac{\partial}{\partial \tau} \Psi(\mathbf{X}, \tau) = \hat{H} \Psi(\mathbf{X}, \tau), \quad (3.102)$$

and the originally oscillating stationary-state solutions are now convergent in the $\tau \rightarrow \infty$ limit

$$\Psi_\alpha(\mathbf{X}, \tau) = \psi_\alpha(\mathbf{X}) e^{-E_\alpha \tau}, \quad (3.103)$$

where

$$\hat{H}\psi_\alpha(\mathbf{X}) = E_\alpha\psi_\alpha(\mathbf{X}). \quad (3.104)$$

Let us assume these energy eigenstates are ordered such that $E_0 \leq E_1 \leq E_2 \leq \dots$. By expanding the general wave function in terms of the transformed stationary states,

$$\Psi(\mathbf{X}, \tau) = \sum_{\alpha=0}^{\infty} c_\alpha \Psi_\alpha(\mathbf{X}, \tau) = \sum_{\alpha=0}^{\infty} c_\alpha \psi_\alpha(\mathbf{X}) e^{-E_\alpha \tau}, \quad (3.105)$$

we find that any state that is not orthogonal to the ground state, $c_0 \neq 0$, will evolve to the ground state

$$\lim_{\tau \rightarrow \infty} \Psi(\mathbf{X}, \tau) = c_0 \psi_0(\mathbf{X}) e^{-E_0 \tau}, \quad (3.106)$$

in the large τ limit. By introducing a constant offset energy E_T to Eq. 3.102, where E_T is as close to the ground state energy E_0 as possible, the large τ limit can be kept finite, while leaving the eigenstates themselves unchanged. Then by ignoring spin momentarily, $\mathbf{X} \rightarrow \mathbf{R}$, the shifted imaginary-time Schrödinger equation then reads

$$\frac{\partial}{\partial \tau} \Psi(\mathbf{R}, \tau) = (D \nabla \cdot \nabla - (V(\mathbf{R}) - E_T)) \Psi(\mathbf{R}, \tau), \quad (3.107)$$

where we have again used the notation $D = \frac{\hbar^2}{2m}$ and $\nabla \equiv (\nabla_1, \nabla_2, \dots, \nabla_N)$. At this point, the above equation should look familiar, as it has a similar form to the Fokker-Planck equation we introduced in Eq. 3.27. In fact, if we make the substitution $\Psi(\mathbf{R}, \tau) =$

$f(\mathbf{R}, \tau)/\Psi_T(\mathbf{R})$, take all the necessary derivatives, and multiply by $\Psi_T(\mathbf{R})$, the equation we obtain for $f(\mathbf{R}, \tau)$, has nearly identical form to the Fokker-Planck equation

$$\frac{\partial}{\partial \tau} f(\mathbf{R}, \tau) = D \nabla \cdot (\nabla - \mathbf{F}(\mathbf{R})) f(\mathbf{R}, \tau) - (E(\mathbf{R}) - E_T) f(\mathbf{R}, \tau), \quad (3.108)$$

except for the last term, where

$$E(\mathbf{R}) = \frac{1}{\Psi_T(\mathbf{R})} \hat{H} \Psi_T(\mathbf{R}). \quad (3.109)$$

This modified form of the Schrödinger equation can be solved for $f(\mathbf{R}, \tau)$ through standard Green's function methods, but as we already mentioned before at the end of Sec. 3.1.1, the Green's function solution to the Fokker-Planck equation coincides with the proposal density $Q(\mathbf{R}'|\mathbf{R}, \Delta t)$ we inferred from the corresponding Langevin equation. Without the extra term in Eq. 3.108, the Green's function solution would be exactly the same as before, as introduced in Eq. 3.34,

$$G_D(\mathbf{R}'|\mathbf{R}, \Delta \tau) = \frac{1}{(4\pi D \Delta \tau)^{Nd/2}} e^{-(\mathbf{R}' - \mathbf{R} - D \Delta \tau \mathbf{F}(\mathbf{R}))^2 / 4D \Delta \tau}, \quad (3.110)$$

for small $\Delta \tau$. The subscript D stands for drift. However, the inclusion of the last term modifies the Green's function, which can be approximated by

$$G_B(\mathbf{R}'|\mathbf{R}, \Delta \tau) = e^{\Delta \tau (E(\mathbf{X}') + E(\mathbf{X}) - 2E_T)/2}, \quad (3.111)$$

where B standard for branching. Thus the total Green's function is approximated as

$$G(\mathbf{R}'|\mathbf{R}, \Delta\tau) \approx G_D(\mathbf{R}'|\mathbf{R}, \Delta\tau)G_B(\mathbf{R}'|\mathbf{R}, \Delta\tau). \quad (3.112)$$

The factorization of the Green's function suggests a way we can evolve to the ground state in practice. First, we initialize an ensemble of MCMC walkers according to the probability distribution of a good estimation of the wave function $|\Psi_T(\mathbf{R})|^2$. Next, by using the importance sampling method we covered in Sec. 3.1.1, we move the walkers as usual. Then, to account for the extra branching term in the Green's function, we make copies of the walkers with a probability given by Eq. 3.111 . Proceeding in this way is the spatial realization of the imaginary-time propagation operator $e^{-\tau(\hat{H}-E_T)}$.

Unlike the variational Monte Carlo method, where walkers are sampled from a positive semi-definite distribution given by $|\Psi_{\theta}(\mathbf{X})|^2$, a diffusion Monte Carlo calculation employs the walkers' distribution itself to represent the wave function. For systems of bosons, this poses no issue, but for fermions, it is difficult to maintain the antisymmetry of the wave function. Therefore, when we reintroduce spin into our formulation, we must add the extra step of killing the walkers if they attempt to cross a node. This solution is known as the fixed-node approximation. The nodes themselves are determined by a previous VMC calculation, and consequently, DMC calculations carry a residual dependence on the VMC calculation.

Chapter 4

Machine Learning

Machine learning encompasses a wide range of computational models and algorithms designed to learn from data rather than rely on explicit programming. It is a subset of artificial intelligence, a rapidly evolving field that aims to endow computers with human-like cognition, such as the ability to interpret complicated data, adapt to new situations, and make decisions autonomously. These ambitious goals have driven the development of machine learning models and have transformed them into the highly versatile tools we witness today.

In this chapter, we introduce fundamental concepts in machine learning, in the hopes of building a broad perspective that extends beyond specific physics applications. We will explore different approaches to learning, providing context for the artificial neural networks we will cover afterwards. In particular, we place emphasis on the mathematical aspects of artificial neural networks, aiming to clarify their relevance and utility as neural-network quantum states.

4.1 The Curse of Dimensionality

When dealing with high-dimensional data, be it in machine learning or any data other analysis endeavor, a recurring challenge arises. It is known as the "curse of dimensionality," wherein the performance of algorithms decreases as the number of dimensions increases. The basis of the curse of dimensionality is geometric: given a fixed number of data points, the

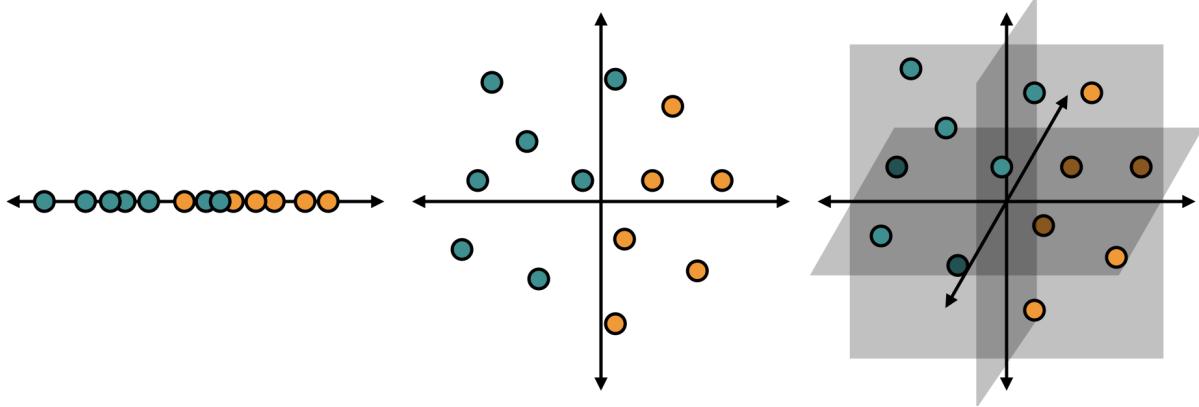


Figure 4.1: Cartoon of the curse of dimensionality; as the dimension increases, the density of a fixed number of data points decreases exponentially.

density of those points in space decreases exponentially with the number of dimensions. To counteract the sparsity of the data, is it common to focus considerable effort and resources on finding a lower-dimensional representation of the data before applying any given algorithm.

Is it reasonable to assume such a lower-dimensional representation exists? The short answer is yes—pure geometric scaling only truly applies if our data set contains no information at all, i.e. it is completely random. Real-life data sets, ones that actually contain *information*, have structure. This is reminiscent of the exponential scaling of the Hilbert space dimension with respect to the size of the quantum system. As we have discussed in Sec. 2.2, the effective Hilbert space dimension decreases if we restrict ourselves only to states that accurately represent real, physical systems.

Despite the challenges posed by the curse of dimensionality, machine learning algorithms have proven exceptional in discovering lower-dimensional representations of data. Among them, artificial neural networks have emerged as unparalleled tools for handling large and intricate datasets, where the underlying structure is difficult to predict a priori. Consequently, neural networks, particularly deep ones, have established themselves as the dominant approach for tackling some of the most formidable problems across various domains. Notably,

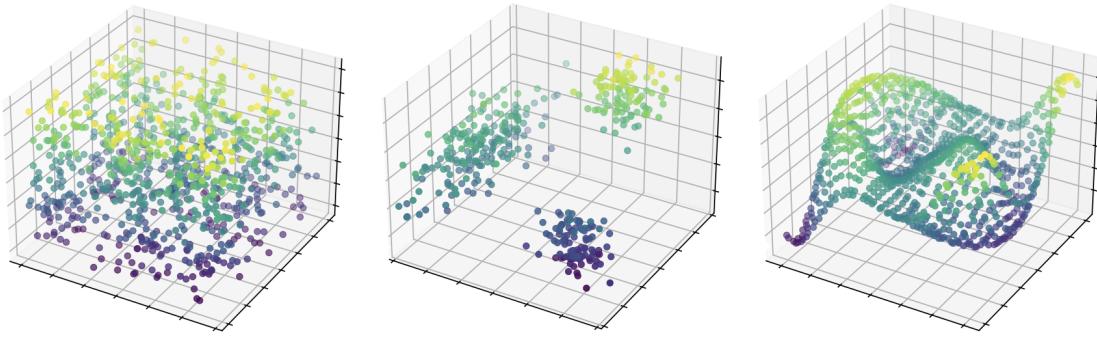


Figure 4.2: Comparison of a random data set (left) with no information, a data set with clear clusters (middle), and a data set with a clear lower-dimensional manifold (right). The latter two are examples of real-life data sets that contain information, i.e. structure.

they have been instrumental in solving complex tasks such as protein folding, game strategizing, natural language processing, image and voice recognition, and autonomous robotics.

4.2 Cost Functions

The goals of a machine learning problem are often narrow in scope, usually involving one or two tasks that can be expressed mathematically through a so-called "cost function". The terms "cost function" and "loss function" are often used interchangeably, but the cost function typically refers to the overall measure of error across the entire dataset, while the loss function calculates the error for individual data points. If the objective function needs to be maximized rather than minimized, it is commonly called a "reward function" in the context of machine learning.

While not all machine learning problems are explicitly framed as optimization problems, many involve minimizing or maximizing certain quantities. In some cases, this optimization can be achieved in a single step, leading to simple and highly efficient models, such as Kernel

Ridge Regression and Gaussian Processes. However, more generally, optimization problems require iterative methods, as discussed in the Sec. 3.2.2. For the purposes of our discussion, we will assume that all machine learning problems have an underlying cost function, even if it may not be obvious at first glance.

Determining the structure of a particular problem ultimately depends on the training data that is, or can be made, available. In the following sections, we will aim to give an intuitive understanding of machine learning and why exactly it has been proven to be so powerful across a wide variety of applications. When faced with a brand new problem, the reader should be able to zoom out, and quickly identify its core elements and interdependencies. Most importantly, the reader should be able to break a large, complicated problem into small, digestible components that are each suitable for a simple machine learning problem.

4.3 Supervised Learning

The goal of a supervised learning problem is to identify a mapping between two spaces, where only examples of the endpoints in each space are known. If we denote \mathbf{x} as the data points in our input space, \mathbf{y} as the data points in our output space, and $\hat{\mathbf{y}}$ as the predictions of our model in the output space, then the cost function of any supervised learning problem can be written as $C(\mathbf{y}(\mathbf{x}), \hat{\mathbf{y}}(\mathbf{x}))$, where C measures the discrepancy between $\hat{\mathbf{y}}$ and \mathbf{y} . The elements within the output space are often referred to as "labels." However, we will refrain from using this terminology to prevent any potential misconception that the output space must necessarily be discrete in nature.

Supervised learning problems can be separated into two major categories: regression and classification. The former involves predicting a continuous numerical value, while the latter

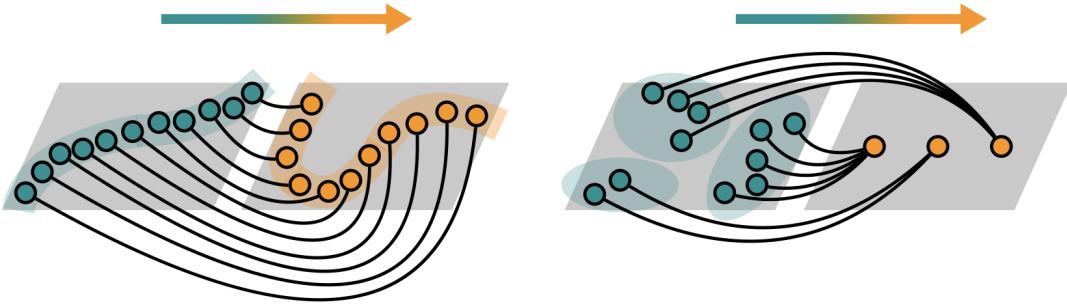


Figure 4.3: The two main types of supervised learning: regression (left) and classification (right).

involves predicting discrete values corresponding to distinct categories. Python packages like `scikit-learn` and `Keras` provide a wide range of pre-implemented supervised learning algorithms that are readily available for use "out-of-the-box," including decision trees, support vector machines, linear and logistic regression, and neural networks, to name a few.

Choosing an appropriate model depends strongly on the desired task, the complexity of the data, and the number of data points. If the number of data points is limited compared to the flexibility of a certain model, it is likely to overfit the data, leading to poor predictive power when new data is introduced. On the other hand, models that are too simple can fail at the task all together. Balancing the risk of overfitting and underfitting is a reflection of the bias-variance trade-off, the clash between a model's sensitivity to fluctuations in the input data and the model's underlying assumptions about the mapping.

To manage these challenges, supervised learning techniques commonly involve dividing a fixed dataset into three distinct parts: a training set, a validation set, and a test set. The training set is utilized to optimize the model, as its name implies. The validation set is used to estimate the bias and variance of the model throughout the training process, helping to determine the ideal hyperparameters. Finally, the model's performance is evaluated on the

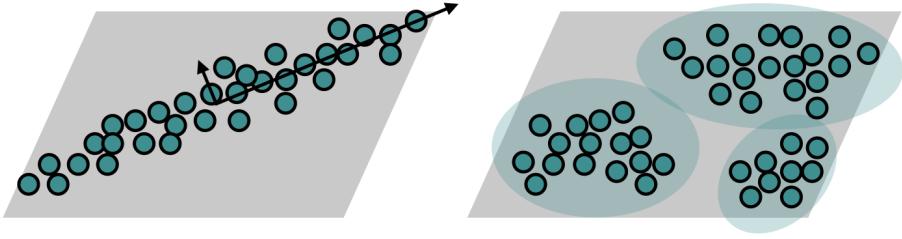


Figure 4.4: A depiction of principal component analysis (left) and clustering (right), two types of unsupervised learning.

test set once training is completed.

4.4 Unsupervised Learning

Unlike supervised learning, where the model learns a mapping from an input space to an output space, unsupervised learning models aim to learn the underlying structure of the input data set itself. Therefore, the associated cost functions solely depend on the inputs, $C(\mathbf{x})$. Unsupervised learning problems can be divided into three main categories: clustering, dimensionality reduction, and anomaly detection. Each of these categories can be used as a preprocessing step in conjunction with supervised or reinforcement learning problems to improve their overall performance.

In unsupervised clustering tasks, the model assumes the input data can be mapped to some latent space in which similar data points are physically closer together. Specifically, the assumed structure of the points in the n -dimensional latent space are (roughly) n -spheres, where n is at most the number of dimensions in the original input data. The key ingredient of such models is the quantification of "distance" between points. In simple problems, the input space itself can act as the latent space, as the data points may already exist in localized

groups.

Clustering and dimensionality reduction are not necessarily distinct tasks, as both aim to find some lower-dimensional representation of the data. However, the latter typically refers to problems in which the underlying structure of the dataset consists of hyperplanes rather than hyperspheres. Commonly used algorithms include Principal Component Analysis and Singular Value Decomposition, both of which seek to find the most relevant degrees of freedom in the data.

Dimensionality reduction problems may also include finding more abstract representations of the data. For instance, autoencoders are neural networks with an architecture that maps to a very small latent space before mapping back to the original input space. After training the autoencoder to reproduce each element of the data set, the latent space representation contains only the most essential and informative features of the data. This style of dimensionality reduction is particularly useful for file or image compression, as the relationships between data points can be highly non-intuitive.

Both clustering and dimensionality reduction go hand in hand with anomaly detection, because it focuses on identifying data points that deviate significantly from the expected behavior. Anomaly detection plays an important role in enhancing security in the digital age by identifying fraudulent behavior and detecting malware at a pace that surpasses human capabilities. Furthermore, anomaly detection finds practical utility in experimental sciences, where it aids in discerning genuine events from background noise. This capability enables efficient allocation of memory resources to statistically significant events, ensuring optimal utilization of computational resources.

4.5 Reinforcement Learning

Reinforcement learning is a specific branch of machine learning in which datasets rely on the model itself. If we let $\hat{\mathbf{y}}_t$ denote the state of the model at iteration t , and $\mathbf{x}_t(\hat{\mathbf{y}}_t)$ denote the data collected based on that state, then the cost function of a reinforcement learning problem may be written schematically as

$$C(\mathbf{x}_t(\hat{\mathbf{y}}_t(\mathbf{x}_{t-1}(\hat{\mathbf{y}}_{t-1}(\cdots(\mathbf{x}_0(\hat{\mathbf{y}}_0))\cdots))))). \quad (4.1)$$

In other words, starting from some initial model $\hat{\mathbf{y}}_0$, the data $\mathbf{x}_0(\hat{\mathbf{y}}_0)$ are generated by letting the model attempt the task at hand. Its performance is evaluated and used to update the model for the next iteration, continuing to collect new data based on its previous experiences. This notation for the cost function is simply meant to reflect the inherent cyclic nature of the problem; it is never explicitly written in this form in practice.

Reinforcement learning is commonly framed in terms of agent-based modeling, because the model learns by interacting with its environment and receiving positive or negative feedback, mirroring the way humans learn through trial and error in their everyday experiences (outside of the classroom). This learning approach is widely employed to train robots in game-playing and navigation tasks, which is impossible to accomplish through supervised learning methods as it would require an immense amount of example data.

Let us now rephrase variational Monte Carlo explicitly in the agent-based modeling framework. We begin by initializing the state of our agent, i.e. setting random parameters in our wave function. Then we allow our agent to interact with the environment, collecting data along the way. In VMC, the environment is the many-body Hilbert space and the data are the sampled configurations \mathbf{X} . Then we evaluate the performance of our agent, by calculat-

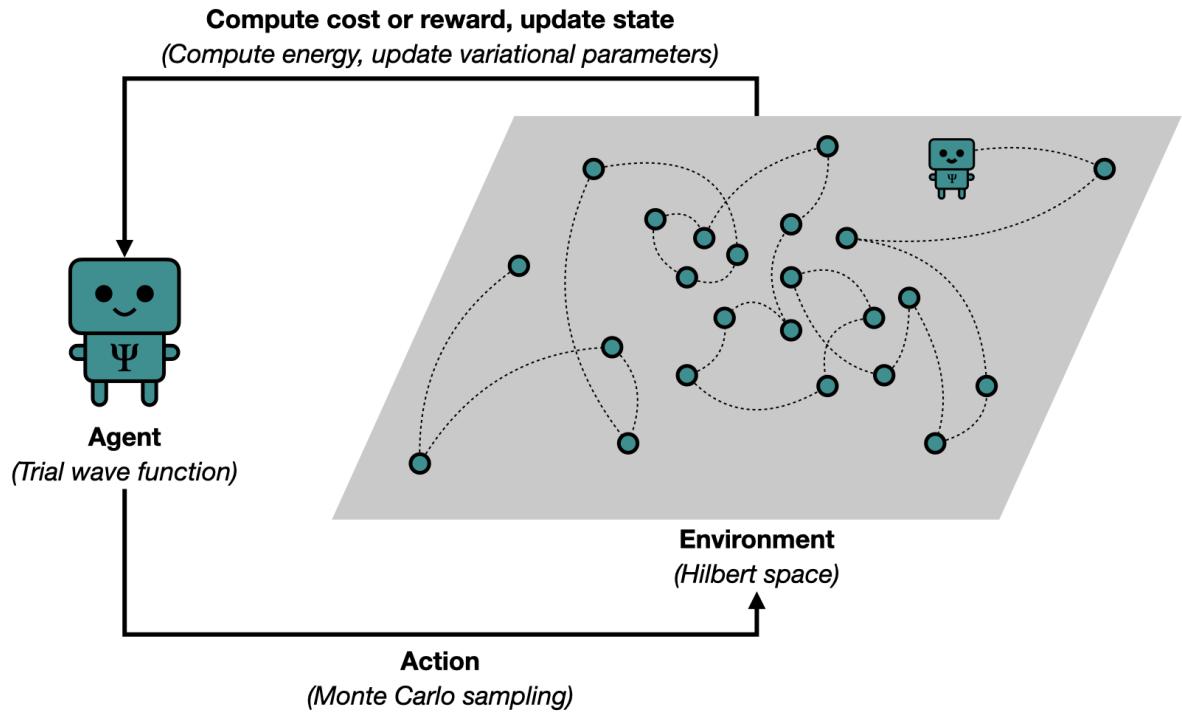


Figure 4.5: A depiction of the variational Monte Carlo algorithm in the context of agent-based modeling, a commonly employed reinforcement learning framework.

ing the average local energy based on the samples. This process continues until convergence is reached.

4.6 Transfer Learning

Transfer learning is a powerful technique that applies the knowledge gained from solving one problem to a more challenging and sometimes even unrelated problem. For example, researchers [8] have found that pre-training a convolutional neural network on ImageNet data provides immediate benefits for training the network on particle detection data. Interestingly, a fixed, pretrained network can be appended with additional trainable layers, while still achieving accelerated and stabilized training. This approach effectively reduces the number of trainable parameters while preserving the entire network’s representation power. Consequently, this finding suggests that all image recognition models rely on a shared set of core capabilities.

When using neural-network quantum states for variational Monte Carlo, we will use transfer learning to handle particularly challenging interaction potentials. Pre-training the neural-network quantum states on softer potentials allows the training on harder potentials to proceed in a smooth, controlled manner. These potentials may contain singularities, which we renormalize according to some newly defined hyperparameter, or they may simply be hard by nature. We will discuss transfer learning in these applications on a case-by-case basis.

4.7 Artificial Neural Networks

We are now in a suitable position to discuss our primary focus: the mathematics of artificial neural networks (ANNs). This class of machine learning models is inspired by the structure of the human brain, consisting of interconnected nodes and nonlinear connections between them. While the term is often used interchangeably with feedforward neural networks, it is important to note that the latter represents just one specific type within the broader ANN framework. We will exclusively refer to the broader class as artificial neural networks to ensure clarity.

In the following sections, we will formulate our artificial neural networks with the eventual goal of applying them as trial wave functions for variational Monte Carlo. The networks will be trained using the same optimization algorithms we use for traditional variational Monte Carlo calculations, discussed in Sec. 3.2.2. Consequently, our discussion may prioritize certain aspects that differ from those typically encountered in standard machine learning applications. Nonetheless, we will strive to maintain a level of generality throughout our discourse, ensuring that the concepts and techniques discussed herein can be applied to a wide range of problems.

4.7.1 Boltzmann Machines

Boltzmann machines are generative stochastic artificial neural networks that are primarily used for unsupervised learning tasks. They have two types of nodes: visible nodes, which represent the inputs to the network, and hidden nodes, which serve as latent variables capturing higher-level features or patterns in the data. In a general Boltzmann machine, all the nodes are fully connected, regardless of type. Additionally, the connections are undirected, unlike

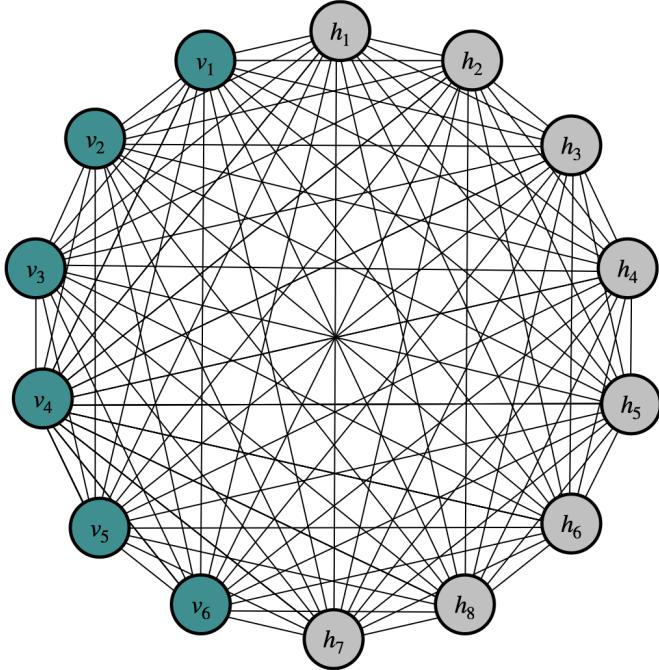


Figure 4.6: A general Boltzmann machine in which all visible and hidden nodes are fully-connected.

other well-known models such as feedforward neural networks. Unconstrained Boltzmann machines have limited practical utility for machine learning and inference. However, their optimization efficiency can be greatly enhanced by appropriately restraining the connectivity. This improvement allows for better applicability in solving real-world problems.

Restricted Boltzmann machines (RBMs) are a variant of Boltzmann machines that organize the visible and hidden nodes into two parallel layers. The layers are joined by undirected connections, but there are no connections between individual nodes in a given layer, hence they are "restricted". This architectural arrangement improves training efficiency by leveraging highly optimized linear algebra operations.

The objective of an RBM is to reveal the underlying probability distribution of the input data, enabling applications such as dimensionality reduction and the generation of samples resembling the data. It is precisely their ability to model probability distributions that

makes them a natural choice for implementing them as neural-network quantum states. Furthermore, RBMs can be stacked to construct deep belief networks, which excel in capturing hierarchical representations of complex datasets. These networks have exhibited remarkable success in addressing challenges in image and speech recognition domains, but training deep belief networks can present greater challenges compared to deep feedforward neural networks, which will be introduced in Sec. 4.7.2.

The most common type of RBM is a binary-binary RBM, in which all the visible and hidden nodes are discrete and constrained to the values of 0 or 1. Since we will exclusively work with real-valued inputs, our visible nodes will be taken to be Gaussian units, each with their own mean and variance. The hidden units will remain binary. The resulting formulation is called a Gaussian-binary restricted Boltzmann machine, which we will delve into shortly. In the subsequent sections, we will develop different variations of continuous restricted Boltzmann machines.

4.7.1.1 Gaussian-Binary Restricted Boltzmann Machines

We will begin by presenting the most common type of Gaussian-binary restricted Boltzmann machine (GB-RBM). The first step is to define an energy-like quantity $\mathcal{E}(\mathbf{v}, \mathbf{h})$ that describes the interplay between the visible nodes $\mathbf{v} \in \mathbb{R}^V$ and hidden nodes $\mathbf{h} \in \mathbb{R}^H$ within the network. This quantity should not be conflated with the physical energy of our quantum system, which is why it has been assigned a distinct name \mathcal{E} instead of E . The energy of the GB-RBM is defined as

$$\mathcal{E}(\mathbf{v}, \mathbf{h}) = \sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} - \sum_{j=1}^H b_j h_j - \sum_{i=1}^V \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j, \quad (4.2)$$

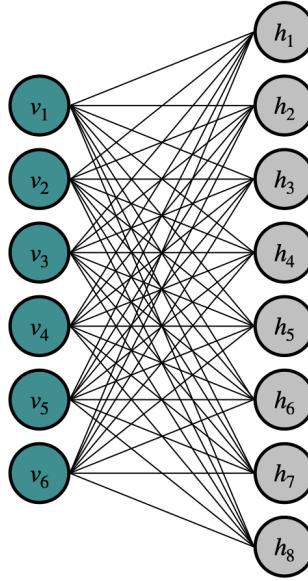


Figure 4.7: A standard restricted Boltzmann machine, where the hidden nodes are binary and the visible nodes can be either binary or Gaussian. There are no connections between nodes within the same layer.

where a_i is the bias or mean of the i th visible node, σ_i^2 is the variance of the i th visible node, b_j is the bias for the j th hidden node, and W_{ij} is the weight between the i th visible node and j th hidden node. In many machine learning applications, it is sufficient to simply treat the variances as a constant hyperparameter $\sigma_i^2 = \sigma^2$. However, training the variances individually will be important in our application, so we reparameterize them as

$$\frac{1}{\sigma_i^2} = \exp(s_i), \quad (4.3)$$

which improves their training and ensures they are always positive. The trainable parameters are collectively denoted as $\boldsymbol{\theta} = (\mathbf{a}, \mathbf{s}, \mathbf{b}, W)$, where $\mathbf{a} \in \mathbb{R}^V$, $\mathbf{s} \in \mathbb{R}^V$, $\mathbf{b} \in \mathbb{R}^H$, and $W \in \mathbb{R}^{V \times H}$.

By taking the Boltzmann distribution, we can define a corresponding joint probability

distribution

$$\mathcal{P}(\mathbf{v}, \mathbf{h}) = \frac{1}{\mathcal{Z}} e^{-\mathcal{E}(\mathbf{v}, \mathbf{h})}, \quad (4.4)$$

where the normalization constant is

$$\mathcal{Z} = \sum_{\mathbf{h}} \int d\mathbf{v} e^{-\mathcal{E}(\mathbf{v}, \mathbf{h})}. \quad (4.5)$$

This particular step justifies the inclusion of "Boltzmann" in the name of the model. We are specifically interested in the marginal probability distribution of the visible nodes, so we sum over the possible values of the hidden nodes

$$\begin{aligned} \mathcal{P}(\mathbf{v}) &= \sum_{\mathbf{h}} \mathcal{P}(\mathbf{v}, \mathbf{h}) = \frac{1}{\mathcal{Z}} \sum_{\mathbf{h}} e^{-\mathcal{E}(\mathbf{v}, \mathbf{h})} \\ &= \frac{1}{\mathcal{Z}} \sum_{\mathbf{h}} \exp \left(-\sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{j=1}^H b_j h_j + \sum_{i=1}^V \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j \right) \\ &= \frac{1}{\mathcal{Z}} \exp \left(-\sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} \right) \prod_{j=1}^H \sum_{h_j=0}^1 \exp \left(b_j h_j + \sum_{i=1}^V \frac{v_i}{\sigma_i^2} W_{ij} h_j \right) \quad (4.6) \\ &= \frac{1}{\mathcal{Z}} \exp \left(-\sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} \right) \prod_{j=1}^H \left(1 + \exp \left(b_j + \sum_{i=1}^V \frac{v_i}{\sigma_i^2} W_{ij} \right) \right) \\ &= \frac{1}{\mathcal{Z}} \exp \left(-\sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} \right) \prod_{j=1}^H \left(1 + \exp(z_j(\mathbf{v})) \right), \end{aligned}$$

where we have defined a new vector $\mathbf{z}(\mathbf{v}) \in \mathbb{R}^H$ with the elements

$$z_j(\mathbf{v}) = b_j + \sum_{i=1}^V \frac{v_i}{\sigma_i^2} W_{ij}, \quad (4.7)$$

for convenience. One of the notable strengths of a GB-RBM lies in its foundation on intuitive

physical laws. Unlike various other neural network models, the GB-RBM avoids being labeled as a "black-box," as it is fairly straightforward to understand how changes in the trainable parameters $\boldsymbol{\theta}$ causes changes in the learned distributions. For example, we can interpret the marginal probability distribution of the visible nodes in Eq. 4.6 as the product of a non-interacting Gaussian part and a nonlinear interaction term.

We can similarly derive the marginal probability distribution of the hidden nodes by integrating over the visible nodes

$$\begin{aligned}\mathcal{P}(\mathbf{h}) &= \int d\mathbf{v} \mathcal{P}(\mathbf{v}, \mathbf{h}) = \frac{1}{\mathcal{Z}} \int d\mathbf{v} e^{-\mathcal{E}(\mathbf{v}, \mathbf{h})} \\ &= \frac{1}{\mathcal{Z}} \int d\mathbf{v} \exp \left(-\sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{j=1}^H b_j h_j + \sum_{i=1}^V \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j \right) \\ &= \frac{1}{\mathcal{Z}} \exp \left(\sum_{j=1}^H b_j h_j \right) \prod_{i=1}^V \int dv_i \exp \left(-\frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j \right).\end{aligned}\quad (4.8)$$

Defining another vector $\mathbf{c}(\mathbf{h}) \in \mathbb{R}^V$ with the elements

$$c_i(\mathbf{h}) = a_i + \sum_j^H W_{ij} h_j, \quad (4.9)$$

and completing the square on the right-hand side of Eq. 4.8

$$\begin{aligned}-\frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j &= -\frac{1}{2\sigma_i^2} \left((v_i - a_i)^2 + \sum_{j=1}^H v_i W_{ij} h_j \right) \\ &= -\frac{1}{2\sigma_i^2} \left(v_i^2 - 2 \left(a_i + \sum_{j=1}^H W_{ij} h_j \right) v_i + a_i^2 \right) \\ &= -\frac{1}{2\sigma_i^2} \left(v_i^2 - 2c_i(\mathbf{h})v_i + a_i^2 \right) \\ &= -\frac{1}{2\sigma_i^2} \left((v_i - c_i(\mathbf{h}))^2 + a_i^2 - c_i(\mathbf{h})^2 \right)\end{aligned}\quad (4.10)$$

yields

$$\begin{aligned}\mathcal{P}(\mathbf{h}) &= \frac{1}{\mathcal{Z}} \exp \left(\sum_{j=1}^H b_j h_j \right) \prod_{i=1}^V \int dv_i \exp \left(-\frac{(v_i - c_i(\mathbf{h}))^2 + a_i^2 - c_i(\mathbf{h})^2}{2\sigma_i^2} \right) \\ &= \frac{1}{\mathcal{Z}} \exp \left(\sum_{j=1}^H b_j h_j \right) \prod_{i=1}^V \exp \left(\frac{c_i(\mathbf{h})^2 - a_i^2}{2\sigma_i^2} \right) \sqrt{2\pi} \sigma_i.\end{aligned}\quad (4.11)$$

The marginal distribution of the hidden nodes is yet again another product of non-interacting and interacting parts. From here, we can obtain the conditional probabilities of the visible nodes $\mathcal{P}(\mathbf{v}|\mathbf{h})$ and hidden nodes $\mathcal{P}(\mathbf{h}|\mathbf{v})$ by dividing the joint probability distribution $\mathcal{P}(\mathbf{v}, \mathbf{h})$ by $\mathcal{P}(\mathbf{h})$ and $\mathcal{P}(\mathbf{v})$, respectively. For the former, we find a product of normalized Gaussian distributions,

$$\begin{aligned}\mathcal{P}(\mathbf{v}|\mathbf{h}) &= \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{h})} \\ &= \frac{\exp \left(-\sum_{i=1}^V \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{i=1}^V \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j \right)}{\prod_{i=1}^V \exp \left(\frac{c_i(\mathbf{h})^2 - a_i^2}{2\sigma_i^2} \right) \sqrt{2\pi} \sigma_i} \\ &= \prod_{i=1}^V \frac{1}{\sqrt{2\pi} \sigma_i} \exp \left(-\frac{(v_i - c_i(\mathbf{h}))^2 + a_i^2 - c_i(\mathbf{h})^2}{2\sigma_i^2} - \frac{c_i(\mathbf{h})^2 - a_i^2}{2\sigma_i^2} \right) \\ &= \prod_{i=1}^V \frac{1}{\sqrt{2\pi} \sigma_i} \exp \left(-\frac{(v_i - c_i(\mathbf{h}))^2}{2\sigma_i^2} \right) \\ &= \prod_{i=1}^V \mathcal{N}(v_i; c_i(\mathbf{h}), \sigma_i^2),\end{aligned}\quad (4.12)$$

with their individual means $c_i(\mathbf{h})$ and variances of σ_i^2 . For the latter we have

$$\mathcal{P}(\mathbf{h}|\mathbf{v}) = \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{v})} = \frac{\exp \left(\sum_{j=1}^H b_j h_j + \sum_{i=1}^V \sum_{j=1}^H \frac{v_i}{\sigma_i^2} W_{ij} h_j \right)}{\prod_{j=1}^H \left(1 + e^{z_j(\mathbf{v})} \right)} = \prod_{j=1}^H \frac{e^{z_j(\mathbf{v}) h_j}}{1 + e^{z_j(\mathbf{v})}}, \quad (4.13)$$

which implies that the conditional probability of the j th hidden node being active is given by

$$\mathcal{P}(h_j = 1 | \mathbf{v}) = \frac{e^{z_j(\mathbf{v})}}{1 + e^{z_j(\mathbf{v})}} = \frac{1}{1 + e^{-z_j(\mathbf{v})}} = \sigma(z_j(\mathbf{v})), \quad (4.14)$$

where

$$\sigma(x) \equiv \frac{1}{1 + e^{-x}} \quad (4.15)$$

is the sigmoid function. In a sense, the binary nodes act as mediators for the interaction, switching between "on" and "off" states with a probability determined by the configuration of the visible nodes.

In standard unsupervised learning problems, the parameters $\boldsymbol{\theta}$ of a Restricted Boltzmann machine are trained by maximizing the log-likelihood $\log P(\mathbf{v})$. The gradients of the log-likelihood with respect to the parameters can be approximated using the Contrastive Divergence method and Gibbs sampling, increasing the efficiency of the training process significantly. However, in our application, the parameters will be trained by minimizing the expectation value of the physical energy of our quantum system, a major departure from traditional machine learning problems. Nonetheless, it will be useful to compute the gradients of the log-likelihood.

In terms of the parameters $\boldsymbol{\theta} = (\mathbf{a}, \mathbf{s}, \mathbf{b}, W)$, the log-likelihood reads

$$\log \mathcal{P}(\mathbf{v}) = -\frac{1}{2} \sum_{i=1}^V \exp(s_i)(v_i - a_i)^2 + \sum_{j=1}^H f(z_j(\mathbf{v})) - \log \mathcal{Z}, \quad (4.16)$$

where $z_j(\mathbf{v})$ was defined in Eq. (4.7) and is dependent on \mathbf{s}, \mathbf{b} , and W , while $f(x)$ is the "softplus" function

$$f(x) = \log(1 + \exp(x)). \quad (4.17)$$

It will become clear why we write the log-likelihood in terms of $f(x)$ when we develop generalizations of the GB-RBM. The derivatives of the log-likelihood with respect to the parameters are then given by

$$\frac{\partial}{\partial a_i} \log \mathcal{P}(\mathbf{v}) = \exp(s_i)(v_i - a_i), \quad (4.18)$$

$$\frac{\partial}{\partial s_i} \log \mathcal{P}(\mathbf{v}) = -\frac{1}{2} \exp(s_i)(v_i - a_i)^2 + \exp(s_i)v_i \sum_{j=1}^H W_{ij}f'(z_j(\mathbf{v})), \quad (4.19)$$

$$\frac{\partial}{\partial b_j} \log \mathcal{P}(\mathbf{v}) = f'(z_j(\mathbf{v})), \quad (4.20)$$

$$\frac{\partial}{\partial W_{ij}} \log \mathcal{P}(\mathbf{v}) = \exp(s_i)v_i f'(z_j(\mathbf{v})), \quad (4.21)$$

where $f'(x)$ is the same as the sigmoid function $\sigma(x)$ defined in Eq. 4.15.

4.7.1.2 Multivariate Gaussian-Binary Restricted Boltzmann Machines

We will now move beyond the standard formulation of the Gaussian-Binary restricted Boltzmann machine, and introduce correlations among the Gaussian visible nodes. The resulting network is a Multivariate Gaussian-Binary restricted Boltzmann machine (mGB-RBM), which is halfway between the standard GB-RBM and a general RBM. We will still label this network as "restricted" because there are no connections between the hidden nodes. The energy of the mGB-RBM can be written as

$$\mathcal{E}(\mathbf{v}, \mathbf{h}) = \frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a}) - \mathbf{b}^T \mathbf{h} - \mathbf{v}^T \Sigma^{-1} W \mathbf{h}, \quad (4.22)$$

where we have used the covariance matrix

$$\Sigma = \text{Cov}(\mathbf{v}, \mathbf{v}) = \begin{bmatrix} \sigma_1^2 & \sigma_1\sigma_2 & \cdots & \sigma_1\sigma_V \\ \sigma_2\sigma_1 & \sigma_2^2 & \cdots & \sigma_2\sigma_V \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_V\sigma_1 & \sigma_V\sigma_2 & \cdots & \sigma_V^2 \end{bmatrix} \in \mathbb{R}^{V \times V}, \quad (4.23)$$

in addition to the usual biases $\mathbf{a} \in \mathbb{R}^V$, $\mathbf{b} \in \mathbb{R}^H$ and weights $W \in \mathbb{R}^{V \times H}$. Following the previous reparameterizing of the variances in Eq. 4.3, we write the diagonal elements of the inverse covariance matrix as

$$\Sigma_{ii}^{-1} = \exp(S_{ii}), \quad (4.24)$$

and the offdiagonal elements as,

$$\Sigma_{ij}^{-1} = S_{ij}, \text{ for } i \neq j, \quad (4.25)$$

for some matrix $S \in \mathbb{R}^{V \times V}$. This choice is intended to avoid explicitly computing inverses of matrices and to guarantee the invertibility of the inverse covariance matrix Σ^{-1} . However, there is still a small risk that the variances corresponding to this reparameterization become negative, so we bias against this possibility by randomly initializing the diagonal elements S_{ii} to small positive values, and the off diagonal elements S_{ij} , $j \neq i$, to small negative values. This initialization results in a covariance matrix Σ close to a diagonal matrix with positive covariances. The covariance matrix (and its inverse) are symmetric, so only the upper triangular part of the matrix S , denoted as $\text{triu}(S)$, is required. Then the trainable parameters are collectively represented as $\boldsymbol{\theta} = (\mathbf{a}, \text{triu}(S), \mathbf{b}, W)$.

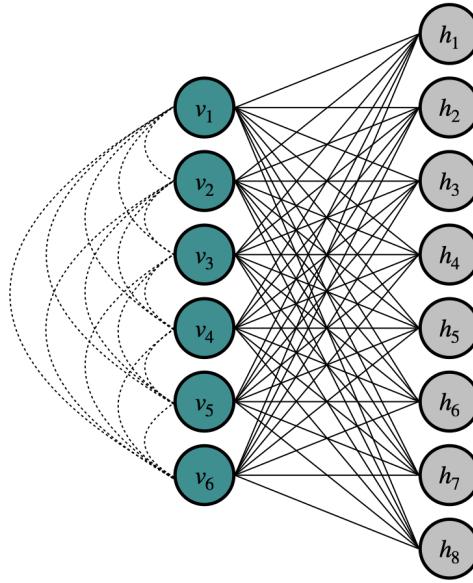


Figure 4.8: A multivariate Gaussian-binary restricted Boltzmann machine, a step between a standard Gaussian-binary restricted Boltzmann machine and a general Boltzmann machine.

The calculations of the joint, marginal, and conditional probability distributions will follow the same process as before. Having previously discussed the process extensively, we will simply list the final results. The joint probability distribution corresponding to the energy in Eq. 4.22 is written as

$$\mathcal{P}(\mathbf{v}, \mathbf{h}) = \exp \left(-\frac{1}{2} (\mathbf{v} - \mathbf{a})^T \Sigma^{-1} (\mathbf{v} - \mathbf{a}) + \mathbf{b}^T \mathbf{h} + \mathbf{v}^T \Sigma^{-1} W \mathbf{h} \right). \quad (4.26)$$

Summing over the hidden nodes yields the marginal probability distribution of the visible nodes

$$\mathcal{P}(\mathbf{v}) = \frac{1}{Z} \exp \left(-\frac{1}{2} (\mathbf{v} - \mathbf{a})^T \Sigma^{-1} (\mathbf{v} - \mathbf{a}) \right) \prod_{j=1}^H (1 + \exp(z_j(\mathbf{v}))), \quad (4.27)$$

where

$$\mathbf{z}(\mathbf{v}) = \mathbf{b} + W^T \Sigma^{-1} \mathbf{v} \in \mathbb{R}^H. \quad (4.28)$$

Unlike the standard GB-RBM, the distribution above does not contain a non-interacting part, as the Gaussian nodes are not independent (unless Σ is diagonal). Integrating the joint distribution over the visible nodes gives the marginal distribution over the hidden nodes

$$\mathcal{P}(\mathbf{h}) = \frac{1}{\mathcal{Z}} \exp \left(\mathbf{b}^T \mathbf{h} + \frac{1}{2} \mathbf{c}^T \Sigma \mathbf{c} - \frac{1}{2} \mathbf{a}^T \Sigma \mathbf{a} \right) \sqrt{(2\pi)^V |\Sigma|}, \quad (4.29)$$

where

$$\mathbf{c}(\mathbf{h}) = \mathbf{a} + W\mathbf{h}. \quad (4.30)$$

The conditional probability distribution of the visible nodes is no longer a product of independent Gaussians, but a V -dimensional multivariate Gaussian

$$\mathcal{P}(\mathbf{v}|\mathbf{h}) = \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{h})} = \mathcal{N}(\mathbf{v}; \mathbf{c}(\mathbf{h}), \Sigma). \quad (4.31)$$

Meanwhile, the conditional probability distribution of the hidden nodes is the same

$$\mathcal{P}(\mathbf{h}|\mathbf{v}) = \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{v})} = \prod_{j=1}^H \frac{e^{z_j(\mathbf{v}) h_j}}{1 + e^{z_j(\mathbf{v})}}, \quad \mathcal{P}(h_j = 1|\mathbf{v}) = \sigma(z_j(\mathbf{v})), \quad (4.32)$$

as long as we use our new definition of $\mathbf{z}(\mathbf{v})$ from Eq. 4.28. Finally, the log-likelihood is given by

$$\log \mathcal{P}(\mathbf{v}) = -\frac{1}{2} (\mathbf{v} - \mathbf{a})^T \Sigma^{-1} (\mathbf{v} - \mathbf{a}) + \sum_{j=1}^H f(z_j(\mathbf{v})) - \log \mathcal{Z}, \quad (4.33)$$

and its derivatives with respect to the parameters are

$$\frac{\partial}{\partial \mathbf{a}} \log \mathcal{P}(\mathbf{v}) = \Sigma^{-1}(\mathbf{v} - \mathbf{a}), \quad (4.34)$$

$$\frac{\partial}{\partial S} \log \mathcal{P}(\mathbf{v}) = \left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})(\mathbf{v} - \mathbf{a})^T + W \mathbf{f}'(\mathbf{z}(\mathbf{v})) \mathbf{v}^T \right) \circ \frac{\partial \Sigma^{-1}}{\partial S}, \quad (4.35)$$

$$\frac{\partial}{\partial \mathbf{b}} \log \mathcal{P}(\mathbf{v}) = \mathbf{f}'(\mathbf{z}(\mathbf{v})), \quad (4.36)$$

$$\frac{\partial}{\partial W} \log \mathcal{P}(\mathbf{v}) = \mathbf{v} \Sigma^{-1} (\mathbf{f}'(\mathbf{z}(\mathbf{v})))^T, \quad (4.37)$$

where $\mathbf{f}'(\mathbf{z}(\mathbf{v}))$ denotes the element-wise application of the derivative of the softplus function (Eq. (4.17)), also known as the sigmoid function, and \circ denotes element-wise multiplication. In addition, $\frac{\partial \Sigma^{-1}}{\partial S}$ is a matrix filled with ones except for the diagonal, which is filled with $\exp(S_{ii})$.

There are clearly many similarities between the standard GB-RBM and the multivariate GB-RBM. However, the inclusion of covariances in the latter enhances its flexibility compared to the former, and makes computing and writing down the various probabilities simpler. Additionally, the multivariate GB-RBM offers the advantage of performing most computations through matrix-matrix multiplication or matrix-vector multiplication, enabling more efficient implementations.

Before moving onto another step in generalization, let us consider one slight variation of the Multivariate Gaussian-Binary restricted Boltzmann machine, in which the binary nodes are allowed to take the values of -1 and 1, instead of 0 and 1. We can take the same energy function as Eq. 4.22, so that our joint probability distribution $\mathcal{P}(\mathbf{v}, \mathbf{h})$ is also the same. However, when we sum over the hidden nodes to compute the marginal distribution of the

visible nodes, we obtain

$$\begin{aligned}\mathcal{P}(\mathbf{v}) &= \frac{1}{\mathcal{Z}} \exp \left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a}) \right) \prod_{j=1}^H \left(\exp(z_j(\mathbf{v})) + \exp(-z_j(\mathbf{v})) \right) \\ &= \frac{1}{\mathcal{Z}} \exp \left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a}) \right) 2^H \prod_{j=1}^H \cosh(z_j(\mathbf{v})),\end{aligned}\tag{4.38}$$

instead of Eq. 4.27. Then the conditional probability distribution for the hidden nodes becomes

$$\begin{aligned}\mathcal{P}(\mathbf{h}|\mathbf{v}) &= \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{v})} = \frac{1}{2^H} \prod_{j=1}^H \frac{\exp(z_j(\mathbf{v})h_j)}{\cosh(z_j(\mathbf{v}))}, \\ \mathcal{P}(h_j = 1|\mathbf{v}) &= \frac{1}{2} \left(\tanh(z_j(\mathbf{v})) + 1 \right).\end{aligned}\tag{4.39}$$

In Figure 4.9, we compare the form of $P(h_j = 1|\mathbf{v})$ given above, with the one obtained in Eq. 4.32. When the hidden nodes were allowed to take the values of 0 or 1, we found that $P(h_j = 1|\mathbf{v})$ was given by the sigmoid function. When we instead allowed the hidden nodes to take the values of -1 or 1 , we obtained the hyperbolic tangent function, normalized to be between 0 and 1. The results are surprisingly intuitive.

To finish off our calculations, we write down the new log-likelihood

$$\log \mathcal{P}(\mathbf{v}) = -\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \exp(S)(\mathbf{v} - \mathbf{a}) + \sum_{j=1}^H f(z_j(\mathbf{v})) - \log \mathcal{Z}',\tag{4.40}$$

where \mathcal{Z}' contains the extra factor of 2^H from Eq. (4.38), and $f(x)$ from Eq. (4.17) is replaced by

$$f(x) = \log(\cosh(x)),\tag{4.41}$$

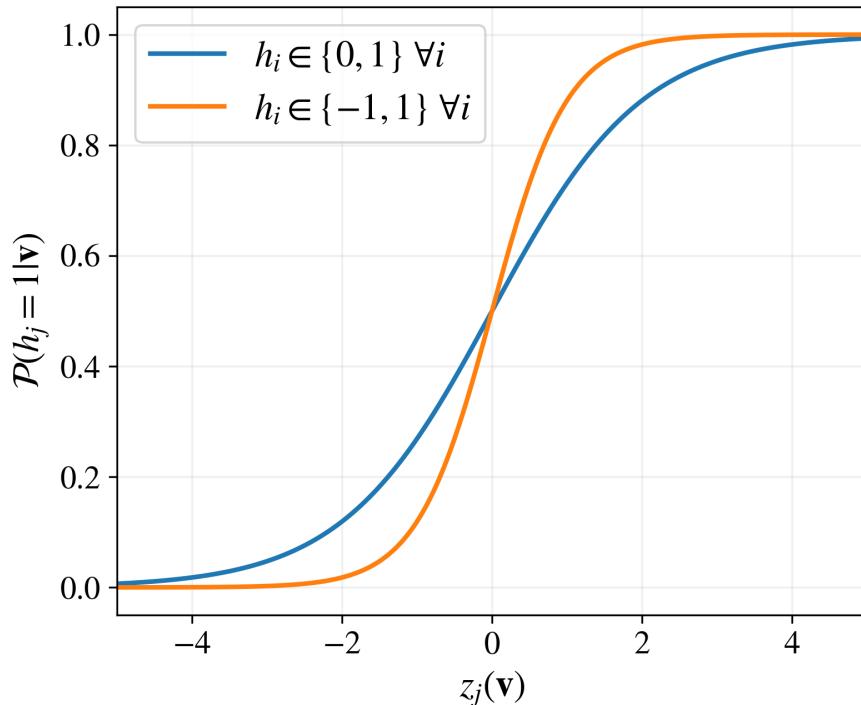


Figure 4.9: The conditional probability of a single binary hidden node h_j activating, given the state of the visible nodes \mathbf{v} , for a multivariate Gaussian-binary restricted Boltzmann machine. The blue line represents the probability if the hidden nodes are allowed to take values of 0 and 1 (Eq. 4.32), while the orange line represents the probability if the hidden nodes are allowed to take values of -1 and 1 (Eq. 4.39). The x-axis is a transformation of the visible nodes given by Eq. 4.28.

which we call the "log-cosh" function¹. With this redefinition of $f(x)$, the derivatives of the log-likelihood take the same form as before, in Eq. (4.37).

4.7.1.3 Multivariate Gaussian-Uniform Restricted Boltzmann Machines

Another variation of the continuous RBM involves allowing the hidden nodes to be continuous rather than discrete. We will name this variation the multivariate Gaussian-Uniform Restricted Boltzmann Machine (mGU-RBM), as we let the hidden nodes take any value between 0 and 1. Using the same form of the energy as Eq. (4.22), the marginal distribution of the visible nodes becomes

$$\begin{aligned}\mathcal{P}(\mathbf{v}) &= \int_0^1 d\mathbf{h} e^{-\mathcal{E}(\mathbf{v}, \mathbf{h})} \\ &= \frac{1}{\mathcal{Z}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a})\right) \prod_{j=1}^H \int_0^1 dh_j \exp(z_j(\mathbf{v})h_j) \\ &= \frac{1}{\mathcal{Z}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a})\right) \prod_{j=1}^H \frac{\exp(z_j(\mathbf{v})) - 1}{z_j(\mathbf{v})},\end{aligned}\quad (4.42)$$

while the marginal distribution of the hidden nodes remains the same as before. Then for the conditional probabilities for the hidden units, we obtain

$$\begin{aligned}\mathcal{P}(\mathbf{h}|\mathbf{v}) &= \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{v})} = \prod_{j=1}^H \frac{z_j(\mathbf{v}) \exp(z_j(\mathbf{v})h_j)}{\exp(z_j(\mathbf{v})) - 1}, \\ \mathcal{P}(h_j = 1|\mathbf{v}) &= \frac{z_j(\mathbf{v})}{1 - \exp(-z_j(\mathbf{v}))}.\end{aligned}\quad (4.43)$$

¹Just as the softplus function, Eq. (4.17), is a smoothly differentiable alternative to the Rectified Linear Unit (ReLU), a non-continuously differentiable activation function commonly used for feedforward neural networks, the log-cosh function, Eq. (4.41), is a smoothly differentiable alternative to the mean-absolute-error (MAE), a non-continuously differentiable cost function that tends to prevent overfitting.

The derivatives of the log-likelihood can remain in the same form as Eq. (4.37) by redefining $f(x)$ as

$$f(x) = \log \frac{\exp(x) - 1}{x}, \quad (4.44)$$

which we will name the "leaky-softplus" function. Even though the limit of this function exists mathematically as $|x| \rightarrow 0$, we may encounter instabilities near $x = 0$ when implemented computationally. Therefore, when $|x|$ is small, we replace $f(x)$ and $f'(x)$ with their Taylor expansions up to first order in x .

Now for the final variation of the mGU-RBM, we allow the hidden nodes to take any value between -1 and 1,

$$\begin{aligned} \mathcal{P}(\mathbf{v}) &= \int_0^1 d\mathbf{h} e^{-\mathcal{E}(\mathbf{v}, \mathbf{h})} \\ &= \frac{1}{\mathcal{Z}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a})\right) \prod_{j=1}^H \int_{-1}^1 dh_j \exp(z_j(\mathbf{v})h_j) \\ &= \frac{1}{\mathcal{Z}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a})\right) \prod_{j=1}^H \frac{\exp(z_j(\mathbf{v})) - \exp(-z_j(\mathbf{v}))}{z_j(\mathbf{v})} \\ &= \frac{1}{\mathcal{Z}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a})\right) 2^H \prod_{j=1}^H \frac{\sinh(z_j(\mathbf{v}))}{z_j(\mathbf{v})}. \end{aligned} \quad (4.45)$$

The conditional probabilities for the hidden units are given by

$$\begin{aligned} \mathcal{P}(\mathbf{h}|\mathbf{v}) &= \frac{\mathcal{P}(\mathbf{v}, \mathbf{h})}{\mathcal{P}(\mathbf{v})} = \frac{1}{2^H} \prod_{j=1}^H \frac{z_j(\mathbf{v}) \exp(z_j(\mathbf{v})h_j)}{\sinh(z_j(\mathbf{v}))}, \\ \mathcal{P}(h_j = 1|\mathbf{v}) &= \frac{1}{2} z_j(\mathbf{v}) \left(\coth(z_j(\mathbf{v})) + 1 \right), \end{aligned} \quad (4.46)$$

where $\mathcal{P}(h_j = 1|\mathbf{v})$ for the two variations of the mGU-RBM will be plotted alongside the softplus function in Fig. 4.10. Notice that both variations of the mGU-RBM involve a con-

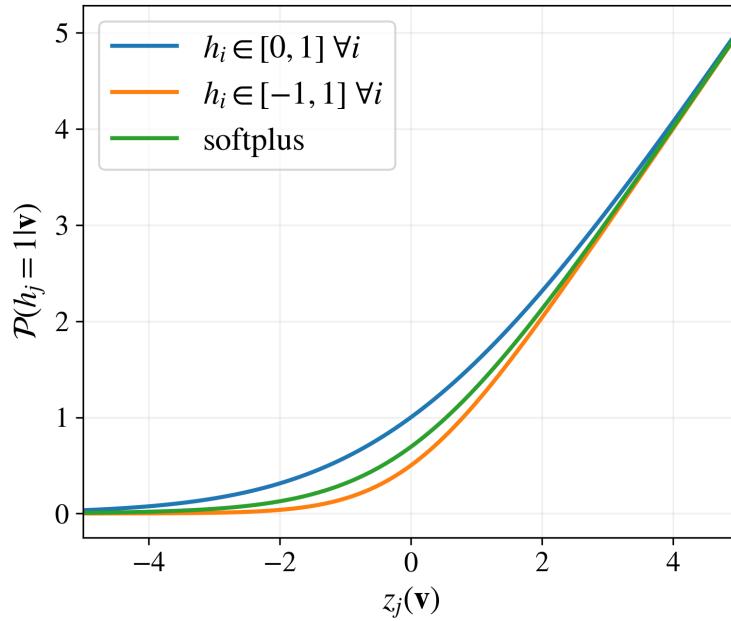


Figure 4.10: The conditional probability of a single uniform node h_j activating, given the state of the visible nodes \mathbf{v} , for a multivariate Gaussian-uniform restricted Boltzmann machine. The blue line represents the probability if the uniform nodes are allowed to take values between 0 and 1 (Eq. 4.43), while the orange line represents the probability if the hidden nodes are allowed to take values of -1 and 1 (Eq. 4.46). The x-axis is a transformation of the visible nodes given by Eq. 4.28. We also plot the softplus function (Eq. 4.17) for comparison, as its exhibits similar character. Probability values greater than 1 are assumed to mean the hidden node h_j is always activated.

ditional probability that is unbounded from above, as opposed to the mGB-RBMs, which produce a conditional probability that is bounded between 0 and 1. To interpret the conditional probabilities we obtained for the mGU-RBMs as true probabilities, let us assume all values are capped at 1.

Again, the derivatives of the log-likelihood can remain in the same form as Eq. (4.37) by redefining $f(x)$ as

$$f(x) = \log \frac{\sinh(x)}{x}, \quad (4.47)$$

which we will name the "log-sinhc" function, inspired by the "sine cardinal" function defined

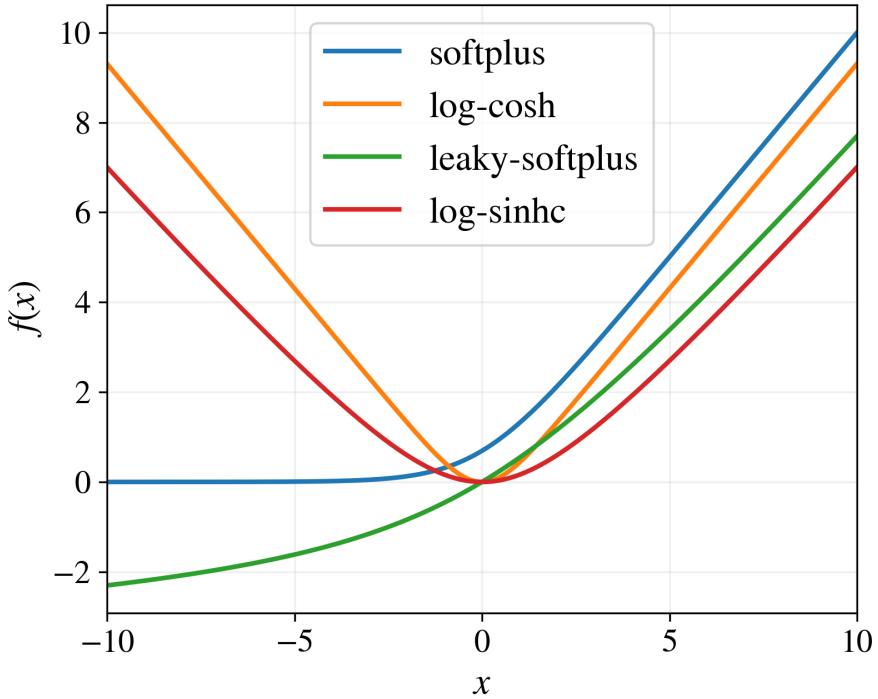


Figure 4.11: Comparison of the four functions we called $f(x)$ during our derivations of the multivariate Gaussian-binary and Gaussian-uniform restricted Boltzmann machines. These functions appear in the log-likelihoods $\log P(\mathbf{v})$.

as $\text{sinc}(x) \equiv \sin(x)/x$. As before, we handle the singularity at $x = 0$ by replacing $f(x)$ and $f'(x)$ by their first-order Taylor expansions.

To end our discussion on the various continuous restricted Boltzmann machines, let us compare the different definitions of $f(x)$ that appear in the log-likelihood. In Fig. 4.11, we plot the softplus, log-cosh, leaky-softplus, and log-sinhc functions, defined in Eqs. (4.17), (4.41), (4.44), and (4.47). The leaky-softplus function is the only function that permits negative values, hence the label "leaky". Both softplus and log-cosh scale as x for large, positive x , while leaky-softplus and log-sinhc scale as $x - \log(x)$. Finally, log-cosh and log-sinhc are both even functions, just like the interval for which the hidden nodes were defined during their construction.

4.7.2 Feedforward Neural Networks

Feedforward neural networks (FNNs) are widely recognized as one of the most prominent forms of artificial neural networks, primarily due to their training simplicity, scalability, and generality. As a result, they have been used to solve a diverse set of complex problems within all three branches of machine learning. However, in contrast to more interpretable models like restricted Boltzmann machines, feedforward neural networks are often characterized as black boxes, as their inner workings are difficult to interpret. FNNs are so successful, in fact, that it is recommended to use them only when low-bias models are required and a large data set is available. Otherwise, they can be prone to overfitting, especially in supervised learning problems.

A deep feedforward neural network is constructed by alternating compositions of affine transformations and simple, nonlinear transformations. It consists of an input layer, at least one hidden layer, and an explicit output layer, all of which are densely-connected to the adjacent layers only. As the name implies, the information flows only in one direction, starting from the input layer, and ending at the output layer. For an FNN with L layers (including the output layer, but excluding the input layer), we construct the nodes as

$$\mathbf{h}^{(0)} = \mathbf{v}, \quad (4.48)$$

$$\mathbf{h}^{(\ell)} = \mathbf{f}_\ell \left(W^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)} \right), \quad \text{for } \ell = 1, \dots, L, \quad (4.49)$$

where \mathbf{v} are the visible nodes for the input layer $\mathbf{h}^{(0)}$, $\mathbf{h}^{(\ell)}$ are the hidden layers for $\ell = 1, \dots, L - 1$, $\mathbf{h}^{(L)}$ is the output layer, and \mathbf{f}_ℓ represent the element-wise application of the nonlinear activation function f_ℓ . Each of the hidden layers have dimension H_ℓ , meaning $W^{(\ell)}$ is an $H_\ell \times H_{\ell-1}$ weight matrix and $\mathbf{b}^{(\ell)}$ is an H_ℓ -dimensional bias vector. The activation

functions f_ℓ introduce nonlinearity to the network, enabling it to learn and approximate complex functions. Without at least one nonlinear activation function, the neural network reduces to a simple linear model.

The output layer of an FNN is commonly distinguished from the hidden layers for two main reasons. Firstly, the output dimension is typically determined by the problem's structure, rather than being a hyperparameter that can be freely chosen. Secondly, the choice of activation function for the output layer is often different from that of the hidden layers. To maintain simplicity in notation, we adopt the indexing notation f_ℓ , where ℓ corresponds to the layer, to represent the activation functions rather than explicitly singling out the output layer. In regression problems, the final activation function f_L is typically linear, while binary-classification problems commonly employ sigmoid or hyperbolic tangent functions. When it comes to the hidden layers, popular choices for activation functions include the sigmoid function, hyperbolic tangent, and the Rectified Linear Unit (ReLU). However, our specific applications require unbounded, twice-continuously differentiable functions due to the evaluation of the local energy. Consequently, we prefer activation functions such as the Gaussian Error Linear Unit (GELU) and the Softplus function.

4.7.2.1 Backpropagation

One of the most attractive features of a deep feedforward neural network is its remarkable ease of training. Since we have written the forward propagation steps in Eq. 4.49 recursively, we can write the backpropagation steps using the chain rule. The derivatives of the output

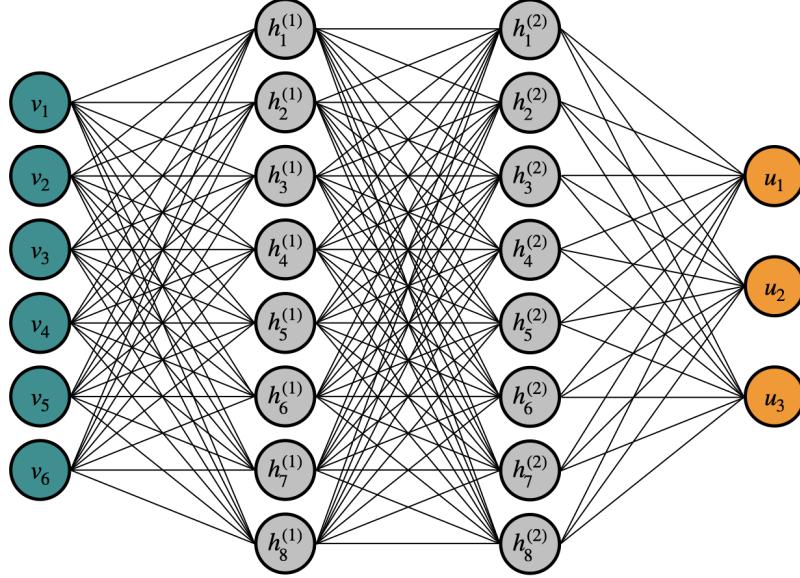


Figure 4.12: A deep feedforward neural network with six inputs (blue) and three outputs (orange). There are two hidden layers with eight nodes each (gray).

layer with respect to the weights and biases in layer ℓ are

$$\frac{\partial \mathbf{h}^{(L)}}{\partial W^{(\ell)}} = \frac{\partial \mathbf{h}^{(L)}}{\partial \mathbf{h}^{(L-1)}} \frac{\partial \mathbf{h}^{(L-1)}}{\partial \mathbf{h}^{(L-2)}} \cdots \frac{\partial \mathbf{h}^{(\ell+1)}}{\partial \mathbf{h}^{(\ell)}} \frac{\partial \mathbf{h}^{(\ell)}}{\partial W^{(\ell)}}, \quad (4.50)$$

$$\frac{\partial \mathbf{h}^{(L)}}{\partial \mathbf{b}^{(\ell)}} = \frac{\partial \mathbf{h}^{(L)}}{\partial \mathbf{h}^{(L-1)}} \frac{\partial \mathbf{h}^{(L-1)}}{\partial \mathbf{h}^{(L-2)}} \cdots \frac{\partial \mathbf{h}^{(\ell+1)}}{\partial \mathbf{h}^{(\ell)}} \frac{\partial \mathbf{h}^{(\ell)}}{\partial \mathbf{b}^{(\ell)}}. \quad (4.51)$$

To simplify notation, we define the intermediate vector

$$\mathbf{g}^{(\ell)} \equiv W^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)}, \quad (4.52)$$

which represents the linear transformation only. Then evaluating each term in Eqs. 4.50 and 4.51, we have

$$\frac{\partial h_i^{(\ell)}}{\partial h_j^{(\ell-1)}} = f'_\ell(g_i^{(\ell)}) W_{ij}^{(\ell)}, \quad (4.53)$$

$$\frac{\partial h_i^{(\ell)}}{W_{jk}^{(\ell)}} = f'_\ell(g_i^{(\ell)}) \delta_{ij} h_k^{(\ell-1)}, \quad (4.54)$$

$$\frac{\partial h_i^{(\ell)}}{b_j^{(\ell)}} = f'_\ell(g_i^{(\ell)}) \delta_{ij}, \quad (4.55)$$

where δ_{ij} is the Kronecker delta. Proceeding in this way, we can propagate the error from the cost function $\frac{\partial C}{\partial \mathbf{h}^{(L)}}$ all the way back to the input layer, updating all the weights and biases along the way.

4.7.2.2 Universal Approximation Theorem

Another contributing factor to the widespread adoption and success of feedforward neural networks is the Universal Approximation Theorem. In simple terms, the theorem states that a feedforward neural network with a single hidden layer, also known as a shallow neural network, can approximate any continuous function on a compact subset of Euclidean space, given there are a sufficient number of hidden nodes and the activation function satisfies certain conditions. This theorem was originally proven for the sigmoid activation function, and since then, there have been efforts to prove the theorem for a number of different activation functions, network depths, and even different types of networks entirely.

We will accept the Universal Approximation Theorem as a foundational principle without proof, although we should exercise caution in its application. While the theorem establishes the existence of a satisfactory approximation for any continuous function using a shallow neu-

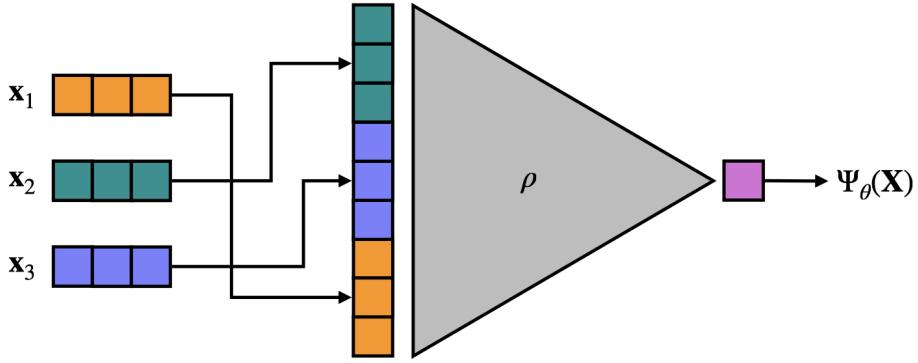


Figure 4.13: An example of a permutation-invariant Jastrow wave function constructed from the single output of a feedforward neural network. To enforce permutation-invariance, the input vectors can be sorted according to some rule before concatenating them into a single vector.

ral network, it does not guarantee that the required number of hidden nodes scales efficiently. Therefore, applying the theorem in practice is not as straightforward as simply increasing the size of our network. We must also consider other factors that can enhance training stability, capture additional correlations, and improve precision. While the Universal Approximation Theorem has not been formally proven for more complex network architectures, in practical applications, the effectiveness and versatility of these architectures are widely acknowledged nonetheless.

4.7.3 Deep Sets

One of the limitations of a standard feedforward neural network is the requirement to organize inputs into a single vector. This poses a challenge when dealing with inputs such as pixels in a 2D image, as flattening the image inherently compromises the spatial structure of the data². Similarly, standard feedforward neural networks encounter challenges when

²This challenge leads to the construction of convolutional neural networks (CNNs), a popular choice of neural network that leverages spatial correlations in grid-like data. As CNNs are specifically designed for grid-like data, we will not apply them as neural-network quantum states for continuous-space systems. Consequently, we will not cover them in this chapter, but we will discuss their generalizations, graph neural

attempting to effectively learn sets that exhibit permutation invariance. Without special treatment, the network would require at least $N!$ times more data points to learn a set of N elements compared to learning on an equivalently sized data set without permutation invariance, in order to effectively capture the inherent redundancies in the data. Even after learning, the network’s representation of permutation invariance is only approximate. This becomes problematic when constructing neural-network quantum states, as exact permutation invariance is often required in advance.

Approaches to address this issue can be categorized into two main strategies: data pre-processing and architecture design. The former approach involves sorting the set according to a specified rule, thereby transforming the set into a fixed-dimensional sequence. However, this preprocessing step still results in an input vector size that scales linearly with the number of elements in the set, making it challenging to handle large sets. In this section, we focus on the architecture-level approach, where the inherent set structure is ingrained within the neural network architecture. This approach is commonly known as the Deep Set architecture, originally introduced by Zaheer et al. in Ref. [9].

Suppose we wish to learn on sets of N elements $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$, where each element \mathbf{v}_i is a v -dimensional vector. Formally, a function $\mathbf{f} : (\mathbb{R}^v)^N \rightarrow \mathbb{R}^O$ is a permutation-invariant function of the set $\{\mathbf{v}_i\}$ if it satisfies

$$\mathbf{f}(\{\mathbf{v}_i\}) = \mathbf{f}(\hat{P}\{\mathbf{v}_i\}), \quad (4.56)$$

for any permutation $\hat{P} \in S_N$. First, we recognize that the function \mathbf{f} is permutation-

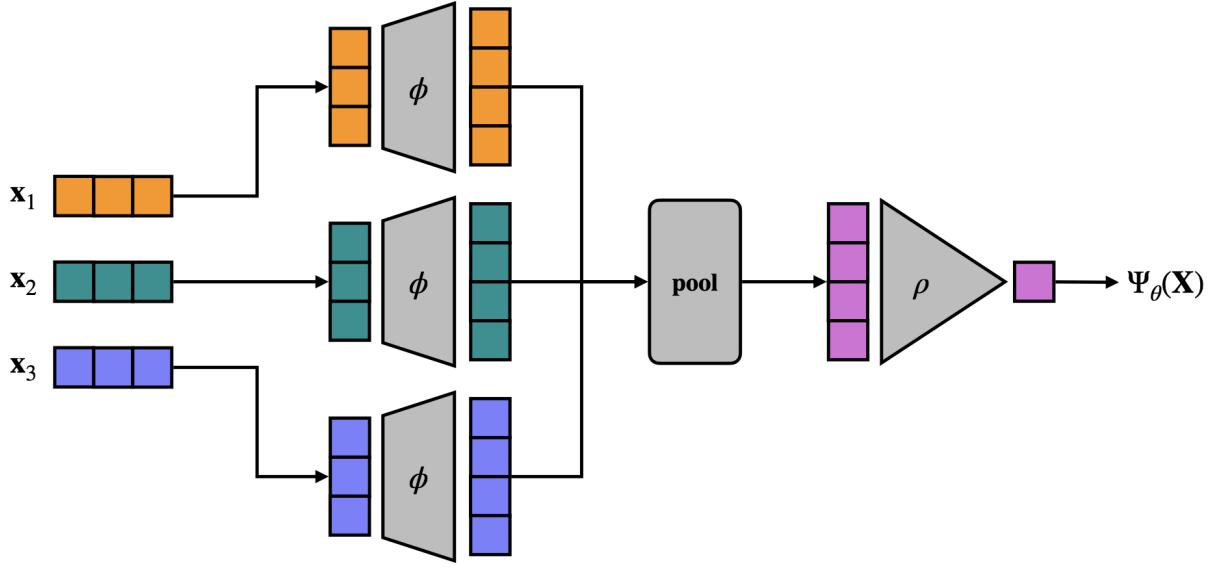


Figure 4.14: A depiction of a permutation-invariant Jastrow wave function constructed from the single output of a Deep Set. As Deep Sets enforce permutation invariance at the architecture level, no prior sorting of the inputs is required.

invariant if and only if it can be decomposed as

$$f(\{v_i\}) = \rho \left(\sum_{i=1}^N \phi(v_i) \right), \quad (4.57)$$

for some functions $\phi : \mathbb{R}^v \rightarrow \mathbb{R}^L$ and $\rho : \mathbb{R}^L \rightarrow \mathbb{R}^O$. The function ϕ maps each element of the set to an L -dimensional latent space. Then, the summation destroys the ordering of the set, resulting in an L -dimensional representation of the set. Then the function ρ maps the set representation into an O -dimensional output space.

A Deep Set is realized by replacing each of the functions ϕ and ρ with feedforward neural networks. Since feedforward neural networks are universal function approximators, the result of this substitution is a universal approximator for permutation-invariant functions. The summation operation in Eq. 4.57 is the key ingredient in this construction, as it aggregates

the latent space representations of the elements into a collective set representation. Since other operations can play the same role, we often write our Deep Sets with the more generic form

$$\mathbf{f}(\{\mathbf{v}_i\}) = \rho(\mathbf{pool}(\{\phi(\mathbf{v}_i)\})), \quad (4.58)$$

where **pool** can be any pooling operation that collects the set of vectors it acts upon, destroys their ordering, and returns a vector of equal size to a single element in the set.

The general idea above can be extended to construct generic permutation-*equivariant* functions as well. The function $\mathbf{f} : (\mathbb{R}^v)^N \rightarrow (\mathbb{R}^O)^N$ is formally permutation-equivariant if

$$\hat{P}\mathbf{f}(\{\mathbf{v}_i\}) = \mathbf{f}(\hat{P}\{\mathbf{v}_i\}) \quad (4.59)$$

holds for all permutations $\hat{P} \in S_N$, with $\mathbf{f}(\{\mathbf{v}_i\}) \equiv (\mathbf{f}_1(\{\mathbf{v}_i\}), \mathbf{f}_2(\{\mathbf{v}_i\}), \dots, \mathbf{f}_N(\{\mathbf{v}_i\}))$. One of the simplest ways to satisfy this condition is to define

$$\mathbf{f}_i(\{\mathbf{v}_j\}) = \rho(\mathbf{pool}(\{\phi(\mathbf{v}_j) | j \neq i\})), \quad (4.60)$$

where ρ and ϕ play the same role as before. Another option includes

$$\mathbf{f}_i(\{\mathbf{v}_j\}) = (\mathbf{v}_i, \mathbf{g}(\{\mathbf{v}_j\})), \quad (4.61)$$

where \mathbf{g} is a permutation-invariant function such as a standard Deep Set.

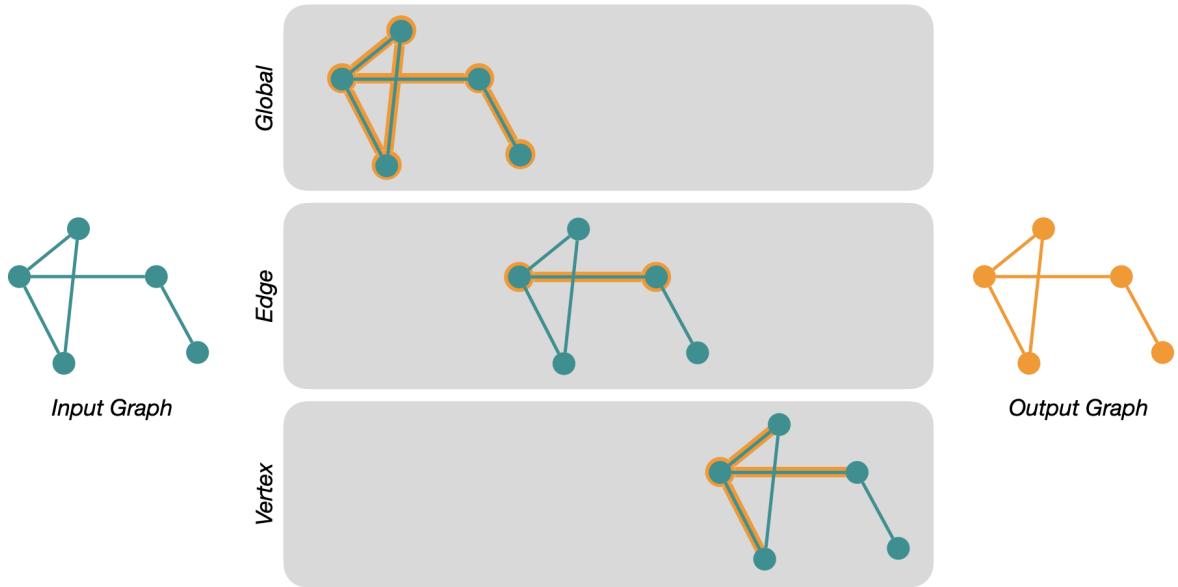


Figure 4.15: One layer of a graph neural network that produces an output graph that preserves the structure of the original input graph. Three types of transformations can be considered: global (top), edge (middle), and vertex (bottom) transformations.

4.7.4 Graph Neural Networks

Graph neural networks (GNNs) can be seen as a generalization of the popular convolutional neural networks, as they extend the concept of convolutions from grid-like data, such as images, to more general graph-structured data. A graph is a collection of nodes, or vertices, connected by edges. The idea of a graph neural network, specifically the graph-to-graph variant, is to take some input graph and iteratively embed correlations between the edges and nodes to produce a new graph with the same connectivity.

For demonstration purposes, we will construct a generic version of a graph-to-graph GNN acting on three different types of information: vertex features $\mathbf{v}_i^{(t)} \in \mathbb{R}^V$, edge features $\mathbf{e}_{ij}^{(t)} \in \mathbb{R}^E$, and global features $\mathbf{g}^{(t)} \in \mathbb{R}^G$, where t denotes the current iteration of the GNN. Then global features $\mathbf{g}^{(t)}$ can be constructed by taking in all three types of information from

the previous graph,

$$\mathbf{g}^{(t)} = \mathbf{G}_t \left(\mathbf{g}^{(t-1)}, \boldsymbol{\rho}_t(\{\mathbf{v}_i^{(t-1)}\}), \boldsymbol{\phi}_t(\{e_{ij}^{(t-1)}\}) \right), \quad (4.62)$$

where \mathbf{G}_t is a feedforward neural network with an output dimension of G , and $\boldsymbol{\rho}_t$ and $\boldsymbol{\phi}_t$ are each Deep Sets acting on the set of all previous vertices and edges, respectively. To preserve the original connectivity of the graph, the edge features can be updated as

$$e_{ij}^{(t)} = \mathbf{E}_t \left(e_{ij}^{(t-1)}, \mathbf{v}_i^{(t-1)}, \mathbf{v}_j^{(t-1)}, \mathbf{g}^{(t)} \right), \quad (4.63)$$

where \mathbf{E}_t is another feedforward neural network with output dimension E . Notice that this update includes the global feature $\mathbf{g}^{(t)}$ from the current step, rather than the previous step. In addition, the update includes contributions from the two vertices connected by the particular edge. Finally, we can update the vertex features as

$$\mathbf{v}_i^{(t)} = \mathbf{V}_t \left(\mathbf{v}_i^{(t-1)}, \mathbf{g}^{(t)}, \boldsymbol{\eta}_t(\{e_{ij}^{(t)} \mid j \in \mathcal{N}(i)\}) \right), \quad (4.64)$$

where \mathbf{V}_t is a feedforward neural network with output dimension V and $\boldsymbol{\eta}_t$ is a Deep Set acting on the set of the current edges connected to the particular node. In the notation above, $\mathcal{N}(i)$ denotes the set of all nodes in the *neighborhood* of node i .

By updating the global, edge, and vertex features in this way, we can preserve the original structure of the graph and embed highly non-linear correlations between all of the features in only a few iterations of this procedure. In practice, there are many ways of constructing a graph-to-graph GNN, but the overall concept remains the same. We leverage the permutation invariance of Deep Sets and the universal approximation properties of feedfor-

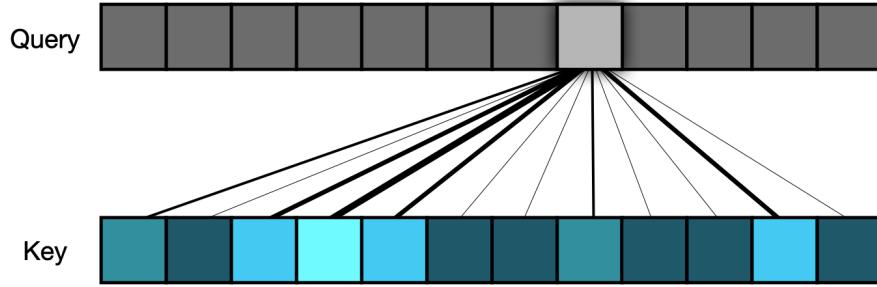


Figure 4.16: A visualization of the attention scores between a single element of a query vector and all the elements of a key vector. Larger attention scores are depicted with darker, bolder lines.

ward neural networks to exchange information between the vertex, edge, and global features. Consequently, neural-networks of this type are commonly called "message-passing" neural networks.

4.7.5 Attention Mechanisms

Attention mechanisms are components originally developed in the field of natural language processing to enable neural networks to selectively focus on specific segments of an input sentence. Their purpose is to enhance the network's ability to comprehend the correlations between individual words and derive the overall meaning. These mechanisms proved instrumental in addressing the challenge posed by word ambiguity, which arises when a word can have many different meanings depending on how it is used. By leveraging attention, neural networks are able to effectively capture contextual clues similar to how humans interpret language.

To build an attention mechanism, we begin by organizing our N input sequences $\mathbf{x}_i \in \mathbb{R}^d$

column-wise into a matrix $X \in \mathbb{R}^{d \times N}$. Then we define query W_Q , key W_K , and value W_V weight matrices all with dimensions $D \times d$. We apply each of these weight matrices to our data to obtain corresponding queries, keys, and values

$$Q = W_Q X, \quad K = W_K X, \quad V = W_V X, \quad (4.65)$$

where $Q, K, V \in \mathbb{R}^{D \times N}$. Then the attention scores are computed as

$$A = \text{softmax} \left(\frac{QK^T}{\sqrt{D}} \right) \in \mathbb{R}^{D \times D}, \quad (4.66)$$

with **softmax** denoting the row-wise application of the function

$$\text{softmax}(\mathbf{a})_i = \frac{e^{a_i}}{\sum_{a=1}^D e^{a_i}}, \quad (4.67)$$

where \mathbf{a} denotes a D -dimensional vector. Applying the **softmax** function computes the normalized overlap between the queries and keys, determining which values to pay attention to, and scaling the denominator by \sqrt{D} helps prevent vanishing gradients. Finally, the output of the attention layer is given by

$$Y = AV \in \mathbb{R}^{D \times N}, \quad (4.68)$$

The attention mechanism above is specifically known as self-attention, since the queries, keys, and values are all derived from the input data X . Other variations of attention include cross-attention, in which the queries are either constructed from an entirely different external data source X' or treated as a completely trainable matrix. In addition, multiple attention-

mechanisms can be implemented in parallel, leading to multi-headed attention.

Chapter 5

Implementation

Throughout the rest of this dissertation, we will reference two different implementations of neural-network quantum states for variational Monte Carlo, one in C++ and another in Python. As their objectives and approaches significantly differ, we will organize this chapter into two sections, each dedicated to its respective implementation.

5.1 NeuralAnsatz: C++ Software for Localized Systems

The first implementation, named `NeuralAnsatz`, is an object-oriented C++ code that leverages distributed-memory programming with `Open MPI` for efficient utilization of multiple CPUs. The code is built from scratch using `Eigen`, a lightweight, header-only linear algebra library, eliminating the need for separate linking during compilation. The primary features and goals for this implementation include: modularity, reproducibility, low memory costs, and exact gradient calculations. The source code can be found in this Github repository or at this url: <https://github.com/kim-jane/NeuralAnsatz>.

5.1.1 Trainer

`NeuralAnsatz` consists of several components, called `Objects`, that are passed to a single universal `Trainer`. The `Object` class is an abstract base class that simply contains relevant

```
System Wave Function Hamiltonian Sampler Optimizer Cost Function Density
```

Objects

Trainer

Figure 5.1: The highest-level structure of the `NeuralAnsatz` code, an object-oriented C++ software for neural-network quantum states.

names, abbreviations, and simulation information so that the `Trainer` can easily handle files and compile reports. This way, it is straightforward to refer back to a past run, reset the simulation parameters, and reproduce the results. In order to facilitate interdependence between the `Objects` beyond the scope of the `Trainer`, each `Object` must be instantiated as a `std::shared_ptr`. This choice enables other `Objects` to assume ownership when necessary and ensures safe deallocation of the pointers when the `Objects` go out of scope.

```
// main.cpp
std::shared_ptr<Object> pObject = std::make_shared<Object>(...);
```

It is essential to start every driver file (e.g., `main.cpp`) by initializing the `MPI_COMM_WORLD` communicator and setting a unique seed on each parallel process, prior to instantiating any `Objects`. These lines only need to be placed at the start of the file, in this particular order, and there is no explicit need to finalize MPI.

```
// main.cpp
mpi::initialize();
rng::initialize();
```

Then after constructing all the desired `Objects`, the trainer can be created as follows:

```
// main.cpp
Trainer T(pSystem, pWaveFunc, pSampler, pOptimizer, pCostFunc, pDensity, filename)
;
T.train();
```

Notice that the `Trainer` accepts a pointer to a `CostFunction` rather than a `Hamiltonian`. This design choice is deliberate because while variational Monte Carlo involves minimizing the energy as a reinforcement learning problem, there are situations where supervised learning can be valuable, in pretraining the `WaveFunction` for example. Our `Trainer` is versatile enough to handle both supervised and reinforcement learning problems.

One of the main jobs of the `Trainer` is to handle all opening, closing, and writing of files. For a given simulation, the `Trainer` sends four different outputs to the `data/` folder:

- `data/history/filename.txt`: This file contains the entire history of the energy, the error in the energy, the norm of the gradient, the norm of the parameter vector, the average positions, the variance of the positions, the acceptance rates for the positions and spins, and the time per iteration. If more than 10000 training epochs are requested, the `Trainer` will skip printing some iterations to prevent the file from becoming too large.
- `data/log/filename.txt`: This file contains snapshots of the contents in `data/history/filename.txt` for a broad, easy-to-interpret overview of the training. Any errors or additional print statements will be sent to this file, including the total execution time of the program.
- `data/params/filename.txt`: This file contains snapshots of the variational parameters.

ters throughout training. One can use this file to provide pretrained parameters for a wave function or a component of a wave function. If an error occurs, one can revert back to a previous stage in the training by deleting the unwanted variational parameters.

- `data/density/filename.txt`: This file contains snapshots of the density during training, provided by the `Density` object. Monitoring the density is helpful for debugging problems during training and, of course, to visualize the training of the wave function.

All files will contain a summary of the simulation details at the top of the file. Because the wave function tends to evolve the most near the beginning of the training process, the iterations at which snapshots are taken follow the pattern 0, 10, 20, ..., 100, 200, ..., 1000, 2000, ..., until the specified number of total training epochs is reached.

5.1.2 Systems

The `System` abstract class is very simple; it contains only the most basic simulation details that most other `Objects` depend on, such as the number of particles, the number of spatial dimensions, the mass, the value of \hbar , etc. The present version of `NeuralAnsatz` offers support for localized, continuous-space systems of spin-0 bosons and spin-1/2 fermions in arbitrary spatial dimensions. Although support for nucleons is present in the source code, it has not undergone testing with a Hamiltonian that includes both spin and isospin dependence. Ongoing development is underway to incorporate support for periodic systems, elements of which will be discussed herein.

In order for all `Objects` to be compatible with both bosons and fermions, we store all

spatial and spin degrees of freedom in a single matrix,

$$X \equiv \begin{bmatrix} \boldsymbol{x}_1^T \\ \boldsymbol{x}_2^T \\ \vdots \\ \boldsymbol{x}_N^T \end{bmatrix} = \begin{bmatrix} \boldsymbol{r}_1^T & \boldsymbol{s}_1^T \\ \boldsymbol{r}_2^T & \boldsymbol{s}_2^T \\ \vdots & \vdots \\ \boldsymbol{r}_N^T & \boldsymbol{s}_N^T \end{bmatrix} \in \mathbb{R}^{N \times (d+s)} \quad (5.1)$$

where the d left-most columns represent the positions, the s right-most columns represent the spins (and isospins), and the rows correspond to the particles. This construction is simply the matrix form of the notation for \mathbf{X} we introduced in Sec. 2.1. The bolded notation will be used to denote vectors, while the unbolded notation will be used for matrices, but they should be understood as interchangeable when used as inputs to functions, e.g. $f(X) = f(\mathbf{X})$. **Objects** that are designed specifically for **Bosons**, **Fermions**, or **Nucleons** will be specified in the constructor. Otherwise, all **Objects** are functions of any **System**. Due to the heirarchy of this group, functions that operate on a **System** will be compatible with any derived class, and functions that operate on **Fermions** will be compatible with **Nucleons**. Derived classes can override functions from their base classes by providing their own implementation of the function with the same name and signature.

The **Bosons** class is the simplest derived **System** class, as it contains no special attributes that are specific to bosons, other than their name. The **Fermions** class, on the other hand, determines how many distinct spatial orbitals are needed in a Slater Determinant, given a specific number of spin-up/spin-down particles. It also stores a convenient spin configuration with the desired polarization for use during the Monte Carlo sampling of spins, as discussed in Sec. 3.1.1, and projects configurations onto the pre-determined spin components of the single-particle spin-orbitals.

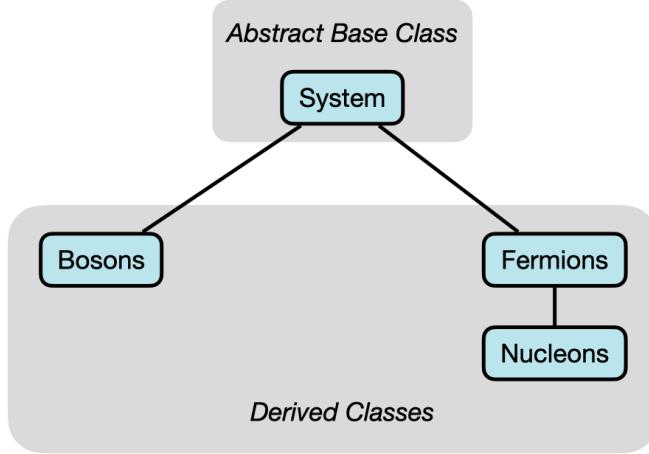


Figure 5.2: A specific example of the hierarchical structure provided by abstract base classes and derived classes. The classes that appear lower on this graph inherit properties from the classes above it.

```
// fermions.cpp
mat Fermions::get_spin_projection(const mat& X)
{
    s_ = X.col(n_dims_);
    for(int i = 0; i < n_parts_; ++i)
    {
        S_.row(i) = 0.5 * (ones_ + spin_orbitals_(i,0) * s_);
    }
    return S_;
}
```

This code snippet highlights a few additional points about our strategy. As we aim to minimize our memory footprint, we declare known array sizes in the constructor of the class, rather than dynamically allocating them on the fly. The names of all member variables are followed by a trailing underscore. In the above example, `S_` is a square matrix of dimension N that holds the spin projections, and `s_` is a vector that holds the spins from the larger matrix `X`. The variable `spin_orbitals_` is an $N \times s$ matrix that stores the spin and isospin quantum numbers corresponding to each single-particle spin-orbital. The function above is overridden in the `Nucleons` class to account for the isospin projections.

5.1.3 Wave Functions

Classes derived from the abstract base class `WaveFunction` must provide definitions for the following capabilities:

- Compute the logarithm of the trial wave function $\log \Psi_{\boldsymbol{\theta}}(\mathbf{X})$, or more generally, $\log |\Psi_{\boldsymbol{\theta}}(\mathbf{X})|$ and $\text{sgn} (\Psi_{\boldsymbol{\theta}}(\mathbf{X}))$.
- Compute the local gradient of the trial wave function with respect to the coordinates of all particles $\nabla \log \Psi_{\boldsymbol{\theta}}(\mathbf{X}) \equiv \left(\nabla_1 \log \Psi_{\boldsymbol{\theta}}(\mathbf{X}), \nabla_2 \log \Psi_{\boldsymbol{\theta}}(\mathbf{X}), \dots, \nabla_N \log \Psi_{\boldsymbol{\theta}}(\mathbf{X}) \right)$.
- Compute the sum of the local Laplacians with respect to the coordinates of each particles $\sum_{i=1}^N \frac{\nabla_i^2 \Psi_{\boldsymbol{\theta}}(\mathbf{X})}{\Psi_{\boldsymbol{\theta}}(\mathbf{X})}$. This will allow the `Hamiltonian` class to compute the local energy $E_{\boldsymbol{\theta}}(\mathbf{X}) = \frac{\hat{H} \Psi_{\boldsymbol{\theta}}(\mathbf{X})}{\Psi_{\boldsymbol{\theta}}(\mathbf{X})}$, as defined in Eq. (3.47).
- Compute the local gradient of the trial wave function with respect to the trainable parameters $\mathbf{O}_{\boldsymbol{\theta}}(\mathbf{X}) = \nabla_{\boldsymbol{\theta}} \log \Psi_{\boldsymbol{\theta}}(\mathbf{X})$, as defined in Eq. (3.48).
- Flatten and unflatten the trainable parameters $\boldsymbol{\theta}$.

As implied by the requirements listed above, computations will be written in terms of $\log \Psi_{\boldsymbol{\theta}}(\mathbf{X})$ instead of $\Psi_{\boldsymbol{\theta}}(\mathbf{X})$, which helps to mitigate the risk of numerical overflow or underflow. In addition, the abstract base class allows snapshots of the parameters to be written to file during training, and conversely, allows pretrained parameters to be loaded from file. All computations are written exactly in terms of the variational parameters $\boldsymbol{\theta}$ and the input data \mathbf{X} , or more specifically, its matrix form X . We will first discuss the various options for inputs to the neural networks, then the implementations of specific neural-network quantum states.

5.1.3.1 Inputs

The purpose of an `Input` object is to format the raw data X into a single vector \mathbf{v} which will be passed into an artificial neural network. The matrix $X \in \mathbb{R}^{N \times (d+s)}$ contains all the spatial and spin degrees of freedom of all N particles, where d and s are the numbers of spatial and spin degrees of freedom per particle, respectively. For instance, nucleons have $d = 3$ and $s = 2$, with $\mathbf{x}_i = (\mathbf{r}_i, s_i, t_i)$ containing the coordinates, spin projections on the z -axis, and the isospin projections on the z -axis. Separating the `Input` object from the `WaveFunction` object allows for increased modularity, as we can interchange the different types of inputs to the network without affecting the source code of the neural-network quantum states.

One-Body. If the one-body features are the inputs to the neural network, then the input vector \mathbf{v} can be divided into N equal parts

$$\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N) \in \mathbb{R}^V \quad (5.2)$$

with $\mathbf{v}_i \in \mathbb{R}^v$, $Nv \equiv V$. To compute the local kinetic energy, we will need the expression for the gradient of \mathbf{v} with respect to all of the particles $\mathbf{R} = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) \in \mathbb{R}^{Nd}$, which can be written in block matrix form as

$$\nabla \mathbf{v} = \begin{bmatrix} \nabla_1 \mathbf{v}_1 & \nabla_2 \mathbf{v}_1 & \cdots & \nabla_N \mathbf{v}_1 \\ \nabla_1 \mathbf{v}_2 & \nabla_2 \mathbf{v}_2 & \cdots & \nabla_N \mathbf{v}_2 \\ \vdots & \vdots & \ddots & \vdots \\ \nabla_1 \mathbf{v}_N & \nabla_2 \mathbf{v}_N & \cdots & \nabla_N \mathbf{v}_N \end{bmatrix} \in \mathbb{R}^{Nv \times Nd}, \quad (5.3)$$

with the blocks

$$\nabla_j \mathbf{v}_i = \frac{\partial \mathbf{v}_i}{\partial \mathbf{r}_j} \in \mathbb{R}^{v \times d}. \quad (5.4)$$

The Laplacian of each visible node with respect to each particle can be organized into a similar block matrix

$$\nabla^2 \mathbf{v} = \begin{bmatrix} \nabla_1^2 \mathbf{v}_1 & \nabla_2^2 \mathbf{v}_1 & \cdots & \nabla_N^2 \mathbf{v}_1 \\ \nabla_1^2 \mathbf{v}_2 & \nabla_2^2 \mathbf{v}_2 & \cdots & \nabla_N^2 \mathbf{v}_2 \\ \vdots & \vdots & \ddots & \vdots \\ \nabla_1^2 \mathbf{v}_N & \nabla_2^2 \mathbf{v}_N & \cdots & \nabla_N^2 \mathbf{v}_N \end{bmatrix} \in \mathbb{R}^{Nv \times N}, \quad (5.5)$$

with

$$\nabla_j^2 \mathbf{v}_i = \frac{\partial}{\partial \mathbf{r}_j} \cdot \frac{\partial \mathbf{v}_i}{\partial \mathbf{r}_j} \in \mathbb{R}^{v \times 1}. \quad (5.6)$$

Writing the Laplacian in this way makes it possible to use the chain rule efficiently later on.

Unsorted, Nonperiodic. In the simplest case, we can take \mathbf{v} to be the row-wise, flattened form of X . This type of input vector is suitable for localized systems and Deep Sets, a neural network that enforces permutation symmetry at the architecture level. Then we simply have $v = d + s$ and

$$\mathbf{v}_i = \mathbf{x}_i, \quad (5.7)$$

$$\nabla_j \mathbf{v}_i = \delta_{i,j} \begin{bmatrix} \mathbf{I}_{d \times d} \\ \mathbf{0}_{s \times d} \end{bmatrix}, \quad (5.8)$$

$$\nabla_j^2 \mathbf{v}_i = \mathbf{0}. \quad (5.9)$$

Sorted, Nonperiodic. Another way to enforce permutation symmetry is to sort the particles according to some scalar value. For example, one-dimensional particles can be sorted by their position, while two- and three-dimensional particles can be sorted by their Euclidean distance from the origin $r_i = \|\mathbf{r}_i\|$. Since permutation symmetry is enforced at the data-processing level, this type of input vector is suitable for any neural network. Let $Z_N = \{1, 2, \dots, N\}$ denote the set of all integers from 1 and N , and $m : Z_N \rightarrow Z_N$ denote the bijective map that sorts the particles according to the chosen rule. Then the result is very similar to before, with $v = d + s$ and

$$\mathbf{v}_i = \mathbf{x}_{m(i)}, \quad (5.10)$$

$$\nabla_j \mathbf{v}_i = \delta_{m(i),j} \begin{bmatrix} \mathbf{I}_{d \times d} \\ \mathbf{0}_{s \times d} \end{bmatrix}, \quad (5.11)$$

$$\nabla_j^2 \mathbf{v}_i = \mathbf{0}. \quad (5.12)$$

Unsorted, Periodic. For periodic systems, it is crucial to transform X so that \mathbf{v} is explicitly periodic, otherwise the neural networks can learn to break the symmetry, leading to erroneous results. To enforce L -periodicity, we transform the coordinates as

$$\mathbf{r}_i \mapsto \tilde{\mathbf{r}}_i = \left(\cos\left(\frac{2\pi \mathbf{r}_i}{L}\right), \sin\left(\frac{2\pi \mathbf{r}_i}{L}\right) \right) \in \mathbb{R}^{2d}, \quad (5.13)$$

such that the components of the input vector are

$$\mathbf{v}_i = (\tilde{\mathbf{r}}_i, \mathbf{s}_i) \in \mathbb{R}^{2d+s}, \quad (5.14)$$

implying $v = 2d + s$. Then the derivatives are

$$\nabla_j \mathbf{v}_i = \frac{2\pi}{L} \delta_{i,j} \begin{bmatrix} -\text{diag} \left[\sin \left(\frac{2\pi \mathbf{r}_i}{L} \right) \right] \\ \text{diag} \left[\cos \left(\frac{2\pi \mathbf{r}_i}{L} \right) \right] \\ \mathbf{0}_{s \times d} \end{bmatrix}, \quad (5.15)$$

$$\nabla_j^2 \mathbf{v}_i = - \left(\frac{2\pi}{L} \right)^2 \delta_{i,j} \begin{bmatrix} \cos \left(\frac{2\pi \mathbf{r}_i}{L} \right) \\ \sin \left(\frac{2\pi \mathbf{r}_i}{L} \right) \\ \mathbf{0} \end{bmatrix} = - \left(\frac{2\pi}{L} \right)^2 \delta_{i,j} \begin{bmatrix} \tilde{\mathbf{r}}_i \\ \mathbf{0} \end{bmatrix}, \quad (5.16)$$

where $\text{diag} [\mathbf{r}]$ refers to a diagonal matrix with \mathbf{r} on the diagonal, and the trigonometric functions are applied element-wise.

Sorted, Periodic. While the inputs themselves must be periodic in this case, the sorting map does not necessarily have to sort according to periodic positions or distances. We therefore use the same sorting map m as the nonperiodic case. The resulting transformation is similar to the previous case, except with the indices permuted

$$\mathbf{v}_i = (\tilde{\mathbf{r}}_{m(i)}, \mathbf{s}_{m(i)}) \in \mathbb{R}^{2d+s}, \quad (5.17)$$

$$\nabla_j \mathbf{v}_i = \frac{2\pi}{L} \delta_{m(i),j} \begin{bmatrix} -\text{diag} \left[\sin \left(\frac{2\pi \mathbf{r}_{m(i)}}{L} \right) \right] \\ \text{diag} \left[\cos \left(\frac{2\pi \mathbf{r}_{m(i)}}{L} \right) \right] \\ \mathbf{0}_{s \times d} \end{bmatrix}, \quad (5.18)$$

$$\nabla_j^2 \mathbf{v}_i = - \left(\frac{2\pi}{L} \right)^2 \delta_{m(i),j} \begin{bmatrix} \cos \left(\frac{2\pi \mathbf{r}_{m(i)}}{L} \right) \\ \sin \left(\frac{2\pi \mathbf{r}_{m(i)}}{L} \right) \\ \mathbf{0} \end{bmatrix} = - \left(\frac{2\pi}{L} \right)^2 \delta_{m(i),j} \begin{bmatrix} \tilde{\mathbf{r}}_{m(i)} \\ \mathbf{0} \end{bmatrix}, \quad (5.19)$$

with $v = 2d + s$.

Two-Body. If the inputs to the neural network are two-body features, then \mathbf{v} can be divided into $P = N(N - 1)/2$ equal parts

$$\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_P) \in \mathbb{R}^V, \quad (5.20)$$

where $\mathbf{v} \in \mathbb{R}^v$ and $Pv \equiv V$. The gradient of \mathbf{v} with respect to all of the particles $\mathbf{r} = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N) \in \mathbb{R}^{ND}$ can be written in block matrix form as

$$\nabla \mathbf{v} = \begin{bmatrix} \nabla_1 \mathbf{v}_1 & \nabla_2 \mathbf{v}_1 & \cdots & \nabla_N \mathbf{v}_1 \\ \nabla_1 \mathbf{v}_2 & \nabla_2 \mathbf{v}_2 & \cdots & \nabla_N \mathbf{v}_2 \\ \vdots & \vdots & \ddots & \vdots \\ \nabla_1 \mathbf{v}_P & \nabla_2 \mathbf{v}_P & \cdots & \nabla_N \mathbf{v}_P \end{bmatrix} \in \mathbb{R}^{Pv \times Nd}, \quad (5.21)$$

with the blocks

$$\nabla_k \mathbf{v}_p = \frac{\partial \mathbf{v}_p}{\partial \mathbf{r}_k} \in \mathbb{R}^{v \times d}. \quad (5.22)$$

The Laplacian of each visible node with respect to each particle can also be organized into a block matrix

$$\nabla^2 \mathbf{v} = \begin{bmatrix} \nabla_1^2 \mathbf{v}_1 & \nabla_2^2 \mathbf{v}_1 & \cdots & \nabla_N^2 \mathbf{v}_1 \\ \nabla_1^2 \mathbf{v}_2 & \nabla_2^2 \mathbf{v}_2 & \cdots & \nabla_N^2 \mathbf{v}_2 \\ \vdots & \vdots & \ddots & \vdots \\ \nabla_1^2 \mathbf{v}_P & \nabla_2^2 \mathbf{v}_P & \cdots & \nabla_N^2 \mathbf{v}_P \end{bmatrix} \in \mathbb{R}^{Pv \times N}, \quad (5.23)$$

with

$$\nabla_k^2 \mathbf{v}_p = \frac{\partial}{\partial \mathbf{r}_k} \cdot \frac{\partial \mathbf{v}_p}{\partial \mathbf{r}_k} \in \mathbb{R}^{v \times 1}. \quad (5.24)$$

In the following, we will assume that the two-body features are symmetric, meaning the feature for the pair i, j is the same as the feature for pair j, i . We store the mapping between a particular pair (i, j) and their pair index p as $i(p)$ and $j(p)$, always assuming $i < j$. To simplify notation, we will use

$$\mathbf{z}_p = \mathbf{r}_{i(p)} - \mathbf{r}_{j(p)} \quad (5.25)$$

to denote the separation vector between particles $i(p)$ and $j(p)$ corresponding to the pair p , and

$$z_p = \|\mathbf{z}_p\| \quad (5.26)$$

to denote the distance between the pair.

Unsorted, Nonperiodic. The components of the input vector are

$$\mathbf{v}_p = (z_p, \mathbf{s}_{i(p)} * \mathbf{s}_{j(p)}) \in \mathbb{R}^{1+s}, \quad (5.27)$$

where $*$ denotes an element-wise product of the spin degrees of freedom and $v = 1 + s$. The corresponding derivatives are

$$\nabla_k \mathbf{v}_p = \left(\delta_{i(p), k} - \delta_{j(p), k} \right) \frac{1}{z_p} \begin{bmatrix} \mathbf{z}_p^T \\ \mathbf{0}_{Q \times D} \end{bmatrix}, \quad (5.28)$$

$$\nabla_k^2 \mathbf{v}_p = \left(\delta_{i(p), k} + \delta_{j(p), k} \right) \frac{2}{z_p} \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}. \quad (5.29)$$

Sorted, Nonperiodic. Let $m : Z_P \longrightarrow Z_P$ be the mapping that sorts the pairs according to the distance between them \varkappa_p . Then,

$$\mathbf{v}_p = (\varkappa_{m(p)}, \mathbf{s}_{i(m(p))} * \mathbf{s}_{j(m(p))}), \quad (5.30)$$

$$\nabla_k \mathbf{v}_p = \left(\delta_{i(m(p)),k} - \delta_{j(m(p)),k} \right) \frac{1}{\varkappa_{m(p)}} \begin{bmatrix} \boldsymbol{\varkappa}_{m(p)}^T \\ \mathbf{0}_{Q \times D} \end{bmatrix}, \quad (5.31)$$

$$\nabla_k^2 \mathbf{v}_p = \left(\delta_{i(m(p)),k} + \delta_{j(m(p)),k} \right) \frac{2}{\varkappa_{m(p)}} \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}. \quad (5.32)$$

Unsorted, Periodic. The distances between pairs are transformed to their periodic counterparts,

$$\varkappa_p \longmapsto \tilde{\varkappa}_p = \left\| \sin \left(\frac{\pi \boldsymbol{\varkappa}_p}{L} \right) \right\|, \quad (5.33)$$

making the components of the input vector become

$$\mathbf{v}_p = (\tilde{\varkappa}_p, \mathbf{s}_{i(p)} * \mathbf{s}_{j(p)}) \in \mathbb{R}^{1+Q}, \quad (5.34)$$

where $\mathbf{s}_i \in \mathbb{R}^Q$ is a vector containing the spin degrees of freedom of the i -th particle and $*$ denotes an element-wise product.

$$\nabla_k \mathbf{v}_p = \left(\delta_{i(p),k} - \delta_{j(p),k} \right) \frac{\pi}{2L\tilde{\varkappa}_p} \begin{bmatrix} \sin \left(\frac{2\pi \boldsymbol{\varkappa}_p}{L} \right)^T \\ \mathbf{0}_{Q \times D} \end{bmatrix} \in \mathbb{R}^{(1+Q) \times D} \quad (5.35)$$

$$\nabla_k^2 \mathbf{v}_p = \left(\delta_{i(p),k} + \delta_{j(p),k} \right) \frac{2}{\tilde{\mathbf{z}}_p} \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{1+Q} \quad (5.36)$$

Sorted, Periodic. Let $m : Z_P \longrightarrow Z_P$ be the mapping that sorts the pairs according to the periodic distance between them $\tilde{\mathbf{z}}_k$. Then,

$$\mathbf{v}_p = (\tilde{\mathbf{z}}_{m(p)}, \mathbf{s}_{i(m(p))} * \mathbf{s}_{j(m(p))}) \in \mathbb{R}^{1+Q}, \quad (5.37)$$

$$\nabla_k \mathbf{v}_p = \left(\delta_{i(m(p)),k} - \delta_{j(m(p)),k} \right) \frac{\pi}{2L\tilde{\mathbf{z}}_{m(p)}} \begin{bmatrix} \sin\left(\frac{2\pi\mathbf{z}_{m(p)}}{L}\right)^T \\ \mathbf{0}_{Q \times D} \end{bmatrix} \in \mathbb{R}^{(1+Q) \times D} \quad (5.38)$$

$$\nabla_k^2 \mathbf{v}_p = \left(\delta_{i(m(p)),k} + \delta_{j(m(p)),k} \right) \frac{2}{\tilde{\mathbf{z}}_{m(p)}} \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{1+Q} \quad (5.39)$$

To compute \mathbf{v} , $\nabla\mathbf{v}$, and $\nabla^2\mathbf{v}$ in terms of the inputs X inside a particular `WaveFunction` class, we call the `format` function using a pointer to the `Input` object. There are three versions of the same `format` function since \mathbf{v} , $\nabla\mathbf{v}$, and $\nabla^2\mathbf{v}$ are not necessary for every computation. All derived classes, e.g. `OneBodyInputs` and `TwoBodyInputs`, must provide definitions for these functions.

```
// input.hpp
virtual void format(const mat& X, vec& v) = 0;
virtual void format(const mat& X, vec& v, mat& dv) = 0;
virtual void format(const mat& X, vec& v, mat& dv, mat& d2v) = 0;
```

5.1.3.2 Continuous Restricted Boltzmann Machines

The various versions of multivariate Gaussian-Binary and Gaussian-Uniform restricted Boltzmann machines are implemented under a single class called `ContinuousRBM`. In Sec. 4.7.1, we defined a new function called $f(x)$ that allowed all the log-likelihoods and their gradients with respect to the parameters to have the same form. Since we already need to define activation functions to implement feedforward neural networks, we implement $f(x)$ as activation functions as well, in order to streamline the `ContinuousRBM` class. More information on the implementation of activation functions will be provided in the next section.

We will use `ContinuousRBM` as a symmetric, positive-definite Jastrow wave function for bosons, or as a Jastrow correlator for fermions, combined with an antisymmetric Slater determinant. Recall the form of the log-likelihood for a continuous RBM written in terms of the trainable parameters $\boldsymbol{\theta} = (\mathbf{a}, \text{triu}(S), \mathbf{b}, W)$,

$$\log \mathcal{P}(\mathbf{v}) = -\frac{1}{2}(\mathbf{v} - \mathbf{a})^T \exp(S)(\mathbf{v} - \mathbf{a}) + \sum_{j=1}^H f(z_j(\mathbf{v})) + C, \quad (5.40)$$

$$\mathbf{z}(\mathbf{v}) = \mathbf{b} + W^T \exp(S)\mathbf{v}. \quad (5.41)$$

where C is a constant. Since $\mathcal{P}(\mathbf{v})$ is the marginal probability distribution of the visible nodes \mathbf{v} , it is natural to define our trial wave function as

$$\log \Psi_{\boldsymbol{\theta}}(\mathbf{X}) \equiv \log \mathcal{P}(\mathbf{v}(\mathbf{X})). \quad (5.42)$$

In fact, we could make the relationship between the probability amplitude $\Psi_{\boldsymbol{\theta}}(\mathbf{X})$ and the probability distribution more exact by writing $\log |\Psi_{\boldsymbol{\theta}}(\mathbf{X})| \equiv \frac{1}{2} \log \mathcal{P}(\mathbf{v}(\mathbf{X}))$, but due to the presence of all the trainable parameters, the two forms are practically equivalent.

Computing the local gradient with respect to $\mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_N)$ becomes straightforward because of the way the `Input` object organizes $\nabla \mathbf{v}(\mathbf{X})$ into blocks. Namely,

$$\nabla \log \Psi_{\theta}(\mathbf{X}) = \left[-(\mathbf{v} - \mathbf{a})^T + \mathbf{f}'(\mathbf{z}(\mathbf{v}(\mathbf{X})))^T W^T \right] \exp(S) \nabla \mathbf{v}(\mathbf{X}). \quad (5.43)$$

To compute the sum of local Laplacians, first notice that

$$\frac{1}{\Psi_{\theta}(\mathbf{X})} \nabla_i^2 \Psi_{\theta}(\mathbf{X}) = \nabla_i^2 \log \Psi_{\theta}(\mathbf{X}) + (\nabla_i \log \Psi_{\theta}(\mathbf{X}))^2. \quad (5.44)$$

If we again, organize our Laplacians into blocks, we can simplify the computation to be in terms of matrix and array operations only, eliminating the need for any for-loops.

```
// continuous_rbm.cpp
long double ContinuousRBM::get_local_laplacian(const mat& X)
{
    pInput_->format(X, v_, dv_, d2v_);

    K_ = exp(S_);
    va_ = v_ - a_;
    dA_ = -va_.transpose() * K_ * dv_;
    d2A_ = - (va_.transpose() * K_ * d2v_).sum()
            + - (dv_.transpose() * K_ * dv_).trace();

    z_ = b_ + W_.transpose() * K_ * v_;
    dz_ = W_.transpose() * K_ * dv_;
    d2z_ = W_.transpose() * K_ * d2v_;

    df_ = pFunc_->get_df(z_);
    d2f_ = pFunc_->get_d2f(z_);

    dB_ = df_.transpose() * dz_;
    d2B_ = df_.dot(d2z_.rowwise().sum())
           + d2f_.dot(dz_.rowwise().squaredNorm());

    return d2A_ + d2B_ + (dA_ + dB_).squaredNorm();
}
```

In the above code snippet, `pInput_` is the member pointer to the `Input` object and `pFunc_` is the member pointer to the (activation) function $f(x)$. Finally, because we defined our trial wave function as the marginal probability distribution of the visible nodes, the derivatives of the log-likelihood in Eq. (4.37) coincide with the derivatives of $\log \Psi_{\theta}(\mathbf{X})$.

In general, the variational parameters of any restricted Boltzmann machine can be initialized to small random values near zero. However, because RBMs are highly interpretable models, we can provide much better guesses for the initial parameters based on our physical intuition. As a reminder, the variational parameters of an mGB-RBM and an mGU-RBM are $\theta = (\mathbf{a}, \text{triu}(S), \mathbf{b}, W)$, where \mathbf{a} are the biases of the visible nodes, S is used to reparameterize the inverse covariance matrix

$$\Sigma_{ii}^{-1} = \exp(S_{ii}), \text{ and } \Sigma_{ij}^{-1} = S_{ij} \text{ for } j \neq i,$$

\mathbf{b} are the biases of the hidden nodes, and W is a weight matrix connecting the visible and hidden nodes. If \mathbf{b} and W are both identically zero, the marginal probability distribution of both RBMs become a multivariate Gaussian. Thus setting \mathbf{b} and W to small random values near zero represents small deviations from the multivariate Gaussian. For small \mathbf{b} and W , the individual peaks of the Gaussian nearly coincide with the biases \mathbf{a} . Therefore, we can set the initial values of \mathbf{a} based on our guesses for the most likely positions of the particles.

The initialization of the matrix S is the most important, as choosing convenient values can prevent the multivariate Gaussian from becoming multi-modal early on during training, leading to the collapse of the sampling process. In addition, the above reparametrization guarantees the invertibility of Σ^{-1} but it does not guarantee that the diagonals of Σ are positive. Since negative variances are unphysical, we bias our network to provide positive

variances at the beginning of training. Then the RBM naturally learns to maintain positive variances on its own. This is accomplished by choosing the diagonal values of S as small positive values, and the off diagonal values as small negative values.

5.1.3.3 Feedforward Neural Networks

Deep feedforward neural networks are implemented under a class called `DeepFNN`. This class assumes the trial wave function is given by

$$\log \Psi_{\theta}(\mathbf{X}) \equiv \rho(\mathbf{v}(\mathbf{X})), \quad (5.45)$$

where ρ is a feedforward neural network with a single output, and $\mathbf{v}(\mathbf{X})$ is constructed by the `format` function belonging to the `Input` object. In addition, ρ is assumed to be rectangular in structure, meaning all the hidden layers have the same number of nodes.

`DeepFNN` consists of a sequence of `DenseLayer` and `ActivationLayer` objects, which represent the affine transformations and the nonlinear, element-wise transformations, respectively. Both `DenseLayer` and `ActivationLayer` are derived from the abstract base class `Layer`. In addition to simple functions that facilitate the flattening and unflattening of the parameters, the `Layer` objects have the following key member functions:

```
// layer.hpp
virtual void forward(vec& h) = 0;
virtual void forward(vec& h, mat& dh) = 0;
virtual void forward(vec& h, mat& dh, mat& d2h) = 0;
virtual void backward(int& iter, vec& gradient, vec& du_dh) = 0;
```

The `forward` functions are used during the computations of the (log of the) wave function, the local gradient with respect to the coordinates, and the sum of local Laplacians. The

`backward` function is used during the computation of the local gradient of the trial wave function with respect to the variational parameters.

Even though activation functions are not formally considered layers of the network themselves, it is convenient to implement them as layers because the forward-passing and backward-passing steps can be written extremely simply. More specifically, if we store all the `DenseLayer` and `ActivationLayer` objects into a single `std::vector` of pointers to each layer, we can write the wave function computation as the following:

```
// deepfnn.cpp
long double DeepFNN::get_logpsi(const mat& X)
{
    pInput_->format(X, h_);
    for(int l = 0; l < n_layers_; ++l) pLayers_[l]->forward(h_);
    return h_(0);
}
```

Similarly, the gradient with respect to the parameters is computed as the following:

```
// deepfnn.cpp
vec DeepFNN::get_param_gradient(const mat& x)
{
    // forward pass to store intermediate values
    pInput_->format(x, h_);
    for(int l = 0; l < n_layers_; ++l) pLayers_[l]->forward(h_);

    // backward pass to calculate gradient
    iter_ = 0;
    du_dh_ = vec::Ones(1);
    for(int l = n_layers_-1; l >= 0; --l) pLayers_[l]->backward(iter_,
                                                                param_grad_,
                                                                du_dh_);
    return param_grad_;
}
```

Notice that the above gradient calculation includes a forward pass. This is because the

gradient of the output depends on the input for both `DenseLayer` and `ActivationLayer`.

Thus, during the forward step, a copy of the input vector is stored for later use.

The forward pass for the `DenseLayer` object is exceptionally simple, as it consists of only a linear transformation.

```
// dense.cpp
void DenseLayer::forward(vec& h, mat& dh, mat& d2h)
{
    h = W_ * h + b_;
    dh = W_ * dh;
    d2h = W_ * d2h;
}
```

Similarly, the backward pass can be written simply as well.

```
// dense.cpp
void DenseLayer::backward(int& iter, vec& gradient, vec& du_dh)
{
    du_dW_ = du_dh * h_in_.transpose();
    du_db_ = du_dh;
    du_dh = W_.transpose() * du_dh;

    gradient.segment(iter, n_params_) = concat(flatten(du_dW_), du_db_);
    iter += n_params_;
}
```

In the above, the character `u` represents the single output of the entire network, written as ρ in Eq. (5.45). The first three lines compute the gradient of the output with respect to the weights and biases, as well as the gradient with respect to the hidden nodes. Then the gradients with respect to the parameters are flattened and stored in their respective locations in the incoming `gradient` vector, which stores the gradients with respect to *all* weights and biases in the network. These weights and biases are initialized using the Glorot normal scheme and the biases are initialized to zero.

On the other hand, the forward step for the `ActivationLayer` is slightly more complicated because of the nonlinearity and the block organization of the gradient and Laplacian.

```
// activation.cpp
void Activation::forward(vec& h, mat& dh, mat& d2h)
{
    for(int i = 0; i < h.size(); ++i)
    {
        d2h.row(i) *= get_df(h(i));
        for(int p = 0; p < n_parts_; ++p)
        {
            d2h(i,p) += get_d2f(h(i))
                        * dh.row(i).segment(p*n_dims_, n_dims_).squaredNorm();
        }
        dh.row(i) *= get_df(h(i));
        h(i) = get_f(h(i));
    }
}
```

Here, the functions `get_f`, `get_df`, and `get_d2f` are pure virtual functions representing $f(x)$, $f'(x)$ and $f''(x)$, respectively. They must be defined in each class derived from `ActivationLayer`. Luckily, the backward step is even simpler than for the `DenseLayer` because there are no trainable parameters.

```
// activation.cpp
void Activation::backward(int& iter, vec& gradient, vec& du_dh)
{
    for(int i = 0; i < du_dh.size(); ++i) du_dh(i) *= get_df(h_in_(i));
}
```

Again, the backward pass involves the original input vector, which was saved in the member `h_in_`.

As a final note, DeepFNN is often accompanied by an enveloping `FixedGaussian`, but we will omit discussion on the latter as it is straightforward to implement and contains

no variational parameters. Products of DeepFNN and FixedGaussian are generated by `ProductWaveFunction`, discussed in Sec. 5.1.3.6.

5.1.3.4 Deep Sets

The present version of `NeuralAnsatz` also provides a simple implementation of Deep Sets in a class called `DeepSet`. In this case, the trial wave function is given by

$$\log \Psi_{\theta}(\mathbf{X}) \equiv \rho \left(\text{pool} \left(\{\phi(\mathbf{v}_i(\mathbf{X}))\} \right) \right), \quad (5.46)$$

where ρ is a feedforward neural network with a single output, ϕ is feedforward neural network with a vector output, and `pool` is a pooling operation specified by a `PoolingLayer`. Additionally, the index i in $\mathbf{v}_i(\mathbf{X})$ may run over the total number of particles N or the total number of pairs $N(N - 1)/2$, depending on which `Input` object is chosen. The overall structure of `DeepSet` is very similar to `DeepFNN` except that there are now *two* `std::vector` members of pointers to the `Layer` objects corresponding to ρ and ϕ , and an additional pointer to the `PoolingLayer`. For example, the local laplacian is implemented as the following:

```
// deepset.cpp
long double DeepSet::get_local_laplacian(const mat& X)
{
    pInput_->format(X, h_, dh_, d2h_);
    for(int l = 0; l < n_layers1_; ++l) pLayers1_[l]->forward(h_, dh_, d2h_);
    pPooling_->forward(h_, dh_, d2h_);
    for(int l = 0; l < n_layers2_; ++l) pLayers2_[l]->forward(h_, dh_, d2h_);
    return dh_.squaredNorm() + d2h_.sum();
}
```

As a side note, the elements of `pLayers1_` include pointers to objects derived from `DenseLayer`, called `DenseSet`, that passes elements of the set forward and backward in parallel rather than

the entire input at once.

The more important change to discuss is the inclusion of the `PoolingLayer`. As a simple illustration, see the implementation of `SumPooling` below.

```
// sum.cpp
void SumPooling::forward(vec& h)
{
    h_out_.setZero();
    for(int s = 0; s < n_set_; ++s) h_out_ += h.segment(s * n_out_, n_out_);
    h = h_out_;
}
```

Since the ordering of the set elements is destroyed in the forward process, and each set element is given equal weight to the eventual output, the backward direction can be completed by essentially creating `n_set_` copies of the gradient and passing it back to the `DenseSet` sequence. Then the `DenseSet` layers average the contributions to the gradient from each element of the set. If the set elements have unequal weights in the eventual output, the inputs need to be stored to compute the corresponding weights in the backward direction.

5.1.3.5 Slater Determinants

The implementation of `SlaterDeterminant` in `NeuralAnsatz` is limited by the fact that our sampling procedure involves perturbing all particles in a single Monte Carlo step, rather than perturbing one particle at a time. In addition, we include the option to sample spins as well, where the number of spin-up and spin-down particles are not necessarily equal. Consequently, we cannot take advantage of commonly-employed tricks involving updating a single row or single column of the determinant matrix, nor is it straightforward to split the determinant into a product of two determinants corresponding to the spin-up and spin-down

particles. Then, the analytical computations of the gradients and Laplacians of the Slater determinant with respect to the particle coordinates require taking the determinant *and* the inverse of the matrix. According to `Eigen`'s documentation, both of these operations are only guaranteed to be stable for at most 4×4 matrices. Nonetheless, we describe the implementation of `SlaterDeterminant` here, but we will exercise caution for systems larger than $N = 4$.

A Slater determinant is given by

$$\Psi(\mathbf{X}) \equiv \frac{1}{\sqrt{N!}} \det [\varphi_\alpha(\mathbf{x}_i)] = \frac{1}{\sqrt{N!}} \det \begin{bmatrix} \varphi_1(\mathbf{x}_1) & \varphi_1(\mathbf{x}_2) & \cdots & \varphi_1(\mathbf{x}_N) \\ \varphi_2(\mathbf{x}_1) & \varphi_2(\mathbf{x}_2) & \cdots & \varphi_2(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_N(\mathbf{x}_1) & \varphi_N(\mathbf{x}_2) & \cdots & \varphi_N(\mathbf{x}_N) \end{bmatrix}, \quad (5.47)$$

where φ_α are single-particle spin-orbitals indexed by α , and we have omitted the subscript θ on the ansatz $\Psi_\theta(\mathbf{X})$ because we do not consider trainable spin-orbitals in this implementation. As mentioned in Sec. 5.1.2, when a system of `Fermions` is constructed, the number of distinct spatial orbitals that is required in the Slater determinant is automatically stored. The member `orbital_index_` acts as the mapping between orbital index α and the quantum numbers of the spatial orbitals. At the same time, the spin components of the spin-orbitals are determined and stored in the member `spin_orbitals_`. Together, `orbital_index_` and `spin_orbitals_` determine each unique spin-orbital.

The orbitals are computed by a `Basis` object which takes the quantum numbers of the spatial orbitals stored in `orbital_index_` and evaluates the orbitals for all particles. For the spin components, the `Fermions` class computes the spin projections using the information stored in `spin_orbitals_`. One can simply take the element-wise product to fill the elements

of the matrix, and then take the determinant.

```
// slaterdeterminant.cpp
long double SlaterDeterminant::get_logpsi(const mat& X)
{
    D_ = (pBasis_->get_states(X)).cwiseProduct(pFermions_->get_spin_projection(X))
        ;
    det_ = D_.determinant();
    return log1(fabs1(det_)) + lognorm_;
}
```

Storing a copy of the determinant in `det_` allows the `Hamiltonian` to access the sign of the determinant if the interaction is spin dependent.

To write the local gradient with respect to a particle's coordinates, we first define the matrix

$$D(\mathbf{X}) \equiv [\varphi_\alpha(\mathbf{x}_i)], \quad (5.48)$$

for brevity, and use Jacobi's formula,

$$\frac{1}{\Psi(\mathbf{X})} \nabla_i \Psi(\mathbf{X}) = \text{tr} \left(D(\mathbf{X})^{-1} \nabla_i D(\mathbf{X}) \right), \quad (5.49)$$

where $\nabla_i D(\mathbf{X})$ means to evaluate the derivative for each spatial dimension of particle i , and the subsequent trace is taken separately for each direction. For the local Laplacian, we also use that $\nabla_i D(\mathbf{X})^{-1} = -D(\mathbf{X})^{-1} \nabla_i D(\mathbf{X}) D(\mathbf{X})^{-1}$,

$$\begin{aligned} \frac{1}{\Psi(\mathbf{X})} \nabla_i^2 \Psi(\mathbf{X}) &= \text{tr} \left(D(\mathbf{X})^{-1} \nabla_i^2 D(\mathbf{X}) \right) + \text{tr}^2 \left(D(\mathbf{X})^{-1} \nabla_i D(\mathbf{X}) \right) \\ &\quad - \text{tr} \left(D(\mathbf{X})^{-1} \nabla_i D(\mathbf{X}) D(\mathbf{X})^{-1} \nabla_i D(\mathbf{X}) \right). \end{aligned} \quad (5.50)$$

The computation of the spatial parts of $\nabla_i D(\mathbf{X})$ and $\nabla_i^2 D(\mathbf{X})$ are handled by the `Basis`.

5.1.3.6 Products

The wrapper class `ProductWaveFunction`, derived from `WaveFunction`, enables the creation of products involving any desired number F of `WaveFunction` factors,

$$\Psi_{\boldsymbol{\theta}}(\mathbf{X}) \equiv \prod_{f=1}^F \Psi_{\boldsymbol{\theta}_f}(\mathbf{X}), \quad (5.51)$$

where the parameters are concatenated together $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_F)$. It is straightforward to show that the local gradient with respect to a particle's coordinates is given by

$$\frac{1}{\Psi_{\boldsymbol{\theta}}(\mathbf{X})} \nabla_i \Psi_{\boldsymbol{\theta}}(\mathbf{X}) = \sum_{f=1}^F \frac{1}{\Psi_{\boldsymbol{\theta}_f}(\mathbf{X})} \nabla_i \Psi_{\boldsymbol{\theta}_f}(\mathbf{X}), \quad (5.52)$$

and the local Laplacians are given by

$$\begin{aligned} \frac{1}{\Psi_{\boldsymbol{\theta}}(\mathbf{X})} \nabla_i \Psi_{\boldsymbol{\theta}}(\mathbf{X}) &= \sum_{f=1}^F \frac{1}{\Psi_{\boldsymbol{\theta}_f}(\mathbf{X})} \nabla_i^2 \Psi_{\boldsymbol{\theta}_f}(\mathbf{X}) \\ &\quad + 2 \sum_{f < g}^F \left(\frac{1}{\Psi_{\boldsymbol{\theta}_f}(\mathbf{X})} \nabla_i \Psi_{\boldsymbol{\theta}_f}(\mathbf{X}) \right) \cdot \left(\frac{1}{\Psi_{\boldsymbol{\theta}_g}(\mathbf{X})} \nabla_i \Psi_{\boldsymbol{\theta}_g}(\mathbf{X}) \right). \end{aligned} \quad (5.53)$$

The utilization of the `ProductWaveFunction` becomes necessary when working with a `SlaterDeterminant` or when employing a `FixedGaussian` as an enveloping function. Furthermore, different neural-network quantum states can be combined to achieve the best of both (or all) worlds.

5.1.4 Hamiltonian

The `Hamiltonian` class is an abstract base class that takes `System` and `WaveFunction` objects as input to its constructor. By default, it is the free Hamiltonian which just contains the

kinetic energy operator. Its virtual member functions `get_external` and `get_interaction` both return zero, unless they are overridden by derived classes.

```
// hamiltonian.cpp
long double Hamiltonian::get_local_energy(const mat& X)
{
    return get_local_kinetic(X) + get_external(X) + get_interaction(X);
}

long double Hamiltonian::get_local_kinetic(const mat& X)
{
    return -hbar2m_ * pWaveFunc_->get_local_laplacian(X);
}
```

If the interaction potential contains a singularity, it is strongly-recommended to regularize the singularity by introducing a hyperparameter r_0 and transforming as

$$\frac{1}{r^k} \mapsto \frac{f^k(r/r_0)}{r^k}, \quad (5.54)$$

where $f(r)$ is a function that scales as $\sim r$ for small r and approaches 1 for large r . By doing so, the long-range behavior of the original interaction is preserved, but at $r = 0$, the potential is finite, taking a value of $(1/r_0)^k$. We typically choose $f(r) \equiv \tanh(r)$. Pretraining the neural-network quantum states on a softer potential before lowering the value of r_0 to tackle harder potentials has been found to significantly stabilize and accelerate performance. In addition, this allows the neural-network quantum states to discover the cusp condition on its own, rather than providing it from the very beginning, demonstrating their representation power.

It should be noted that adding new `Hamiltonian` objects is extremely straightforward, as one only needs to specify the external and interaction potentials for a single configu-

ration sample X . Since the neural-network quantum states are entirely separate from the Hamiltonian, one can swap different interactions without affecting the states.

5.1.5 Sampler

`Metropolis` and `Importance` are classes derived from the abstract `Sampler` class. The `Sampler` class also provides the option to sample the spins as well. The implementations follow the descriptions provided in Sec. 3.1.1.

In addition to pointers to the `System` and `WaveFunction` objects, the class constructor takes `n_eq_init`, `n_eq`, `n_void`, `sigma0`, `step` as arguments. The integer `n_void` specifies the number of samples to skip between each effective sample, constituting one Monte Carlo cycle. At the very beginning of the simulation, initial positions are drawn from a Gaussian with standard deviation `sigma0`, and the sampler is equilibrated using `n_eq_init` Monte Carlo cycles. Instead of equilibrating the sampler from scratch at the beginning of each training epoch, we recover the last configuration and equilibrate using only `n_eq` cycles. Typically, `n_eq` can be one or two orders of magnitude smaller than `n_eq_init`. This construction assumes the learning rate is not so large that the wave function changes significantly in a single epoch. The parameter `step` determines the size of the spatial perturbations.

5.1.6 Optimizer

The abstract `Optimizer` class takes a pointer to the `WaveFunction` object, the initial learning rate `eta0`, and the total number of epochs `n_epochs` as input. Derived classes only need to provide a definition for the following pure virtual function.

```
// optimizer.hpp
```

```
virtual void update_params(vec gradient) = 0;
```

This function simply takes a gradient vector and updates the parameters of the `WaveFunction` object. The specific optimization algorithms were discussed in Sec. 3.2.2.

5.1.7 Cost Functions

The `CostFunction` computes both the cost and the gradient of the cost with respect to the parameters. Derived classes can be designed for reinforcement learning or for supervised learning. In variational Monte Carlo, we use the `Energy` as the cost, which has the following constructor signature.

```
// energy.hpp
Energy(int n_samples,
       long double epsilon,
       std::shared_ptr<System> pSystem,
       std::shared_ptr<WaveFunction> pWaveFunc,
       std::shared_ptr<Sampler> pSampler,
       std::shared_ptr<Hamiltonian> pHamiltonian);
```

If the parameter `epsilon` is larger than zero, stochastic reconfiguration is used to precondition the gradient. To minimize the memory footprint of the program, samples are taken on the fly and expectation values are computed as running averages. Once all parallel processes have completed the computations of their local expectation values, Open MPI is used to communicate between processes and aggregate the results. Since the `Energy` requires the `Hamiltonian` object as input, there is no need for the `Trainer` to have explicit ties to the `Hamiltonian`. The gradient calculation for the `Energy` is shown below.

```
// energy.cpp
```

```

vec Energy::get_gradient()
{
    reset_gradient_estimator();
    pSampler_->equilibrate();
    for(int n = 0; n < n_samples_loc_; ++n)
    {
        X_ = pSampler_->get_next_config();
        E_ = pHamiltonian_->get_local_energy(X_);
        O_ = pWaveFunc_->get_param_gradient(X_);
        R_ = X_.leftCols(n_dims_);

        E_avg_loc_ += E_;
        E2_avg_loc_ += E_ * E_;
        O_avg_loc_ += O_;
        EO_avg_loc_ += E_ * O_;
        r_avg_loc_ += R_.sum();
        r2_avg_loc_ += R_.squaredNorm();

        if(use_stochastic_reconfig_) O2_avg_loc_ += O_ * O_.transpose();
        if(pSampler_->positions_accepted()) r_acc_loc_ += 1;
        if(pSampler_->spins_accepted()) s_acc_loc_ += 1;
    }
    finalize_gradient_estimator();
    return gradient_;
}

```

If supervised learning of the `WaveFunction` is desired, the `CostFunction` needs to take a target `WaveFunction` object as input in place of the `Hamiltonian`. For instance, the constructor for the `Overlap` between two wave functions has the following signature.

```

// overlap.hpp
Overlap(int n_samples,
         long double sigma,
         std::shared_ptr<System> pSystem,
         std::shared_ptr<WaveFunction> pWaveFunc,
         std::shared_ptr<Sampler> pSampler,
         std::shared_ptr<WaveFunction> pTargetWaveFunc);

```

The parameter `sigma` specifies Gaussian probability distribution from which samples are drawn.

5.1.8 Density

The `Density` object takes snapshots of the spatial distribution during training. If the `WaveFunction` is a product, it is possible to take snapshots of each individual factor as well. The constructor for the abstract class is given by the following.

```
// density.hpp
Density(int n_samples,
         int n_bins,
         long double bounds,
         std::shared_ptr<System> pSystem,
         std::shared_ptr<WaveFunction> pWaveFunc,
         std::shared_ptr<Sampler> pSampler,
         bool sample_factors = false);
```

Since obtaining an accurate density requires a large number of samples compared to the number of samples required to calculate the cost, the `Density` object takes its own `n_samples` parameter. While taking these samples, the `OneBodyDensity` object bins them according to their position in one-dimension, or their distance from the origin in two- and three-dimensions. Alternatively, the `TwoBodyDensity` pair distribution functions can be computed by binning the distances between all pairs.

5.1.9 Running NeuralAnsatz

First, clone a copy of the Github repository at this url: <https://github.com/kim-jane/NeuralAnsatz>. Ensure that both `make` and `Open MPI` are installed. On Ubuntu, this can be accomplished by running the following commands in the terminal:

```
sudo apt-get update
sudo apt install build-essential
sudo apt install make
```

```
sudo apt install openmpi-bin
```

On MacOS, one can instead run the following:

```
brew update  
brew upgrade  
brew install gcc  
brew install make  
brew install open-mpi
```

Alternatively, the appropriate modules can be loaded on a high-performance computing cluster. Next, enter repository and compile the `NeuralAnsatz` library.

```
cd ~/NeuralAnsatz  
make
```

Driver files in the `examples/` folder will be compiled into corresponding executables in the `bin/` folder. For example, if there is a driver file named `examples/myprogram.cpp`, the code can be run by typing the command below.

```
mpirun -np <number of parallel tasks> ./bin/myprogram
```

Once completed, remove all binaries and build files.

```
make clean
```

5.2 Python Software for Periodic Systems

In this section, we introduce a different implementation of neural-network quantum states that leverages the flexibility and high-performance computing capability of **JAX**. This code was initially developed by Alessandro Lovato for nuclei and has subsequently been adapted by Bryce Fore and myself for periodic systems. Since some of the over-arching details of the simulation are similar to the `NeuralAnsatz` code, we will focus only on the key differences here. In addition, we will discuss the basic ideas required to implement the neural-network quantum states relevant to the work in Chapters 9 and 7.

5.2.1 JAX: High-Performance Array Computing

JAX is an open-source software that combines the flexibility and ease of use of `numpy` with the efficiency of compiled languages like C++. The key features of **JAX** include:

- Compute gradients with `grad`: Easily compute gradients of arbitrary functions by creating a new function.
- Just-in-time compilation with `jit`: Functions decorated with `jit` will be dynamically compiled and cached upon first encounter.
- Vectorization with `vmap`: Apply functions elementwise to one or more arrays for efficient batch operations.
- Parallelization with `pmap`: Execute functions in parallel across multiple devices, maximizing hardware resources.

5.2.1.1 Building Neural Networks

The functionalities provided by JAX allows for far more freedom in network structure, as one only needs to worry about computing the forward direction. We begin by importing a few required modules.

```
import jax
import jax.numpy as jnp
from jax.example_libraries import stax
```

All layer construction functions within `stax` return a pair of functions (`init_func`, `apply_func`) that represent the initialization of a layer with random weights and the application of a layer, respectively. All initialization functions, `init_func`, take the input shape and a `jax.random.key`, and subsequently returns the new output shape and the initialized parameters. All application functions, `apply_func`, take a set of parameters and input data, then evaluates the forward direction of the layer. Multiply layers can be joined together by using `stax.serial`.

```
init_func, apply_func = stax.serial(*layers)
```

If one defines the list `layers` as an alternating sequence of `stax.Dense` layers and activation layers, e.g. `stax.Gelu`, then `stax.serial` generates the constructor functions of a standard feedforward neural network.

The separation of the `init_func` from the `apply_func` is beneficial if many copies of identical neural network structures are needed, where each copy is initialized with different parameters. For instance, the single-particle orbitals of a Slater determinant can have one common `apply_func`, but N different sets of unique parameters by simply calling `init_func`

N times.

To construct a Deep Set, on the other hand, one must create two sets of constructor functions, one that maps to the latent space and another that maps to the output space.

```
lat_init, lat_apply = stax.serial(stax.Dense(n), stax.Gelu, stax.Dense(n))
out_init, out_apply = stax.serial(stax.Dense(n), stax.Gelu, 1)
```

In the above, the integer `n` is the number of nodes in the layer. Computing the forward direction of the Deep Set amounts to calling `lat_apply` and `out_apply` with an appropriate pooling function in between. Below is an example of a Deep Set used as a Jastrow factor, where `phi_params` and `rho_params` refer to the variational parameters in the mappings to and from the latent space.

```
phi = self.lat_apply(phi_params, X)
phi = jnp.sum(phi, axis=0)
logpsi = self.out_apply(rho_params, phi)
```

Virtually all artificial neural networks, no matter how complicated the connections, can be constructed by appropriately combining feedforward neural networks, concatenation operations, and pooling operations.

5.2.2 Symmetries and Boundary Conditions

Now that we have established the basics of building a network, we must discuss the specific requirements the network must satisfy if it is to be treated as a wave function. If either of these elements are broken, the wave function provide erroneous results. For periodic systems, we must transform the coordinate or distances such that obey strict periodicity at all times.

Some options for periodic inputs with L -periodicity include:

```
# periodic positions
cos_i = jnp.cos( 2 * jnp.pi * r / L)
sin_i = jnp.sin( 2 * jnp.pi * r / L )

# periodic separation vectors
r_ij = r[ ip[k] ] - r[ jp[k] ]
r_ij = r_ij - jnp.rint( r_ij / self.L ) * self.L
cos_ij = jnp.cos( 2 * jnp.pi * r_ij / self.L )
cos_ij = jnp.sin( 2 * jnp.pi * r_ij / self.L )

# periodic distance
d_ij = jnp.sin( jnp.pi * r_ij / self.L )
d_ij = jnp.sum( d_ij**2 )
```

In addition to the periodicity and the antisymmetry of the fermionic wave function, which is constrained by a Slater determinant or a Pfaffian, discrete symmetries such as parity and time-reversal can be enforced. Below is an example of parity and time-reversal enforcement for a system with even parity and odd time-reversal symmetry

```
@jit
def logsignpsi_PT(params, r, sz):

    r_PT = jnp.stack((r, -r, r, -r))
    sz_PT = jnp.stack((sz, sz, -sz, -sz))

    sign, logpsi = vmap(logsignpsi, in_axes=(None, 0, 0))(params, r_PT, sz_PT)
    sign = sign.at[2:].set(-sign[2:])

    logpsi_PT, sign_PT = jax.scipy.special.logsumexp(a=logpsi,
                                                    b=sign,
                                                    return_sign=True)

    return sign_PT, logpsi_PT
```

5.2.3 Computing the Pfaffian

Unfortunately, JAX does not provide built-in implementation of the Pfaffian of a matrix, nor does it provide convenient decomposition schemes useful for the Pfaffian. For our work on ultra-cold Fermi gases, we implemented the Pfaffian ourselves, complete with pivoting. In the code below, we assume A is even dimensional and symmetric. The algorithm is implemented according to Ref. [10].

```
@jit
def pfaffian(self, A):

    def pivot(kp, k, A, pf_A):

        # exchange rows
        temp = A.at[k+ 1, :].get()
        A = A.at[k + 1, :].set(A.at[kp, :].get())
        A = A.at[kp, :].set(temp)

        # exchange columns
        temp = A.at[:, k + 1].get()
        A = A.at[:, k + 1].set(A.at[:, kp].get())
        A = A.at[:, kp].set(temp)

        # flip sign
        pf_A *= -1.0
        return A, pf_A

    def no_pivot(kp, k, A, pf_A):
        return A, pf_A

    pf_A = 1.0
    for k in range(0, A.shape[0]-1, 2):

        # find pivot index
        kp = k + 1 + jnp.argmax(jnp.abs(A[k + 1:, k]))

        # pivot if necessary
        A, pf_A = cond(kp != k + 1, pivot, no_pivot, kp, k, A, pf_A)

        # update pfaffian
        pf_A *= A[k, k + 1]
```

```

# update matrix
mu = A.at[k, :].get() / A[k, k+1]
nu = A.at[:, k+1].get()
A = A + jnp.outer(mu, nu) - jnp.outer(nu, mu)

return jnp.sign(pf_A), jnp.log(jnp.abs(pf_A))

```

5.2.4 Stochastic Reconfiguration with RMSprop Regularization

One of the unique components of the JAX-based VMC code is the alternative regularization scheme used for the Fisher-information matrix. As discussed in Sec. 3.2.2.5, the inversion of the Fisher-information matrix S is typically stabilized by offsetting the diagonal elements by a small ϵ , usually chosen between $10^{-5} - 10^{-2}$. However, by choosing a constant shift, the relative magnitudes of the variational parameters are ignored, ultimately leading to an inefficient optimization process. In the alternative scheme based on RMSprop, a running average of the square of the gradients are stored to normalize the diagonal shift according to each parameter. When the optimizer is initialized, \mathbf{v} is set to zero. Then during each parameter update,

$$\mathbf{v} \leftarrow \beta\mathbf{v} + (1 - \beta)(\nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle)^2, \quad (5.55)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta (S + \text{diag}(\sqrt{\mathbf{v}} + \epsilon))^{-1} \nabla_{\boldsymbol{\theta}} \langle E_{\boldsymbol{\theta}}(\mathbf{X}) \rangle. \quad (5.56)$$

Chapter 6

The Calogero-Sutherland Model

To evaluate the performance of the C++ software discussed in Sec. 5.1, we benchmark against the Calogero-Sutherland model, an exactly-solvable model of interacting bosons trapped in a one-dimensional harmonic oscillator well. This model provides an excellent basis to assess our neural-network quantum states (NQS) for several reasons. Firstly, the interaction contains a singularity. Because the bosons are trapped, the effects of the singularity will be particularly important, as the likelihood of two bosons becoming close to one another is high. Learning how to handle this singularity will prove to be great practice for more realistic systems. Secondly, one-dimensional systems are easy to visualize, giving us the opportunity to gain an intuitive understanding of the two types of neural-network quantum states we will test. Lastly, the availability of exact analytical solutions enables us to validate the accuracy of our implementation.

6.1 Hamiltonian

The Hamiltonian of the Calogero-Sutherland model is given by

$$\hat{H}^{\text{cs}} = \sum_{i=1}^N \left(-\frac{1}{2} \frac{\partial^2}{\partial x_i^2} + \frac{1}{2} \omega^2 x_i^2 \right) + \sum_{i < j}^N \frac{\beta(\beta-1)}{x_{ij}^2}, \quad (6.1)$$

where β is an interaction parameter, $x_{ij} = |x_i - x_j|$ is the distance between particles i and j , and we have set $\hbar = m = 1$ for convenience. During our discussion on Kato's cusp condition in Sec. 3.2.1.1, we showed that requiring the local energy to be finite as any two particles approach each other led us to the exact ground state of the system

$$\Psi_0^{\text{cs}}(\mathbf{X}) = \exp\left(-\sum_{i=1}^N \frac{1}{2}\omega x_i^2\right) \prod_{j < k} x_{ij}^\beta, \quad (6.2)$$

with the corresponding energy

$$E_0^{\text{cs}} = \frac{1}{2}N\omega + \frac{1}{2}\beta N(N-1)\omega = \frac{1}{2}N\omega(1 + \beta(N-1)). \quad (6.3)$$

Throughout this work, we will set $N = 6$ and $\beta = 2$, resulting in an exact interacting energy that is 11 times greater than the noninteracting harmonic oscillator ground state energy.

To prevent numerical instabilities caused by the singularity in Eq. 6.1, we regularize the two-body potential as

$$\frac{1}{x^2} \mapsto \frac{\tanh^2(x/x_0)}{x^2}, \quad (6.4)$$

where x_0 is a regularization hyperparameter that smoothly controls the height of the potential, while preserving the long-range behavior. As $x_0 \rightarrow 0$, the regularized potential converges to the original potential containing the singularity. See Fig. 6.1 to compare the regularized potential using different values of x_0 and $\beta = 2$.

Our general strategy is to first train the neural-network quantum states using a fairly large value of the regularization parameter x_0 . Then after the training has stabilized, we progressively decrease the value of x_0 until the energy converges. Another possibility could be to fix x_0 to a small value, then gradually increase β from a value slightly greater than 1

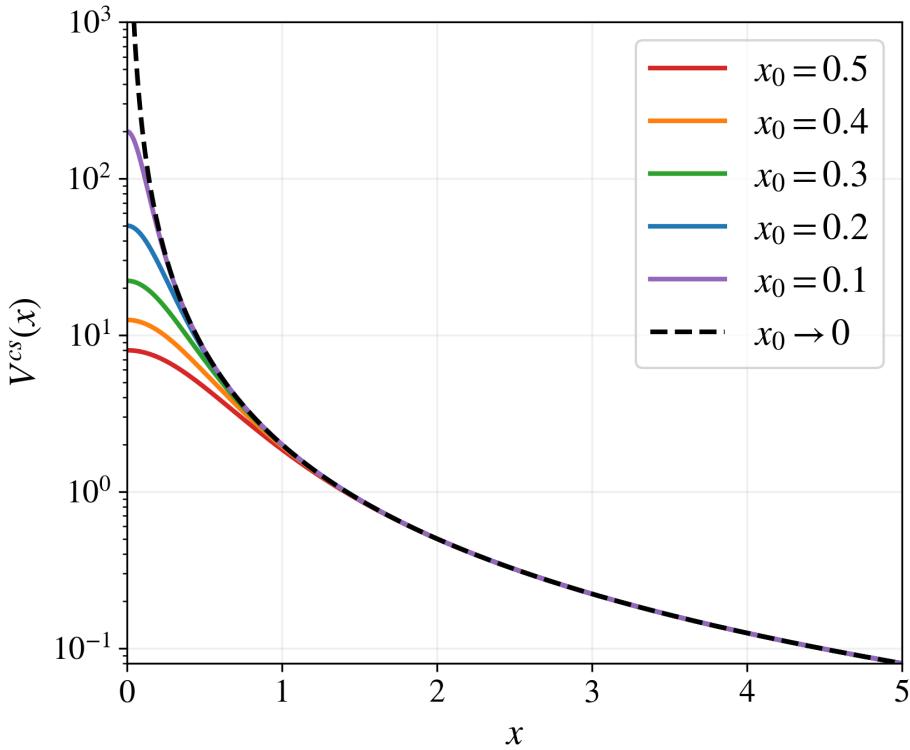


Figure 6.1: The regularized two-body potential for the Calogero-Sutherland model, where x_0 is the regularization parameter (Eq. 6.4). The long-range behavior of the original potential (black, dashed line) is preserved, but the height of the potential decreases with larger values of x_0 . Here, we have used an interaction parameter value of $\beta = 2$.

to our desired value of 2. Both of these strategies are examples of transfer learning, as they involve pretraining the NQS on easier tasks before transitioning to more challenging ones.

6.2 Neural-Network Quantum States

This work focuses on a comparative analysis of two distinct types of neural-network quantum states: feedforward neural networks and continuous restricted Boltzmann machines. The inputs for both networks comprise the positions of all particles, with permutation invariance guaranteed by sorting the positions. Sorting the inputs was discussed in depth in Sec. 5.1.3.1. For brevity, we will use the notation $\mathbf{v}(\mathbf{X})$ to imply the positions $\mathbf{X} = (x_1, x_2, \dots, x_N)$ have

been sorted and serve as the visible nodes \mathbf{v} of the networks.

6.2.1 Feedforward Neural Networks

Inspired by Saito in Ref. [11], we take our positive-definite Jastrow ansatz to be

$$\Psi_{\theta}^{\text{fnn}}(\mathbf{X}) = \exp(\rho(\mathbf{v}(\mathbf{X}))) \exp\left(-\frac{\mathbf{X}^2}{2\sigma^2}\right), \quad (6.5)$$

where $\rho(\mathbf{v}(\mathbf{X}))$ is the single output of a feedforward neural network, which takes sorted positions as input. Since neural networks are initialized with small random values around zero, ρ initially takes a value around zero as well. To ensure we sample from a localized distribution rather than a uniform distribution of infinite range, our ansatz includes a Gaussian envelope with a variance σ^2 large enough to not significantly affect the training of the feedforward neural network. In this work with $N = 6$, we have set $\sigma^2 = 16^2 = 256$. Alternatively, one could remove the Gaussian envelope and pretrain the ansatz to a Gaussian instead, either by turning off the interaction entirely, leaving a harmonic oscillator potential, or through supervised training of the wave function. We opt for the Gaussian envelope as it removes the need for pretraining, but the boundary conditions imposed by the alternative approaches are generally upheld by the network during the reinforcement learning process.

6.2.2 Continuous Restricted Boltzmann Machines

We also employ neural-network quantum states parameterized by the various types of multivariate Gaussian-binary and Gaussian-uniform restricted Boltzmann machine, discussed in Sec. 4.7.1.2 and Sec. 4.7.1.3. Since they are all positive-definite by default, we can simply

write

$$\Psi_{\boldsymbol{\theta}}^{\text{rbm}}(\mathbf{X}) = \mathcal{P}(\mathbf{v}(\mathbf{X})), \quad (6.6)$$

where $\mathcal{P}(\mathbf{v}(\mathbf{X}))$ denotes the marginal probability distribution of the visible nodes, with $\mathbf{v}(\mathbf{X})$ representing the sorted positions. When the variational parameters are set to small random values, the initial state of the RBMs are similar to a Gaussian. We can bias our initial state to more closely resemble the eventual wave function we expect by initializing the visible biases \mathbf{a} across a much wider range. Moreover, we sort the visible biases in order to bring them closer to the mean of the corresponding sorted positions.

6.3 Results

We begin our analysis by testing four commonly employed activation functions for a fixed feedforward neural network (FNN) architecture. For $N = 6$, we choose to use a network with two hidden layers and 12 hidden nodes each. The regularization parameter is initially set to $x_0 = 0.5$ in order to guarantee fast and stable training. Additionally, we utilize the stochastic reconfiguration (SR) algorithm, with a diagonal offset of $\epsilon = 0.001$, which accelerates training significantly. As the SR algorithm manipulates the gradient such that the curvature of the energy landscape is encoded, we have found that a simple stochastic gradient descent optimizer with a constant learning rate of $\eta = 0.001$ provides the most stable training.

In Fig. 6.2, we validate the efficient training of the feedforward neural networks. For each activation function, about 200 iterations were sufficient to reach convergence. Notice that the converged value of the energy is well below 0, an artifact of the regularization of the potential. At this initial stage, the performances of all activation functions are essentially

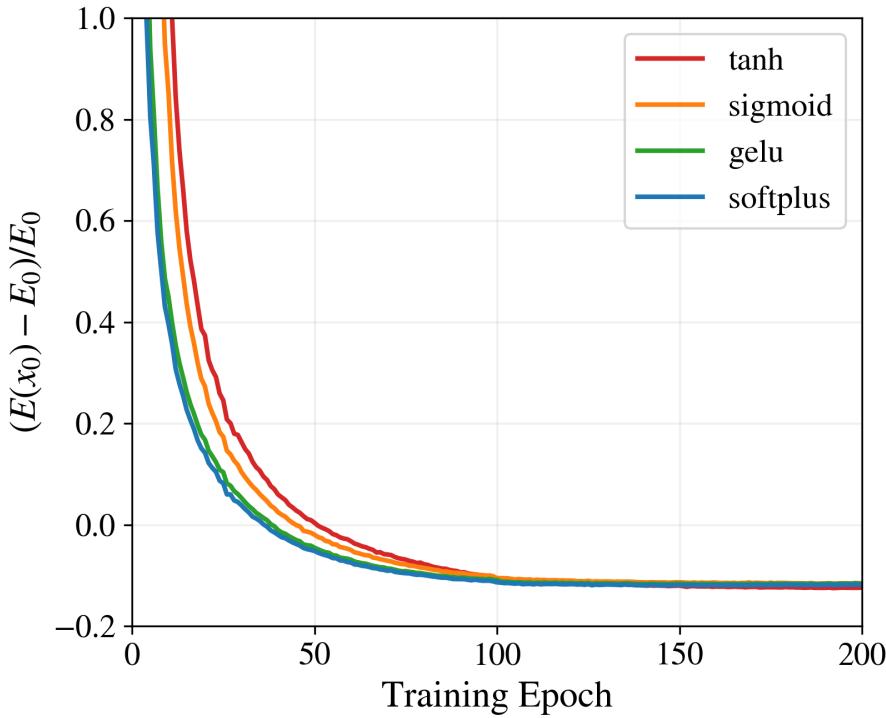


Figure 6.2: Training curves for feedforward neural networks with different activation functions. In this initial pretraining stage, $x_0 = 0.5$.

equal, with the sigmoid activation function yielding an energy slightly greater hyperbolic tangent, the Gaussian error linear unit (gelu), and softplus.

We also validate the training of the Gaussian-binary and Gaussian-uniform restricted Boltzmann machines, where we consider both non-symmetric and symmetric hidden node values for each type. Since the computed marginal probability distribution of the visible nodes for each RBM type only differed by the definition of a function we simply called " $f(x)$," one could interpret these functions as activation functions. In Fig. 6.3, we show the training curves for all four variants, each with 6 hidden nodes. Similar to the FNN case, we use the SR algorithm with $\epsilon = 0.001$ to compute the preconditioned gradient. However, for the RBMs we use a smaller learning rate of $\eta = 0.0005$.

While the RBMs yield an initial energy that is much closer to the eventual converged value, training still required about 3000 iterations, an order of magnitude more than the FNNs. The converged value of the energy still did not reach quite as low as the FNNs, reaching a minimum percent difference around $(E(x_0) - E_0)/E_0 = -0.0941$ for all RBM types. Here, E_0 is the exact ground state energy in the limit $x_0 \rightarrow 0$. In contrast, the FNNs reached a minimum between $(E(x_0) - E_0)/E_0 = -0.113 \sim -0.119$. In addition, notice that the training curves for the RBMs with symmetrically-valued hidden nodes (shown in orange and blue) show slight deviations from the nonsymmetrically-valued hidden nodes (red and green) after about ~ 250 training iterations. This may suggest that both needed to overcome a local minimum before proceeding further to the ground state. During this initial pretraining stage, the multivariate Gaussian-Uniform restricted Boltzmann machine with hidden node values between 0 and 1 had the most stable training curve.

In order to better predict the ground state of the full Calogero-Sutherland Hamiltonian, we progressively decrease the value of the regularization parameter x_0 from a value of 0.5 to 0.01. In Figures 6.4 and 6.5, we show the converged energies as a function of x_0 , both the FNNs and RBMs, respectively. For the FNNs, most of the activation functions exhibit similar behavior as x_0 is decreased, except for the sigmoid activation function. This difference in performance is somewhat expected, as the sigmoid activation function is known to be prone to the vanishing-gradient problem. The other three activation functions perform similarly as x_0 is decreased, with the softplus function consistently producing the lowest energies. The tanh, gelu, and softplus functions all seem to converge to the exact ground state as expected.

On the other hand, the different restricted Boltzmann machines exhibit almost identical behavior as x_0 is decreased, with slight deviations occurring for the smallest values of x_0 . In particular, both of the Gaussian-uniform RBMs perform slightly better than the Gaussian-

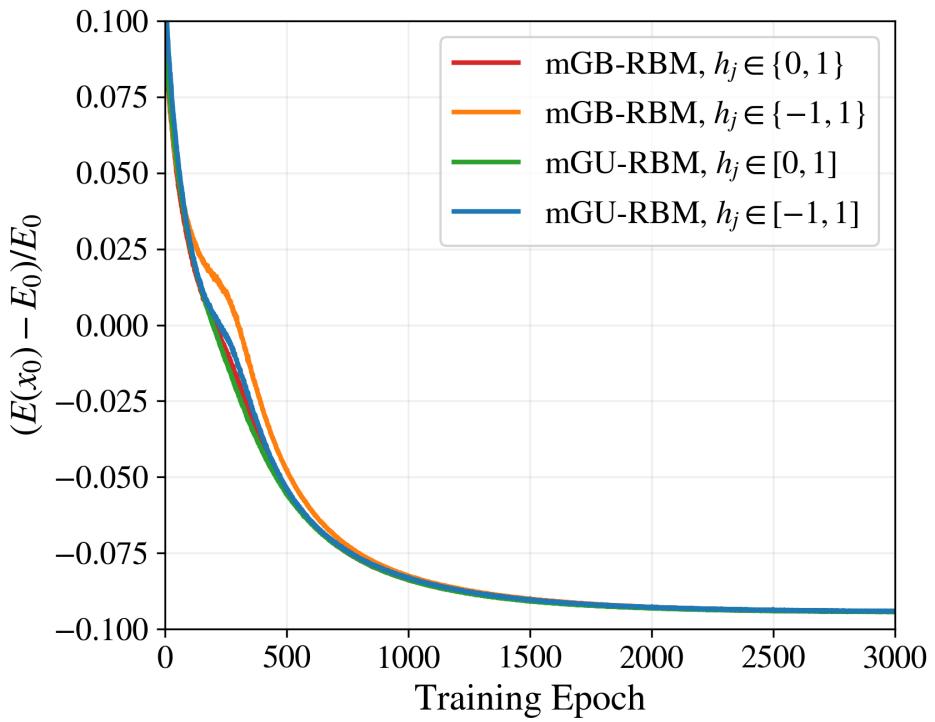


Figure 6.3: Training curves for the different variants of Gaussian-binary and Gaussian-uniform restricted Boltzmann machines, with $x_0 = 0.5$.

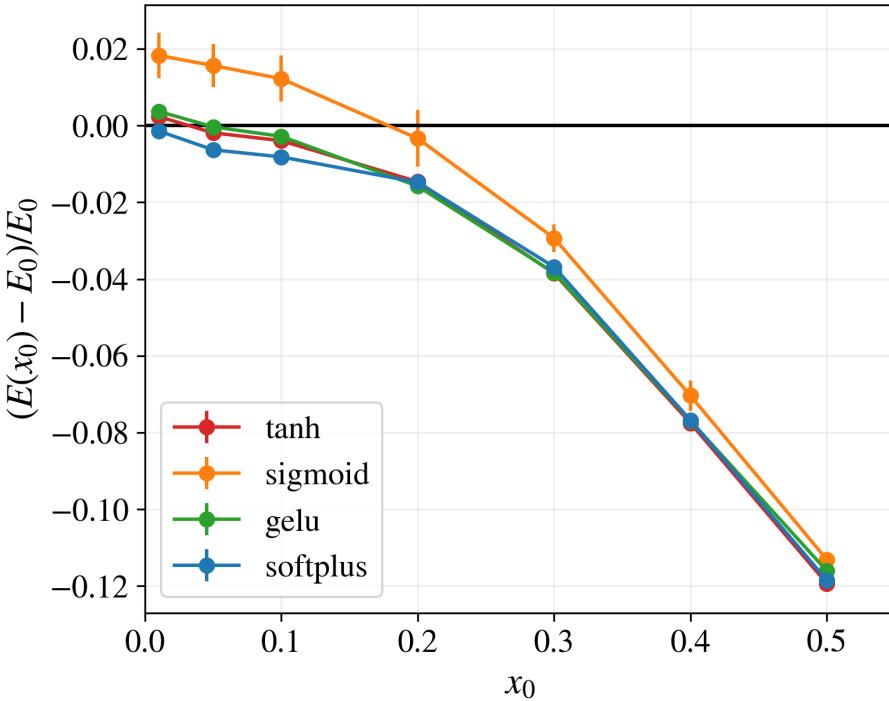


Figure 6.4: The converged energy as a function of the regularization parameter x_0 for feed-forward neural networks with different activation functions.

binary RBMs for $x_0 = 0.01$. However, unlike the FNNs which usually converge to the exact ground state

Finally, in Fig. 6.6, we plot the one-body densities for the symmetric multivariate Gaussian-Uniform restricted Boltzmann machine and the feedforward neural network with the soft-plus activation function. The top row shows the initial state for both networks. The RBM naturally provides localized initial state with some distinguished peaks, demonstrating the highly-biased nature of the ansatz. As x_0 is decreased, the location of the peaks match well with the exact ground state, but there are small deviations in the shape of the peaks and in the tails of the distribution. The FNN is a low-bias model, and consequently, we used a Gaussian envelope to localize the distribution around 0. After the initial pretraining stage

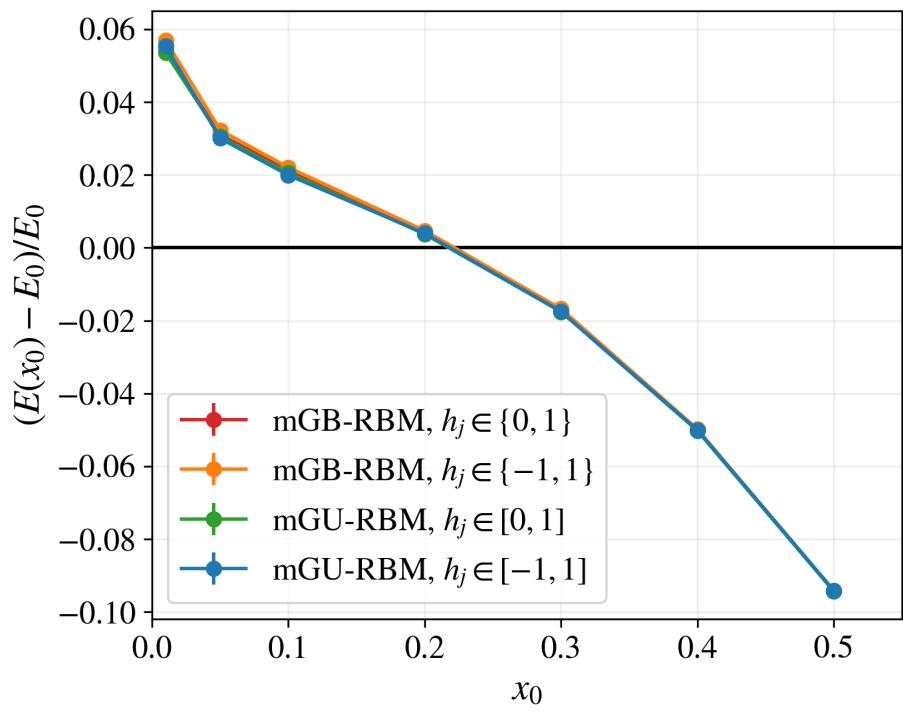


Figure 6.5: The converged energy as a function of the regularization parameter x_0 for the various Gaussian-binary and Gaussian-uniform restricted Boltzmann machines.

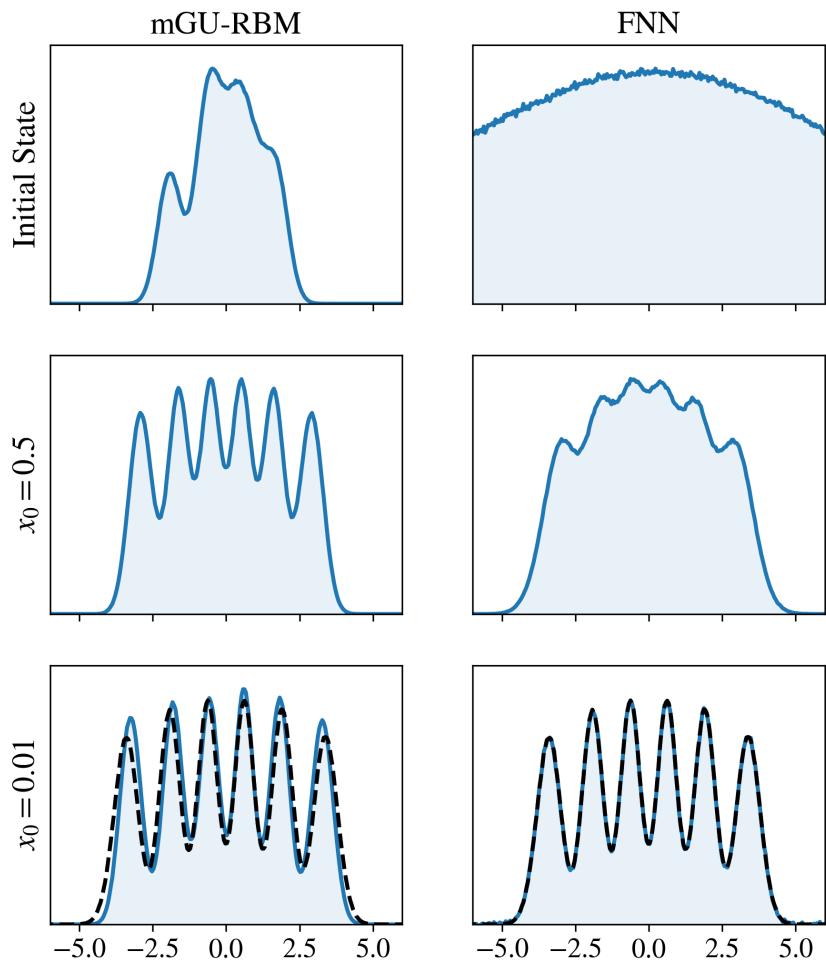


Figure 6.6: One-body densities during three different points during training: the initial stage, after pretraining with $x_0 = 0.5$, and after training with $x_0 = 0.01$. The black dashed line represents the one-body distribution from the exact ground state wave function.

using $x_0 = 0.5$ the wave function learned by the FNN has significantly different shape to the one learned by the RBM, as the peaks are far less distinguished in the former. After decreasing the regularization parameter to $x_0 = 0.05$, the FNN ansatz matches the exact ground state almost perfectly.

6.4 Conclusions

In this illustrative example of neural-network quantum states, we demonstrate a transfer learning technique beneficial for handling singularities in the potential. By comparing different activation functions, we confirm that the sigmoid function suffers from a vanishing gradient, but nevertheless, still performs better than all the RBMs tested near the $x_0 \rightarrow 0$ limit.

The RBMs perform similarly to one another, with the Gaussian-uniform variants performing slightly better as x_0 decreases. However, due to the high-biased nature of the model, RBMs are limited in being able to accurately capture all behaviors in the wave function. We visualize this fact by investigating the evolution of the one-body distributions throughout the training process. Overall, a feedforward neural network proves to be the more accurate model and is therefore recommended for use as a neural-network quantum state, even if it requires comparatively more variational parameters.

Initial investigations of products of RBMs and FNNs show promise, as this eliminates the need for a Gaussian envelope and reduced the overall number of parameters required in the FNN. We leave a more thorough report on these products for future works.

Chapter 7

Dilute Neutron Matter

Low-density neutron matter is characterized by fascinating emergent quantum phenomena, such as the formation of Cooper pairs and the onset of superfluidity. We model this density regime by capitalizing on the expressivity of the hidden-nucleon neural-network quantum states combined with variational Monte Carlo and stochastic reconfiguration techniques. Our approach is competitive with the auxiliary-field diffusion Monte Carlo method at a fraction of the computational cost. Using a leading-order pionless effective field theory Hamiltonian, we compute the energy per particle of infinite neutron matter and compare it with those obtained from highly realistic interactions. In addition, a comparison between the spin-singlet and triplet two-body distribution functions indicates the emergence of pairing in the 1S_0 channel.

The following paper titled '**Dilute neutron star matter from neural-network quantum states**' represents a collaborative effort involving Bryce Fore¹, myself², Giuseppe Carleo³, Morten Hjorth-Jensen^{2,4}, Alessandro Lovato^{5,6,7}, and Maria Piarulli^{8,9}.

¹*Physics Division, Argonne National Laboratory, Argonne, IL 60439, USA*

²*Department of Physics and Astronomy and Facility for Rare Isotope Beams, Michigan State University, East Lansing, MI 48824, USA*

³*Institute of Physics, École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland*

⁴*Department of Physics and Center for Computing in Science Education, University of Oslo, N-0316 Oslo, Norway*

⁵*Physics Division, Argonne National Laboratory, Argonne, IL 60439, USA*

⁶*Computational Science Division, Argonne National Laboratory, Argonne, IL 60439, USA*

⁷*INFN-TIFPA Trento Institute for Fundamental Physics and Applications, 38123 Trento, Italy*

⁸*Physics Department, Washington University, St Louis, MO 63130, USA*

⁹*McDonnell Center for the Space Sciences at Washington University in St. Louis, MO 63130, USA*

7.1 Introduction

Multi-messenger astronomy has opened new windows into the state of matter at densities and isospin asymmetries that cannot be directly probed by terrestrial experiments [12, 13, 14, 15]. Concurrently, nuclear many-body theory has made considerable progress in computing the nucleonic-matter equation of state at densities corresponding to the inner core of neutron stars starting from realistic Hamiltonians [16, 17, 18, 19, 20, 20, 21]. Comparisons between theoretical predictions and astrophysical observation pose stringent constraints on models of nuclear dynamics, particularly three-nucleon forces [22].

In this work, we focus on lower densities, $\rho \lesssim 0.04 \text{ fm}^{-3}$, which are relevant to the phenomenology of the stellar inner crust and outer core. In this region, both conditions

for superfluidity — strong Fermi degeneracy and an attractive interaction between neutron pairs in the 1S_0 channel — are believed to be met [23, 24, 25]. In addition to lowering the system’s energy, the formation of Cooper pairs plays a critical role in neutrino emission [26, 27], and the phenomenology of glitches [28]. Pairing is also relevant in modeling neutron-rich nuclei, which are the subject of intense experimental activities [29].

Quantum Monte Carlo approaches [30], and in particular the auxiliary-field diffusion Monte Carlo (AFDMC) method [31] have been extensively applied to accurately compute neutron-matter properties [18, 17, 32]. In the low-density regime, AFDMC calculations have convincingly shown a depletion of the superfluid gap with respect to the Bardeen–Cooper–Schrieffer theory [33, 34]. However, because of the fermion sign problem, AFDMC predictions depend upon the starting variational wave function. For instance, the superfluid phase must be assumed *a priori* by using pfaffian wave functions [35].

Neural-network quantum states [1] (NQS) have gained popularity in solving the Schrödinger equation of atomic nuclei both in real space [36, 37, 38, 39, 40] and in the occupation-number formalism [41]. In this work, we introduce a periodic NQS suitable to model both the normal and superfluid phases of neutron matter. The ansatz is based on the “hidden-nucleon” architecture, which can model the ground-state wave functions of nuclei up to ${}^{16}\text{O}$ with high accuracy [39]. Inspired by chemistry applications [42, 43], we further improve the expressivity of the hidden-nucleon NQS using generalized backflow correlations, which generalize both the pfaffian and the spin-dependent backflow of Ref. [44].

Our model of nuclear dynamics is the leading-order pionless effective field theory ($\not\!\text{EFT}$) Hamiltonian of Ref. [45], which qualitatively reproduces the binding energies of nuclei with up to $A = 90$ nucleons. Arguments based on the expansion around the unitary limit [46], and Brueckner-Hartree-Fock calculations of infinite nuclear matter [47], indicate that $\not\!\text{EFT}$

should provide accurate energies of dilute neutron matter. We quantitatively address this point by carrying out variational Monte Carlo calculations based on NQS that are specifically designed to model wavefunctions of nuclear systems in the presence of spatial periodicity. We compare the $\not{\pi}$ EFT energy per particle against the phenomenological Argonne v_{18} [48] plus Urbana IX [49] (AV18+UIX) Hamiltonian used in the Akmal-Pandharipande-Ravenhall (APR) [50] equation of state. We additionally consider the local, Δ -full chiral-EFT potentials that include tritium β -decay in the fitting procedure and do not yield self-bound neutron matter [17, 32].

To better quantify the role of dynamical correlations, we evaluate the two-body spatial distribution functions, separating the spin-triplet and spin-singlet channels. We analyze the self-emergence of pairing correlations, not explicitly included in the NQS ansatz, as a function of neutron-matter density.

7.2 Method

We model the interactions among neutrons through the leading-order $\not{\pi}$ EFT Hamiltonian “o” of Ref. [45]. The two-body contact potential is designed to reproduce the np scattering lengths and effective ranges in the $S/T = 0/1$ and $1/0$ channels. Thus, it yields a neutron-neutron scattering length of $a_{nn} = -22.5$ fm, slightly larger than the experimental value of $-18.9(4)$ fm, see [51] and references therein, while the effective range is well reproduced. The Hamiltonian also contains a repulsive three-body force that ensures the stability of nuclei with $A \leq 3$ nucleons. As for the latter, we take the parameterization with $R_3 = 1.0$ fm, as it reproduces the binding energies of nuclei in the $A \leq 90$ mass range reasonably well. For benchmarking purposes, we also consider the leading-order $\not{\pi}$ EFT Hamiltonian “a” with

$R_3 = 1.0$ fm, which describes the trend of binding energies of light and medium-mass nuclei.

We approximate the ground-state solution of the nuclear many-body problem with an NQS ansatz that belongs to the hidden-fermion family [52], recently generalized to continuum Hilbert spaces and applied to atomic nuclei in Ref. [39]. In addition to the visible spatial and spin coordinates of the A neutrons, $R = \{\mathbf{r}_1 \dots \mathbf{r}_A\}$ and $S = \{s_1^z \dots s_A^z\}$, the Hilbert space contains fictitious A_h hidden-nucleon degrees of freedom. In this work we use $A_h = A = 14$ so that the system is as flexible as possible, but in practice we have also found using as few as 8 hidden nucleons gives very similar results. The wave function can be conveniently expressed in a block matrix form as

$$\Psi_{HN}(R, S) \equiv \det \begin{bmatrix} \phi_v(R, S) & \phi_v(R_h, S_h) \\ \chi_h(R, S) & \chi_h(R_h, S_h) \end{bmatrix}. \quad (7.1)$$

As in Ref. [39], $\phi_v(R, S)$ is the $A \times A$ matrix representing visible single-particle orbitals computed on the visible coordinates while the $A_h \times A_h$ matrix $\chi_h(R_h, S_h)$ yields the amplitudes of hidden orbitals evaluated on the coordinates of the A_h hidden nucleons. Finally, $\chi_h(R, S)$ and $\phi_v(R_h, S_h)$ are $A_h \times A$ and $A \times A_h$ matrices giving the amplitudes of hidden orbitals on visible coordinates and visible orbitals on hidden coordinates, respectively. All the above matrices are expressed in terms of deep neural networks with differentiable activation functions — see Ref. [39] for additional details. To respect the Pauli principle, the coordinates of the hidden nucleons must be permutation-invariant functions of the visible ones. We enforce this symmetry by using a Deep-Sets architecture [9, 53] with *logsumexp* pooling. Additionally, the discrete parity and time reversal symmetries, are enforced in the same manner as Ref. [39].

Inspired by recent developments in quantum-chemistry NQS [54, 42, 43], we augment

the flexibility of the ansatz by performing a generalized backflow transformation to the visible coordinates of the hidden-nucleon matrix: $(R, S) \rightarrow (\tilde{R}, \tilde{S})$. We use the Deep-Sets architecture again to enforce fermion anti-symmetry

$$(\tilde{\mathbf{r}}_i, \tilde{s}_i^z) = \rho_{\text{bf}} \left(\mathbf{r}_i, s_i^z, \log \left(\sum_j \exp(\phi_{\text{bf}}(\mathbf{r}_j, s_j^z)) \right) \right). \quad (7.2)$$

To further augment the expressivity, separate ρ_{bf} and ϕ_{bf} neural networks are used for each of the A visible orbitals.

We simulate infinite neutron matter using 14 particles in a box with periodic boundary conditions. Following Ref. [55], the latter are imposed by mapping the spatial coordinates onto periodic functions by

$$\mathbf{r}_i \rightarrow \left(\sin \left(\frac{2\pi \mathbf{r}_i}{L} \right), \cos \left(\frac{2\pi \mathbf{r}_i}{L} \right) \right) \quad (7.3)$$

which ensures the wave function is continuous and differentiable at the box boundary. Here L is the size of the simulation periodic box, and the sin and cos functions are applied element-wise to \mathbf{r}_i . Finite-size effects due to the tail corrections of two- and three-body potentials are accounted for by summing the contributions given by neighboring cells to the simulation box [56].

Evaluating the expectation values of quantum mechanical operators, including the Hamiltonian, requires carrying out multi-dimensional integration over the spatial and spin coordinates of the neutrons. To this aim, we exploit Monte Carlo quadrature and sample R and S from $|\Psi_{HN}(R, S)|^2$ using the Metropolis-Hastings algorithm [57] — additional details can be found in the supplemental material of Ref. [37]. The best variational parameters

defining the NQS are found by minimizing the system’s energy, which we carry out using the R(oot)M(ean)S(quared)Prop(agation)-enhanced version of the stochastic-reconfiguration optimization method introduced in Ref. [39].

7.3 Results and Discussion

We first benchmark the expressivity of the hidden-nucleon NQS for periodic systems by comparing the energy per particle of infinite neutron matter against “conventional” variational Monte Carlo (VMC), and both constrained-path and AFDMC results. The variational wave function used in state-of-the-art neutron-matter studies, see for example [18, 32], contains a spin-independent Jastrow factor that multiplies a Slater determinant augmented by spin-dependent backflow correlations. The constrained-path approximation, commonly employed to alleviate the AFDMC fermion-sign problem [30], brings about a bias in the ground-state energy estimate [17, 32]. Exact results can be obtained by performing unconstrained propagations, but the statistical error grows exponentially with the imaginary time.

As shown in Fig. 7.1 for $\rho = 0.04 \text{ fm}^{-3}$, after $\simeq 2000$ stochastic-reconfiguration steps, the NQS ansatz converges to the *virtually exact* unconstrained AFDMC energy, using a fraction of its computing time: about 100 hours on NVIDIA-A100 GPUs vs approximately 1.2 million hours on Intel-KNL CPUs. Notice that the constrained-path approximation violates the variational principle. In contrast, variational Monte Carlo calculations based on the NQS never yield energies below that of the Hamiltonian’s ground state. Comparing with the Hartree-Fock approximation, it appears that the hidden-nucleon ansatz captures the overwhelming majority of the correlation energy.

In Fig. 7.2, we compare the π EFT energies obtained with the NQS ansatz against

AFDMC calculations of 14 particles with periodic-box boundary conditions, so that finite-size effects are the same in both approaches. The AFDMC takes as input the AV18+UIX Hamiltonian used in the celebrated APR equation of state [50], and the local, Δ -full chiral-EFT potentials that include tritium β -decay in the fitting procedure and do not make neutron-matter collapse, i.e. models NV2+3-Ia*/b*, and NV2+3-IIb* [17, 32]. Since for all the densities we consider, the AV18+UIX, NV2+3-Ia*/b*, and NV2+3-IIb* are in excellent agreement, they are collectively displayed by the “ π -full” band, stressing that all these interactions explicitly retain pion-exchange terms.

The $\not\pi$ EFT Hamiltonian “o” is in excellent agreement with the π -full models — both providing energies much below the non-interacting Fermi gas (not shown in the Figure). These minor differences are likely because model “o” yields a slightly larger nn scattering length than the experimental value and, therefore, more attraction in neutron matter. For benchmark purposes, we also consider the $\not\pi$ EFT model “a” of Ref. [45], which provides a slightly stiffer equation of state than model “o”. By checking the individual expectation value of the two and three-body potentials, we find that this behavior is primarily due to the three-body force contribution that is more repulsive in model “a” than model “o”, which arises from a more bound ^3H when the two-body force alone is employed.

Once trained on the system’s energy, the NQS can be used to accurately evaluate a variety of quantum-mechanical observables, such as the spin-singlet and triplet two-body distribution functions defined in Ref. [58]. Figure 7.3 shows these distributions at $\rho = 0.01 \text{ fm}^{-3}$ (panel a) and $\rho = 0.04 \text{ fm}^{-3}$ (panel b). The significant increase in the spin-singlet channel compared to the non-interacting Fermi Gas indicates that the NQS wave function can capture the emergence of the 1S_0 neutron pairing, despite not being explicitly encoded in the ansatz. Consistent with the behavior of the pairing gap [33, 24], the enhancement

is more prominent at $\rho = 0.01 \text{ fm}^{-3}$ than $\rho = 0.04 \text{ fm}^{-3}$. On the other hand, at these densities, no pairing correlations are present in the spin-triplet channel.

7.4 Conclusion

In this work, we have put forward an NQS suitable to model the normal and superfluid phases of infinite neutron matter in a unified fashion. We improve the expressivity of the hidden-nucleon ansatz of Ref. [39] by adding state-dependent generalized backflow correlations, whose inclusion has proven beneficial in condensed-matter applications [42, 43]. Periodic-box boundary conditions are imposed by mapping the spatial coordinates of the neutrons onto periodic functions.

Combined with Monte Carlo techniques to sample the Hilbert space and the stochastic-reconfiguration algorithm to optimize the variational parameters, the NQS yields energies per particle of low-density neutron matter that are in excellent agreement with unconstrained AFDMC calculations at a fraction of the computational cost. In contrast, the computationally-inexpensive AFDMC constrained-path approximation brings about appreciable violations of the variational principle.

We have shown that π EFT yields a low-density neutron matter equation of state that is remarkably close to those obtained from AFDMC calculations that take as input highly-realistic phenomenological and chiral-EFT interactions [50, 17, 32]. This finding paves the way for more systematic comparisons between dilute neutron matter and Fermi gas around the unitary limit. In addition, it enables studies of phenomena relevant to understand the inner crust and the outer core of neutron stars, such as pairing and superfluidity, using relatively simple models of nuclear dynamics.

Finally, we have analyzed the possible onset of Cooper pairing in the neutron medium. Specifically, the NQS two-body distribution functions corresponding to pairs of neutrons in the spin-singlet 1S_0 channel exhibit a clear enhancement at small inter-particle distances with respect to the non-interacting case, which is absent in the spin-triplet channel. Consistent with pairing-gap calculations [33, 24, 34], this behavior is more prominent at smaller densities. Note that this feature has not been encoded in the NQS; rather, it is a self-emerging quantum mechanical phenomenon.

As a future development, we plan on including more sophisticated interactions, including highly-realistic phenomenological ones, including AV18+UIX and the local, Δ -full chiral-EFT potentials of Ref. [59, 17, 32]. The flexibility of the NQS ansatz will also be tested in isospin-asymmetric nucleonic matter at low densities, where strong clustering is expected to occur [60].

7.5 Acknowledgements

We are grateful to O. Benhar, S. Gandolfi, A. Kievsky, and B. Wiringa for many illuminating discussions. A.L. and B.F. are supported by the U.S. Department of Energy, Office of Science, Office of Nuclear Physics, under contracts DE-AC02-06CH11357, by the 2020 DOE Early Career Award number ANL PRJ1008597, by the NUCLEI SciDAC program, and Argonne LDRD awards. M.P. is supported by the 2021 Early Career Award number DE-SC0022002 and the FRIB Theory Alliance award DE-SC0013617. J.M.K. and M.H.-J. are supported by the U.S. National Science Foundation Grants No. PHY-1404159 and PHY-2013047. Numerical calculations were performed using resources of the Laboratory Computing Resource Center at Argonne National Laboratory, and the computers of the Ar-

gonne Leadership Computing Facility via the ALCC grant “Short Range Correlations from a Quantum Monte Carlo perspective.”

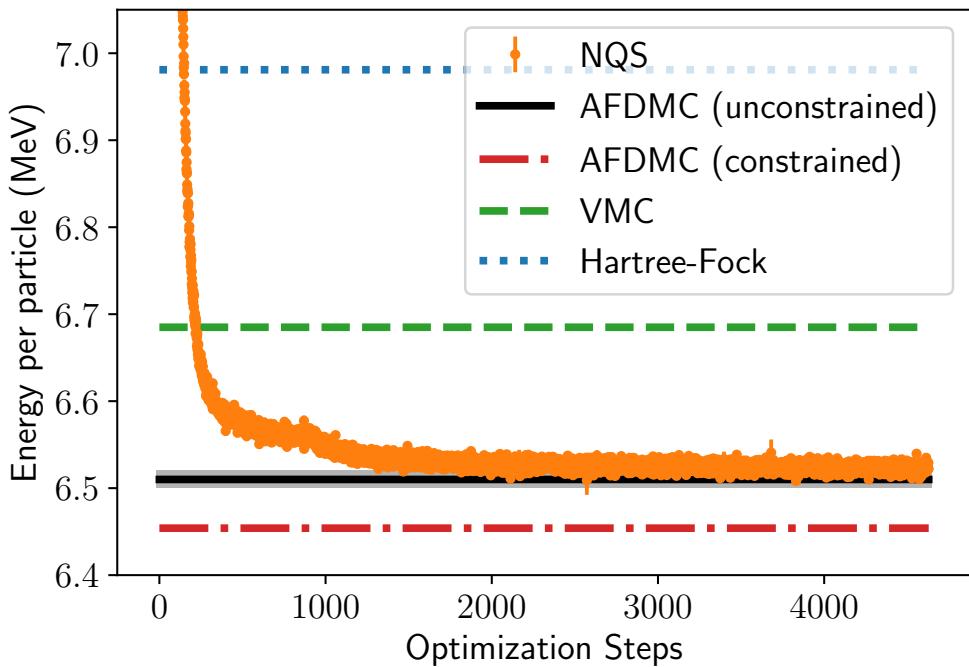


Figure 7.1: NQS training data in neutron matter at $\rho = 0.04 \text{ fm}^{-3}$ (data points) compared with Hartree-Fock (dotted line), conventional VMC (dashed line), constrained-path ADMC (dash-dotted line) and unconstrained-path ADMC results (solid line).

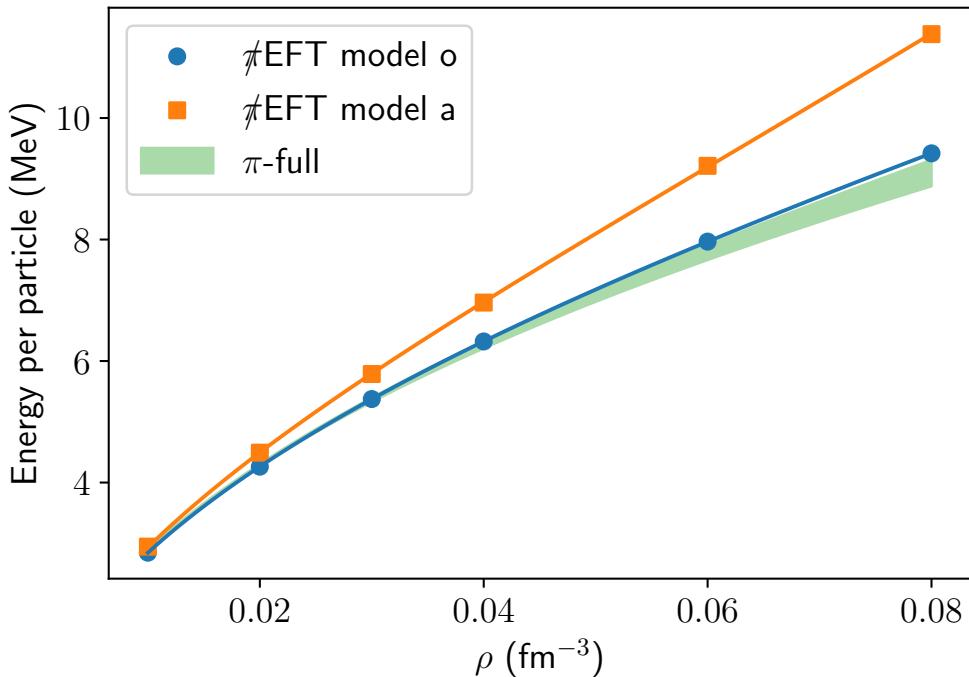


Figure 7.2: Low-density neutron-matter $\not{\pi}\text{EFT}$ equation of state as obtained with the hidden-nucleon NQS for $\not{\pi}\text{EFT}$ potential "o" (blue circles) and $\not{\pi}\text{EFT}$ potential "a" (orange squares) compared with interactions which retain pion-exchange terms (green band). We see that the "o" potential is in excellent agreement with the π -full interactions while the "a" potential has a slightly stiffer equations of state due primarily to a more repulsive three-body force.

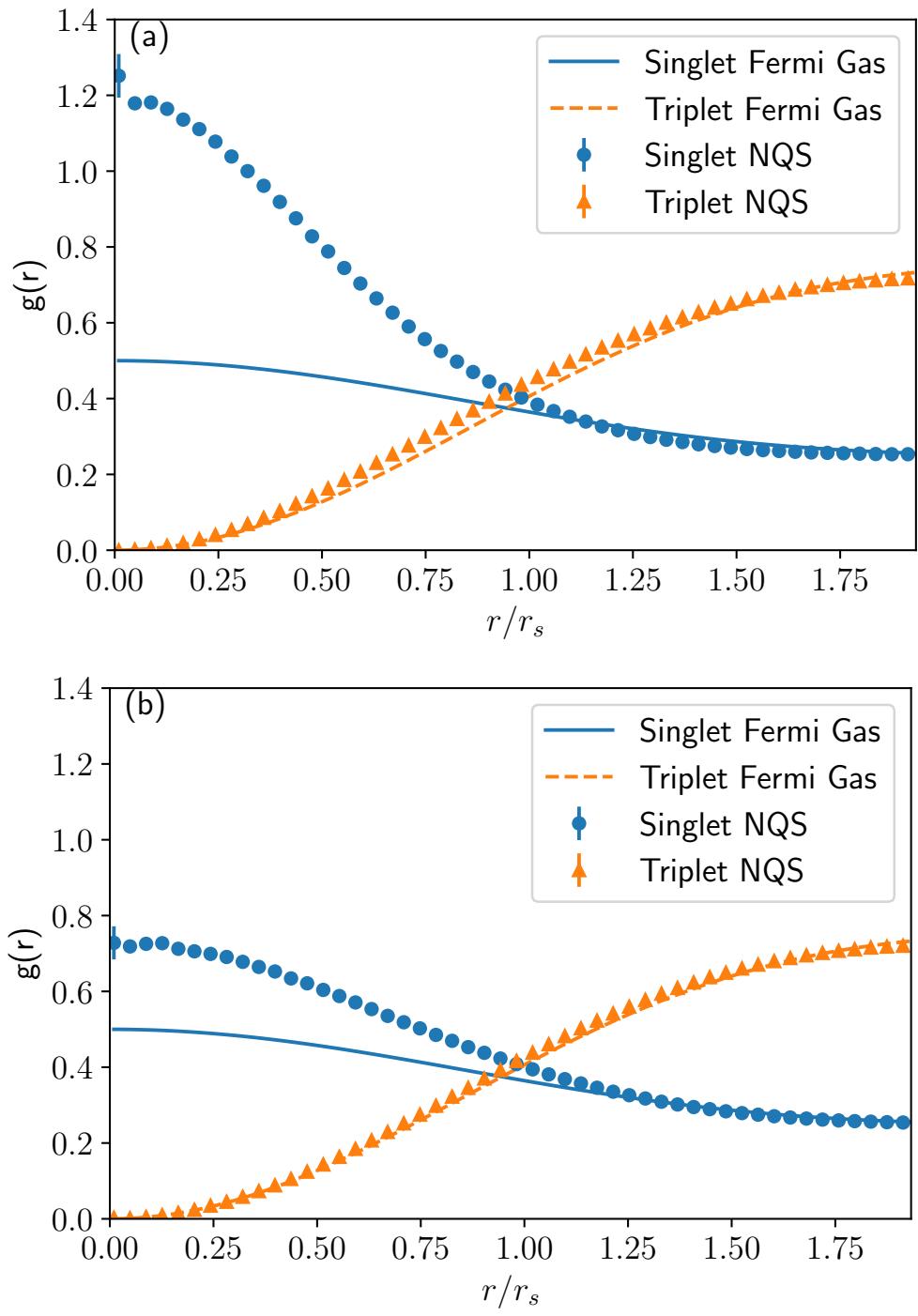


Figure 7.3: Spin-singlet and triplet two-body distribution functions at $\rho = 0.01 \text{ fm}^{-3}$ (panel a) and $\rho = 0.04 \text{ fm}^{-3}$ (panel b). The NQS calculations (solid symbols) are compared with non-interacting Fermi Gas results (solid lines).

Chapter 8

Homogeneous Electron Gas

We introduce a message-passing-neural-network-based wave function Ansatz to simulate extended, strongly interacting fermions in continuous space. Symmetry constraints, such as continuous translation symmetries, can be readily embedded in the model. We demonstrate its accuracy by simulating the ground state of the homogeneous electron gas in three spatial dimensions at different densities and system sizes. With orders of magnitude fewer parameters than state-of-the-art neural-network wave functions, we demonstrate better or comparable ground-state energies. Reducing the parameter complexity allows scaling to $N = 128$ electrons, previously inaccessible to neural-network wave functions in continuous space, enabling future work on finite-size extrapolations to the thermodynamic limit. We also show the Ansatz's capability of representing qualitative different phases of matter.

The following paper titled '**Message-Passing Neural Quantum States for the Homogeneous Electron Gas**' represents a collaborative effort involving Gabriel Pescia^{1,2}, Jannes Nys^{1,2}, myself³, Alessandro Lovato^{4,5,6} and Giuseppe Carleo^{1,2}.

¹*Institute of Physics, École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland*

²*Center for Quantum Science and Engineering, École Polytechnique Fédérale de Lausanne*

(EPFL), CH-1015 Lausanne, Switzerland

³Department of Physics and Astronomy and Facility for Rare Isotope Beams, Michigan State University, East Lansing, Michigan 48824, USA

⁴Physics Division, Argonne National Laboratory, Argonne, Illinois 60439, USA

⁵Computational Science Division, Argonne National Laboratory, Argonne, Illinois 60439, USA

⁶INFN-TIFPA Trento Institute for Fundamental Physics and Applications, 38123 Trento, Italy

8.1 Introduction

Predicting emergent physical phenomena and system properties from the *ab-initio* description of the system's constituents is notoriously difficult [61, 62]. Fermionic systems, in particular, can exhibit strong correlations among the particles, leading to collective phenomena in the form of exotic phases of matter, e.g. superconductivity and superfluidity [63, 64]. In recent years, progress in numerical simulations of strongly correlated systems was triggered by the development of increasingly precise machine-learning approximation techniques. Most notably, artificial neural network (NN) architectures, in combination with Variational Monte Carlo (VMC), have shown great promise to represent ground states of quantum spin systems, especially in more than one spatial dimension [1, 65, 66, 67, 68, 69, 70, 71]. Due to the universal approximation property of NNs, neural-network quantum states (NQS) can, in theory, accurately represent any quantum many-body state [72, 73]. NQS have been extended to fermionic degrees of freedom in a discrete basis [74, 75, 76], by incorporating the indistinguishability of quantum particles. More recently, advancements to ground and

excited state searches for fermionic and bosonic continuous degrees of freedom with open [3, 77, 78] and periodic boundary conditions (PBCs) [79, 80, 81], have been introduced.

The flexibility of NQS, compared to more traditional models, allows representing multiple phases of matter and even different physical systems with a single Ansatz. To exemplify this point, we refer to NQS studies on the ground-state of molecular systems [3, 77], solutions to effective field theory Hamiltonians describing atomic nuclei [82, 83, 4], bulk studies of fermionic and bosonic extended systems [79, 84, 81, 85], as well as NQS simulations of low-density neutron matter found in neutron stars [86]. The downside of this flexibility, especially in continuous space, is that NQS typically need a significant amount of variational parameters to reach a given accuracy. This makes optimization challenging and costly, preventing the usage of refined optimization schemes, e.g. second order optimization procedures [87, 88]. As a result, the accessible system sizes are limited to a few tens of particles. However, studying larger system sizes is of utmost importance to estimate physical properties in the thermodynamic limit [89, 90, 91, 92]. To remedy the situation, novel NQS architectures must be developed that significantly reduce the parameter complexity while retaining high accuracy.

This work introduces a neural-network wave function suitable for simulating strongly interacting fermionic quantum systems in continuous space with *one to two* orders of magnitude fewer parameters than current state-of-the-art NQS. The general form of the Ansatz is motivated by an analytical argument, relating the exact ground-state wave function to a many-body coordinate transformation of the electronic coordinates. It uses a permutation-equivariant message-passing architecture on a graph, inherently implementing the indistinguishability of same-species quantum particles [93]. As an application, we study the Homogeneous Electron Gas (HEG) in three spatial dimensions without explicitly breaking any of

the fundamental symmetries of the system, such as translations and spin-inversion symmetry. This allows characterizing, from first principles, the onset of Wigner crystallization at low densities.

8.2 Exact Backflow Transformations

Throughout this work, we consider a non-relativistic Hamiltonian of identical particles with mass m in d spatial dimensions:

$$H = -\frac{\hbar^2}{2m} \sum_i^N \nabla_i^2 + V(\mathbf{X}), \quad (8.1)$$

where the potential and interaction energy, V , is assumed to be diagonal in position representation, defined by the particle coordinates $\mathbf{X} = (\mathbf{r}_1, \dots, \mathbf{r}_N)$, $\mathbf{r}_i \in \mathbb{R}^d$. In the following, we derive an analytic functional form of the ground-state wave function and relate it to our variational Ansatz.

We use a suitable reference state $|\Phi_0\rangle$, as initial condition for the imaginary-time (τ) evolution induced by the Hamiltonian: $\Phi_\tau(\mathbf{X}) = \langle \mathbf{X} | e^{-\tau H} | \Phi_0 \rangle$. The exact ground-state is obtained in the large imaginary-time limit: $\lim_{\tau \rightarrow \infty} \Phi_\tau(\mathbf{X}) \propto \Psi_0(\mathbf{X})$, provided $|\Phi_0\rangle$ is non-orthogonal to the exact ground state, $|\langle \Psi_0 | \Phi_0 \rangle| > 0$. For fermions, non-orthogonality implies that the wave function must be antisymmetric w.r.t. the exchange of two particles i.e. $\Phi_0(\mathcal{P}_{ij}(\mathbf{X})) = -\Phi_0(\mathbf{X})$ (\mathcal{P}_{ij} permutes particles i and j). Applying the mean-value theorem to the imaginary-time evolved state, yields:

$$\Phi_\tau(\mathbf{X}) = \int_{\Omega} d\mathbf{X}' G_\tau(\mathbf{X}, \mathbf{X}') \Phi_0(\mathbf{X}') \quad (8.2)$$

$$= \text{Vol}(\Omega) \times G_\tau(\mathbf{X}, \mathbf{Y}(\mathbf{X})) \Phi_0(\mathbf{Y}(\mathbf{X})), \quad (8.3)$$

where Ω is the integration domain for the positional degrees of freedom, and $G_\tau(\mathbf{X}, \mathbf{X}') = \langle \mathbf{X} | e^{-\tau H} | \mathbf{X}' \rangle$ is the matrix element of the imaginary-time propagator. In (8.3), we introduced the mean-value point $\mathbf{Y}(\mathbf{X}) = (\mathbf{y}_1(\mathbf{X}), \dots, \mathbf{y}_N(\mathbf{X})) \in \Omega$, depending parametrically on the coordinates \mathbf{X} .

For non-relativistic Hamiltonians, Eq. (8.1), we have $G_\tau(\mathbf{X}, \mathbf{X}') \geq 0$, for all \mathbf{X}, \mathbf{X}' . Moreover, $G_\tau(\mathbf{X}, \mathbf{X}')$ is invariant under the exchange of particle coordinates: $G_\tau(\mathcal{P}_{ij}(\mathbf{X}), \mathcal{P}_{ij}(\mathbf{X}')) = G_\tau(\mathbf{X}, \mathbf{X}')$. In the fermionic case, the latter implies that $\mathbf{Y}(\mathbf{X})$ must be equivariant under particle exchange, $\mathbf{Y}(\mathcal{P}_{ij}(\mathbf{X})) = \mathcal{P}_{ij}(\mathbf{Y}(\mathbf{X}))$, to ensure antisymmetry of the total wavefunction. Eq. (8.3) therefore yields the product between a permutation symmetric, positive semi-definite function $J(\mathbf{X}) = G_\tau(\mathbf{X}, \mathbf{Y}(\mathbf{X})) \times \text{Vol}(\Omega)$ and a reference state computed at modified coordinates $\mathbf{Y}(\mathbf{X})$:

$$\Phi_\tau(\mathbf{X}) = J(\mathbf{X}) \times \Phi_0(\mathbf{Y}(\mathbf{X})). \quad (8.4)$$

Identification of the mean-value point $Y(X)$ with a many-body coordinate transformation gives an alternative justification for the backflow transformations [94] of single-particle coordinates. With a Slater determinant of given spin orbitals $\phi_\mu(\mathbf{r}_i)$ as initial state, $\Phi_0(\mathbf{X}) = \det \phi_\mu(\mathbf{r}_i)/\sqrt{N!}$, Eq. (8.4) is structurally related to the heuristic Jastrow-Backflow variational form [95, 96]. We remark that the symmetric contribution, $J(X)$, can be incorporated to the determinant:

$$\Phi_\tau(\mathbf{X}) = \mathcal{K} \times \det \varphi_\mu(\mathbf{y}_i(\mathbf{X})), \quad (8.5)$$

with $\varphi_\mu(\mathbf{y}_i(\mathbf{X})) = \phi_\mu(\mathbf{y}_i(\mathbf{X})) \times \sqrt[N]{J(\mathbf{X})}$, and \mathcal{K} a normalization constant.

The functional form (8.4),(8.5) is exact, provided that the symmetric factor $J(\mathbf{X})$ and the mean-value coordinates $\mathbf{Y}(\mathbf{X})$ satisfy Eq. (8.3), and the reference state is not orthogonal to the exact ground state. An approximate but explicit form for the coordinate transformation $\mathbf{Y}(\mathbf{X})$ can be obtained by repeatedly applying the imaginary-time propagator to the reference state in the limit of small τ . This process gives rise to the iterative backflow transformation, as introduced in Ref. [97, 98].

8.3 Message-Passing Neural-Network Quantum States

Motivated by Eq. (8.5), we use single-particle orbitals, $\{\phi_\mu\}_{\mu=1}^N$, evaluated at many-body backflow coordinates, $\mathbf{Y}(\mathbf{X})$, to construct the variational Ansatz. The backflow transformation is parameterized with permutation-equivariant message-passing NNs (MPNN) [93], hence we name it Message-Passing Neural Quantum State (MP-NQS).

In the MPNN, an all-to-all connected graph, encoding effective particle positions (nodes) $\mathbf{x}_i^{(t)} = [\mathbf{x}_i^{(0)}, \mathbf{h}_i^{(t)}] \in \mathbb{R}^{D_1}$ and their interactions (edges) $\mathbf{x}_{ij}^{(t)} = [\mathbf{x}_{ij}^{(0)}, \mathbf{h}_{ij}^{(t)}] \in \mathbb{R}^{D_2}$, is updated iteratively (we have introduced a time-step $t \geq 0$, an initial graph $(\mathbf{x}_i^{(0)}, \mathbf{x}_{ij}^{(0)})$, successively updated auxiliary variables known as hidden states $(\mathbf{h}_i^{(t)} \in \mathbb{R}^{D_1^h}, \mathbf{h}_{ij}^{(t)} \in \mathbb{R}^{D_2^h})$, suitably chosen *feature* dimensions $D_1^{(h)}$ and $D_2^{(h)}$, and we denoted concatenation with $[\cdot, \cdot]$). Construction of the initial graph $(\mathbf{x}_i^{(0)}, \mathbf{x}_{ij}^{(0)})$ is system dependent and will be discussed in detail later.

The hidden states are updated using equivariant *messages*, obtained from an attention mechanism [99]. The messages are given by weighted transformations of the edges $\mathbf{x}_{ij}^{(t)}$: $\mathbf{m}_{ij}^{(t+1)} = \boldsymbol{\omega}_{ij}^{(t)}(\mathbf{x}_{ij}^{(t)}) \odot \boldsymbol{\phi}(\mathbf{x}_{ij}^{(t)})$, where \odot represents element-wise multiplication along the feature dimension, and $\boldsymbol{\omega}_{ij}^{(t)} \in \mathbb{R}^{D_2}$ are weight vectors. The weights are obtained using

query/key matrices given by $\mathbf{Q}_{ij}^{(t)} = W_Q^{(t)} \cdot \mathbf{x}_{ij}^{(t)}$ and $\mathbf{K}_{ij}^{(t)} = W_K^{(t)} \cdot \mathbf{x}_{ij}^{(t)}$, with weight matrices $W_Q^{(t)}, W_K^{(t)} \in \mathbb{R}^{D_2 \times D_2}$. Applying an element-wise GELU non-linearity [100] to the overlap between queries and keys along the *particle dimension* (as opposed to the feature dimension [99, 101]), results in permutation-equivariant weights

$$\omega_{ij}^{(t)} = \text{GELU} \left(\sum_l \mathbf{Q}_{il}^{(t)} \mathbf{K}_{lj}^{(t)} \right). \quad (8.6)$$

This attention mechanism compares environments of particles i and j , and effectively increases the order of correlations that can be embedded in a single iteration of the network. This is crucial to reduce the total number of network iterations (parameters) and capture many-body effects. The hidden states are updated using the current graph and messages:

$$\mathbf{h}_i^{(t+1)} = \mathbf{f} \left(\mathbf{x}_i^{(t)}, \sum_{j \neq i} \mathbf{m}_{ij}^{(t+1)} \right) \quad (8.7)$$

$$\mathbf{h}_{ij}^{(t+1)} = \mathbf{g} \left(\mathbf{x}_{ij}^{(t)}, \mathbf{m}_{ij}^{(t+1)} \right) \quad (8.8)$$

where the functions ϕ , \mathbf{f} , and \mathbf{g} are parameterized by Multilayer Perceptrons (MLPs). The updated graph has then the same structure as the former: $\mathbf{x}_i^{(t+1)} = [\mathbf{x}_i^{(0)}, \mathbf{h}_i^{(t+1)}]$, $\mathbf{x}_{ij}^{(t+1)} = [\mathbf{x}_{ij}^{(0)}, \mathbf{h}_{ij}^{(t+1)}]$. Inclusion of the initial inputs, referred to as a “skip connection” in ML literature [102], mitigates the vanishing gradient problem and allows a more efficient capture of correlations.

The final backflow coordinates are constructed as $\mathbf{y}_i(\mathbf{X}) = \mathbf{r}_i + \delta\mathbf{r}_i(\mathbf{X})$, where the displacements, $\delta\mathbf{r}_i(\mathbf{X})$, are obtained via a linear transformation of the final node states to d dimensions, i.e. $\delta\mathbf{r}_i(\mathbf{X}) = W \cdot \mathbf{x}_i^{(T)}$ with $W \in \mathbb{C}^{d \times D_1}$. The complex-valued backflow transformation allows changing the degree of localization, determined by the chosen single-particle

orbitals, and representing complex-valued wave functions in general.

Following (8.5), we further augment the orbitals with a permutation-invariant factor J , yielding:

$$\Psi(\mathbf{X}) = \det \varphi_\mu(\mathbf{y}_i(\mathbf{X})), \quad (8.9)$$

with $\varphi_\mu(\mathbf{y}_i) = \exp[J(\mathbf{Y}, \mu)] \times \phi_\mu(\mathbf{y}_i)$.

8.4 Hamiltonian

We now study the case of the homogeneous electron gas (HEG) in $d = 3$ spatial dimensions, a prototypical model for the electronic structure in solids. It includes Coulomb interactions among the solids' electrons while treating its positively charged ions as uniform, static, positive background [103]. Despite this simplification, the HEG exhibits different phases of matter and captures properties of real solids, particularly of Alkali metals. The Hamiltonian (in units of Hartree), for a system of N electrons with uniform density $n = \frac{N}{V}$, is given by:

$$H = -\frac{1}{2r_s^2} \sum_i^N \nabla_i^2 + \frac{1}{r_s} \sum_{i < j}^N \frac{1}{\|\mathbf{r}_i - \mathbf{r}_j\|} + \text{const.} \quad (8.10)$$

where we introduced the Wigner-Seitz radius $r_s = \sqrt[3]{3/(4\pi n)}$, and a constant arising from the electron-background interaction [95]. The conditionally convergent series of pairwise Coulomb interactions is evaluated using the Ewald summation technique, as is standard for extended systems in QMC [104, 105, 106]. We'll assume a fixed spin-polarization with $N = N_\uparrow + N_\downarrow$, where $N_{\uparrow/\downarrow}$ denotes the number of up/down spins. $s_i \in \{\uparrow, \downarrow\}$ denotes the i-th electron's spin. We equip the cubic simulation cell of side length L with periodic boundary

conditions (PBCs) in all spatial directions to access the bulk of the system.

As in (8.5), we use a single Slater determinant as a reference state, $\Phi_0(\mathbf{X})$. For the translation-invariant HEG, plane-wave orbitals are a natural and physically-motivated choice:

$$\phi_{\mathbf{k}}(\mathbf{r}) = \exp(i\mathbf{k} \cdot \mathbf{r}) \quad (8.11)$$

with $\mathbf{k} = \frac{2\pi}{L}\mathbf{n}$, where $\mathbf{n} \in \mathbb{Z}^d$. To account for the spin s of a particle located at \mathbf{r} , we use spin-orbitals $\phi_{\mu}(\mathbf{r}, s) = \phi_{\mathbf{k}_{\mu}}(\mathbf{r})\delta_{s_{\mu}, s}$, where each spin-orbital is characterized by the quantum numbers $\mu = (\mathbf{k}_{\mu}, s_{\mu})$. These orbitals allow modeling of translation invariant systems at fixed total momentum $\mathbf{k}_{\text{tot}} = \sum_{i=1}^N \mathbf{k}_i$. Furthermore, the determinant factorizes into a product of determinants of up and down spin orbitals.

We further specialize the MP-NQS to the HEG by defining initial feature vectors. Respecting the spin inversion and translation symmetries of the HEG requires us to ignore single-particle positions and spins. We, therefore, initialize the nodes to a learnable embedding vector $\mathbf{e} \in \mathbb{R}^{D_1}$, that does not depend on the particle index i . For the edge features, we use the translation invariant particle-distances $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$ and their norm. Same- and opposite-spin pairs are distinguished using products of the form $s_i \cdot s_j = \pm 1$ in the edge features. Overall, we obtain the following initial features:

$$\mathbf{x}_i^{(0)} = \mathbf{e}, \quad \mathbf{x}_{ij}^{(0)} = (\mathbf{r}_{ij}, \|\mathbf{r}_{ij}\|, s_i \cdot s_j). \quad (8.12)$$

Notice that this choice preserves the spin quantum number of each particle.

The PBCs of the simulation box are incorporated by mapping the components of vectors $\mathbf{r} \in \mathbb{R}^d$ (where $\mathbf{r} = \mathbf{r}_i$ or $\mathbf{r} = \mathbf{r}_{ij}$) to a Fourier basis $\mathbf{r} \mapsto \left[\sin\left(\frac{2\pi}{L}\mathbf{r}\right), \cos\left(\frac{2\pi}{L}\mathbf{r}\right) \right] \in$

\mathbb{R}^{2d} , and the norm of the distance between two particles, $\|\mathbf{r}_{ij}\|$, to a periodic surrogate $\|\mathbf{r}_{ij}\| \mapsto \|\sin(\frac{\pi}{L}\mathbf{r}_{ij})\|$, as in Ref. [79]. In sum, our Ansatz is translation- and spin-inversion-invariant and possesses fixed total momentum, \mathbf{k}_{tot} . Its number of parameters is system-size independent (here ~ 19000) and, using Stochastic Reconfiguration (SR) [88], only $\mathcal{O}(10^3)$ optimization steps are needed to reach convergence. A comparison to other NQS approaches is given in the Supplemental Material.

8.5 Results

We study the fully spin-polarized and unpolarized HEG in different density regimes $r_s \in [1, 200]$ and up to system sizes of $N = 128$ electrons. We benchmark our ground-state energies against a variety of methods, including transcorrelation augmented full configuration interaction method (FCI), distinguishable clusters with doubles (DCD) method [91] for large densities, and *Diffusion Monte Carlo* (DMC) with backflow (BF-DMC) [107, 108] for small densities. We also compare to state-of-the-art NQS architectures – FermiNet [81] and WAP-Net [80] – available for small system sizes $N \in \{14, 19\}$. The effect of our backflow transformation on the nodal surface is studied by comparing it to fixed-node DMC (FN-DMC) results (see Supplemental Material). We use an energy of 1.5 mHa per particle (chemical accuracy) to assess the significance of energy differences between the different methods. An overview of our results and benchmarks for the various system sizes and densities is provided in the Supplemental Material.

8.5.1 Energy Benchmarks for Small Systems

For $N = 14$, FCI provides the lowest available ground-state energy for the HEG for $r_s \leq 5$. The energy difference between the MP-NQS and FCI results is lower than 1.5 mHa per particle. State-of-the-art NQS architectures perform comparably to the MP-NQS: The unrestricted FermiNet performs slightly better ($\mathcal{O}(10^{-5})$ Ha/N) than both MP-NQS and WAPNet for $r_s \leq 2$, while the MP-NQS and WAPNet improve over this version of FermiNet for $r_s = 5$. The restricted FermiNet yields worse ground-state energies than the MP-NQS over all probed densities [81] (see Fig. 8.1). For $r_s \geq 5$ we compare to results obtained with WAPNet and DCD. We find slightly better performance than WAPNet for all of the reported densities. Furthermore, both neural-network-based methods outperform the DCD method. Nevertheless, all results lie within a range of 1.5 mHa per particle. A similar pattern is seen for $N = 19$: the MP-NQS obtains slightly higher energies than WAPNet for large densities ($r_s \leq 5$) and marginally lower ones at smaller densities ($r_s \geq 5$). The differences are lower than 1.5 mHa per particle.

8.5.2 Energy Benchmarks for Large Systems

For $N = 54$ particles, we compare to FCI at high density $r_s = 1$, where our method achieves the same accuracy. The difference between our results and FCI decreases with increasing system size, yielding statistically indistinguishable results for $N = 54$. We obtain significantly better ground-state energies than BF-DMC, especially at high densities. This is in stark contrast to the (FermiNet-based) architecture of [84] (dubbed LiNet in the following), which does not improve upon BF-DMC energies over the whole density regime (see Figure 8.1, bottom). Comparison to DCD data shows that the MP-NQS can consistently improve

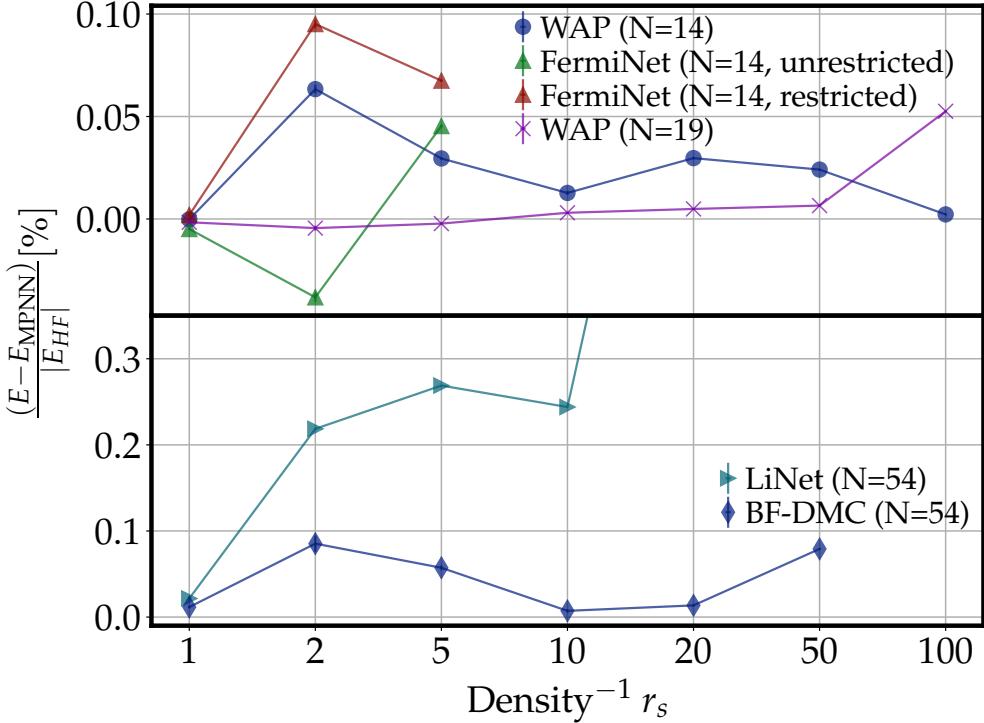


Figure 8.1: Energy differences between ground-state energies, obtained with other methods and with the MP-NQS, in units of the thermodynamic Hartree Fock energy, for various densities, polarizations, and system sizes. (Top) $N = 14, 19$, (Bottom): $N = 54$ particles. Error bars are too small to be visible for most densities. The corresponding numerical data can be found in the Supplemental Material.

the ground-state energies for all examined densities, with energy differences exceeding 1.5 mHa per particle. This indicates a deterioration of DCD’s accuracy with increasing system size, while our method appears to maintain its precision independent of system size. As expected, we reach higher accuracy compared to FN-DMC due to the optimization of the nodal surface. The difference is ~ 1 mHa per particle at the largest density. Still, it decreases with decreasing density because the nodal surface does not contribute as much to the ground-state energy in this regime.

A liquid-crystal phase transition is expected for the HEG, as a function of the density n . The dominating kinetic energy in Eq. (8.10) ($\sim 1/r_s^2$) for large n leads to the well-known Fermi liquid behavior. For small n the potential energy ($\sim 1/r_s$) dominates and

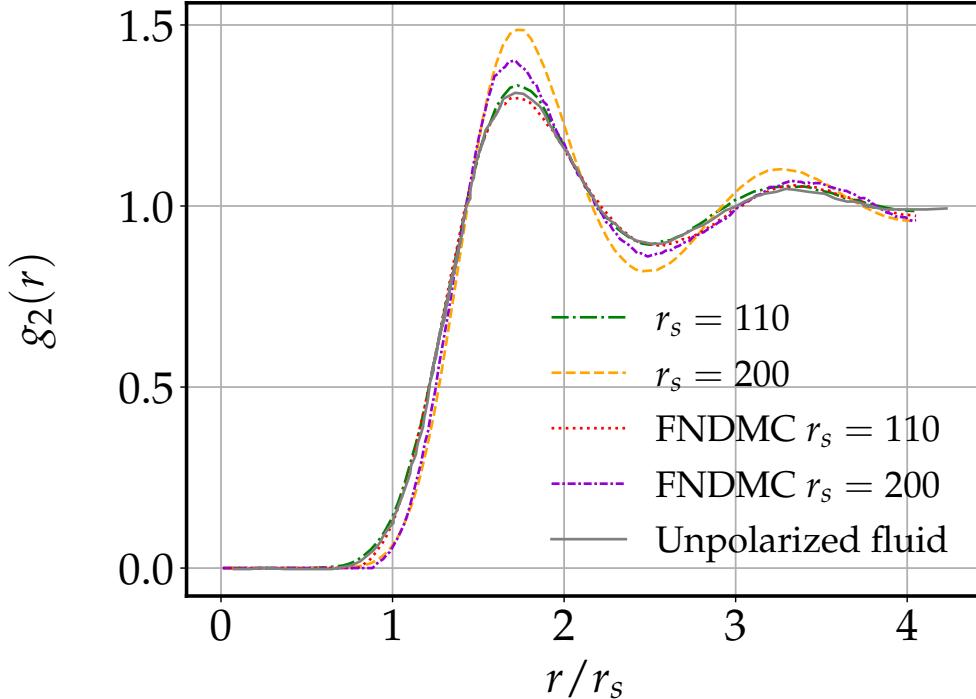


Figure 8.2: Spin-averaged radial distribution function for the HEG system with $N = 128$ electrons at low densities ($r_s = 110, 200$). Error bars are smaller than the symbols. The unpolarized fluid is obtained from [110] for $r_s = 110$.

enforces a crystalline structure among the electrons, known as Wigner crystal. The Wigner crystal, expected to be of BCC type [109], is translation invariant, and the resulting crystal structure is called a *floating* crystal. A homogeneous single-particle density distribution and a crystalline two-body radial distribution function characterize the latter. Detection of the translation invariant transition has not been realized up to now, to the best of our knowledge. In previous QMC simulations of the Wigner crystal phase (e.g., in the works of Refs. [110, 92, 81]), the transition is investigated by comparison of energies obtained with a crystalline and liquid variational Ansatz, explicitly breaking translation symmetry. Consequently, these works find a pinned BCC lattice. Crystallization is further favored by Gaussian orbitals centered around BCC lattice sites, a primitive BCC simulation cell, or both. We use the MP-NQS as unbiased, translation invariant variational Ansatz to investigate the transition

in a conventional BCC cell, i.e., a simple cubic simulation cell. We study a system size of $N = 128$ particles at $r_s = 110, 200$. We display the radial distribution functions for different densities in Figure 8.2. We compare results from a liquid FN-DMC and BF-DMC calculation to results obtained with the MP-NQS at the predicted transition density around $r_s = 110$ [110]. We observe remarkable agreement between the three correlators, implying that the MP-NQS favors the fluid behavior at this density. For the lower density of $r_s = 200$, we observe increasing oscillations up to larger distances, departing from the liquid FN-DMC result. This shows the capability of the MP-NQS to describe qualitatively different phases.

N	r_s	MP-NQS	WAP [80]	FermiNet [81]	FCI*/DCD** [91]
14	1	0.568967(6)	0.568965(1)	0.568904(1)	0.56861(1)*
	2	-0.008391(1)	-0.0083310(3)	-0.008427(1)	-0.00868(2)*
	5	-0.0798544(4)	-0.0798360(1)	-0.079821(1)	-0.08002(2)*
	10	-0.0552126(6)	-0.05520380(3)	N/A	-0.05509**
	20	-0.0324553(2)	-0.0324434(1)	N/A	-0.03201**
	50	-0.01462631(6)	-0.01462211(4)	N/A	-0.01384**
	100	-0.00773018(3)	-0.007729980(2)	N/A	N/A

Table 8.1: Total energy per particle in Hartree for unpolarized system of $N = 14$ particles. WAPNet and FermiNet are alternative NQS architectures optimized via VMC. We include FCI and DCD results as benchmarks from quantum chemistry.

N	r_s	MP-NQS	LiNet [84]	FN-DMC	BF-DMC [107, 108]	FCI*/DCD** [91]
54	1	0.52973(1)	0.530019(1)	0.53094(2)	0.52989(4)	0.52973(3)*
	2	-0.014046(8)	-0.013840(1)	-0.01326(2)	-0.013966(2)	-0.01379**
	5	-0.079090(2)	-0.0788354(2)	-0.07867(1)	-0.079036(3)	-0.07837**
	10	-0.054448(1)	-0.0542785(1)	-0.054269(8)	-0.054443(2)	-0.05322**
	20	-0.0320524(5)	-0.0316886(1)	-0.031976(8)	-0.032047(2)	-0.03113**
	50	-0.0145015(1)	N/A	-0.01387(2)	-0.0144877(1)	-0.01281**
	100	-0.0076793(1)	N/A	-0.007674(3)	N/A	N/A

Table 8.2: Total energy per particle in Hartree for the unpolarized system of $N = 54$ particles. We include FCI and DCD results as benchmarks from quantum chemistry.

8.6 Conclusions

We have introduced MP-NQS, a novel NQS architecture, leveraging MPNNs to build highly expressive backflow coordinates. We demonstrate its power on the HEG system, reducing the number of parameters by orders of magnitudes compared to state-of-the-art NQS in continuous space while reaching at par or better accuracy. We also show improvement upon state-of-the-art BF-DMC results on large systems. The favorable scaling allows us to accurately simulate large periodic electronic systems, previously inaccessible to state-of-the-art NQS models. We increase the available system sizes from $N = 27$ and $N = 54$ electrons in periodic systems [84, 81, 80] to $N = 128$ electrons in this work. We hence open the door to extrapolation methods to the thermodynamic limit for extended systems.

The simulation of the HEG phase transition using a translation invariant Ansatz in an unbiased simulation cell shows the MP-NQS' capability to represent different phases of matter. We reproduce the liquid phase of the HEG up to around $r_s = 110$, while observing pronounced density fluctuations, potentially compatible with Wigner crystallization, at around $r_s = 200$. We leave it to further, more specialized, investigation to study the nature of the floating phase found at small densities, potentially resolving the tension with existing predictions [110, 92] (based on translation-symmetry breaking wave function states), on the location and nature of the phase transition.

In addition to the numerical results, we introduced an analytical argument, justifying commonly adopted backflow transformations. Our argument shows that a backflow transformation over a reference state is sufficient to obtain the exact ground-state wave function. It will be of particular interest to characterize the geometrical properties of these transformations and understand in what cases neural-network parameterizations can efficiently

describe them.

8.7 Acknowledgements

We acknowledge stimulating discussions with Markus Holzmann and Francesco Pederiva. The numerical tools used in this work are based on the open-source software NetKet[111, 112, 113] version 3.0. The present research is supported by the Swiss National Science Foundation under Grant No. 200021_200336, Microsoft Research, by the U.S. Department of Energy, Office of Science, Office of Nuclear Physics, under contracts DE-AC02-06CH11357 (A.L.), by the 2020 DOE Early Career Award program (A.L.), the NUCLEI SciDAC program (A.L.), and Argonne LDRD awards (A. L.). J. K. is supported by the U.S. Department of Energy, Office of Science, office of Nuclear Physics under grant No. DE-SC0021152 and U.S. National Science Foundation Grants No. PHY-1404159 and PHY-2013047.

Chapter 9

Ultra-cold Fermi Gases

Ultra-cold Fermi gases exhibit a rich array of quantum mechanical properties, including the transition from a fermionic superfluid BCS state to a bosonic superfluid BEC state, which can be precisely probed experimentally. However, accurately describing these properties poses significant theoretical challenges due to strong pairing correlations and the non-perturbative nature of particle interactions. Here, we introduce a Pfaffian-Jastrow neural-network quantum state with a message-passing architecture to encode pairing and other quantum correlations efficiently. Our novel approach surpasses existing Slater-Jastrow frameworks and outperforms state-of-the-art diffusion Monte Carlo methods, as evidenced by lower ground-state energies. We observe the emergence of strong pairing correlations by analyzing the opposite-spin pair distribution functions. Additionally, we demonstrate that transfer learning enhances the training of neural-network wave functions, facilitating the exploration of the BCS-BEC crossover region near unitarity. Our findings highlight the potential of neural-network quantum states as a promising strategy for investigating ultra-cold Fermi gases.

The following paper titled '**Neural-network quantum states for ultra-cold Fermi gases**' represents a collaborative effort involving myself¹, Gabriel Pescia^{3,4}, Bryce Fore², Jannes Nys^{3,4}, Giuseppe Carleo^{3,4}, Stefano Gandolfi⁵, Morten Hjorth-Jensen^{1,6}, and Alessandro Lovato^{2,7,8}.

¹*Department of Physics and Astronomy and Facility for Rare Isotope Beams, Michigan State University, East Lansing, Michigan 48824, USA*

²*Physics Division, Argonne National Laboratory, Argonne, Illinois 60439, USA*

³*Institute of Physics, École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland*

⁴*Center for Quantum Science and Engineering, École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland*

⁵*Theory Division T-2, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA*

⁶*Department of Physics and Center for Computing in Science Education, University of Oslo, N-0316 Oslo, Norway*

⁷*Computational Science Division, Argonne National Laboratory, Argonne, Illinois 60439, USA*

⁸*INFN-TIFPA Trento Institute of Fundamental Physics and Applications, 38123 Trento, Italy*

9.1 Introduction

The study of ultra-cold Fermi gases has received considerable experimental and theoretical attention in recent years due to their unique properties and potential applications in fields ranging from condensed matter physics to astrophysics. These systems can be created and manipulated in the laboratory with high precision, providing a versatile platform for investigating a wide variety of phenomena. By tuning the *s*-wave scattering length a via external

magnetic fields near a Feshbach resonance, one can smoothly crossover from a fermionic superfluid BCS state ($a < 0$) of long-range Cooper pairs to a bosonic superfluid BEC state ($a > 0$) of tightly-bound, repulsive dimers. Given their diluteness, the behavior of these systems is mainly governed by a and the effective range of the potential r_e , with natural units provided by the Fermi momentum k_F and the Fermi gas energy per particle in the thermodynamic limit $E_{FG} = \frac{3}{5} \frac{\hbar^2}{2m} k_F^2$ (see Ref. [114] and references therein).

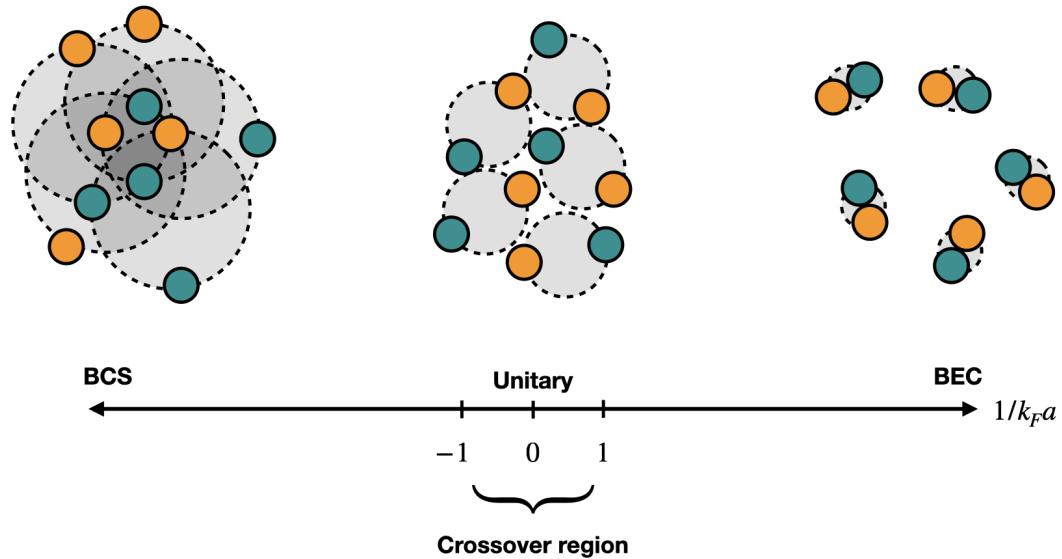


Figure 9.1: A cartoon of the BCS-BEC crossover. Moving from left to right, the attractive interaction between opposite-spin fermions increases. However, in the BEC regime, the attraction binds the pairs so tightly that they behave as weakly repulsive bosons. The region between the weakly attractive Cooper pairs and the weakly repulsive dimers is known as the unitary limit.

The region between the BCS and BEC states, known as the “unitary limit,” is particularly interesting as a diverges and r_e approaches zero. The unitary Fermi gas (UFG) is a strongly-interacting system that exhibits surprisingly stable superfluid behavior. Studying the BCS-BEC crossover near the unitary limit can reveal critical aspects of the underlying mechanism

behind superfluidity in fermionic matter. The UFG is also universal, meaning its properties are independent of the details of the two-body potential. This universality allows for robust comparisons and predictions between seemingly disparate quantum systems. For instance, the UFG is relevant for neutron stars, as they provide a means to study superfluid low-density neutron matter [115, 116], whose properties are crucial for the phenomenology of glitches [28] and the cooling of these stars via neutrino emission [26, 27, 117].

The onset of strong pairing correlations and the non-perturbative nature of the interaction makes the theoretical study of these systems particularly challenging for quantum many-body methods. Among them, quantum Monte Carlo (QMC) has proven to be exceptionally efficient in calculating various properties with high accuracy, including the energy [118], pairing gap [119], and other quantities related to the so-called contact parameter [120]. Diffusion Monte Carlo (DMC), in particular, is an accurate tool for calculating the properties of quantum many-body systems [121]. The fixed-node approximation typically employed in DMC calculations to control the fermion-sign problem provides a rigorous upper bound to the ground-state energy that agrees well with other methods and experiments [122, 123]. Moreover, unlike the Auxiliary Field Quantum Monte Carlo method, DMC can handle broad classes of local interactions, which provides exact results that are sign-problem free but is limited only to unpolarized systems with a purely attractive interaction [118]. However, the fixed-node approximation limits the accuracy of DMC energies, which induces a residual dependence on the starting variational wave function. The latter has a critical role in DMC calculations of expectation values of operators that do not commute with the Hamiltonian, such as spatial and momentum distributions. The analytical form of the variational ansatz is usually tailored to specific problems of interest and biased by the physical intuition of the researchers.

In this work, we overcome these limitations by performing variational Monte Carlo (VMC) calculations of ultra-cold Fermi gases with neural-network quantum states (NQS) [1] that incorporate only the most essential symmetries and boundary conditions. After their initial application to quantum-chemistry problems [42, 124], continuous-space NQS have been successfully employed to study quantum many-body systems in the presence of spatial periodicities, such as interacting quantum gases of bosons [125], the homogeneous electron gas [126, 127], and dilute neutron matter [128]. Recent works have also used NQS to solve the nuclear Schrödinger equation in both real space[36, 129, 130, 131, 132] and the occupation number formalism [133]. When dealing with fermions, the antisymmetry is usually enforced using generalized Slater determinants, the expressivity of which can be augmented with either backflow transformations [134] or by adding “hidden” degrees of freedom [52].

Strong pairing correlations in fermionic systems motivate adopting an antisymmetrized wave function constructed from pairing orbitals rather than single-particle orbitals. One such construction, often called the geminal wave function [135, 136], considers determinants of spin-singlet pairs, while other more general wave functions based on the Pfaffian [137, 5, 138, 139], consider both singlet and triplet contributions. Pfaffian wave functions combined with neural-network Jastrow correlators [140] have successfully modeled lattice fermions, even revealing the existence of a quantum spin liquid phase in the J1-J2 models on two-dimensional lattices [141].

We propose a novel NQS that extends the conventional Pfaffian-Jastrow [5] ansatz by incorporating neural backflow transformations into a fully trainable pairing orbital. The backflow transformations are generated by a message-passing architecture recently introduced to model the homogeneous electron gas [2]. In addition to being a significant departure from generalized Slater determinants, our Pfaffian-Jastrow NQS naturally encodes pairing

in the singlet and triplet channels, without stipulating a particular form for the pairing orbital. In view of this, it is broadly applicable to other strongly-interacting systems with the same symmetries and boundary conditions. We demonstrate the representative power of our NQS by computing ground-state properties of ultra-cold Fermi gases in the BCS-BEC crossover. Our Pfaffian-Jastrow NQS outperforms Slater-Jastrow NQS by a large margin, even when generalized backflow transformations are included in the latter. Most notably, we find lower energies than those obtained with state-of-the-art DMC methods, which start from highly-accurate BCS-like trial wave functions.

The rest of the paper is organized as follows. In Section ??, we introduce the Hamiltonian used to model ultra-cold atomic gases near the unitary limit and the many-body techniques used to solve the Schrödinger equation. In Section 9.6, we compare the Pfaffian-Jastrow NQS with other NQS ansätze and state-of-the-art DMC results. Finally, in Section 9.7, we draw our conclusion and provide future perspectives of this work.

9.2 Hamiltonian

As customary in QMC approaches, we simulate the infinite system using a finite number of fermions N in a cubic simulation cell with side length L , equipped with periodic boundary conditions (PBCs) in all $d = 3$ spatial dimensions. We use $\mathbf{r}_i \in \mathbb{R}^d$ and $s_i \in \{\uparrow, \downarrow\}$ to denote the positions and spin projections on the z -axis of the i -th particle, and the length L can be determined from the uniform density of the system $N/L^3 = k_F^3/(3\pi^2)$. The dynamics of the unpolarized gas is governed by the non-relativistic Hamiltonian

$$H = -\frac{\hbar^2}{2m} \sum_i^N \nabla_i^2 + \sum_{ij}^N v_{ij}, \quad (9.1)$$

where the attractive two-body interaction

$$v_{ij} = (\delta_{s_i, s_j} - 1)v_0 \frac{2\hbar^2}{m} \frac{\mu^2}{\cosh^2(\mu r_{ij})}, \quad (9.2)$$

acts only between opposite-spin pairs, making the interaction mainly in *s*-wave for small values of r_e . In the above equations, ∇_i^2 is the Laplacian with respect to \mathbf{r}_i and $r_{ij} = \|\mathbf{r}_i - \mathbf{r}_j\|$ is the Euclidean distance between particles i and j . The Pöschl-Teller interaction potential of Eq. (9.2) provides an analytic solution of the two-body problem and has been employed in several previous QMC calculations [142, 143, 144, 120]. The parameters v_0 and μ tune the scattering length a and effective range r_e , respectively. In the unitary limit $|a| \rightarrow \infty$, the zero-energy ground state between two particles corresponds to $v_0 = 1$ and $r_e = 2/\mu$. In order to analyze the crossover between the BCS and BEC phases, we will use different combinations of v_0 and μ that correspond to the same effective range. In addition, we will consider various values of μ with fixed $v_0 = 1$ to extrapolate the zero effective range behavior at unitarity.

9.3 Neural-Network Quantum States

We solve the Schrödinger equation associated with the Hamiltonian of Eq. (9.1) using two different families of NQS. All ansätze have the general form

$$\Psi(X) = e^{J(X)} \Phi(X), \quad (9.3)$$

where the Jastrow correlator $J(X)$ is symmetric under particle exchange and $\Phi(X)$ is antisymmetric. Here, we have introduced $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, with $\mathbf{x}_i = (\mathbf{r}_i, s_i)$, to represent

the set of all single-particle positions and spins compactly.

In addition to the antisymmetry of fermionic wave functions, the periodic boundary conditions, and the translational symmetry (which will be discussed in Sec. ??), we also enforce the discrete parity and time-reversal symmetries as prescribed in Ref. [131]. More specifically, we carry out the VMC calculations for the unpolarized gas using $\Psi^{PT}(R, S)$ given by

$$\Psi^P(R, S) = \Psi(R, S) + \Psi(-R, S), \quad (9.4)$$

$$\Psi^{PT}(R, S) = \Psi^P(R, S) + (-1)^n \Psi^P(R, -S), \quad (9.5)$$

where $n = N/2$ and we have used the notation $R = \{\mathbf{r}_1, \dots, \mathbf{r}_N\}$ and $S = \{s_1, \dots, s_N\}$ for the set of all positions and spins, respectively. Enforcing these symmetries has been shown to accelerate the convergence of ground-state energies for both atomic nuclei [131] and dilute neutron matter [128].

9.3.1 Pfaffian-Jastrow-Backflow

The antisymmetric part of the wave function employed in QMC studies of ultra-cold Fermi gases is typically constructed as an antisymmetrized product of BCS spin-singlet pairs [145, 142, 115, 143, 146]. It goes by a variety of names, such as the geminal wave function [135, 136], the singlet pairing wave function [5], and the (number-projected) BCS wave function [146], just to name a few. Although geminal wave functions have demonstrated significant improvements over single-determinant wave functions of single-particle orbitals, the energy gains are typically smaller for partially spin-polarized systems [136], as contributions from the spin-triplet channel are missing. This naturally leads to the singlet-triplet-unpaired

(STU) Pfaffian wave function [137, 5], in which the pairing orbitals are explicitly decomposed into singlet and triplet channels. Then, the STU ansatz is expressed as the Pfaffian of a block matrix, with the singlet, triplet, and unpaired contributions partitioned into separate blocks. When the triplet blocks are zero, the STU wave function reduces to the geminal wave function.

Both the geminal and the STU wave functions rely on fixing the spin ordering of the interacting fermions. Consequently, they are not amenable to potentials that exchange spin, such as those used to model the interaction among nucleons [147]. In neutron-matter calculations, for instance, the pairing orbital for the Pfaffian wave function can be taken as a product of a radial part and a spin-singlet part [138, 139]. The spin-triplet pairing has so far been neglected in neutron-matter calculations, but they can be treated similarly without requiring spin ordering.

To address the limitations of the works mentioned above, we take the most general form of the Pfaffian wave function [137, 5] as the antisymmetric part of our ansatz

$$\Phi_{PJ}(X) = \text{pf} \begin{bmatrix} 0 & \phi(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1, \mathbf{x}_N) \\ \phi(\mathbf{x}_2, \mathbf{x}_1) & 0 & \cdots & \phi(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N, \mathbf{x}_1) & \phi(\mathbf{x}_N, \mathbf{x}_2) & \cdots & 0 \end{bmatrix}, \quad (9.6)$$

where we assume the unpolarized case for this initial investigation. We do not keep the spins fixed, nor do we mandate a specific form for the pairing orbital $\phi(\mathbf{x}_i, \mathbf{x}_j)$. Instead, we capitalize on the universal approximation property of feed-forward neural networks (FNN) [148]

by defining the pairing orbital as

$$\phi(\mathbf{x}_i, \mathbf{x}_j) = \nu(\mathbf{x}_i, \mathbf{x}_j) - \nu(\mathbf{x}_j, \mathbf{x}_i), \quad (9.7)$$

where ν is a dense FNN. The above expression ensures that the Pfaffian is mathematically well-defined, as the matrix is skew-symmetric by construction $\phi(\mathbf{x}_i, \mathbf{x}_j) = -\phi(\mathbf{x}_j, \mathbf{x}_i)$. Since ν takes all the degrees of freedom of a given pair of particles as input, including the spins, our pairing orbital has the capacity to discover the spin-singlet and spin-triplet correlations on its own.

This design leaves our Pfaffian-Jastrow (PJ) ansatz agnostic to any particular form of the interaction and systematically improvable by simply increasing the size of ν . The input dimension of ν only depends on the spatial dimension d and not the total number of particles N , leading to an exceptionally scalable ansatz. Given the generality of our formulation, the Pfaffian ansatz calculation cannot be reduced to a determinant of singlet pairing orbitals as in the geminal wave function. Thus, the efficient computation of the Pfaffian is crucial to the scalability of our approach. To this aim, we implement the Pfaffian computation according to Ref. [10].

We further improve the nodal structure of our PJ ansatz through backflow (BF) transformations [94]. To our knowledge, this is the first time neural BF transformations have been used in a Pfaffian wave function, although they have demonstrated their superiority over traditional BF transformations within the Slater-Jastrow formalism in numerous applications [134, 42, 124]. We replace the original single-particle coordinates \mathbf{x}_i by new ones $\tilde{\mathbf{x}}_i(X)$, such that correlations generated by the presence of all particles are incorporated into the pairing orbital. To ensure that the Pfaffian remains antisymmetric, the backflow trans-

formation must be permutation equivariant with respect to the original \mathbf{x}_i , i.e. $\tilde{\mathbf{x}}_i$ depends on \mathbf{x}_i and is invariant with respect to the set $\{\mathbf{x}_j\}_{j \neq i}$. In Sec. ??, we discuss in detail how the backflow correlations are encoded via a permutation-equivariant message-passing neural network. All calculations labeled as PJ-BF assume that we apply the transformation $\nu(\mathbf{x}_i, \mathbf{x}_j) \rightarrow \nu(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)$ to the FNN in Eq. (9.7).

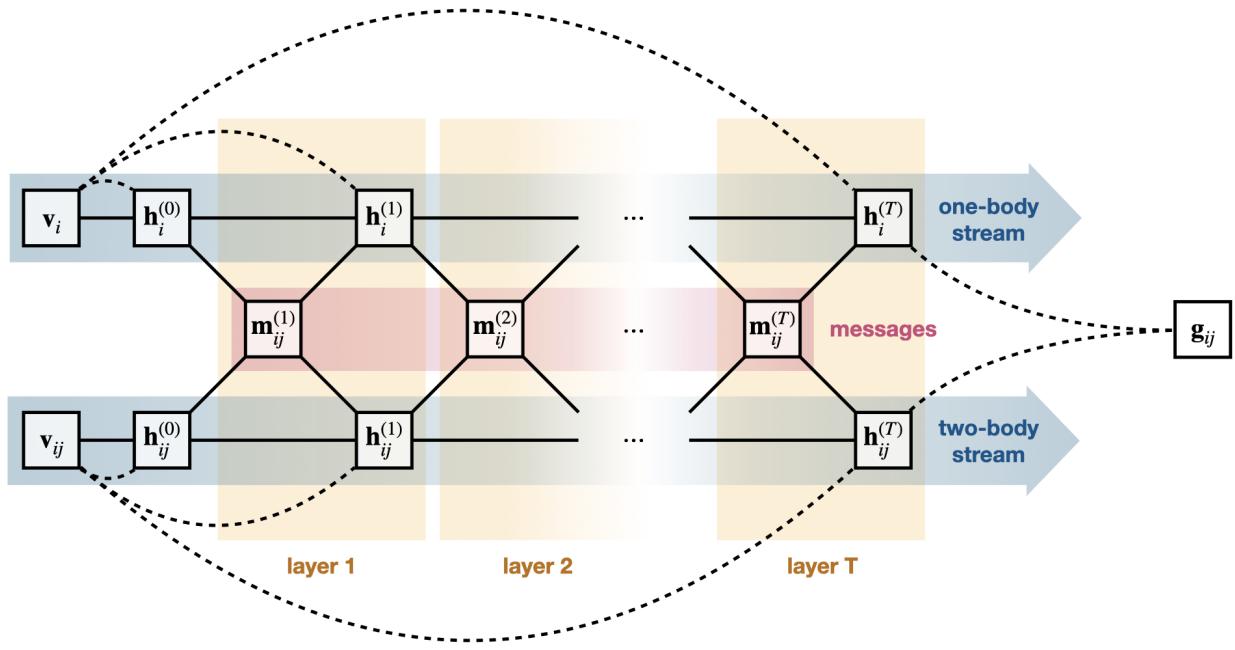


Figure 9.2: Schematic representation of a message-passing neural network with T iterations. Dashed lines represent the concatenation operations, while solid lines represent the parameterized transformations (linear transformations and nonlinear feedforward neural networks). Messages, highlighted in pink, mediate the exchange of information between the one- and two-body streams, in blue. A yellow box indicates a single iteration of the network.

9.3.2 Fixed-Node Slater-Jastrow

For comparison, we will report results obtained using a Slater-Jastrow (SJ) ansatz, which amounts to taking the antisymmetric part of the wave function to be a Slater determinant

of single-particle states

$$\Phi_{SJ}(X) = \det \begin{bmatrix} \phi_1(\mathbf{x}_1) & \phi_1(\mathbf{x}_2) & \cdots & \phi_1(\mathbf{x}_N) \\ \phi_2(\mathbf{x}_1) & \phi_2(\mathbf{x}_2) & \cdots & \phi_2(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_N(\mathbf{x}_1) & \phi_N(\mathbf{x}_2) & \cdots & \phi_N(\mathbf{x}_N) \end{bmatrix}. \quad (9.8)$$

In the fixed-node approximation, the single-particle states are the products of spin eigenstates with definite spin projection on the z -axis s_α and plane wave (PW) orbitals with discrete momenta $\mathbf{k}_\alpha = 2\pi\mathbf{n}_\alpha/L$, $\mathbf{n}_\alpha \in \mathbb{Z}^d$,

$$\phi_\alpha(\mathbf{x}_i) = e^{i\mathbf{k}_\alpha \cdot \mathbf{r}_i} \chi_\alpha(s_i), \quad (9.9)$$

where $\chi_\alpha(s_i) = \delta_{s_\alpha, s_i}$. Here, $\alpha = (\mathbf{k}_\alpha, s_\alpha)$ denotes the quantum numbers characterizing the state. We will label Slater-Jastrow NQS calculations using above plane wave orbitals as SJ-PW.

9.3.3 Slater-Jastrow-Backflow

As in the Pfaffian case, we improve the nodal structure of the above Slater determinant using backflow transformations generated by the message-passing neural network discussed in Sec. ???. We modify the spatial coordinates of Eq. (9.9) as

$$\mathbf{r}_i \rightarrow \mathbf{r}_i + \mathbf{u}_i(X), \quad (9.10)$$

where the complex backflow displacement $\mathbf{u}_i(X) \in \mathbb{C}^d$ allows for changes in both the phases and amplitudes of the spatial part of the single-particle states. We also map the single-particle spinors onto the Bloch sphere as

$$\chi_\alpha(\tilde{\mathbf{x}}_i) = \cos\left(\frac{\theta_i(X)}{2}\right) \delta_{s_\alpha, s_i} + \sin\left(\frac{\theta_i(X)}{2}\right) (1 - \delta_{s_\alpha, s_i}), \quad (9.11)$$

where $\theta_i(X) \in \mathbb{R}$ is the polar angle on the sphere. Both $\mathbf{u}_i(X)$ and $\theta_i(X)$ are permutation-equivariant functions of the original coordinates \mathbf{x}_i , the functional form of which will be discussed in Sec. 9.3.4.

To motivate the form of the backflow transformation in Eqs. (9.10) and (9.11), let us first revisit the original plane wave orbitals in Eq. (9.9). We simulate our system in the basis

$$|\mathbf{x}_i\rangle = |\mathbf{r}_i\rangle \otimes |s_i\rangle, \quad (9.12)$$

where $|\mathbf{r}_i\rangle$ are eigenstates of the position operator, with $\mathbf{r}_i \in \mathbb{R}^d$, and $|s_i\rangle$ are eigenspinors of the S_z operator, with $s_i \in \{\uparrow, \downarrow\}$. In the fixed-node approximation, we take the single-particle states to be products of momentum eigenstates with definite wave vector $\mathbf{k}_\alpha = 2\pi\mathbf{n}_\alpha/L$, $\mathbf{n}_\alpha \in \mathbb{Z}^d$, and eigenspinors with definite spin projection s_α ,

$$|\phi_\alpha\rangle = |\mathbf{k}_\alpha\rangle \otimes |s_\alpha\rangle. \quad (9.13)$$

Omitting overall normalization constants, the probability amplitude of measuring particle i

in state α is

$$\phi_\alpha(\mathbf{x}_i) = \langle \mathbf{x}_i | \phi_\alpha \rangle = \langle \mathbf{r}_i | \mathbf{k}_\alpha \rangle \langle s_i | s_\alpha \rangle = e^{i \mathbf{k}_\alpha \cdot \mathbf{r}_i} \delta_{s_\alpha, s_i}, \quad (9.14)$$

which we call the plane wave orbitals.

Now, let us transform to a new basis with modified position eigenstates and a superposition of eigenspinors

$$|\tilde{\mathbf{x}}_i\rangle = |\tilde{\mathbf{r}}_i\rangle \otimes |\chi_i\rangle. \quad (9.15)$$

While permutation equivariance is the sole essential property required for the backflow transformation $|\mathbf{x}_i\rangle \mapsto |\tilde{\mathbf{x}}_i\rangle$ to preserve the antisymmetry of the fermionic wave function, an additional property is desirable for computational convenience. Specifically, when the transformation depends on certain parameters, we aim to have $|\tilde{\mathbf{x}}_i\rangle = |\mathbf{x}_i\rangle$ when the parameters are identically zero. Then, nonzero parameters signify deviations from the original plane wave orbitals, such that less training is required compared to completely trainable orbitals.

An appropriate spatial transformation is trivial. We simply define new parameters $\mathbf{u}_i \in \mathbb{C}^d$, called the backflow displacement, and shift the coordinates as $\mathbf{r}'_i = \mathbf{r}_i + \mathbf{u}_i$. The backflow displacement is complex, allowing for changes in both the phases and amplitudes of the original plane wave orbitals.

For the spin part of the transformation, we look to spinors on the Bloch sphere for inspiration,

$$|\chi_i\rangle = \cos\left(\frac{\theta_i}{2}\right) |s_i\rangle + \sin\left(\frac{\theta_i}{2}\right) \sigma_i^x |s_i\rangle. \quad (9.16)$$

In the above, we have introduced another backflow variable $\theta_i \in \mathbb{R}$ akin to the polar angle of a Bloch spinor, and we have excluded the relative phase in favor of a completely real-valued

wave function. We also write the superposition in terms of $|s_i\rangle$ and the Pauli X -operator σ_i^x , which flips the spin of the i -th particle, rather than $|\uparrow\rangle$ and $|\downarrow\rangle$. This way, it is obvious that $|\chi_i\rangle = |s_i\rangle$ when $\theta_i = 0$, as desired. The overlap of two spinors is given by

$$\begin{aligned}
\langle \chi_i | \chi_j \rangle &= \cos\left(\frac{\theta_i}{2}\right) \cos\left(\frac{\theta_j}{2}\right) \langle s_i | s_j \rangle \\
&\quad + \sin\left(\frac{\theta_i}{2}\right) \cos\left(\frac{\theta_j}{2}\right) \langle s_i | \sigma_i^{x\dagger} | s_j \rangle \\
&\quad + \cos\left(\frac{\theta_i}{2}\right) \sin\left(\frac{\theta_j}{2}\right) \langle s_i | \sigma_j^x | s_j \rangle \\
&\quad + \sin\left(\frac{\theta_i}{2}\right) \sin\left(\frac{\theta_j}{2}\right) \langle s_i | \sigma_i^{x\dagger} \sigma_j^x | s_j \rangle \\
&= \left[\cos\left(\frac{\chi_i}{2}\right) \cos\left(\frac{\chi_j}{2}\right) + \sin\left(\frac{\chi_i}{2}\right) \sin\left(\frac{\chi_j}{2}\right) \right] \delta_{s_i, s_j} \\
&\quad \left[\sin\left(\frac{\theta_i}{2}\right) \cos\left(\frac{\theta_j}{2}\right) - \cos\left(\frac{\theta_i}{2}\right) \sin\left(\frac{\theta_j}{2}\right) \right] (1 - \delta_{s_i, s_j}) \\
&= \cos\left(\frac{\theta_i - \theta_j}{2}\right) \delta_{s_i, s_j} + \sin\left(\frac{\theta_i - \theta_j}{2}\right) (1 - \delta_{s_i, s_j}),
\end{aligned} \tag{9.17}$$

where the limiting cases are summarized in the following table.

$\langle \chi_i \chi_j \rangle$	$\theta_i = \theta_j$	$\theta_i - \theta_j = \pm\pi$
$s_i = s_j$	1	0
$s_i \neq s_j$	0	± 1

Finally, we can compute the backflow orbitals with the transformed degrees of freedom

$$\begin{aligned}
\phi_\alpha(\tilde{\mathbf{x}}_i) &= \langle \tilde{\mathbf{x}}_i | \phi_\alpha \rangle \\
&= \langle \tilde{\mathbf{r}}_i | \mathbf{k}_\alpha \rangle \langle \chi_i | s_\alpha \rangle \\
&= e^{i\mathbf{k}_\alpha \cdot \tilde{\mathbf{r}}_i} \left(\cos\left(\frac{\theta_i}{2}\right) \langle s_i | s_\alpha \rangle + \sin\left(\frac{\theta_i}{2}\right) \langle s_i | \sigma_i^{x\dagger} | s_\alpha \rangle \right) \\
&= e^{i\mathbf{k}_\alpha \cdot (\mathbf{r}_i + \mathbf{u}_i)} \left(\cos\left(\frac{\theta_i}{2}\right) \delta_{s_\alpha, s_i} + \sin\left(\frac{\theta_i}{2}\right) (1 - \delta_{s_\alpha, s_i}) \right).
\end{aligned} \tag{9.18}$$

In Eqs. (9.10) and (9.11), we use the notation $\mathbf{u}_i(X)$ and $\theta_i(X)$ to emphasize that the backflow “parameters” we define here are not variational parameters, but a function of all other particles. More specifically, they are permutation-equivariant functions of the original \mathbf{x}_i , whose functional forms depend on the outputs of the permutation-equivariant message-passing neural network (MPNN) described in Sec. ??.

In addition to omitting the relative phase in Eq. (9.16) so that the spinors remain real, we map the spatial components of the plane wave and backflow orbitals to the equivalent real ones. For the unpolarized system of fermions, the latter mapping is

$$\begin{Bmatrix} e^{i\mathbf{k}_\alpha \cdot (\mathbf{r}_i + \mathbf{u}_i)} \\ e^{-i\mathbf{k}_\alpha \cdot (\mathbf{r}_i + \mathbf{u}_i)} \end{Bmatrix} \mapsto \begin{Bmatrix} e^{\mathbf{k}_\alpha \cdot \text{Im}(\mathbf{u}_i)} \cos(\mathbf{k}_\alpha \cdot (\mathbf{r}_i + \text{Re}(\mathbf{u}_i))) \\ e^{\mathbf{k}_\alpha \cdot \text{Im}(\mathbf{u}_i)} \sin(\mathbf{k}_\alpha \cdot (\mathbf{r}_i + \text{Re}(\mathbf{u}_i))) \end{Bmatrix}. \quad (9.19)$$

Remember that the spatial inputs to the MPNN are L -periodic, which guarantees that \mathbf{u}_i is as well. Therefore it is less redundant to compute

$$\begin{Bmatrix} e^{\text{Im}(\mathbf{u}_i)} \cos(\mathbf{k}_\alpha \cdot \mathbf{r}_i + \text{Re}(\mathbf{u}_i)) \\ e^{\text{Im}(\mathbf{u}_i)} \sin(\mathbf{k}_\alpha \cdot \mathbf{r}_i + \text{Re}(\mathbf{u}_i)) \end{Bmatrix}, \quad (9.20)$$

instead of the right-hand side of Eq. (9.19). Then the modified phase $\text{Re}(\mathbf{u}_i)$ and amplitude $e^{\text{Im}(\mathbf{u}_i)}$ still have the correct periodicity, but they do not reduce to the trivial case when $\mathbf{k}_\alpha = \mathbf{0}$. In the end, this last step is a small detail and does not affect the final result.

9.3.4 Message-Passing Neural Network

Implementing the aforementioned neural-network quantum states is possible using X as direct inputs to the appropriate FNNs and Deep-Sets [149]. Still, it is advantageous to devise

new inputs that already capture a large portion of the correlations. As in Ref. [2], we employ a permutation-equivariant message-passing neural network (MPNN) to iteratively build correlations into new one-body and two-body features from the original “visible” features. The visible features are chosen to be

$$\mathbf{v}_i = (s_i), \quad (9.21)$$

$$\mathbf{v}_{ij} = (\mathbf{r}_{ij}, \|\mathbf{r}_{ij}\|, s_{ij}), \quad (9.22)$$

with the separation vectors $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$ and distances $\|\mathbf{r}_{ij}\| = r_{ij}$ replaced by their L -periodic surrogates

$$\mathbf{r}_{ij} \mapsto (\cos(2\pi\mathbf{r}_{ij}/L), \sin(2\pi\mathbf{r}_{ij}/L)), \quad (9.23)$$

$$\|\mathbf{r}_{ij}\| \mapsto \|\sin(\pi\mathbf{r}_{ij}/L)\|, \quad (9.24)$$

and the quantity $s_{ij} \equiv 2\delta_{s_i, s_j} - 1$ assigned a value of $+1$ for aligned spins and -1 for anti-aligned spins. Note that we have excluded explicit dependence on the particle positions \mathbf{r}_i in the visible one-body features, thereby enforcing translational invariance in the new features. Linear transformations are applied to and concatenated with each feature to obtain the initial hidden features

$$\mathbf{h}_i^{(0)} = (\mathbf{v}_i, A\mathbf{v}_i), \quad (9.25)$$

$$\mathbf{h}_{ij}^{(0)} = (\mathbf{v}_{ij}, B\mathbf{v}_{ij}). \quad (9.26)$$

The main purpose of the linear transformations is to preprocess the input data. Still, they also help simplify the implementation by keeping the dimension of the hidden features $\mathbf{h}_i^{(t)}$, $\mathbf{h}_{ij}^{(t)}$ constant for all t . In each iteration, $t = 1, \dots, T$ of the MPNN, information is exchanged between the one- and two-body streams through a so-called “message”

$$\mathbf{m}_{ij}^{(t)} = \mathbf{M}_t \left(\mathbf{h}_i^{(t-1)}, \mathbf{h}_j^{(t-1)}, \mathbf{h}_{ij}^{(t-1)} \right). \quad (9.27)$$

For a given particle i , relevant messages are collected and pooled together to destroy the ordering with respect to all other particles $j \neq i$,

$$\mathbf{m}_i^{(t)} = \mathbf{pool} \left(\{\mathbf{m}_{ij}^{(t)} \mid j \neq i\} \right). \quad (9.28)$$

The pooling operation **pool** collapses the order of the elements in the set it acts upon and produces a vector with the same dimension as an individual element. Throughout this work, we use logsumexp-pooling, the smooth variation of max-pooling.

The pairwise messages $\mathbf{m}_{ij}^{(t)}$ and the implied particle messages $\mathbf{m}_i^{(t)}$ are then used to update the hidden features

$$\mathbf{h}_i^{(t)} = \left(\mathbf{v}_i, \mathbf{F}_t \left(\mathbf{h}_i^{(t-1)}, \mathbf{m}_i^{(t)} \right) \right), \quad (9.29)$$

$$\mathbf{h}_{ij}^{(t)} = \left(\mathbf{v}_{ij}, \mathbf{G}_t \left(\mathbf{h}_{ij}^{(t-1)}, \mathbf{m}_{ij}^{(t)} \right) \right). \quad (9.30)$$

The functions \mathbf{M}_t , \mathbf{F}_t , and \mathbf{G}_t are all unique FNNs with the same output dimension as the linear preprocessors A and B . By incorporating concatenated skip connections to the visible features, we guarantee that the signal originating from the raw data remains discernible even as the MPNN depth T increases. Finally, we combine the resulting outputs $\mathbf{h}_i^{(T)}$ and $\mathbf{h}_{ij}^{(T)}$

into pairwise feature vectors

$$\mathbf{g}_{ij} = \left(\mathbf{h}_i^{(T)}, \mathbf{h}_j^{(T)}, \mathbf{h}_{ij}^{(T)} \right) \quad (9.31)$$

to feed into subsequent networks. The flow of information through the MPNN can be visualized in Fig. 9.2. Notice how the hidden features in a given layer depend on the hidden features of the previous layer *and* the original visible features.

For all our NQS, we use a Jastrow correlator based on a Deep-Set [149] to enforce permutation invariance over the set of all pairwise features

$$J(X) = \rho \left(\mathbf{pool} \left(\{\zeta(\mathbf{g}_{ij}) \mid i \neq j\} \right) \right). \quad (9.32)$$

Here, ρ and ζ are FNNs, and the pooling operation is the same as in Eq. (9.28). While many Jastrow functions are typically designed to satisfy Kato's cusp condition [150] for specific systems, we take a different approach and allow our neural networks to learn the cusp fully. The short-range behavior of the ground state is particularly important for the UFG, so leaving our NQS completely unbiased serves as an important test for evaluating the overall capabilities of NQS.

The Slater-Jastrow ansatz with plane wave orbitals (SJ-PW) does not require any additional neural networks beyond ρ and ζ , so it establishes a baseline for the number of trainable parameters in this work. On the other hand, the backflow variables \mathbf{u}_i and θ_i for the Slater-Jastrow ansatz with backflow orbitals (SJ-BF) are the outputs of another Deep-Set

$$(\text{Re}(\mathbf{u}_i), \text{Im}(\mathbf{u}_i), \theta_i) = \boldsymbol{\rho}_{bf} \left(\mathbf{pool} \left(\{\zeta_{bf}(\mathbf{g}_{ij}) \mid j \neq i\} \right) \right), \quad (9.33)$$

which is permutation invariant with respect to all $j \neq i$ by construction. The size of $\boldsymbol{\rho}_{bf}$

and ζ_{bf} determines the number of extra variational parameters present in the SJ-BF ansatz compared to the SJ-PW ansatz. For the PJ-BF ansatz, the pairing orbital ν in Eq. (9.7) simply takes \mathbf{g}_{ij} as input in place of $(\mathbf{x}_i, \mathbf{x}_j)$. Therefore, the number of additional variational parameters in the PJ-BF ansatz relative to the SJ-PW ansatz is determined by the size of ν .

All of the feedforward neural networks mentioned throughout this section have at least two hidden layers with 16 nodes each. The activation function is GELU [151] and the weights/biases are initialized with glorot normal/zeros unless pretrained parameters are used.

It is worth highlighting that the individual feedforward neural networks within our NQS are solely dependent on the spatial dimension d and not the system size N . Therefore, even though this study focuses on benchmarking the $N = 14$ case, the trained NQS can be used as starting points for larger even- N , unpolarized systems without requiring any modifications to the network structure. This is an example of transfer learning, a powerful strategy that involves applying knowledge gained from solving one problem to another, often more challenging problem.

9.4 Variational Monte Carlo and Training

We train our NQS by minimizing the energy

$$E(\mathbf{p}) \equiv \frac{\langle \Psi(\mathbf{p}) | H | \Psi(\mathbf{p}) \rangle}{\langle \Psi(\mathbf{p}) | \Psi(\mathbf{p}) \rangle} \quad (9.34)$$

with respect to the variational parameters \mathbf{p} . To compute the energy and its gradient $\nabla_{\mathbf{p}} E$ using Monte Carlo integration, we sample positions R and spins S from $|\Psi(R, S)|^2$ in a way

that preserves periodicity and total spin projection on the z -axis, as in Refs. [125, 131]. Since the ordering of the spins is not fixed, our ansätze can be immediately applied to any continuous-space Hamiltonian that exchange spin, such as Ref. [152].

A sophisticated optimization technique is critical for achieving an ansatz that is both compact and expressive. In this work, we employ the stochastic reconfiguration [6] (SR) algorithm with regularization based on the RMSprop method, introduced in Ref. [131]. The parameters are updated as

$$\mathbf{p} \leftarrow \mathbf{p} - \eta G^{-1} \nabla_{\mathbf{p}} E, \quad (9.35)$$

where η is a constant learning rate and G is the quantum geometric tensor [7].

Due to the strong and short-range nature of the interaction in Eq. (9.2), it is likely for the optimization process to get trapped in a local minimum when initialized with random parameters. To avoid this problem, we use transfer learning by pretraining the NQS on a softer interaction ($\mu = 5$) before proceeding to harder ones ($\mu = 10, 20, 40$). Not only does this approach improve the final converged energy, but the efficiency of the optimization process overall. The training for lower values of μ can handle a more aggressive learning rate δ and fewer samples. As a general guideline, we reduce δ by a factor of 10 and double the number of samples each time the value of μ is doubled. The number of optimization steps required for training ranges from $\mathcal{O}(10^3)$ to $\mathcal{O}(10^4)$, depending on whether the NQS were pretrained or initialized with random parameters.

9.5 Diffusion Monte Carlo

The fixed-node DMC calculations are performed as described in Ref. [153]. The initial state is prepared using VMC methods with a variational wave function with the same general

form as Eq. (9.3). Note that, within the fixed-node approximation, DMC provides a strict upperbound to the energy of the system. While DMC is a precise method, its accuracy relies on the choice of nodal surface and the quality of the preceding VMC calculation. The symmetric Jastrow factor is given by

$$J(X) = \sum_{ii'}^n u(r_{ii'}), \quad (9.36)$$

$$u(r) = K \tanh(\mu_J r) \cosh(\gamma r)/r, \quad (9.37)$$

where $n = N/2$ and the unprimed and primed indices denote the spin-up and spin-down particles, respectively. The parameters K and γ are adjusted so that $u(d) = 0$ and $u'(d) = 0$, and μ_J and d are variational parameters. Considering that the s -wave channel dominates the interaction, the antisymmetric part is given by the number-projected BCS wave function

$$\Phi_{BCS}(X) = \det \begin{bmatrix} \phi(\mathbf{r}_{11'}) & \phi(\mathbf{r}_{12'}) & \cdots & \phi(\mathbf{r}_{1n'}) \\ \phi(\mathbf{r}_{21'}) & \phi(\mathbf{r}_{22'}) & \cdots & \phi(\mathbf{r}_{2n'}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{r}_{n1'}) & \phi(\mathbf{r}_{n2'}) & \cdots & \phi(\mathbf{r}_{nn'}) \end{bmatrix}, \quad (9.38)$$

with the pairing orbitals

$$\phi(\mathbf{r}) = \tilde{\beta}(r) + \sum_i a(k_i^2) e^{i\mathbf{k}_i \cdot \mathbf{r}}, \quad (9.39)$$

$$\tilde{\beta}(r) = \beta(r) + \beta(L - r) - 2\beta(L/2), \quad (9.40)$$

$$\beta(r) = (1 + cbr) (1 - e^{-dbr}) \frac{e^{-br}}{dbr}. \quad (9.41)$$

The parameters $a(k_i^2)$, b and d are obtained by minimizing the energy, and c is chosen so that the function β has zero slope at the origin. If we instead let $\beta = 0$ and restrict the sum in Eq. (9.39) to momentum states filled up to k_F , the antisymmetric part is equivalent to the Slater determinant with single-particle plane waves as in Eqs. (9.8) and (9.9). Since this approach does not involve pairing, we will refer to the related DMC results as DMC-PW. Conversely, the approach that accounts for pairing will be identified as DMC-BCS.

It should be emphasized that the BCS wave function of Eq. (9.38) is a special case of the generalized Pfaffian of Eq. (9.6). In fact, it can be easily shown [5] that by only retaining the spin-singlet blocks, the calculation of the Pfaffian reduces to the determinant of spin-singlet block.

9.6 Results

9.6.1 Energy

We first compare the performance of the various neural-network quantum states outlined in Sec. 9.3 as the message-passing neural network (MPNN) depth T is varied. As shown in Fig. 9.3, the final converged energies per particle for the Slater-Jastrow ansatz with plane wave orbitals (SJ-PW) decreases monotonically towards the corresponding DMC-PW benchmark, with remarkable agreement at $T = 5$. This behavior echoes the findings of Ref. [154], and demonstrates the impact of the MPNN on the flexibility of our Jastrow. Incorporating backflow correlations into the Slater-Jastrow ansatz (SJ-BF) significantly improves results compared to the fixed-node approach with PW, but more than half of the discrepancy between the two DMC energies remains. Due to the observed weak dependence on T , it is unlikely that further increasing T would yield a substantial improvement in energy. The

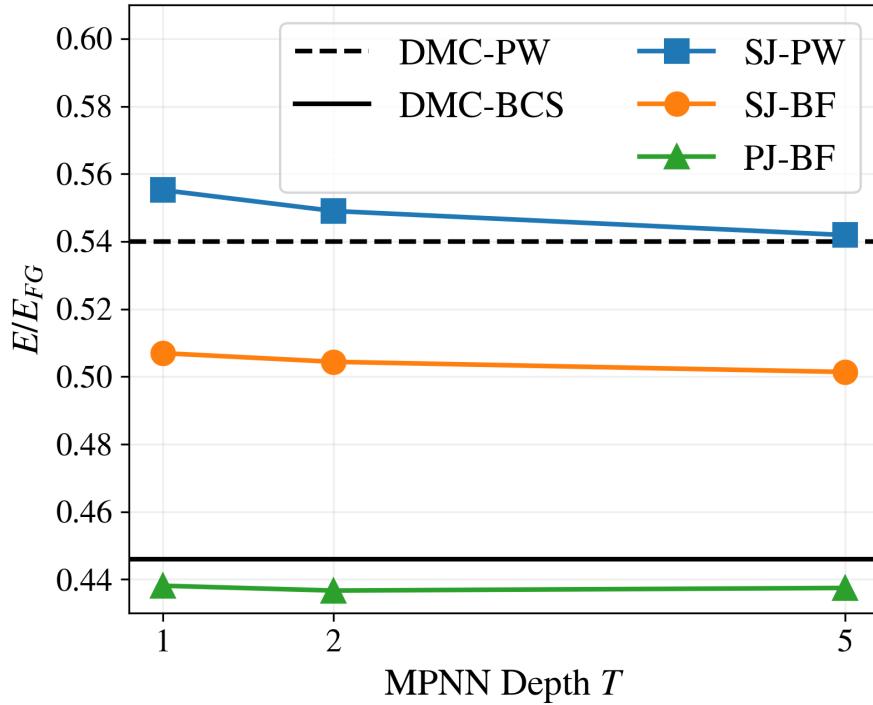


Figure 9.3: Ground-state energies per particle as a function of the MPNN depth T for the SJ-PW (blue squares), SJ-BF (orange circles), and PJ-BF (green triangles) ansätze. The interaction parameters are set to $v_0 = 1$ and $\mu = 5$, corresponding to an effective range of $r_e k_F = 0.4$. The DMC benchmark energies with and without pairing are displayed as solid and dashed lines, respectively.

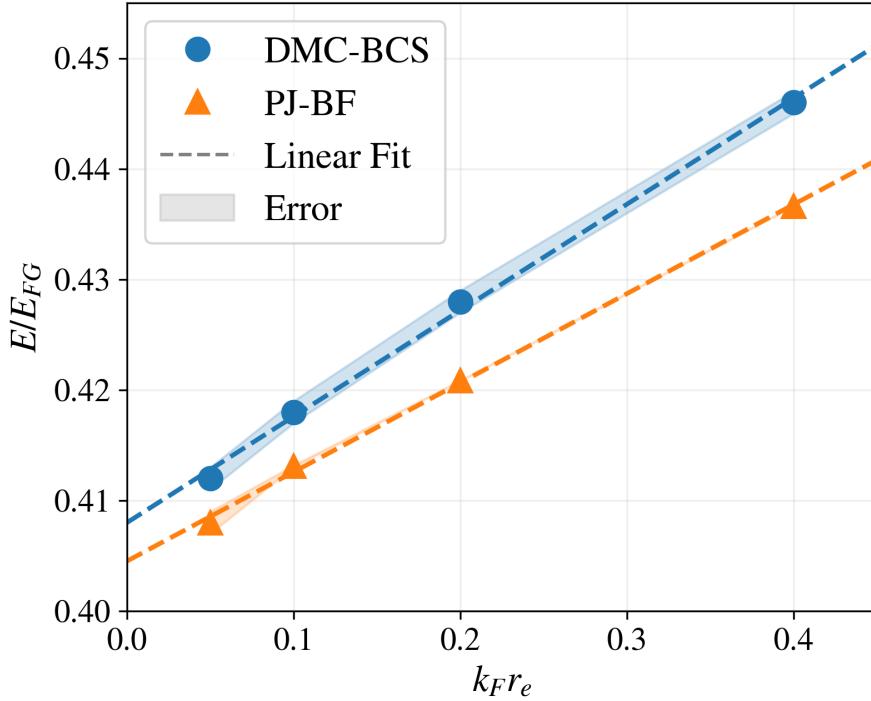


Figure 9.4: Ground-state energies per particle as a function of the effective range. The DMC-BCS benchmark energies (blue circles) and the Pfaffian-Jastrow with backflow (PJ-BF) energies (orange triangles) are extrapolated to zero effective range using linear fits (dashed lines). The shaded regions are the error bands for the DMC-BCS and PJ-BF energies.

SJ-BF ansatz may be able to achieve energies more similar to the DMC-BCS benchmark by increasing the width of the feedforward neural networks. Still, the associated computational expenses are expected to be prohibitively high. Therefore, we turn our attention to the Pfaffian-Jastrow-Backflow (PJ-BF) ansatz. Even with a single MPNN layer, the PJ-BF ansatz easily outperforms DMC-BCS while also possessing fewer parameters (~ 5600 v.s. ~ 6200) than the single-layer SJ-BF ansatz. The overall dependence on the MPNN depth is weak, with $T = 2$ giving slightly lower energy and variance than $T = 5$. For the remainder of our analysis, we will use the PJ-BF ansatz with $T = 2$, which contains about 8500 variational parameters.

μ	$k_{F}r_e$	DMC-BCS	PJ-BF
5	0.4	0.446(1)	0.4366(3)
10	0.2	0.428(1)	0.4208(3)
20	0.1	0.418(1)	0.4131(8)
40	0.05	0.412(1)	0.408(1)
∞	0.0	0.408(1)*	0.405(1)*

Table 9.1: Energy per particle for various values of μ and the corresponding values of r_e . The values with asterisks (*) are extrapolations from the linear fits shown in Fig. 9.4. The parameter $v_0 = 1$ is fixed.

As the unitary limit is characterized by a vanishing effective range, we study how the ground-state energy responds to changing $k_{F}r_e$ in Fig. 9.4. The PJ-BF ansatz gives energies \sim 1-2% lower than DMC-BCS as the effective range is decreased from $k_{F}r_e = 0.4$ to $k_{F}r_e = 0.1$. At $k_{F}r_e = 0.05$, our energy falls below the range of the DMC-BCS error band, suggesting our approach is likely to maintain its superior performance as r_e is decreased further. To estimate the energy at zero effective range, we also perform simple linear fits — See Table 9.1 for the extrapolated values. Note that our results have been obtained by simulating a system of $N = 14$ particles for benchmark purposes. In order to obtain energies closer to the thermodynamic limit, further simulations with more particles will be needed [122, 155].

9.6.2 Pair Distribution Functions

In Fig. 9.5, we show the opposite-spin pair distribution functions at unitarity for $\mu = 5, 10$, and 20. Notice how the peaks of the distributions at $k_Fr = 0$ grow roughly quadratically with μ , demonstrating the presence of strong pairing correlations as we approach the unitary limit $\mu \rightarrow \infty$. Clearly, the short-range character of the distributions are important to capture at unitarity, as they begin to converge around $k_Fr \gtrsim 0.4$.

Fig. 9.6 presents a complementary set of opposite-spin pair distribution functions in the crossover region with fixed effective range of $k_{F}r_e = 0.2$. When we lean towards the BCS

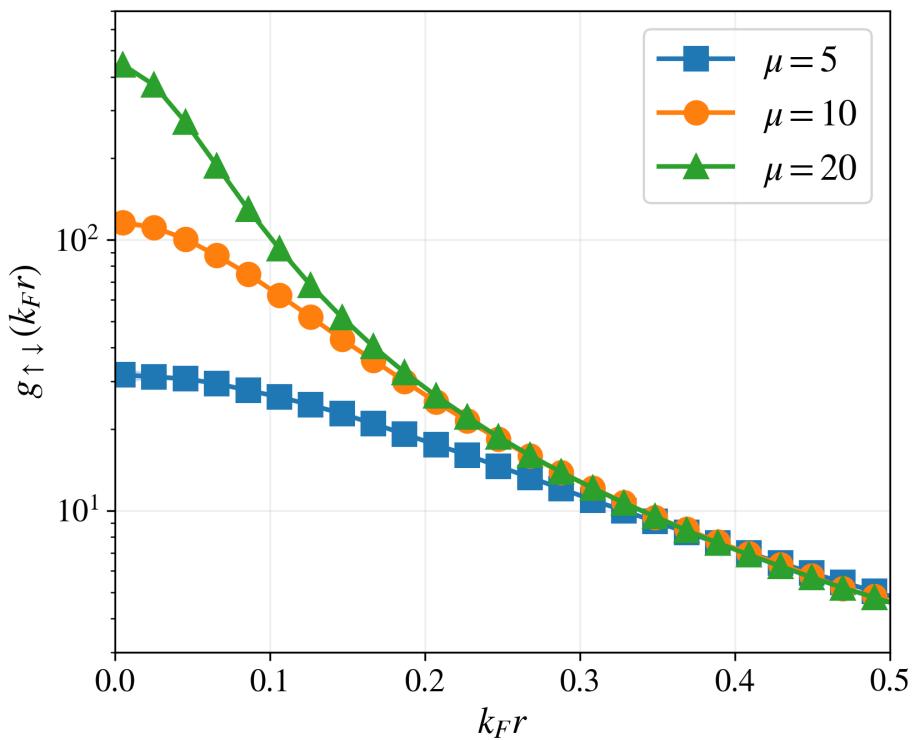


Figure 9.5: Opposite-spin pair densities as a function of small $k_F r$ at unitarity ($v_0 = 1$) and $\mu = 5$ (blue squares), $\mu = 10$ (orange circles), and $\mu = 20$ (green triangles).

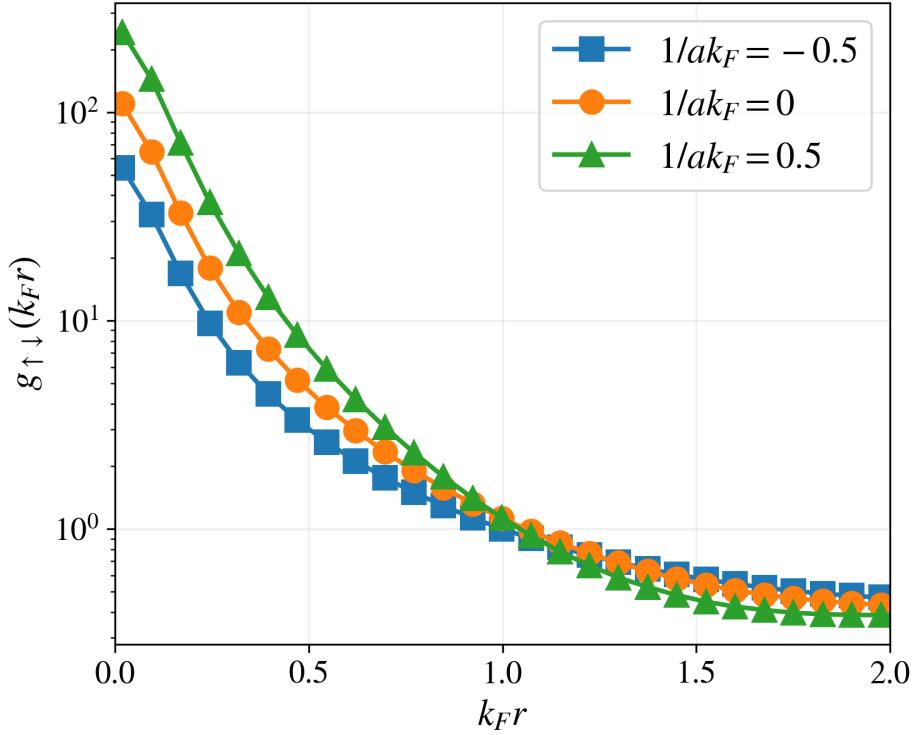


Figure 9.6: Opposite-spin pair densities in the crossover region for the BCS phase $1/ak_F = -0.5$ (blue squares), unitarity $1/ak_F = 0$ (orange circles), and BEC phase $1/ak_F = 0.5$ (green triangles). The effective range of all cases are fixed $k_F r_e = 0.2$. See Table 9.3 for the corresponding values of v_0 and μ .

$1/ak_F$	v_0	μ	DMC-BCS	PJ-BF
-1	0.879214	11.06247	0.801(1)	0.7930(2)
-0.5	0.933216	10.55715	0.705(1)	0.6937(3)
-0.2	0.971423	10.23012	0.578(1)	0.5671(3)
-0.1	0.985366	10.11637	0.510(1)	0.5014(4)
0	1.0	10.0	0.428(1)	0.4208(3)
0.1	1.015388	9.880801	0.328(1)	0.3218(3)
0.2	1.031602	9.758564	0.208(1)	0.2017(3)
0.5	1.086081	9.371025	-0.319(1)	-0.3244(4)
1	1.204354	8.632898	-2.053(1)	-2.0566(6)

Table 9.2: Energies per particle and parameters for the two-body potential in Eq. (9.2) giving different scattering lengths with the same effective range $k_F r_e = 0.2$.

phase $1/ak_F = -0.5$, the long-range tail of the density is enhanced compared to the unitary case $1/ak_F = 0$. On the other hand, the tail is diminished in the BEC phase $1/ak_F = -0.5$, suggesting the initiation of dimer formation. The differences in the peaks of the distributions are not as dramatic as in Fig. 9.5, but they are consistent with the expected behavior in the BCS and BEC regimes near unitarity.

Finally, we explore the BCS-BEC crossover region for a fixed effective range $k_F r_e = 0.2$ in Fig. 9.7. See Table 9.3 for the values of the interaction parameters v_0 and μ , as well as the corresponding DMC-BCS benchmarks and the PJ-BF ansatz results. The cases closer to unitarity were used to pretrain the cases further away. In the BCS regime, our PJ-BF ansatz consistently yields energies $\sim 0.01E_{FG}$ lower than those obtained from DMC-BCS, albeit with slightly inferior performance in the BEC regime. We attribute this effect to the need for increased flexibility in capturing the short-range behavior of pairs in the BEC regime. Simply increasing the size of feedforward neural network ν that defines the pairing orbital should alleviate this discrepancy. In any case, the PJ-BF ansatz gives lower energies than DMC-BCS for all scattering lengths tested.

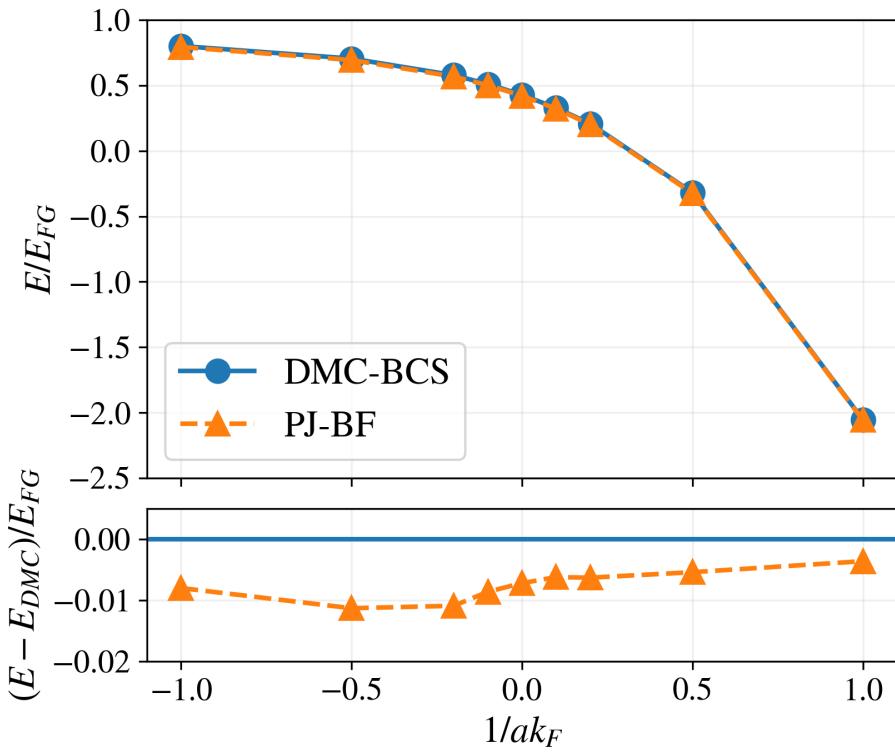


Figure 9.7: Upper panel: Energy per particle in the BCS-BEC crossover region as a function of the scattering length a for a fixed effective range $k_F r_e = 0.2$. Lower panel: Difference between Pfaffian-Jastrow with backflow (PJ-BF) and DMC-BCS benchmark energies. See Table 9.3 for the corresponding values of v_0 and μ .

9.6.3 Pairing Gap

N	DMC-BCS	PJ-BF
14	0.428(1)	0.4208(3)
15	0.4900(2)	0.4766(5)
15	0.5357(2)*	0.5209(5)*
16	0.4240(2)	0.4177(4)

Table 9.3: Energies per particle for different numbers of particles, with $k_F r_e = 0.2$. Values with asterisks (*) indicate translation-invariant calculations.

The pairing gap can be evaluated using

$$\Delta(N) = E(N) - \frac{1}{2} \left(E(N+1) + E(N-1) \right), \quad (9.42)$$

where N is taken to be odd and $E(N)$ denotes the *total* energy for N particles. Since our construction of the Pfaffian-Jastrow-Backflow ansatz remains independent of the system size N (assuming N is even), the trained ansatz for $N = 14$ serves as the initial state for the $N = 16$ calculation. To accommodate odd- N cases, we incorporated an additional neural network to represent the unpaired single-particle orbital. Furthermore, for the $N = 15$ calculation, the pairing orbital from the $N = 16$ calculation was utilized as the starting point.

In Table 9.4, we present two different evaluations of the pairing gap, corresponding to the DMC-BCS results using a single-particle orbital with and without translation invariance. To enforce translation invariance in our odd- N PJ-BF ansatz, we simply take the one-body output of the message-passing neural network $\mathbf{h}_i^{(T)}$ as input to the single-particle pairing orbital. To break translation invariance, we concatenate the positions $(\mathbf{r}_i, \mathbf{h}_i^{(T)})$ with the MPNN output, keeping in mind the spins have already been included in the latter. Due to our lower energies for $N = 15$ compared to DMC-BCS, we predict a pairing gap that

is $\sim 0.10E_{FG}$ smaller in the non-translation invariant case and $\sim 0.12E_{FG}$ smaller in the translation invariant case.

$\Delta(15)/E_{FG}$	DMC-BCS	PJ-BF
Non-translation invariant	0.9620	0.8618
Translation invariant	1.6475	1.5263

Table 9.4: Energies per particle for different numbers of particles, with $k_F r_e = 0.2$. Values with asterisks (*) indicate translation-invariant calculations.

9.7 Conclusions and perspectives

In this study, we propose a novel neural-network quantum state based on the Pfaffian-Jastrow (PJ) framework that utilizes a message-passing neural network (MPNN) to encode pairing and backflow (BF) correlations. We evaluate its performance against comparable Slater-Jastrow (SJ) ansätze with identical MPNN architectures. Our results indicate that increasing the depth of the MPNN systematically improves the performance of the SJ ansätze, but backflow correlations within the single-particle picture are still insufficient in capturing all pairing correlations. However, we demonstrate that a simple and compact PJ-BF ansatz surpasses the DMC-BCS benchmark with ease.

Transfer learning has proven to be an essential tool in this work. It enables the realization of the unitary limit in a controlled manner, mitigating the risk of becoming trapped in local minima. It also allows for the efficient exploration of regions beyond unitarity, unlocking new avenues for studying the BCS-BEC crossover. Transfer learning will remain a crucial part of our training procedure as we move to larger systems. All unpolarized systems can be treated with a single architecture, while the $N \pm 1$ systems can be treated by introducing one additional FNN to represent the unpaired single-particle orbital. This modification is

straightforward to implement, making the calculation of the gap in the thermodynamic limit more accessible and enabling further advancements in our work.

Besides calculating the gap for larger systems, our next steps include a direct comparison with the STU Pfaffian wave function of Ref. [5]. We also plan to perform a more careful extrapolation to the $r_e \rightarrow 0$ limit since we have used relatively large values of $k_F r_e$ for this initial investigation. However, more hyperparameter tuning will be needed, especially about the width of the hidden layers, since the smaller values of r_e will require more flexibility.

Our Pfaffian-Jastrow-Backflow NQS displays immense potential in the study of ultracold Fermi gases. Unlike conventional methods, our PJ-BF ansatz is not subject to biases arising from physical intuition or a lack thereof, as it does not require specifying a particular form for the pairing orbitals. For this reason, it can be readily applied to other strongly-correlated systems, including molecules and other strongly-correlated quantum systems. In stark contrast to the commonly used geminal wave function, our ansatz does not rely on ordering the spin of the interacting fermions, and it is therefore amenable to Hamiltonians that exchange spin, such as those modeling nuclear dynamics. In this regard, we anticipate calculations of atomic nuclei and low-density isospin-asymmetric nucleonic matter and carry out detailed investigations on the nature of nuclear pairing [156].

When the stochastic reconfiguration algorithm and transfer learning techniques are combined with the enforcement of translational, parity, and time-reversal symmetries, highly non-perturbative correlations can be encoded in a small number of parameters by modern standards. This approach will pave the way for future developments in the study of many-body systems, as it offers a powerful tool for encoding correlations in a compact and computationally feasible manner.

Note Added: A work very recently appeared in pre-print [157] introduces neural backflow transformations in a geminal wave function and studies the unitary Fermi gas. We leave systematic comparisons between the two approaches to future works while already observing that the Pfaffian wave function is a strict generalization of the geminal wave function [158, 5].

9.8 Acknowledgments

The work of JK and MHJ is supported by the U.S. National Science Foundation Grants No. PHY-1404159 and PHY-2013047. AL and BF are supported by the U.S. Department of Energy, Office of Science, Office of Nuclear Physics, under contracts DE-AC02-06CH11357, by the 2020 DOE Early Career Award number ANL PRJ1008597, by the NUCLEI SciDAC program, and Argonne LDRD awards. The work of SG is supported by the U.S. Department of Energy, Office of Science, Office of Nuclear Physics, under contract No. DE-AC52-06NA25396, by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Scientific Discovery through Advanced Computing (SciDAC) NUCLEI program, and by the Department of Energy Early Career Award Program. The work of GP, JN, GC is supported by the Swiss National Science Foundation under Grant No. 200021_200336, and by Microsoft Research. Computer time was provided by the Los Alamos National Laboratory Institutional Computing Program, which is supported by the U.S. Department of Energy National Nuclear Security Administration under Contract No. 89233218CNA000001.

Bibliography

- [1] Giuseppe Carleo and Matthias Troyer. “Solving the quantum many-body problem with artificial neural networks”. In: *Science* 355.6325 (Feb. 2017), pp. 602–606. ISSN: 0036-8075, 1095-9203. DOI: 10.1126/science.aag2302. URL: <http://science.sciencemag.org/content/355/6325/602>.
- [2] Gabriel Pescia et al. *Message-Passing Neural Quantum States for the Homogeneous Electron Gas*. 2023. arXiv: 2305.07240 [quant-ph].
- [3] David Pfau et al. “Ab initio solution of the many-electron Schrödinger equation with deep neural networks”. In: *Phys. Rev. Res.* 2 (3 Sept. 2020), p. 033429. DOI: 10.1103/PhysRevResearch.2.033429. URL: <https://link.aps.org/doi/10.1103/PhysRevResearch.2.033429>.
- [4] Alessandro Lovato et al. “Hidden-nucleons neural-network quantum states for the nuclear many-body problem”. In: *Physical Review Research* 4.4 (2022), p. 043178.
- [5] M. Bajdich et al. “Pfaffian pairing and backflow wavefunctions for electronic structure quantum Monte Carlo methods”. In: *Phys. Rev. B* 77 (11 Mar. 2008), p. 115112. DOI: 10.1103/PhysRevB.77.115112. URL: <https://link.aps.org/doi/10.1103/PhysRevB.77.115112>.
- [6] Sandro Sorella. “Wave function optimization in the variational Monte Carlo method”. In: *Phys. Rev. B* 71 (24 June 2005), p. 241103. DOI: 10.1103/PhysRevB.71.241103. URL: <https://link.aps.org/doi/10.1103/PhysRevB.71.241103>.
- [7] James Stokes et al. “Quantum Natural Gradient”. In: *Quantum* 4 (May 2020), p. 269. ISSN: 2521-327X. DOI: 10.22331/q-2020-05-25-269. URL: <http://dx.doi.org/10.22331/q-2020-05-25-269>.
- [8] M. P. Kuchera et al. “Machine Learning Methods for Track Classification in the AT-TPC”. In: *Nucl. Instrum. Meth. Phys. Res. A* 940 (2019), pp. 156–167. DOI: 10.1016/j.nima.2019.05.097.
- [9] Manzil Zaheer et al. “Deep Sets”. In: *arXiv e-prints*, arXiv:1703.06114 (Mar. 2017), arXiv:1703.06114. arXiv: 1703.06114 [cs.LG].
- [10] M. Wimmer. “Algorithm 923: Efficient Numerical Computation of the Pfaffian for Dense and Banded Skew-Symmetric Matrices”. In: *ACM Trans. Math. Softw.* 38.4 (Aug. 2012). ISSN: 0098-3500. DOI: 10.1145/2331130.2331138. URL: <https://doi.org/10.1145/2331130.2331138>.

- [11] Hiroki Saito. “Method to Solve Quantum Few-Body Problems with Artificial Neural Networks”. In: *J. Phys. Soc. Jpn.* 87.7 (2018), p. 074002. DOI: 10.7566/JPSJ.87.074002. URL: <https://doi.org/10.7566/JPSJ.87.074002>.
- [12] B. P. Abbott et al. “GW170817: Observation of Gravitational Waves from a Binary Neutron Star Inspiral”. In: *Phys. Rev. Lett.* 119.16 (2017), p. 161101. DOI: 10.1103/PhysRevLett.119.161101. arXiv: 1710.05832 [gr-qc].
- [13] B. P. Abbott et al. “Multi-messenger Observations of a Binary Neutron Star Merger”. In: *Astrophys. J.* 848.2 (2017), p. L12. DOI: 10.3847/2041-8213/aa91c9. arXiv: 1710.05833 [astro-ph.HE].
- [14] Andrea Sabatucci and Omar Benhar. “Tidal Deformation of Neutron Stars from Microscopic Models of Nuclear Dynamics”. In: *Phys. Rev. C* 101.4 (2020), p. 045807. DOI: 10.1103/PhysRevC.101.045807. arXiv: 2001.06294 [nucl-th].
- [15] Peter Senger. “Probing Dense Nuclear Matter in the Laboratory: Experiments at FAIR and NICA”. In: *Universe* 7.6 (2021). ISSN: 2218-1997. DOI: 10.3390/universe7060171. URL: <https://www.mdpi.com/2218-1997/7/6/171>.
- [16] C. Drischler et al. “Neutron matter from chiral two- and three-nucleon calculations up to N^3LO ”. In: *Phys. Rev. C* 94.5 (2016), p. 054307. DOI: 10.1103/PhysRevC.94.054307. arXiv: 1608.05615 [nucl-th].
- [17] M. Piarulli et al. “Benchmark calculations of pure neutron matter with realistic nucleon-nucleon interactions”. In: *Phys. Rev. C* 101.4 (2020), p. 045801. DOI: 10.1103/PhysRevC.101.045801. arXiv: 1908.04426 [nucl-th].
- [18] D. Lonardoni et al. “Nuclear and neutron-star matter from local chiral interactions”. In: *Phys. Rev. Res.* 2.2 (2020), p. 022033. DOI: 10.1103/PhysRevResearch.2.022033. arXiv: 1912.09411 [nucl-th].
- [19] W. G. Jiang et al. “Accurate bulk properties of nuclei from $A = 2$ to ∞ from potentials with Δ isobars”. In: *Phys. Rev. C* 102.5 (2020), p. 054301. DOI: 10.1103/PhysRevC.102.054301. arXiv: 2006.16774 [nucl-th].
- [20] Francesca Sammarruca and Randy Millerson. “Overview of symmetric nuclear matter properties from chiral interactions up to fourth order of the chiral expansion”. In: *Phys. Rev. C* 104.6 (2021), p. 064312. DOI: 10.1103/PhysRevC.104.064312. arXiv: 2109.01985 [nucl-th].
- [21] H. Heiselberg and M. Hjorth-Jensen. “Phases of dense matter in neutron stars”. In: *Phys. Rep.* 328.5 (2000), p. 237. ISSN: 0370-1573. DOI: [https://doi.org/10.1016/S0370-1573\(99\)00110-6](https://doi.org/10.1016/S0370-1573(99)00110-6). URL: <https://www.sciencedirect.com/science/article/pii/S0370157399001106>.
- [22] Andrea Sabatucci et al. “Sensitivity of neutron star observations to three-nucleon forces”. In: *Phys. Rev. D* 106.8 (2022), p. 083010. DOI: 10.1103/PhysRevD.106.083010. arXiv: 2206.11286 [astro-ph.HE].

- [23] Armen Sedrakian, John W Clark, and Mark Alford. *Pairing in Fermionic Systems*. WORLD SCIENTIFIC, 2006. DOI: 10.1142/6194. eprint: <https://www.worldscientific.com/doi/pdf/10.1142/6194>. URL: <https://www.worldscientific.com/doi/abs/10.1142/6194>.
- [24] Omar Benhar and Giulia De Rosi. “Superfluid Gap in Neutron Matter from a Microscopic Effective Interaction”. In: *J. Low Temp. Phys.* 189.5-6 (2017), pp. 250–261. DOI: 10.1007/s10909-017-1823-x. arXiv: 1705.06607 [nucl-th].
- [25] D. J. Dean and M. Hjorth-Jensen. “Pairing in nuclear systems: from neutron stars to finite nuclei”. In: *Rev. Mod. Phys.* 75 (2003), p. 607. ISSN: 00346861. DOI: 10.1103/RevModPhys.75.607. URL: <https://journals.aps.org/rmp/abstract/10.1103/RevModPhys.75.607>.
- [26] Dima G. Yakovlev and C. J. Pethick. “Neutron star cooling”. In: *Ann. Rev. Astron. Astrophys.* 42 (2004), pp. 169–210. DOI: 10.1146/annurev.astro.42.053102.134013. arXiv: astro-ph/0402143.
- [27] Dany Page et al. “Rapid Cooling of the Neutron Star in Cassiopeia A Triggered by Neutron Superfluidity in Dense Matter”. In: *Phys. Rev. Lett.* 106 (2011), p. 081101. DOI: 10.1103/PhysRevLett.106.081101. arXiv: 1011.6142 [astro-ph.HE].
- [28] C. Monrozeau, J. Margueron, and N. Sandulescu. “Nuclear superfluidity and cooling time of neutron-star crust”. In: *Phys. Rev. C* 75 (2007), p. 065807. DOI: 10.1103/PhysRevC.75.065807. arXiv: nucl-th/0703064.
- [29] Frédéric Nowacki, Alexandre Obertelli, and Alfredo Poves. “The neutron-rich edge of the nuclear landscape: Experiment and theory.” In: *Prog. Part. Nucl. Phys.* 120 (2021), p. 103866. DOI: 10.1016/j.ppnp.2021.103866. arXiv: 2104.06238 [nucl-th].
- [30] J. Carlson et al. “Quantum Monte Carlo methods for nuclear physics”. In: *Rev. Mod. Phys.* 87 (2015), p. 1067. DOI: 10.1103/RevModPhys.87.1067. arXiv: 1412.3081 [nucl-th].
- [31] K. E. Schmidt and S. Fantoni. “A quantum Monte Carlo method for nucleon systems”. In: *Phys. Lett. B* 446 (1999), pp. 99–103. DOI: 10.1016/S0370-2693(98)01522-6.
- [32] A. Lovato et al. “Benchmark calculations of infinite neutron matter with realistic two- and three-nucleon potentials”. In: *Phys. Rev. C* 105.5 (2022), p. 055808. DOI: 10.1103/PhysRevC.105.055808. arXiv: 2202.10293 [nucl-th].
- [33] S. Gandolfi et al. “Equation of state of superfluid neutron matter and the calculation of S(0)-1 pairing gap”. In: *Phys. Rev. Lett.* 101 (2008), p. 132501. DOI: 10.1103/PhysRevLett.101.132501. arXiv: 0805.2513 [nucl-th].
- [34] Stefano Gandolfi et al. “The 1S0 Pairing Gap in Neutron Matter”. In: *Condens. Mat.* 7.1 (2022), p. 19. DOI: 10.3390/condmat7010019. arXiv: 2201.01308 [nucl-th].
- [35] M. Bajdich et al. “Pfaffian pairing wave functions in electronic structure quantum Monte Carlo”. In: *Phys. Rev. Lett.* 96 (2006), p. 130201. DOI: 10.1103/PhysRevLett.96.130201. arXiv: cond-mat/0512327.

- [36] J. W. T. Keeble and A. Rios. “Machine learning the deuteron”. In: *Phys. Lett. B* 809 (2020), p. 135743. DOI: 10.1016/j.physletb.2020.135743. arXiv: 1911.13092 [nucl-th].
- [37] Corey Adams et al. “Variational Monte Carlo Calculations of $A \leq 4$ Nuclei with an Artificial Neural-Network Correlator Ansatz”. In: *Phys. Rev. Lett.* 127.2 (2021), p. 022502. DOI: 10.1103/PhysRevLett.127.022502. arXiv: 2007.14282 [nucl-th].
- [38] Alex Gnech et al. “Nuclei with up to $A = 6$ nucleons with artificial neural network wave functions”. In: *Few Body Syst.* 63.1 (2022), p. 7. DOI: 10.1007/s00601-021-01706-0. arXiv: 2108.06836 [nucl-th].
- [39] A. Lovato et al. “Hidden-nucleons neural-network quantum states for the nuclear many-body problem”. In: (June 2022). arXiv: 2206.10021 [nucl-th].
- [40] Y. L. Yang and P. W. Zhao. “A consistent description of the relativistic effects and three-body interactions in atomic nuclei”. In: *Phys. Lett. B* 835 (2022), p. 137587. DOI: 10.1016/j.physletb.2022.137587. arXiv: 2206.13208 [nucl-th].
- [41] Mauro Rigo et al. “Solving the nuclear pairing model with neural network quantum states”. In: (Nov. 2022). arXiv: 2211.04614 [nucl-th].
- [42] Jan Hermann, Zeno Schätzle, and Frank Noé. “Deep-neural-network solution of the electronic Schrödinger equation”. In: *Nature Chemistry* 12.10 (Sept. 2020), pp. 891–897. DOI: 10.1038/s41557-020-0544-y.
- [43] David Pfau et al. “Ab initio solution of the many-electron Schrödinger equation with deep neural networks”. In: *Physical Review Research* 2.3, 033429 (Sept. 2020), p. 033429. DOI: 10.1103/PhysRevResearch.2.033429. arXiv: 1909.02487 [physics.chem-ph].
- [44] L. Brualla et al. “Spin orbit induced backflow in neutron matter with auxiliary field diffusion Monte Carlo”. In: *Phys. Rev. C* 67 (2003), p. 065806. DOI: 10.1103/PhysRevC.67.065806. arXiv: nucl-th/0304042.
- [45] R. Schiavilla et al. “Two- and three-nucleon contact interactions and ground-state energies of light- and medium-mass nuclei”. In: *Phys. Rev. C* 103.5 (2021), p. 054003. DOI: 10.1103/PhysRevC.103.054003. arXiv: 2102.02327 [nucl-th].
- [46] Sebastian König et al. “Nuclear Physics Around the Unitarity Limit”. In: *Phys. Rev. Lett.* 118.20 (2017), p. 202501. DOI: 10.1103/PhysRevLett.118.202501. arXiv: 1607.04623 [nucl-th].
- [47] A. Kievsky et al. “Correlations imposed by the unitary limit between few-nucleon systems, nuclear matter and neutron stars”. In: *Phys. Rev. Lett.* 121.7 (2018), p. 072701. DOI: 10.1103/PhysRevLett.121.072701. arXiv: 1806.02636 [nucl-th].
- [48] Robert B. Wiringa, V. G. J. Stoks, and R. Schiavilla. “An Accurate nucleon-nucleon potential with charge independence breaking”. In: *Phys. Rev. C* 51 (1995), pp. 38–51. DOI: 10.1103/PhysRevC.51.38. arXiv: nucl-th/9408016.
- [49] B. S. Pudliner et al. “Quantum Monte Carlo calculations of $A \leq 6$ nuclei”. In: *Phys. Rev. Lett.* 74 (1995), pp. 4396–4399. DOI: 10.1103/PhysRevLett.74.4396. arXiv: nucl-th/9502031.

- [50] A. Akmal, V. R. Pandharipande, and D. G. Ravenhall. “The Equation of state of nucleon matter and neutron star structure”. In: *Phys. Rev. C* 58 (1998), pp. 1804–1828. DOI: 10.1103/PhysRevC.58.1804. arXiv: nucl-th/9804027.
- [51] R. Machleidt and D.R. Entem. “Chiral effective field theory and nuclear forces”. In: *Phys. Rep.* 503.1 (2011), p. 1. ISSN: 0370-1573. DOI: <https://doi.org/10.1016/j.physrep.2011.02.001>. URL: <https://www.sciencedirect.com/science/article/pii/S0370157311000457>.
- [52] Javier Robledo Moreno et al. “Fermionic wave functions from neural-network constrained hidden states”. In: *Proceedings of the National Academy of Sciences of the United States of America* 119 (32 Aug. 2022), e2122059119. ISSN: 10916490. DOI: 10.1073/PNAS.2122059119.SUPPL_FILE/PNAS.2122059119.SAPP.PDF. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.2122059119>.
- [53] Edward Wagstaff et al. “On the Limitations of Representing Functions on Sets”. In: *arXiv e-prints*, arXiv:1901.09006 (Jan. 2019), arXiv:1901.09006. arXiv: 1901.09006 [cs.LG].
- [54] Di Luo and Bryan K. Clark. “Backflow Transformations via Neural Networks for Quantum Many-Body Wave Functions”. In: *Phys. Rev. Lett.* 122 (22 June 2019), p. 226401. DOI: 10.1103/PhysRevLett.122.226401. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.122.226401>.
- [55] Gabriel Pescia et al. “Neural-network quantum states for periodic systems in continuous space”. In: *Phys. Rev. Res.* 4.2 (2022), p. 023138. DOI: 10.1103/PhysRevResearch.4.023138. arXiv: 2112.11957 [quant-ph].
- [56] A. Sarsa et al. “Neutron matter at zero temperature with auxiliary field diffusion Monte Carlo”. In: *Phys. Rev. C* 68 (2003), p. 024308. DOI: 10.1103/PhysRevC.68.024308. arXiv: nucl-th/0303035 [nucl-th].
- [57] Nicholas Metropolis et al. “Equation of State Calculations by Fast Computing Machines”. In: *The Journal of Chemical Physics* 21.6 (June 1953), pp. 1087–1092. DOI: 10.1063/1.1699114.
- [58] S. Gandolfi et al. “Equation of state of low-density neutron matter, and the 1S0 pairing gap”. In: *Physical Review C - Nuclear Physics* 80 (4 Oct. 2009), p. 045802. ISSN: 1089490X. DOI: 10.1103/PHYSREVC.80.045802/FIGURES/6/MEDIUM. URL: <https://journals.aps.org/prc/abstract/10.1103/PhysRevC.80.045802>.
- [59] M. Piarulli et al. “Light-nuclei spectra from chiral dynamics”. In: *Phys. Rev. Lett.* 120.5 (2018), p. 052503. DOI: 10.1103/PhysRevLett.120.052503. arXiv: 1707.02883 [nucl-th].
- [60] John W. Negele and D. Vautherin. “Neutron star matter at subnuclear densities”. In: *Nucl. Phys. A* 207 (1973), pp. 298–320. DOI: 10.1016/0375-9474(73)90349-7.
- [61] Brian M Austin, Dmitry Yu Zubarev, and William A Lester Jr. “Quantum Monte Carlo and related approaches”. In: *Chemical reviews* 112.1 (2012), pp. 263–288.
- [62] J Carlson et al. “Quantum Monte Carlo methods for nuclear physics”. In: *Reviews of Modern Physics* 87.3 (2015), p. 1067.

- [63] Anthony J. Leggett. “A theoretical description of the new phases of liquid ^3He ”. In: *Rev. Mod. Phys.* 47 (2 Apr. 1975), pp. 331–414. DOI: 10.1103/RevModPhys.47.331. URL: <https://link.aps.org/doi/10.1103/RevModPhys.47.331>.
- [64] Michael Tinkham. *Introduction to superconductivity*. Courier Corporation, 2004.
- [65] Juan Carrasquilla. “Machine learning for quantum matter”. In: *Advances in Physics: X* 5.1 (2020), p. 1797528.
- [66] Kenny Choo et al. “Symmetries and Many-Body Excitations with Neural-Network Quantum States”. In: *Phys. Rev. Lett.* 121 (16 Oct. 2018), p. 167204. DOI: 10.1103/PhysRevLett.121.167204. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.121.167204>.
- [67] Francesco Ferrari, Federico Becca, and Juan Carrasquilla. “Neural Gutzwiller-projected variational wave functions”. In: *Phys. Rev. B* 100 (12 Sept. 2019), p. 125131. DOI: 10.1103/PhysRevB.100.125131. URL: <https://link.aps.org/doi/10.1103/PhysRevB.100.125131>.
- [68] Mohamed Hibat-Allah et al. “Recurrent neural network wave functions”. In: *Phys. Rev. Res.* 2 (2 June 2020), p. 023358. DOI: 10.1103/PhysRevResearch.2.023358. URL: <https://link.aps.org/doi/10.1103/PhysRevResearch.2.023358>.
- [69] Marin Bukov, Markus Schmitt, and Maxime Dupont. “Learning the ground state of a non-stoquastic quantum Hamiltonian in a rugged neural network landscape”. In: *SciPost Phys.* 10 (2021), p. 147. DOI: 10.21468/SciPostPhys.10.6.147. URL: <https://scipost.org/10.21468/SciPostPhys.10.6.147>.
- [70] Yusuke Nomura and Masatoshi Imada. “Dirac-Type Nodal Spin Liquid Revealed by Refined Quantum Many-Body Solver Using Neural-Network Wave Function, Correlation Ratio, and Level Spectroscopy”. In: *Phys. Rev. X* 11 (3 Aug. 2021), p. 031034. DOI: 10.1103/PhysRevX.11.031034. URL: <https://link.aps.org/doi/10.1103/PhysRevX.11.031034>.
- [71] Nikita Astrakhantsev et al. “Broken-Symmetry Ground States of the Heisenberg Model on the Pyrochlore Lattice”. In: *Phys. Rev. X* 11 (4 Oct. 2021), p. 041021. DOI: 10.1103/PhysRevX.11.041021. URL: <https://link.aps.org/doi/10.1103/PhysRevX.11.041021>.
- [72] Dong-Ling Deng, Xiaopeng Li, and S Das Sarma. “Quantum entanglement in neural network states”. In: *Physical Review X* 7.2 (2017), p. 021021.
- [73] Yoav Levine et al. “Quantum Entanglement in Deep Learning Architectures”. In: *Phys. Rev. Lett.* 122 (6 Feb. 2019), p. 065301. DOI: 10.1103/PhysRevLett.122.065301. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.122.065301>.
- [74] Di Luo and Bryan K. Clark. “Backflow Transformations via Neural Networks for Quantum Many-Body Wave Functions”. In: *Phys. Rev. Lett.* 122 (22 June 2019), p. 226401. DOI: 10.1103/PhysRevLett.122.226401. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.122.226401>.

- [75] Kenny Choo, Antonio Mezzacapo, and Giuseppe Carleo. “Fermionic neural-network states for ab-initio electronic structure”. In: *Nature communications* 11.1 (2020), pp. 1–7.
- [76] Javier Robledo Moreno et al. “Fermionic wave functions from neural-network constrained hidden states”. In: *Proceedings of the National Academy of Sciences* 119.32 (2022), e2122059119.
- [77] Jan Hermann, Zeno Schätzle, and Frank Noé. “Deep-neural-network solution of the electronic Schrödinger equation”. In: *Nature Chemistry* 12.10 (2020), pp. 891–897.
- [78] MT Entwistle et al. “Electronic excited states in deep variational Monte Carlo”. In: *Nature Communications* 14.1 (2023), p. 274.
- [79] Gabriel Pescia et al. “Neural-network quantum states for periodic systems in continuous space”. In: *Physical Review Research* 4.2 (2022), p. 023138.
- [80] Max Wilson et al. “Wave function Ansatz (but Periodic) Networks and the Homogeneous Electron Gas”. In: *arXiv preprint arXiv:2202.04622* (2022).
- [81] Gino Cassella et al. “Discovering Quantum Phase Transitions with Fermionic Neural Networks”. In: *Phys. Rev. Lett.* 130 (3 Jan. 2023), p. 036401. DOI: 10.1103/PhysRevLett.130.036401. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.130.036401>.
- [82] Corey Adams et al. “Variational Monte Carlo Calculations of $A \leq 4$ Nuclei with an Artificial Neural-Network Correlator Ansatz”. In: *Phys. Rev. Lett.* 127 (2 July 2021), p. 022502. DOI: 10.1103/PhysRevLett.127.022502. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.127.022502>.
- [83] Alex Gnech et al. “Nuclei with Up to $\{A=6\}A = 6$ Nucleons with Artificial Neural Network Wave Functions”. In: *Few-Body Systems* 63.1 (2022), pp. 1–10.
- [84] Xiang Li, Zhe Li, and Ji Chen. “Ab initio calculation of real solids via neural network ansatz”. In: *Nature Communications* 13.1 (2022), p. 7895.
- [85] Hao Xie, Linfeng Zhang, and Lei Wang. *\$m^{\wedge} | ast\$ of two-dimensional electron gas: a neural canonical transformation study*. arXiv:2201.03156 [cond-mat, physics:physics]. Jan. 2022. URL: <http://arxiv.org/abs/2201.03156>.
- [86] Bryce Fore et al. “Dilute neutron star matter from neural-network quantum states”. In: *arXiv preprint arXiv:2212.04436* (2022).
- [87] Julien Toulouse and C. J. Umrigar. “Optimization of quantum Monte Carlo wave functions by energy minimization”. In: *The Journal of Chemical Physics* 126.8 (Feb. 2007), pp. 084102–084102. DOI: 10.1063/1.24372150.48550/arXiv.physics/0701039. arXiv: physics/0701039 [physics.chem-ph].
- [88] Sandro Sorella. “Green Function Monte Carlo with Stochastic Reconfiguration”. In: *Phys. Rev. Lett.* 80 (20 May 1998), pp. 4558–4561. DOI: 10.1103/PhysRevLett.80.4558. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.80.4558>.
- [89] C Lin, FH Zong, and David M Ceperley. “Twist-averaged boundary conditions in continuum quantum Monte Carlo algorithms”. In: *Physical Review E* 64.1 (2001), p. 016702.

- [90] ND Drummond and RJ Needs. “Diffusion quantum Monte Carlo calculation of the quasiparticle effective mass of the two-dimensional homogeneous electron gas”. In: *Physical Review B* 87.4 (2013), p. 045131.
- [91] Ke Liao et al. “Towards efficient and accurate ab initio solutions to periodic systems via transcorrelation and coupled cluster theory”. In: *Phys. Rev. Research* 3 (3 July 2021), p. 033072. DOI: 10.1103/PhysRevResearch.3.033072. URL: <https://link.aps.org/doi/10.1103/PhysRevResearch.3.033072>.
- [92] Sam Azadi, ND Drummond, and SM Vinko. “Correlation energy of the paramagnetic electron gas at the thermodynamic limit”. In: *arXiv preprint arXiv:2209.10227* (2022).
- [93] Justin Gilmer et al. “Neural message passing for quantum chemistry”. In: *International conference on machine learning*. PMLR. 2017, pp. 1263–1272.
- [94] R. P. Feynman and Michael Cohen. “Energy Spectrum of the Excitations in Liquid Helium”. In: *Phys. Rev.* 102 (5 June 1956), pp. 1189–1204. DOI: 10.1103/PhysRev.102.1189. URL: <https://link.aps.org/doi/10.1103/PhysRev.102.1189>.
- [95] Yongkyung Kwon, D. M. Ceperley, and Richard M. Martin. “Effects of backflow correlation in the three-dimensional electron gas: Quantum Monte Carlo study”. In: *Phys. Rev. B* 58 (11 Sept. 1998), pp. 6800–6806. DOI: 10.1103/PhysRevB.58.6800. URL: <https://link.aps.org/doi/10.1103/PhysRevB.58.6800>.
- [96] Yongkyung Kwon, DM Ceperley, and Richard M Martin. “Effects of three-body and backflow correlations in the two-dimensional electron gas”. In: *Physical Review B* 48.16 (1993), p. 12037.
- [97] Michele Taddei et al. “Iterative backflow renormalization procedure for many-body ground-state wave functions of strongly interacting normal Fermi liquids”. In: *Physical Review B* 91.11 (2015), p. 115106.
- [98] Michele Ruggeri, Saverio Moroni, and Markus Holzmann. “Nonlinear Network Description for Many-Body Quantum Systems in Continuous Space”. In: *Phys. Rev. Lett.* 120 (20 May 2018), p. 205302. DOI: 10.1103/PhysRevLett.120.205302. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.120.205302>.
- [99] Ashish Vaswani et al. “Attention is all you need”. In: *Advances in neural information processing systems* 30 (2017).
- [100] Dan Hendrycks and Kevin Gimpel. “Gaussian error linear units (gelus)”. In: *arXiv preprint arXiv:1606.08415* (2016).
- [101] Ingrid von Glehn, James S Spencer, and David Pfau. “A Self-Attention Ansatz for Ab-initio Quantum Chemistry”. In: *arXiv preprint arXiv:2211.13672* (2022).
- [102] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [103] D Pines. *Elementary Excitations In Solids*, edited by JD Jackson and D. Pines. 1963.
- [104] Paul P Ewald. “Die Berechnung optischer und elektrostatischer Gitterpotentiale”. In: *Annalen der physik* 369.3 (1921), pp. 253–287.

- [105] Louisa M. Fraser et al. “Finite-size effects and Coulomb interactions in quantum Monte Carlo calculations for homogeneous systems with periodic boundary conditions”. In: *Phys. Rev. B* 53 (4 Jan. 1996), pp. 1814–1832. DOI: 10.1103/PhysRevB.53.1814. URL: <https://link.aps.org/doi/10.1103/PhysRevB.53.1814>.
- [106] Abdulnour Y Toukmaji and John A Board Jr. “Ewald summation techniques in perspective: a survey”. In: *Computer physics communications* 95.2-3 (1996), pp. 73–92.
- [107] P. LPópez Ríos et al. “Inhomogeneous backflow transformations in quantum Monte Carlo calculations”. In: *Phys. Rev. E* 74 (6 Dec. 2006), p. 066701. DOI: 10.1103/PhysRevE.74.066701. URL: <https://link.aps.org/doi/10.1103/PhysRevE.74.066701>.
- [108] Sam Azadi and N. D. Drummond. “Low-density phase diagram of the three-dimensional electron gas”. In: *Phys. Rev. B* 105 (24 June 2022), p. 245135. DOI: 10.1103/PhysRevB.105.245135. URL: <https://link.aps.org/doi/10.1103/PhysRevB.105.245135>.
- [109] E. Wigner. “On the Interaction of Electrons in Metals”. In: *Phys. Rev.* 46 (11 Dec. 1934), pp. 1002–1011. DOI: 10.1103/PhysRev.46.1002. URL: <https://link.aps.org/doi/10.1103/PhysRev.46.1002>.
- [110] ND Drummond et al. “Diffusion quantum Monte Carlo study of three-dimensional Wigner crystals”. In: *Physical Review B* 69.8 (2004), p. 085116.
- [111] Giuseppe Carleo et al. “NetKet: a machine learning toolkit for many-body quantum systems”. In: *SoftwareX* 10 (2019), p. 100311.
- [112] Dion Häfner and Filippo Vicentini. “mpi4jax: Zero-copy MPI communication of JAX arrays”. In: *Journal of Open Source Software* 6.65 (2021), p. 3419.
- [113] Filippo Vicentini et al. “NETKET 3: Machine learning toolbox for many-body quantum systems”. In: *SciPost Physics Codebases* (2022), p. 007.
- [114] S Gandolfi. “Quantum Monte Carlo study of strongly interacting Fermi gases”. In: *Journal of Physics: Conference Series* 529 (Aug. 2014), p. 012011. DOI: 10.1088/1742-6596/529/1/012011. URL: <https://doi.org/10.1088%2F1742-6596%2F529%2F1%2F012011>.
- [115] Alexandros Gezerlis and J. Carlson. “Strongly paired fermions: Cold atoms and neutron matter”. In: *Phys. Rev. C* 77 (3 Mar. 2008), p. 032801. DOI: 10.1103/PhysRevC.77.032801.
- [116] Stefano Gandolfi, Alexandros Gezerlis, and J. Carlson. “Neutron Matter from Low to High Density”. In: *Annual Review of Nuclear and Particle Science* 65.1 (Oct. 2015), pp. 303–328. DOI: 10.1146/annurev-nucl-102014-021957. URL: <https://doi.org/10.1146%2Fannurev-nucl-102014-021957>.
- [117] Wynn C. G. Ho et al. “Tests of the nuclear equation of state and superfluid and superconducting gaps using the Cassiopeia A neutron star”. In: *Phys. Rev. C* 91 (1 Jan. 2015), p. 015806. DOI: 10.1103/PhysRevC.91.015806. URL: <https://link.aps.org/doi/10.1103/PhysRevC.91.015806>.

- [118] J. Carlson et al. “Auxiliary-field quantum Monte Carlo method for strongly paired fermions”. In: *Phys. Rev. A* 84 (6 2011), p. 061602. DOI: 10.1103/PhysRevA.84.061602.
- [119] J. Carlson and Sanjay Reddy. “Superfluid Pairing Gap in Strong Coupling”. In: *Phys. Rev. Lett.* 100.15 (2008), p. 150403. DOI: 10.1103/PhysRevLett.100.150403.
- [120] S. Gandolfi, K. E. Schmidt, and J. Carlson. “BEC-BCS crossover and universal relations in unitary Fermi gases”. In: *Phys. Rev. A* 83.4 (2011), p. 041601. DOI: 10.1103/PhysRevA.83.041601.
- [121] W. M. C. Foulkes et al. “Quantum Monte Carlo simulations of solids”. In: *Rev. Mod. Phys.* 73 (1 Jan. 2001), pp. 33–83. DOI: 10.1103/RevModPhys.73.33.
- [122] Michael McNeil Forbes, Stefano Gandolfi, and Alexandros Gezerlis. “Resonantly Interacting Fermions in a Box”. In: *Phys. Rev. Lett.* 106 (23 June 2011), p. 235303. DOI: 10.1103/PhysRevLett.106.235303. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.106.235303>.
- [123] M. J. H. Ku et al. “Revealing the Superfluid Lambda Transition in the Universal Thermodynamics of a Unitary Fermi Gas”. In: *Science* 335 (Feb. 2012), p. 563. DOI: 10.1126/science.1214987.
- [124] David Pfau et al. “Ab initio solution of the many-electron Schrödinger equation with deep neural networks”. In: *Phys. Rev. Research* 2 (3 Sept. 2020), p. 033429. DOI: 10.1103/PhysRevResearch.2.033429. URL: <https://link.aps.org/doi/10.1103/PhysRevResearch.2.033429>.
- [125] Gabriel Pescia et al. “Neural-network quantum states for periodic systems in continuous space”. In: *Physical Review Research* 4.2 (2022), p. 023138. DOI: 10.1103/PhysRevResearch.4.023138. URL: <https://link.aps.org/doi/10.1103/PhysRevResearch.4.023138> (visited on 07/06/2022).
- [126] Max Wilson et al. “Wave function Ansatz (but Periodic) Networks and the Homogeneous Electron Gas”. In: *arXiv e-prints*, arXiv:2202.04622 (Feb. 2022), arXiv:2202.04622. DOI: 10.48550/arXiv.2202.04622. arXiv: 2202.04622 [physics.chem-ph].
- [127] Gino Cassella et al. “Discovering Quantum Phase Transitions with Fermionic Neural Networks”. In: *Phys. Rev. Lett.* 130.3, 036401 (Jan. 2023), p. 036401. DOI: 10.1103/PhysRevLett.130.036401. arXiv: 2202.05183 [physics.comp-ph].
- [128] Bryce Fore et al. “Dilute neutron star matter from neural-network quantum states”. In: *arXiv e-prints*, arXiv:2212.04436 (Dec. 2022), arXiv:2212.04436. DOI: 10.48550/arXiv.2212.04436. arXiv: 2212.04436 [nucl-th].
- [129] C. Adams et al. “Variational Monte Carlo calculations of $A \leq 4$ nuclei with an artificial neural-network correlator ansatz”. In: *Phys. Rev. Lett.* 127 (2 July 2021), p. 022502. DOI: 10.1103/PhysRevLett.127.022502. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.127.022502>.
- [130] Alex Gnech et al. “Nuclei with up to $A = 6$ nucleons with artificial neural network wave functions”. In: *Few-Body Systems* 63.1 (2021), p. 7. DOI: 10.1007/s00601-021-01706-0. URL: <https://doi.org/10.1007/s00601-021-01706-0>.

- [131] Alessandro Lovato et al. “Hidden-nucleons neural-network quantum states for the nuclear many-body problem”. In: *Phys. Rev. Res.* 4.4 (2022), p. 043178. DOI: 10.1103/PhysRevResearch.4.043178. arXiv: 2206.10021 [nucl-th].
- [132] Y. L. Yang and P. W. Zhao. “A consistent description of the relativistic effects and three-body interactions in atomic nuclei”. In: *Phys. Lett. B* 835 (2022), p. 137587. DOI: 10.1016/j.physletb.2022.137587. arXiv: 2206.13208 [nucl-th].
- [133] Mauro Rigo et al. “Solving the nuclear pairing model with neural network quantum states”. In: *Phys. Rev. E* 107.2 (2023), p. 025310. DOI: 10.1103/PhysRevE.107.025310. arXiv: 2211.04614 [nucl-th].
- [134] Di Luo and Bryan K. Clark. “Backflow Transformations via Neural Networks for Quantum Many-Body Wave Functions”. In: *Phys. Rev. Lett.* 122 (22 June 2019), p. 226401. DOI: 10.1103/PhysRevLett.122.226401. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.122.226401>.
- [135] Michele Casula and Sandro Sorella. “Geminal wave functions with Jastrow correlation: A first application to atoms”. In: *The Journal of Chemical Physics* 119.13 (2003), pp. 6500–6511. DOI: 10.1063/1.1604379.
- [136] Michele Casula, Claudio Attaccalite, and Sandro Sorella. “Correlated geminal wave function for molecules: An efficient resonating valence bond approach”. In: *J. Comput. Phys.* 121.15 (Oct. 2004), pp. 7110–7126. DOI: 10.1063/1.1794632. arXiv: cond-mat/0409644 [cond-mat.other].
- [137] M. Bajdich et al. “Pfaffian Pairing Wave Functions in Electronic-Structure Quantum MonteCarlo Simulations”. In: *Phys. Rev. Lett.* 96.13, 130201 (Apr. 2006), p. 130201. DOI: 10.1103/PhysRevLett.96.130201. arXiv: cond-mat/0512327 [cond-mat.other].
- [138] S. Gandolfi et al. “Equation of state of low-density neutron matter, and the 1S_0 pairing gap”. In: *Phys. Rev. C* 80 (4 Oct. 2009), p. 045802.
- [139] Stefano Gandolfi et al. “Atomic nuclei from quantum Monte Carlo calculations with chiral EFT interactions”. In: *Front. in Phys.* 8 (2020), p. 117. DOI: 10.3389/fphy.2020.00117. arXiv: 2001.01374 [nucl-th].
- [140] Yusuke Nomura et al. “Restricted Boltzmann machine learning for solving strongly correlated quantum systems”. In: *Physical Review B* 96.20 (2017), p. 205152.
- [141] Yusuke Nomura and Masatoshi Imada. “Dirac-type nodal spin liquid revealed by refined quantum many-body solver using neural-network wave function, correlation ratio, and level spectroscopy”. In: *Physical Review X* 11.3 (2021), p. 031034.
- [142] S. Y. Chang et al. “Quantum Monte Carlo studies of superfluid Fermi gases”. In: *Phys. Rev. A* 70.4 (2004), p. 043602. DOI: 10.1103/PhysRevA.70.043602.
- [143] Alexandros Gezerlis et al. “Heavy-Light Fermion Mixtures at Unitarity”. In: *Phys. Rev. Lett.* 103.6 (2009), p. 060403. DOI: 10.1103/PhysRevLett.103.060403.
- [144] Andrew J. Morris, P. López Ríos, and R. J. Needs. “Ultracold atoms at unitarity within quantum Monte Carlo methods”. In: *Phys. Rev. A* 81.3 (2010), p. 033619. DOI: 10.1103/PhysRevA.81.033619.

- [145] J. Carlson et al. “Superfluid Fermi Gases with Large Scattering Length”. In: *Phys. Rev. Lett.* 91 (2003), p. 050401. DOI: 10.1103/PhysRevLett.91.050401. arXiv: physics/0303094.
- [146] Alexander Galea et al. “Diffusion Monte Carlo study of strongly interacting two-dimensional Fermi gases”. In: *Phys. Rev. A* 93 (2 Feb. 2016), p. 023602. DOI: 10.1103/PhysRevA.93.023602. URL: <https://link.aps.org/doi/10.1103/PhysRevA.93.023602>.
- [147] Maria Piarulli and Ingo Tews. “Local Nucleon-Nucleon and Three-Nucleon Interactions Within Chiral Effective Field Theory”. In: *Front. in Phys.* 7 (2020), p. 245. DOI: 10.3389/fphy.2019.00245. arXiv: 2002.00032 [nucl-th].
- [148] Kurt Hornik. “Approximation capabilities of multilayer feedforward networks”. In: *Neural Networks* 4.2 (1991), pp. 251–257. ISSN: 0893-6080. DOI: [https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T). URL: <https://www.sciencedirect.com/science/article/pii/089360809190009T>.
- [149] M. Zaheer et al. “Deep Sets”. In: *arXiv e-prints*, arXiv:1703.06114 (Apr. 2018), arXiv:1703.06114. DOI: 10.48550/arXiv.1703.06114.
- [150] Tosio Kato. “On the eigenfunctions of many-particle systems in quantum mechanics”. In: *Communications on Pure and Applied Mathematics* 10.2 (1957), pp. 151–177. DOI: 10.1002/cpa.3160100201.
- [151] Dan Hendrycks and Kevin Gimpel. “Gaussian error linear units (gelus)”. In: *arXiv preprint arXiv:1606.08415* (2016).
- [152] R. Schiavilla et al. “Two- and three-nucleon contact interactions and ground-state energies of light- and medium-mass nuclei”. In: *Physical Review C* 103.5 (May 2021). DOI: 10.1103/physrevc.103.054003. URL: <https://doi.org/10.1103/physrevc.103.054003>.
- [153] Renato Pessoa et al. “Contact interaction in a unitary ultracold Fermi gas”. In: *Phys. Rev. A* 92 (6 Dec. 2015), p. 063625. DOI: 10.1103/PhysRevA.92.063625. URL: <https://link.aps.org/doi/10.1103/PhysRevA.92.063625>.
- [154] Z. Schätzle, J. Hermann, and F. Noé. “Convergence to the fixed-node limit in deep variational Monte Carlo”. In: *The Journal of Chemical Physics* 154.12 (Mar. 2021), p. 124108. DOI: 10.1063/5.0032836. URL: <https://doi.org/10.1063/5.0032836>.
- [155] Michael McNeil Forbes, Stefano Gandolfi, and Alexandros Gezerlis. “Effective-range dependence of resonantly interacting fermions”. In: *Physical Review A* 86.5 (Nov. 2012). DOI: 10.1103/physreva.86.053603. URL: <https://doi.org/10.1103/physreva.86.053603>.
- [156] D. J. Dean and M. Hjorth-Jensen. “Pairing in nuclear systems: From neutron stars to finite nuclei”. In: *Rev. Mod. Phys.* 75 (2003), pp. 607–656. DOI: 10.1103/RevModPhys.75.607. arXiv: nucl-th/0210033.
- [157] Wan Tong Lou et al. *Neural Wave Functions for Superfluids*. 2023. arXiv: 2305.06989 [cond-mat.quant-gas].

- [158] Claudio Genovese et al. “General Correlated Geminal Ansatz for Electronic Structure Calculations: Exploiting Pfaffians in Place of Determinants”. In: *Journal of Chemical Theory and Computation* 16.10 (2020), pp. 6114–6131. DOI: 10.1021/acs.jctc.0c00165. eprint: <https://pubs.acs.org/doi/pdf/10.1021/acs.jctc.0c00165>. URL: <https://pubs.acs.org/doi/10.1021/acs.jctc.0c00165>.