

RANDOMIZED NUMERICAL LINEAR ALGEBRA FOR LARGE-SCALE MATRIX DATA

A Dissertation

Presented to the Faculty of the Graduate School

of Cornell University

in Partial Fulfillment of the Requirements for the Degree of

Doctor of Philosophy

by

Kun Dong

August 2019

© 2019 Kun Dong

ALL RIGHTS RESERVED

RANDOMIZED NUMERICAL LINEAR ALGEBRA FOR LARGE-SCALE

MATRIX DATA

Kun Dong, Ph.D.

Cornell University 2019

Your abstract goes here. Make sure it sits inside the brackets. If not, your biosketch page may not be roman numeral iii, as required by the graduate school.

BIOGRAPHICAL SKETCH

Kun Dong was born in Shaoxing, Zhejiang Province, China on September 26, 1992 to Jian Dong and Chamei Sang. Kun became interested in puzzles and mathematics in early elementary school, and participated in mathematics competitions until the end of high school.

In 2008, Kun moved to Newmarket, Ontario, Canada and attended Sir William Mulock Secondary School. He spent the next two years learning the new language and culture.

After graduation, Kun attended the University of California, Los Angeles in 2010 to study applied mathematics. During the summers of 2012 and 2013, he participated in the UCLA Applied Math REU program to work on dynamical system and crime modeling. He received tremendous help and encouragement on his way to graduate school from his REU mentors, Scott McCalla and James von Brecht. In 2014, Kun completed the departmental scholar program in mathematics, earning his B.S. and M.A. concurrently. He graduated Summa Cum Laude, receiving the departmental highest honor in applied mathematics and the Daus memorial award for his achievement in mathematics as an undergraduate.

In 2014, Kun was admitted to the Ph.D. program at Center for Applied Mathematics in Cornell University. He soon started working with Professor David Bindel, who later became his advisor. During the summer of 2017, he interned at the Lawrence Berkeley National Laboratory, where he worked under the mentorship of Professor Lin Lin. For the last three years at Cornell, he was partially supported by a National Science Foundation grant. In September 2019, Kun will move to Seattle and become a research scientist at Facebook, Inc.

This document is dedicated to my parents and wife.

ACKNOWLEDGEMENTS

Your acknowledgements go here. Make sure it sits inside the brackets.

TABLE OF CONTENTS

Biographical Sketch	iii
Dedication	iv
Acknowledgements	v
Table of Contents	vi
List of Tables	ix
List of Figures	x
1 Introduction	1
2 Network Density of States	2
2.1 Abstract	2
2.2 Introduction	3
2.3 Background	5
2.3.1 Graph Operators and Eigenvalues	5
2.3.2 Spectral Density (Density of States — DOS)	7
2.4 Methods	9
2.4.1 Kernel Polynomial Method (KPM)	10
2.4.2 Gauss Quadrature and Lanczos (GQL)	12
2.4.3 Nested Dissection (ND)	13
2.4.4 Motif Filtering	15
2.5 Error Analysis	17
2.5.1 KPM Approximation Error	17
2.5.2 Perturbation Analysis	20
2.6 Experiments	21
2.6.1 Gallery of DOS/PDOS	21
2.6.2 Computation time	23
2.6.3 Model Verification	24
2.7 Conclusion	27
3 Scalable Gaussian Processes	28
3.1 Abstract	28
3.2 Introduction	28
3.3 Background	32
3.4 Methods	34
3.4.1 Chebyshev Expansion	35
3.4.2 Gauss Quadrature via Lanczos	36
3.4.3 Diagonal Correction to SKI	37
3.4.4 Estimating higher derivatives	38
3.4.5 Radial Basis Functions	39
3.5 Error Properties	39
3.6 Experiments	41
3.6.1 Natural sound modeling	42

3.6.2	Daily precipitation prediction	43
3.6.3	Hickory data	44
3.6.4	Crime prediction	46
3.6.5	Deep kernel learning	47
3.7	Conclusion	48
4	Scalable Gaussian Processes with Derivatives	49
4.1	Abstract	49
4.2	Introduction	49
4.3	Background	51
4.4	Methods	54
4.4.1	D-SKI	55
4.4.2	D-SKIP	55
4.4.3	Preconditioning	56
4.4.4	Dimensionality reduction	57
4.5	Experiments	58
4.5.1	Approximation Benchmark	58
4.5.2	Dimensionality reduction	60
4.5.3	Rough terrain reconstruction	61
4.5.4	Implicit surface reconstruction	62
4.5.5	Bayesian optimization with derivatives	63
4.6	Conclusion	65
5	Robust Large-Vocabulary Topic Modeling	67
5.1	Abstract	67
5.2	Introduction	67
5.3	Background and Related Work	71
5.4	Low-rank Rectification and Compression	74
5.4.1	Epsilon Non-Negative Rectification (ENN)	74
5.4.2	Proximal Alternating Linearized Minimization Rectification (PALM)	76
5.5	Low-rank Anchor Word Algorithm	77
5.6	Experimental Results	81
5.7	Conclusion	85
6	Weighted K-Means for Electronic Structure Calculation	87
6.1	Abstract	87
6.2	Introduction	88
6.3	Interpolative Separable Density Fitting (ISDF) decomposition	92
6.4	Centroidal Voronoi Tessellation based ISDF decomposition	95
6.5	Numerical results	99
6.5.1	Accuracy: Si_{216} and $\text{Al}_{176}\text{Si}_{24}$	102
6.5.2	Efficiency: Si_{1000}	103
6.5.3	Ab Initio Molecular Dynamics: Si_{64} and $(\text{H}_2\text{O})_{64}$	106

6.6 Conclusion	110
6.7 Acknowledgments	111
7 Conclusion	112

LIST OF TABLES

2.1	Average Computation Time per Chebyshev Moment for Graphs from the SNAP Repository ^a	24
3.1	Prediction Comparison for the Daily Precipitation Data ^a	44
3.2	Hyperparameters Recovered on the Hickory Dataset.	45
3.3	Hyperparameters Recovered, Recovery Time and RMSE for Lanczos and Scaled Eigenvalues on the Chicago Assault Data ^a	47
3.4	Prediction RMSE and Per Training Iteration Runtime.	47
4.1	Relative RMSE on Test Functions Using SKI and Derivatives ^a	59
4.2	Relative RMSE and SMAE prediction error for Welsh ^a	61
4.3	Hyperparameters Recovered, Recovery SMAE, and Recovery Time for SKI and D-SKI on Mountain St. Helens Data ^a	62
6.1	Accuracy of ACE-ISDF Based Hybrid Functional Calculations (HSE06) Obtained by Using the CVT method To Select Interpolation Points, with Varying Rank Parameter c for Semiconducting Si_{216} and Metallic $\text{Al}_{176}\text{Si}_{24}$ Systems ^a	104
6.2	Wall Clock Time (in seconds) Spent in the Components of the ACE-ISDF and ACE Enabled Hybrid DFT Calculations Related to the Exchange Operator, for Si_{1000} on 2002 Edison cores at Different E_{cut} Levels ^a	106

LIST OF FIGURES

2.1	Spectral histogram for the normalized adjacency matrix for the CAIDA autonomous systems graph [87], an Internet topology with 22,965 nodes and 47,193 edges. Blue bars are the real spectrum, and red points are the approximated heights. (a) contains high multiplicity around eigenvalue 0, so (b) zooms in to height between [0, 500].	9
2.2	Common motifs (induced subgraphs) in graph data that result in localized spikes in the spectral density. Each motif generates a specific eigenvalue with locally-supported eigenvectors. Here we uses the normalized adjacency matrix to represent the graph, although we can perform the same analysis for the adjacency, Laplacian, or normalized Laplacian (only the eigenvalues would be different). The eigenvectors are supported only on the labeled nodes.	16
2.3	The improvement in accuracy of the spectral histogram approximation on the normalized adjacency matrix for the High Energy Physics Theory (HepTh) Collaboration Network, as we sweep through spectrum and filter out motifs. The graph has 8,638 nodes and 24,816 edges. Blue bars are the real spectrum, and red points are the approximated heights. Fig. 2.3a- 2.3e use 100 moments and 20 probe vectors. Fig. 2.3f shows the relative L_1 error of the spectral histogram when using no filter, filter at $\lambda = 0$, and all filters.	18
2.4	DOS(top)/PDOS(bottom) histograms for the normalized adjacency of 10 networks from five domains. For DOS, blue bars are the true spectrum, and red points are from KPM (500 moments and 20 Hadamard probes). For PDOS, the spectral histograms of all nodes are aligned vertically. Red indicates high weight around an eigenvalue, and blue indicates low weight. The true spectrum for the California Road Network (2.4j) is omitted, as it is too large to compute exactly (1,965,206 nodes).	22
2.5	Spectral histogram for scale-free model with 5000 nodes and different m . Blue bars are the real spectrum, red points are from KPM (500 moments and 20 probes).	25
2.6	Spectral histograms for small-world model with 5000 nodes and rewiring probability $p = 0.5$, starting with 5000 (2.6a) and 50000 (2.6b) edges. Blue bars are the real spectrum, red points are from KPM (5000 moments and 20 probes).	25
2.7	Comparison of spectral histogram between Erdős Collaboration Network and the BTER model. Both DOS and PDOS are computed with 500 moments and 20 probe vectors.	26

3.1	Sound modeling using 59,306 training points and 691 test points. The intensity of the time series can be seen in (a). Train time for RBF kernel hyperparameters is in (b) and the time for inference is in (c). The standardized mean absolute error (SMAE) as a function of time for an evaluation of the marginal likelihood and all derivatives is shown in (d). Surrogate is (—), Lanczos is (- - -), Chebyshev is (— ◊ —), scaled eigenvalues is (— + —), and FITC is (— o —).	43
3.2	Predictions by exact, scaled eigenvalues, and Lanczos on the Hickory dataset.	45
4.1	An example where gradient information pays off; the true function is on the left. Compare the regular GP without derivatives (middle) to the GP with derivatives (right). Unlike the former, the latter is able to accurately capture critical points of the function.	54
4.2	(Left two images) \log_{10} error in SKI approximation and comparison to the exact spectrum.(Right two images) \log_{10} error in SKIP approximation and comparison to the exact spectrum.	58
4.3	Scaling tests for D-SKI in two dimensions and D-SKIP in 11 dimensions.D-SKIP uses fewer data points for identical matrix sizes.	59
4.4	Fig. 4.4a shows the top 10 eigenvalues of the gradient covariance. Welsh is projected onto the first and second active direction in 4.4b and 4.4c. After joining them together, we see in 4.4d that points of different color are highly mixed, indicating a very spiky surface.	60
4.5	On the left is the true elevation map of Mount St. Helens. In the middle is the elevation map calculated with the SKI. On the right is the elevation map calculated with D-SKI.	62
4.6	(Left) Original surface (Middle) Noisy surface (Right) SKI reconstruction from noisy surface ($s = 0.4$, $\sigma = 0.12$).	63
4.7	In the following experiments, 5D Ackley and 5D Rastrigin are embedded into 50 a dimensional space. We run Algorithm 1, comparing it with BO exact, multi-start BFGS, and random sampling. D-SKI with active subspace learning clearly outperforms the other methods.	65
5.1	Experiment on four datasets. ENN and LR-JSMF mostly agree with AP, whereas PALM has slight inconsistency. The general information of each dataset is above the corresponding row. Recovery, approximation, and runtimes are in \log_{10} scale. Note that ENN and LR-JSMF are almost two orders of magnitude faster than AP. The x -axis indicates the number of clusters K . Lower numbers in y -axis are better except Specificity and Dissimilarity.	82
5.2	As we increase the vocabulary size for four collections, anchor quality and sparsity improve, but running time is stable. The x -axis indicates the vocabulary size N in thousands. Values above 15k will not fit in memory on standard hardware with previous algorithms.	84

5.3	Losses or gains in topic words depending on the vocabulary size. Each row represents a topic from the NeurIPS dataset, with the top 6 topical words shown in the middle column. The red and green cells denote topic words that are lost or gained by shifting the vocabulary size from the default size 5000, respectively. The intensities of the colors indicate the words' contributions towards the specific topic.	86
6.1	Schematic illustration of the CVT procedure in a 2D domain, including (a) initial random choice of centroids and Voronoi tessellation and centroidal Voronoi tessellation generated by the weighted K-Means algorithm. The weight function is given by the linear superposition of 4 Gaussian functions.	99
6.2	The decomposition reaction process of BH_3NH_3 computed with hybrid functional (HSE06) calculations by using the CVT procedure to select interpolation points, including (a) the electron density (yellow isosurfaces), (b) the interpolation points (yellow squares) $\{\hat{\mathbf{r}}_\mu\}_{\mu=1}^{N_\mu}$ ($N_\mu = 8$) selected from the real space grid points $\{\mathbf{r}_i\}_{i=1}^{N_g}$ ($N_g = 100^3$ and $E_{\text{cut}} = 60$ Ha) when the BN distance respectively is 1.3, 1.7 and 2.8 Å and (c) the binding energy as a function of BN distance for BH_3NH_3 in a $10 \text{ \AA} \times 10 \text{ \AA} \times 10 \text{ \AA}$ box. The white, pink and blue pink balls denote hydrogen, boron and nitrogen atoms, respectively.	100
6.3	The accuracy of ACE-ISDF based hybrid functional calculations (HSE06) obtained by using the CVT and QRCP procedures to select the interpolation points, with varying rank parameter c from 4 to 20 for Si_{216} and $\text{Al}_{176}\text{Si}_{24}$, including the error of (a) Hartree-Fock exchange energy ΔE_{HF} (Ha/atom) and (b) total energy ΔE (Ha/atom).	105
6.4	Comparison of the ISDF-CVT method by using either random or QRCP initialization for hybrid DFT AIMD simulations on bulk silicon system Si_{64} and liquid water system $(\text{H}_2\text{O})_{64}$, including (a) the fraction of points what switch cluster in each K-Means iteration and (b) the number of K-Means iterations during each MD step.	108
6.5	Comparison of hybrid HSE06 DFT AIMD simulations by using the ISDF-CVT and ISDF-QRCP methods as well as exact nested two-level SCF iteration procedure as the reference on the bulk silicon Si_{64} , including (a) relatively energy drift and (b) potential energy during MD steps.	109
6.6	The oxygen-oxygen radial distribution functions $g_{\text{OO}}(r)$ of liquid water system $(\text{H}_2\text{O})_{64}$ at $T = 295$ K obtained from hybrid HSE06 + DFT-D2 AIMD simulations with the ISDF-CVT and ISDF-QRCP methods, exact nested two-level SCF iteration procedure (as the reference) as well as previous hybrid PBE0 + TS-vdW calculation [47].	110

CHAPTER 1

INTRODUCTION

CHAPTER 2

NETWORK DENSITY OF STATES

2.1 Abstract

Spectral analysis connects graph structure to the eigenvalues and eigenvectors of associated matrices. Much of spectral graph theory descends directly from spectral geometry, the study of differentiable manifolds through the spectra of associated differential operators. But the translation from spectral geometry to spectral graph theory has largely focused on results involving only a few extreme eigenvalues and their associated eigenvalues. Unlike in geometry, the study of graphs through the overall distribution of eigenvalues — the *spectral density* — is largely limited to simple random graph models. The interior of the spectrum of real-world graphs remains largely unexplored, difficult to compute and to interpret.

In this paper, we delve into the heart of spectral densities of real-world graphs. We borrow tools developed in condensed matter physics, and add novel adaptations to handle the spectral signatures of common graph motifs. The resulting methods are highly efficient, as we illustrate by computing spectral densities for graphs with over a billion edges on a single compute node. Beyond providing visually compelling fingerprints of graphs, we show how the estimation of spectral densities facilitates the computation of many common centrality measures, and use spectral densities to estimate meaningful information about graph structure that cannot be inferred from the extremal eigenpairs alone.

2.2 Introduction

Spectral theory is a powerful analysis tool in graph theory [43, 37, 36], geometry [33], and physics [90]. One follows the same steps in each setting:

- Identify an object of interest, such as a graph or manifold;
- Associate the object with a matrix or operator, often the generator of a linear dynamical system or the Hessian of a quadratic form over functions on the object;
- Connect spectral properties of the matrix or operator to structural properties of the original object.

In each case, the *complete* spectral decomposition is enough to recover the original object; the interesting results relate structure to *partial* spectral information.

Many spectral methods use extreme eigenvalues and associated eigenvectors. These are easy to compute by standard methods, and are easy to interpret in terms of the asymptotic behavior of dynamical systems or the solutions to quadratic optimization problems with quadratic constraints. Several network centrality measures, such as PageRank [136], are expressed via the stationary vectors of transition matrices, and the rate of convergence to stationarity is bounded via the second-largest eigenvalue. In geometry and graph theory, Cheeger's inequality relates the second-smallest eigenvalue of a Laplacian or Laplace-Beltrami operator to the size of the smallest bisecting cut [34, 127]; in the graph setting, the associated eigenvector (the Fiedler vector) is the basis for spectral algorithms for graph partitioning [141]. Spectral algorithms for graph coordinates and clustering use the first few eigenvectors of a transition matrix or (normalized) adjacency or Laplacian [23, 133]. For a survey of such approaches in network science, we refer to [36].

Mark Kac popularized an alternate approach to spectral analysis in an expository article [92] in which he asked whether one can determine the shape of a physical object (Kac used a drum as an example) given the spectrum of the Laplace operator; that is, can one “hear” the shape of a drum? One can ask a similar question in graph theory: can one uniquely determine the structure of a network from the spectrum of the Laplacian or another related matrix? Though the answer is negative in both cases [67, 43], the spectrum is enormously informative even without eigenvector information. Unlike the extreme eigenvalues and vectors, eigenvalues deep in the spectrum are difficult to compute and to interpret, but the overall distribution of eigenvalues — known as the spectral density or density of states — provides valuable structural information. For example, knowing the spectrum of a graph adjacency matrix is equivalent to knowing $\text{tr}(A^k)$, the number of closed walks of any given length k . In some cases, one wants *local* spectral densities in which the eigenvalues also have positive weights associated with a location. Following Kac, this would give us not only the frequencies of a drum, but also amplitudes based on where the drum is struck. In a graph setting, the local spectral density of an adjacency matrix at node j is equivalent to knowing $(A^k)_{jj}$, the number of closed walks of any given length k that begin and end at the node.

Unfortunately, the analysis of spectral densities of networks has been limited by a lack of scalable algorithms. While the normalized Laplacian spectra of Erdős-Rényi random graphs have an approximately semicircular distribution [173], and the spectral distributions for other popular scale-free and small-world random graph models are also known [53], there has been relatively little work on computing spectral densities of large “real-world” networks. Obtaining the full eigendecomposition is $O(N^3)$ for a graph with N nodes, which is prohibitive for graphs of more than a few thousand nodes. In prior work, researchers have employed methods, such as thick-restart Lanczos, that still do not scale to very large graphs [53], or heuristic approximations with no convergence

analysis [16]. It is only recently that clever computational methods were developed simply to *test* for hypothesized power laws in the spectra of large real-world matrices by computing only *part* of the spectrum [52].

In this paper, we show how methods used to study densities of states in condensed matter physics [169] can be used to study spectral densities in networks. We study these methods for both the *global* density of states and for *local* densities of states weighted by specific eigenvector components. We adapt these methods to take advantage of graph-specific structure not present in most physical systems, and analyze the stability of the spectral density to perturbations as well as the convergence of our computational methods. Our methods are remarkably efficient, as we illustrate by computing densities for graphs with billions of edges and tens of millions of nodes on a single cloud compute node. We use our methods for computing these densities to create compelling visual fingerprints that summarize a graph. We also show how the density of states reveals graph properties that are not evident from the extremal eigenvalues and eigenvectors alone, and use it as a tool for fast computation of standard measures of graph connectivity and node centrality. This opens the door for the use of complete spectral information as a tool in large-scale network analysis.

2.3 Background

2.3.1 Graph Operators and Eigenvalues

We consider weighted, undirected graphs $G = (V, E)$ with vertices $V = \{v_1, \dots, v_N\}$ and edges $E \subseteq V \times V$. The weighted adjacency matrix $A \in \mathbb{R}^{N \times N}$ has entries $a_{ij} > 0$ to give the weight of an edge $(i, j) \in E$ and $a_{ij} = 0$ otherwise. The degree matrix $D \in \mathbb{R}^{N \times N}$ is

the diagonal matrix of weighted node degrees, i.e. $D_{ii} = \sum_j a_{ij}$. Several of the matrices in spectral graph theory are defined in terms of D and A . We describe a few of these below, along with their connections to other research areas. For each operator, we let $\lambda_1 \leq \dots \leq \lambda_N$ denotes the eigenvalues in ascending order.

Adjacency Matrix: A . Many studies on the spectrum of A originate from random matrix theory where A represents a random graph model. In these cases, the limiting behavior of eigenvalues as $N \rightarrow \infty$ is of particular interest. Besides the growth of extremal eigenvalues [37], Wigner's semicircular law is the most renowned result about the spectral distribution of the adjacency matrix [173]. When the edges are i.i.d. random variables with bounded moments, the density of eigenvalues within a range converges to a semicircular distribution. One famous graph model of this type is the Erdős-Rényi graph, where $a_{ij} = a_{ji} = 1$ with probability $p < 1$, and 0 with probability $1 - p$. Farkas et al. [53] has extended the semicircular law by investigating the spectrum of scale-free and small-world random graph models. They show the spectra of these random graph models relate to geometric characteristics such as the number of cycles and the degree distribution.

Laplacian Matrix: $L = D - A$. The Laplace operator arises naturally from the study of dynamics in both spectral geometry and spectral graph theory. The continuous Laplace operator and its generalizations are central to the description of physical systems including heat diffusion [122], wave propagation [105], and quantum mechanics [51]. It has infinitely many non-negative eigenvalues, and Weyl's law [171] relates their asymptotic distribution to the volume and dimension of the manifold. On the other hand, the discrete Laplace matrix appears in the formulation of graph partitioning problems. If $f \in \{\pm 1\}^N$ is an indicator vector for a partition $V = V_+ \cup V_-$, then $f^T L f / 4$ is the number of edges between V_+ and V_- , also known as the cut size. L is a positive-semidefinite ma-

trix with the vector of all ones as a null vector. The eigenvalue λ_2 , called the *algebraic connectivity*, bounds from below the smallest bisecting cut size; $\lambda_2 = 0$ if and only if the graph is disconnected. In addition, eigenvalues of L also appear in bounds for vertex connectivity (λ_2) [44], minimal bisection (λ_2) [48], and maximum cut (λ_N) [163].

Normalized Laplacian Matrix: $\bar{L} = I - D^{-1/2}AD^{-1/2}$. We will also mention the normalized adjacency matrix $\bar{A} = D^{-1/2}AD^{-1/2}$ and graph random walk matrix $P = D^{-1}A$ here, because these matrices have the same eigenvalues as \bar{L} up to a shift. The connection to some of the most influential results in spectral geometry is established in terms of eigenvalues and eigenvectors of normalized Laplacian. A prominent example is the extension of Cheeger's inequality to the discrete case, which relates the set of smallest conductance $h(G)$ (the Cheeger constant) to the second smallest eigenvalue of the normalized Laplacian, $\lambda_2(\bar{L})$ [128]:

$$\lambda_2(\bar{L})/2 \leq h(G) = \min_{S \subset V} \frac{|\{(i, j) \in E, i \in S, j \notin S\}|}{\min(\text{vol}(S), \text{vol}(V \setminus S))} \leq \sqrt{2\lambda_2(\bar{L})}, \quad (2.1)$$

where $\text{vol}(T) = \sum_{i \in T} \sum_{j=1}^N a_{ij}$. Cheeger's inequality offers crucial insights and powerful techniques for understanding popular spectral graph algorithms for partitioning [123] and clustering [133]. It also plays a key role in analyzing the mixing time of Markov chains and random walks on a graph [126, 156]. For all these problems, extremal eigenvalues again emerge from relevant optimization formulations.

2.3.2 Spectral Density (Density of States — DOS)

Let $H = \mathbb{R}^{N \times N}$ be any symmetric graph matrix with an eigendecomposition $H = Q\Lambda Q^T$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$ and $Q = [q_1, \dots, q_N]$ is orthogonal. The spectral density

induced by H is the generalized function

$$\mu(\lambda) = \frac{1}{N} \sum_{i=1}^N \delta(\lambda - \lambda_i), \quad \int f(\lambda) \mu(\lambda) = \text{tr}(f(H)) \quad (2.2)$$

where δ is the Dirac delta function and f is any analytic test function. The spectral density μ is also referred to as the *density of states* (DOS) in the condensed matter physics literature [169], as it describes the number of states at different energy levels.

For any vector $u \in \mathbb{R}^N$, the local density of states (LDOS) is

$$\mu(\lambda; u) = \sum_{i=1}^N |u^T q_i|^2 \delta(\lambda - \lambda_i), \quad \int f(\lambda) \mu(\lambda; u) = u^T f(H) u. \quad (2.3)$$

Most of the time, we are interested in the case $u = e_k$ where e_k is the k th standard basis vector—this provides the spectral information about a particular node. We will write $\mu_k(\lambda) = \mu(\lambda; e_k)$ for the pointwise density of states (PDOS) for node v_k . It is noteworthy $|e_k^T q_i| = |q_i(k)|$ gives the magnitude of the weight for v_k in the i -th eigenvector, thereby the set of $\{\mu_k\}$ encodes the entire spectral information of the graph up to sign differences. These concepts can be easily extended to directed graphs with asymmetric matrices, for which the eigenvalues are replaced by singular values, and eigenvectors by left/right singular vectors.

Naively, to obtain the DOS and LDOS requires computing all eigenvalues and eigenvectors for an N -by- N matrix, which is infeasible for large graphs. Therefore, we turn to algorithms that approximate these densities. Since the DOS is a generalized function, it is important we specify how the estimation is evaluated. One choice is to treat μ (or μ_k) as a distribution, and measure its approximation error with respect to a chosen function space \mathcal{L} . For example, when \mathcal{L} is the set of Lipschitz continuous functions taking the value 0 at 0, the error for estimated $\tilde{\mu}$ is in the Wasserstein distance (a.k.a. earth-mover distance) [93]

$$W_1(\mu, \tilde{\mu}) = \sup \left\{ \int (\mu(\lambda) - \tilde{\mu}(\lambda)) f(\lambda) d\lambda : \text{Lip}(f) \leq 1 \right\}. \quad (2.4)$$

This notion is particularly useful when μ is integrated against in applications such as computing centrality measures.

On the other hand, we can regularize μ with a mollifier K_σ (i.e., a smooth approximation of the identity function):

$$(K_\sigma * \mu)(\lambda) = \int_{\mathbb{R}} \sigma^{-1} K\left(\frac{\lambda - v}{\sigma}\right) \mu(v) dv \quad (2.5)$$

A simplified approach is numerically integrating μ over small intervals of equal size to generate a spectral histogram. The advantage is the error is now easily measured and visualized in the L_∞ norm. For example, Figure 2.1 shows the exact and approximated spectral histogram for the normalized adjacency matrix of an Internet topology.

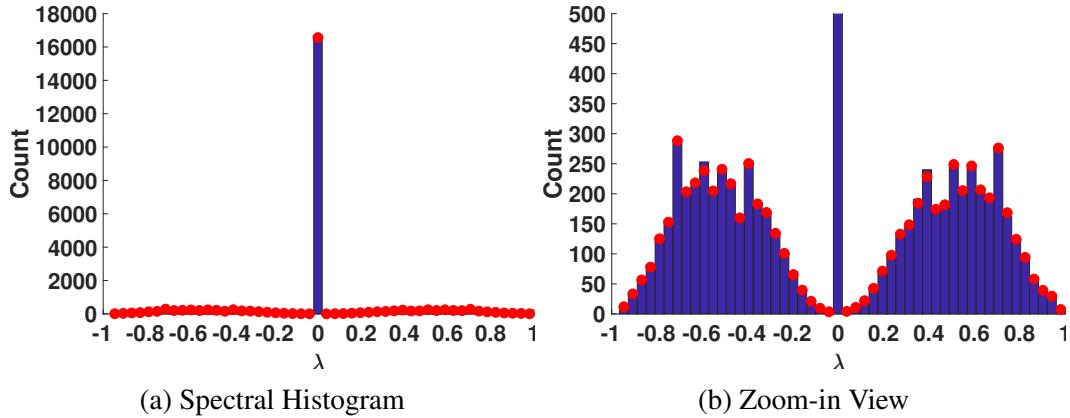


Figure 2.1: Spectral histogram for the normalized adjacency matrix for the CAIDA autonomous systems graph [87], an Internet topology with 22,965 nodes and 47,193 edges. Blue bars are the real spectrum, and red points are the approximated heights. (a) contains high multiplicity around eigenvalue 0, so (b) zooms in to height between [0, 500].

2.4 Methods

The density of states plays a significant role in understanding electronic band structure in solid state physics, and so several methods have been proposed in that litera-

ture to estimate spectral densities. We review two such methods: the kernel polynomial method (KPM) which involves a polynomial expansion of the DOS/LDOS, and the Gauss Quadrature via Lanczos iteration (GQL). These methods have not previously been applied in the network setting, though Cohen-Steiner et al. [38] have independently invented an approach similar to KPM for the global DOS alone, albeit using a less numerically stable polynomial basis (the power basis associated with random walks). We then introduce a new direct nested dissection method for LDOS, as well as new graph-specific modifications to improve the convergence of the KPM and GQL approaches.

Throughout this section, H denotes any symmetric matrix.

2.4.1 Kernel Polynomial Method (KPM)

The Kernel Polynomial Method (KPM) [169] approximates the spectral density through an expansion in the dual basis of an orthogonal polynomial basis. Traditionally, the Chebyshev basis $\{T_m\}$ is used because of its connection to the best polynomial interpolation. Chebyshev approximation requires the spectrum to be supported on the interval $[-1, 1]$ for numerical stability. However, this condition can be satisfied by any graph matrix after shifting and rescaling:

$$\widetilde{H} = \frac{2H - (\lambda_{\max}(H) + \lambda_{\min}(H))}{\lambda_{\max}(H) - \lambda_{\min}(H)}. \quad (2.6)$$

We can compute these extremal eigenvalues efficiently for our sparse matrix H , so the pre-computation is not an issue [137].

The Chebyshev polynomials $T_m(x)$ satisfy the recurrence

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{m+1}(x) = 2xT_m(x) - T_{m-1}(x). \quad (2.7)$$

They are orthogonal with respect to $w(x) = 2 / \left[(1 + \delta_{0n})\pi \sqrt{1 - x^2} \right]$:

$$\int_{-1}^1 w(x) T_m(x) T_n(x) dx = \delta_{mn}. \quad (2.8)$$

(Here and elsewhere, δ_{ij} is the Kronecker delta: 1 if $i = j$ and 0 otherwise.) Therefore, $T_m^*(x) = w(x)T_m(x)$ also forms the dual Chebyshev basis. Using Eq. 2.8, we can expand our DOS $\mu(\lambda)$ as

$$\mu(\lambda) = \sum_{m=0}^{\infty} d_m T_m^*(\lambda), \quad (2.9)$$

$$d_m = \int_{-1}^1 T_m(\lambda) \mu(\lambda) d\lambda = \frac{1}{N} \sum_{i=1}^N T_m(\lambda_i) = \frac{1}{N} \text{tr}(T_m(H)). \quad (2.10)$$

Here, $T_m(H)$ is the m th Chebyshev polynomial of the matrix H . The last equality comes from the spectral mapping theorem, which says that taking a polynomial of H maps the eigenvalues by the same polynomial. Similarly, we express the PDOS $\mu_k(\lambda)$ as

$$d_{mk} = \int_{-1}^1 T_m(\lambda) \mu_k(\lambda) d\lambda = \sum_{i=1}^N |q_i(k)|^2 T_m(\lambda_i) = T_m(H)_{kk}. \quad (2.11)$$

We want to efficiently extract the diagonal elements of the matrices $\{T_m(H)\}$ without forming them explicitly; the key idea is to apply the stochastic trace/diagonal estimation, proposed by Hutchinson [88] and Bekas et al. [21]. Given a random probe vector z such that z_i 's are i.i.d. with mean 0 and variance 1,

$$\mathbb{E}[z^T H z] = \sum_{i,j} H_{ij} \mathbb{E}[z_i z_j] = \text{tr}(H) \quad (2.12)$$

$$\mathbb{E}[z \odot H z] = \text{diag}(H) \quad (2.13)$$

where \odot represents the Hadamard (elementwise) product. Choosing N_z independent probe vectors Z_j , we obtain the unbiased estimator

$$\text{tr}(H) = \mathbb{E}[z^T H z] \approx \frac{1}{N_z} \sum_{j=1}^{N_z} Z_j^T H Z_j \quad (2.14)$$

and similarly for the diagonal. Avron and Toledo [13] review many possible choices of probes for Eqs. (2.12) and (2.13); a common choice is vectors with independent standard normal entries. Using the Chebyshev recurrence (Eq. (2.7)), we can compute the sequence $T_j(H)z$ for each probe at a cost of one matrix-vector product per term, for a total cost of $O(|E|N_z)$ time per moment $T_m(H)$.

In practice, we only use a finite number of moments rather than an infinite expansion. The number of moments required depends on the convergence rate of the Chebyshev approximation for the class of functions DOS/LDOS is integrated with. For example, the approximation error decays exponentially for test functions that are smooth over the spectrum [162], so only a few moments are needed. On the other hand, such truncation leads to Gibbs oscillations that cause error in the interpolation [162]. However, to a large extent, we can use smoothing techniques such as Jackson damping to resolve this issue [89] (we will formalize this in Theorem 1).

2.4.2 Gauss Quadrature and Lanczos (GQL)

Golub and Meurant developed the well-known Gauss Quadrature and Lanczos (GQL) algorithm to approximate bilinear forms for smooth functions of a matrix [66]. Using the same stochastic estimation from § 2.4.1, we can also apply GQL to compute DOS.

For a starting vector z and graph matrix H , Lanczos iterations after M steps produce a decomposition

$$HZ_M = Z_M^T \Gamma_M + r_M e_M^T, \quad (2.15)$$

where $Z_M^T Z_M = I_M$, $Z_M^T r_M = 0$, and Γ_M tridiagonal. GQL approximates $z^T f(H)z$ with

$\|z\|^2 e_1^T f(T_M) e_1$, implying

$$z^T f(H) z = \sum_{i=1}^N |z^T q_i|^2 f(\lambda_i) \approx \|z\|^2 \sum_{i=1}^M |p_{i1}|^2 f(\tau_i), \quad (2.16)$$

where $(\tau_1, p_1), \dots, (\tau_M, p_M)$ are the eigenpairs of Γ_M . Consequently,

$$\|z\|^2 \sum_{i=1}^M |p_{i1}|^2 \delta(\lambda - \tau_i) \quad (2.17)$$

approximates the LDOS $\mu(\lambda; z)$.

Building upon the stochastic estimation idea and the invariance of probe vectors under orthogonal transformation, we have

$$\mathbb{E}[\mu(\lambda; z)] = \sum_{i=1}^N \delta(\lambda - \lambda_i) = N\mu(\lambda). \quad (2.18)$$

Hence

$$\mu(\lambda) \approx \sum_{i=1}^M |p_{i1}|^2 \delta(\lambda - \tau_i). \quad (2.19)$$

The approximate generalized function is exact when applied to polynomials of degree $\leq 2M + 1$. Furthermore, if we let $z = e_k$ then GQL also provides an estimation for the PDOS $\mu_k(\lambda)$. Estimation from GQL can also be converted to Chebyshev moments if needed.

2.4.3 Nested Dissection (ND)

The estimation error via Monte Carlo method intrinsically decays at the rate $O(1/\sqrt{N_z})$, where N_z is the number of random probing vectors. Hence, we have to tolerate the higher variance when increasing the number of probe vectors becomes too expensive. This is particularly problematic when we try to compute the PDOS for all nodes using the stochastic diagonal estimator. Therefore, we introduce an alternative divide-and-conquer method, which computes more accurate PDOS for any set of nodes at a cost comparable to the stochastic approach in practice.

Suppose the graph can be partitioned into two subgraphs by removal of a small vertex separator. Permuting the vertices so that the two partitions appear first, followed by the separator vertices. Up to vertex permutations, we can rewrite H in block form as

$$H = \begin{bmatrix} H_{11} & 0 & H_{13} \\ 0 & H_{22} & H_{23} \\ H_{13}^T & H_{23}^T & H_{33} \end{bmatrix}, \quad (2.20)$$

where the indices indicate the groups identities. Leveraging this structure, we can update the recurrence relation for Chebyshev polynomials to become

$$T_{m+1}(H)_{11} = 2H_{11}T_m(H)_{11} - T_{m-1}(H)_{11} + 2H_{13}T_m(H)_{31}. \quad (2.21)$$

Recurising on the partitioning will lead to a nested dissection, after which we will use direct computation on sufficiently small sub-blocks. We denote the indexing of each partition with $I_p^{(t)} = I_s^{(t)} \cup I_\ell^{(t)} \cup I_r^{(t)}$, which represents all nodes in the current partition, the separators, and two sub-partitions, respectively. For the separators, Eq. (2.21) leads to

$$\begin{aligned} T_{m+1}(H)(I_p^{(t)}, I_s^{(t)}) &= 2H(I_p^{(t)}, I_p^{(t)})T_m(H)(I_p^{(t)}, I_s^{(t)}) \\ &\quad - T_{m-1}(H)(I_p^{(t)}, I_s^{(t)}) + 2 \sum_{t' \in S_t} H(I_p^{(t)}, I_s^{(t')})T_m(H)(I_s^{(t')}, I_s^{(t)}), \end{aligned} \quad (2.22)$$

where S_t is the path from partition t to the root; and for the leaf blocks, $I_s^{(t)} = I_p^{(t)}$ in Eq. (2.22). The result is Algorithm 1.

The multilevel nested dissection process itself has a well-established algorithm by Karypis and Kumar, and efficient implementation is available in **METIS** [94]. Note that this approach is only viable when the graph can be partitioned with a separator of small size. Empirically, we observe this assumption to hold for many real-world networks. The biggest advantage of this approach is we can very efficiently obtain PDOS estimation for a subset of nodes with much better accuracy than KPM.

Algorithm 1: Nested Dissection for PDOS Approximation

Input: Symmetric graph matrix H with eigenvalues in $[-1, 1]$
Output: $C \in \mathbb{R}^{N \times M}$ where c_{ij} is the j -th Chebyshev moment for i -th node.

begin

- Obtain partitions $\{I_p^{(t)}\}$ in a tree structure through multilevel nested dissection.
- for** $m = 1$ **to** M **do**
- Traverse partition tree in pre-order:
- Compute the separator columns with Eq. (2.22).
- if** $I_p^{(t)}$ is a leaf block **then**
- Compute the diagonal entries with equation (2.22).
- end**
- end**
- end**

2.4.4 Motif Filtering

In many graphs, there are large spikes around particular eigenvalues; for example, see Fig. 2.1. This phenomenon affects the accuracy of DOS estimation in two ways. First, the singularity-like behavior means we need many more moments to obtain a good approximation in polynomial basis. Secondly, due to the equi-oscillation property of Chebyshev approximation, error in irregularities (say, at a point of high concentration in the spectral density), spreads to other parts of the spectrum. This is a problem in our case, as the spectral density of real-world networks are far from uniform.

High multiplicity eigenvalues are typically related to local symmetries in a graph. The most prevalent example is two dangling nodes attached to the same neighbor as shown in Fig. 2.2a, which accounts for most eigenvalues around 0 for (normalized) adjacency matrix with a localized eigenvector taking value +1 on one node and -1 on the other. In addition, we list a few more motifs in Fig. 2.2 that appear most frequently in real-world graphs. All of them can be associated with specific eigenvalues, and we include the corresponding ones in normalized adjacency matrix for our example.

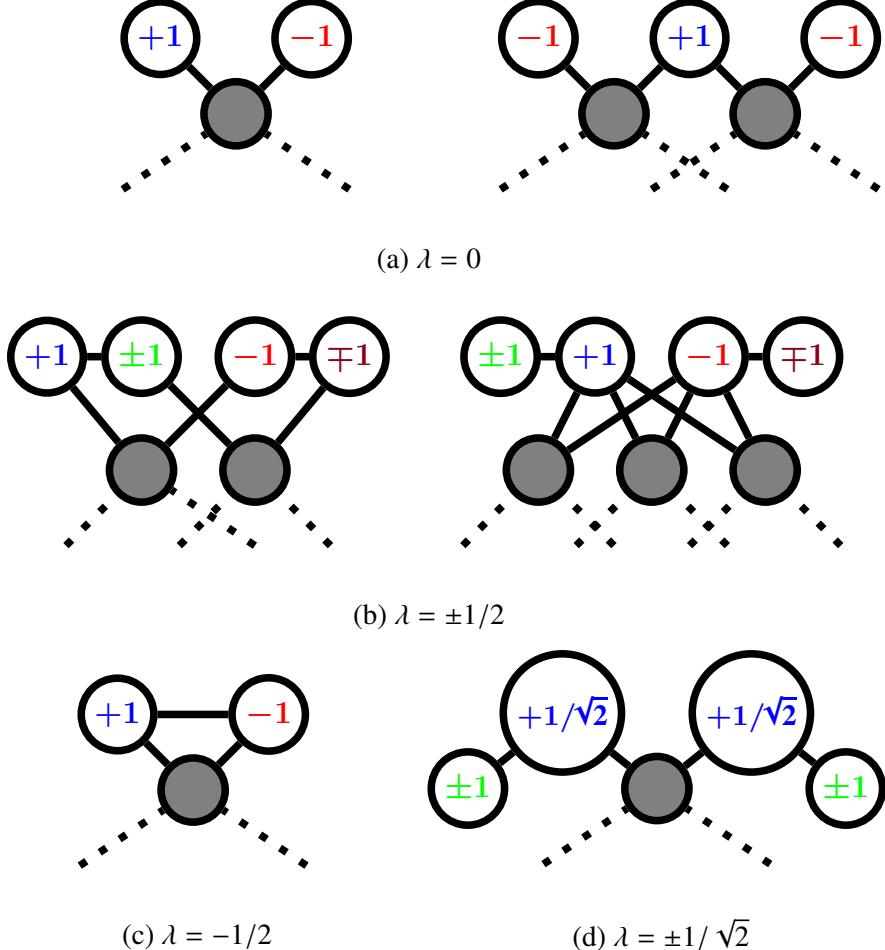


Figure 2.2: Common motifs (induced subgraphs) in graph data that result in localized spikes in the spectral density. Each motif generates a specific eigenvalue with locally-supported eigenvectors. Here we use the normalized adjacency matrix to represent the graph, although we can perform the same analysis for the adjacency, Laplacian, or normalized Laplacian (only the eigenvalues would be different). The eigenvectors are supported only on the labeled nodes.

To detect these motifs in large graphs, we deploy a randomized hashing technique. Given a random vector z , the hashing weight $w = Hz$ encodes all the neighborhood information of each node. To find node copies (left in Figure 2.2a), we seek pairs (i, j) such that $w_i = w_j$; with high probability, this only happens when v_i and v_j share the same neighbors. Similarly, all motifs in Figure 2.2 can be characterized by union and intersection of neighborhood lists.

After identifying motifs, we need only approximate the (relatively smooth) density of the remaining spectrum. The eigenvectors associated with these remaining non-motif eigenvalues must be constant across cycles in the canonical decomposition of the associated permutations. Let $P \in \mathbb{R}^{N \times r}$ denote an orthonormal basis for the space of such vectors formed from columns of the identity and (normalized) indicators for nodes cyclically permuted by the motif. The matrix $H_r = P^T HP$ then has identical eigenvalues to H , except with all the motif eigenvalues omitted. We may form H_r explicitly, as it has the same sparsity structure as H but with a supernode replacing the nodes in each instance of a motif cycle; or we can achieve the same result by replacing each random probe Z with the projected probe $Z_r = PP^T Z$ at an additional cost of $O(N_{\text{motif}})$ per probe, where N_{motif} is the number of nodes involved in motifs.

The motif filtering method essentially allow us to isolate the spiky components from the spectrum. As a result, we are able to obtain a more accurate approximation using fewer Chebyshev moments. Figure 2.3 demonstrates the improvement on the approximation as we procedurally filter out motifs at $0, -1/3, -1/2$, and $-1/4$. The eigenvalue $-1/m$ can be generated by an edge attached to the graph through $m - 1$ nodes, similar to motif (2.2c).

2.5 Error Analysis

2.5.1 KPM Approximation Error

This section provides an error bound for our regularized DOS approximation $K_\sigma * \mu$. We will start with the following theorem.

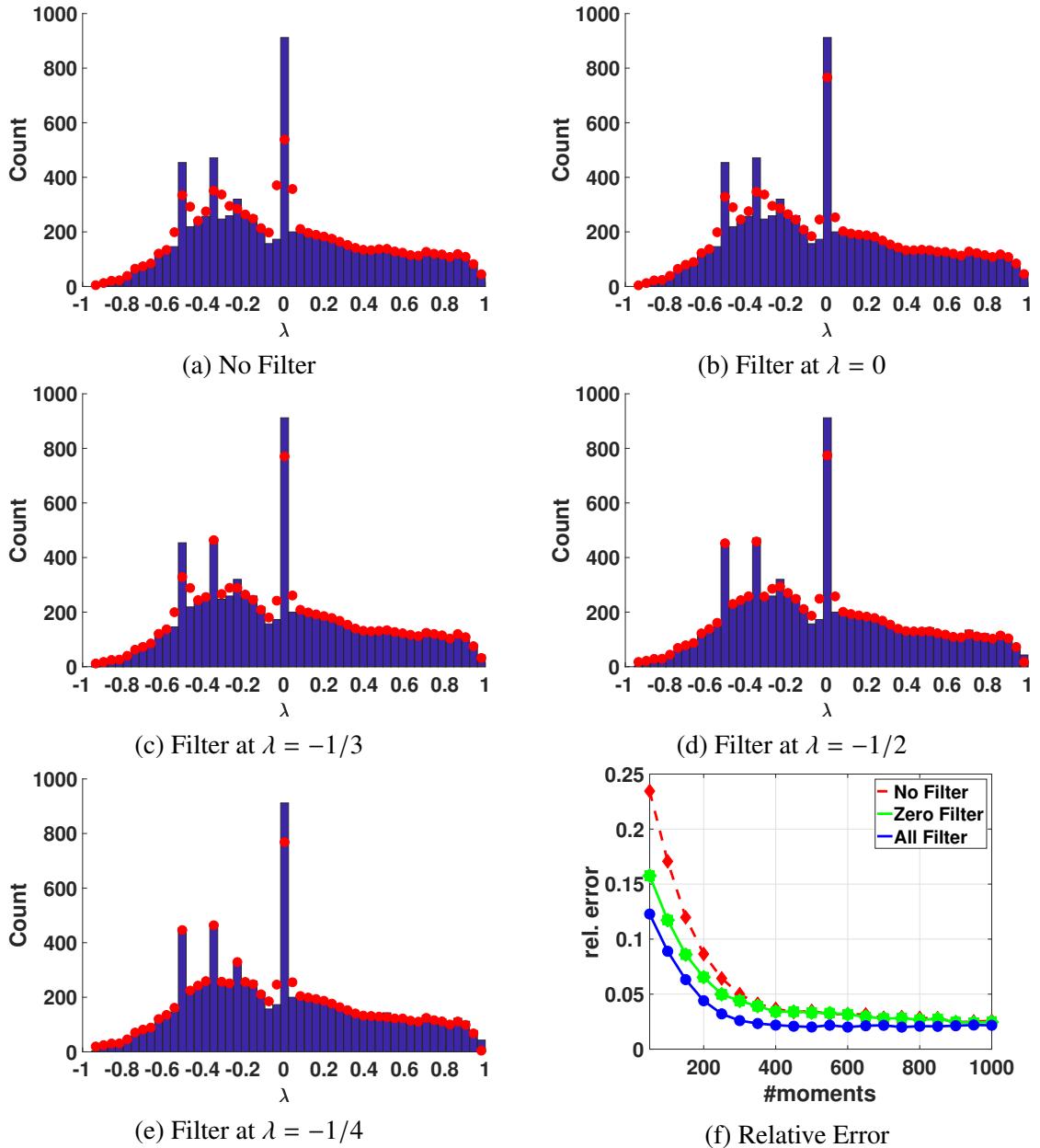


Figure 2.3: The improvement in accuracy of the spectral histogram approximation on the normalized adjacency matrix for the High Energy Physics Theory (HepTh) Collaboration Network, as we sweep through spectrum and filter out motifs. The graph has 8,638 nodes and 24,816 edges. Blue bars are the real spectrum, and red points are the approximated heights. Fig. 2.3a- 2.3e use 100 moments and 20 probe vectors. Fig. 2.3f shows the relative L_1 error of the spectral histogram when using no filter, filter at $\lambda = 0$, and all filters.

Theorem 1 (Jackson's Theorem [89]). If $f : [-1, 1] \rightarrow \mathbb{R}$ is Lipschitz continuous with constant L , its best degree M polynomial approximation \hat{f}^M has an L_∞ error of at most $6L/M$. The approximation can be constructed as

$$\hat{f}^M = \sum_{m=0}^M J_m c_m T_m(x), \quad (2.23)$$

where J_m are Jackson smoothing factors and c_m are the Chebyshev coefficients.

We can pick a smooth mollifier K with $\text{Lip}(K) = 1$. For any $\nu \in \mathbb{R}$ and $\lambda \in [-1, 1]$ there exists a degree M polynomial such that

$$\left| K_\sigma(\nu - \lambda) - \widehat{K}_\sigma^M(\nu - \lambda) \right| < \frac{6L}{M\sigma}. \quad (2.24)$$

Define $\hat{\mu}^M = \sum_{m=0}^M J_m d_m \phi_m$ to be the truncated DOS series,

$$\int_{-1}^1 \hat{f}^M(\lambda) \mu(\lambda) d\lambda = \int_{-1}^1 f(\lambda) \hat{\mu}^M(\lambda) d\lambda = \sum_{m=0}^M J_m c_m d_m. \quad (2.25)$$

Therefore,

$$\begin{aligned} \|K_\sigma * (\mu - \hat{\mu}^M)\|_\infty &= \max_\nu \left| \int_{-1}^1 K_\sigma(\nu - \lambda) (\mu(\lambda) - \hat{\mu}^M(\lambda)) d\lambda \right| \\ &\leq \max_\nu \int_{-1}^1 \left| K_\sigma(\nu - \lambda) - \widehat{K}_\sigma^M(\nu - \lambda) \right| \mu(\lambda) d\lambda \\ &\leq \frac{6L}{M\sigma}. \end{aligned}$$

Consider $\tilde{\mu}^M$ to be the degree M approximation from KPM,

$$\|K_\sigma * (\mu - \tilde{\mu}^M)\|_\infty \leq \|K_\sigma * (\mu - \hat{\mu}^M)\|_\infty + \|K_\sigma\|_\infty \|\hat{\mu}^M - \tilde{\mu}^M\|_1. \quad (2.26)$$

If we use a probe z with independent standard normal entries for the trace estimation,

$$\tilde{\mu}(\lambda) = \sum_{i=1}^N w_i^2 \delta(\lambda - \lambda_i) \quad (2.27)$$

where $w = Q^T z$ is the weight for z in the eigenbasis. Hence

$$\|\hat{\mu}^M - \tilde{\mu}^M\|_1 \leq \sum_{i=1}^N |1 - w_i^2|. \quad (2.28)$$

Finally,

$$\mathbb{E}[\|K_\sigma * (\mu - \tilde{\mu}^M)\|] \leq \frac{1}{\sigma} \left(\frac{6L}{M} + \|K\|_\infty \mathbb{E}[|1 - w_1^2|] \right). \quad (2.29)$$

If we take N_z independent probe vectors, then $N_z w_1^2 \sim \chi^2(N_z)$, which means the expectation decays asymptotically like $\sqrt{2/(\pi N_z)}$.

2.5.2 Perturbation Analysis

In this section, we limit our attention to symmetric graph matrix H . Extracting graph information using DOS, whether as a distribution for functions on a graph or as a direct feature in the form of spectral moments, requires stability under small perturbations. In the case of removing/adding a few number of nodes/edges, the Cauchy Interlacing Theorem [119] gives a bound on each individual new eigenvalue by the old ones. For example, if we remove $r \ll N$ nodes to get a new graph matrix \tilde{H} , then

$$\lambda_i(H) \leq \lambda_i(\tilde{H}) \leq \lambda_{i+r}(H) \quad \text{for } i \leq N - r. \quad (2.30)$$

However, this bound may not be helpful when the impact of the change is not reflected by its size. Hence, we provide a theorem that relates the Wasserstein distance (see Eq. (2.4)) change and the Frobenius norm of the perturbation. Without loss of generality, we assume the eigenvalues of H lie in $[-1, 1]$ already.

Theorem 2. Suppose $\tilde{H} = H + \delta H$ is the perturbed graph matrix with spectral density $\tilde{\mu}$, then

$$W_1(\mu, \tilde{\mu}) \leq \|\delta H\|_F$$

Proof. Let \mathcal{L} be the space of Lipschitz functions with $f(0) = 0$.

$$\begin{aligned} W_1(\mu, \tilde{\mu}) &= \sup_{f \in \mathcal{L}, \text{Lip}(f)=1} \int f(\lambda)(\mu(\lambda) - \tilde{\mu}(\lambda))d\lambda \\ &= \frac{1}{N} \sup_{f \in \mathcal{L}, \text{Lip}(f)=1} \text{tr}(f(H) - f(\tilde{H})) \\ &\leq \sup_{f \in \mathcal{L}, \text{Lip}(f)=1, \|v\|=1} v^T(f(H) - f(\tilde{H}))v. \end{aligned}$$

By Theorem 3.8 from Higham [78], the perturbation on $f(H)$ is bounded by the Fréchet derivative,

$$\|f(H) - f(\tilde{H})\|_2 \leq \text{Lip}(f)\|\delta H\|_F + o(\|\delta H\|_F). \quad (2.31)$$

□

2.6 Experiments

2.6.1 Gallery of DOS/PDOS

We first present our spectral histogram approximation from DOS/ PDOS on a wide variety of graphs, including collaboration networks, online social networks, road networks and autonomous systems (dataset details are in the appendix). For all examples, we apply our methods to the normalized adjacency matrices using 500 Chebyshev moments and 20 Hadamard probe vectors. Afterwards, the spectral density is integrated into 50 histogram bins. In Fig. 2.4, the DOS approximation is on the first row, and the PDOS approximation is on the second. When a spike exists in the spectrum, we apply motif filtering, and DOS is zoomed appropriately to show the remaining part. For PDOS, we stack the spectral histograms for all nodes vertically, sorted by their projected weights on the leading left singular vector. Red indicates that a node has high weight at certain parts of the spectrum, whereas blue indicates low weight.

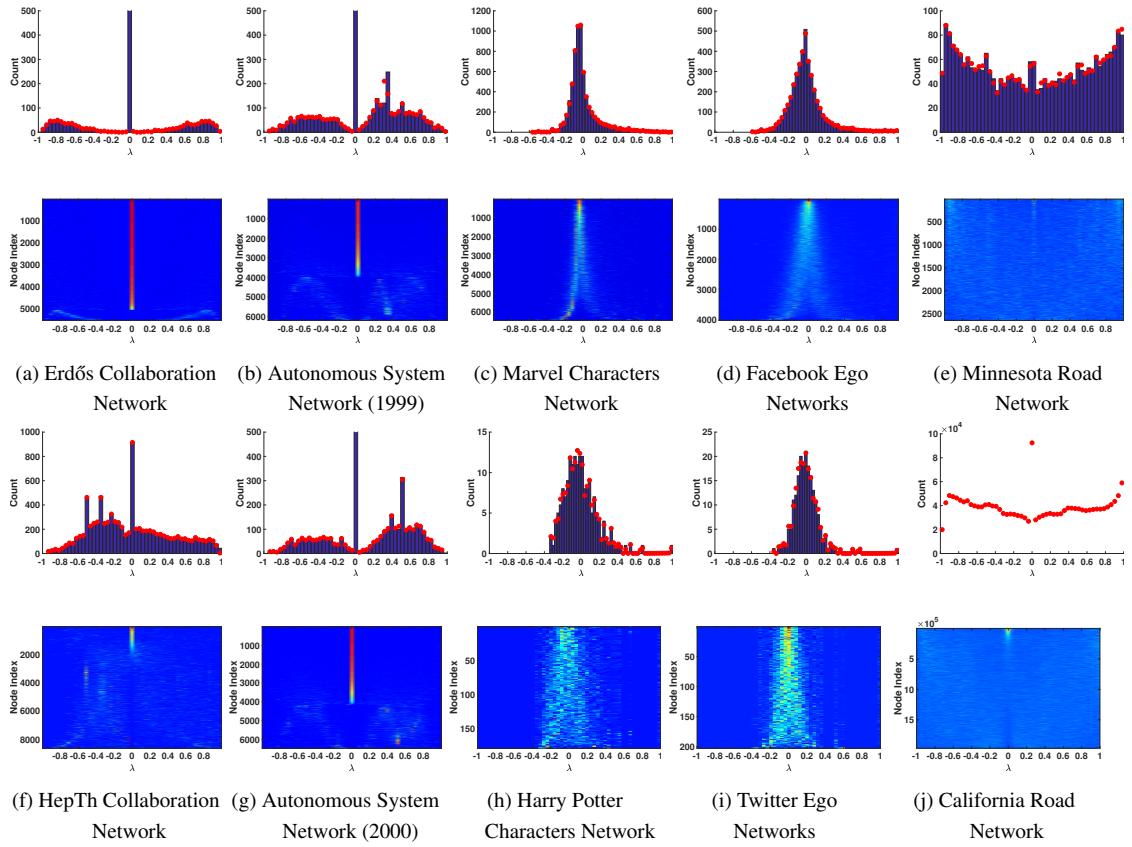


Figure 2.4: DOS(top)/PDOS(bottom) histograms for the normalized adjacency of 10 networks from five domains. For DOS, blue bars are the true spectrum, and red points are from KPM (500 moments and 20 Hadamard probes). For PDOS, the spectral histograms of all nodes are aligned vertically. Red indicates high weight around an eigenvalue, and blue indicates low weight. The true spectrum for the California Road Network (2.4j) is omitted, as it is too large to compute exactly (1,965,206 nodes).

We observe many distinct shapes of spectrum in our examples. The eigenvalues of denser graphs, such as the Marvel characters network (2.4c: average degree 52.16) and Facebook union of ego networks (2.4d: average degree 43.69), exhibit decay similar to the power-law around $\lambda = 0$. There has been study on the power-law distribution in the eigenvalues of the adjacency and the Laplacian matrix, but it only focuses on the leading eigenvalues rather than the entire spectrum [52] for large real-world datasets. Relatively sparse graphs (2.4a: average degree 3.06; 2.4b: average degree 4.13) often possess spikes, especially around $\lambda = 0$, which reflect a larger set of loosely-connected

boundary nodes. It is much more evident in the PDOS spectral histograms, which allow us to pick out the nodes with dominant weights at $\lambda = 0$ and those that contribute most to local structures. Finally, though the road network is quite sparse (ave. deg 2.50), its regularity results in a lack of special features, and most nodes contribute evenly to the spectrum according to PDOS.

2.6.2 Computation time

In this experiment, we show the scaling of our methods by applying them to graphs of varying size of nodes, edges, and sparsity patterns. Rather than computation power, the memory cost of loading a graph with 100M-1B edges is more often the constraint. Hence, we report runtimes for a Python version on a Google Cloud instance with 200GB memory and an Intel Xeon E5 v3 CPU at 2.30GHz.

The datasets we use are obtained from the SNAP repository [104]. For each graph, we compute the first 10 Chebyshev moments using KPM with 20 probe vectors. Most importantly, the cost for each moment is independent of the total number of moments we compute. Table 2.1 reports number of nodes, number of edges, average degree of nodes, and the average runtime for computing each moment. We can observe that the runtime is in accordance with the theoretical complexity $O(N_z(|V| + |E|))$. For the Friendster social network with about 1.8 billion edges, computing each moment takes about 1000 seconds to compute, which means we could obtain a rough approximation to its spectrum within a day. As the dominant cost is matrix-matrix multiplication and we use several probe vectors, our approach has ample opportunity for parallel computation.

Table 2.1: Average Computation Time per Chebyshev Moment for Graphs from the SNAP Repository^a.

Network	# Nodes	# Edges	Avg. Deg.	Time (s)
Facebook	4,039	88,234	43.69	0.007
AstroPh	18,772	198,110	21.11	0.028
Enron	36,692	183,831	10.02	0.046
Gplus	107,614	13,673,453	254.12	1.133
Amazon	334,863	925,872	5.53	0.628
Neuron	1,018,524	24,735,503	48.57	9.138
RoadNetCA	1,965,206	2,766,607	2.82	2.276
Orkut	3,072,441	117,185,083	76.28	153.7
LiveJournal	3,997,962	34,681,189	17.35	14.52
Friendster	65,608,366	1,806,067,135	55.06	1,017

^a 20 probe vectors are used throughout the experiment. The runtime is averaged over 5 moments.

2.6.3 Model Verification

In this experiment, we investigate the spectrum for some of the popular graph models, and whether they resemble the behavior of real-world data. Two of the most popular models used to describe real-world graphs are the scale-free model [18] and the small-world model [167]. Farkas et al. [53] has analyzed the spectrum of the adjacency matrix; we instead consider the normalized adjacency.

The scale-free model grows a random graph with the preferential attachment process, starting from an initial seed graph and adding one node and m edges at every step. Fig. 2.5 shows spectral histograms for this model with 5000 nodes and different choices of m . When $m = 1$, the generated graph has abundant local motifs like many sparse real-world graphs. By searching in PDOS for the nodes that have high weight at the two spikes, we find node-doubles ($\lambda = 0$) and singly-attached chains ($\lambda = \pm 1/\sqrt{2}$). When $m = 5$, the graph is denser, without any particular motifs, resulting in an approximately semicircular spectral distribution.

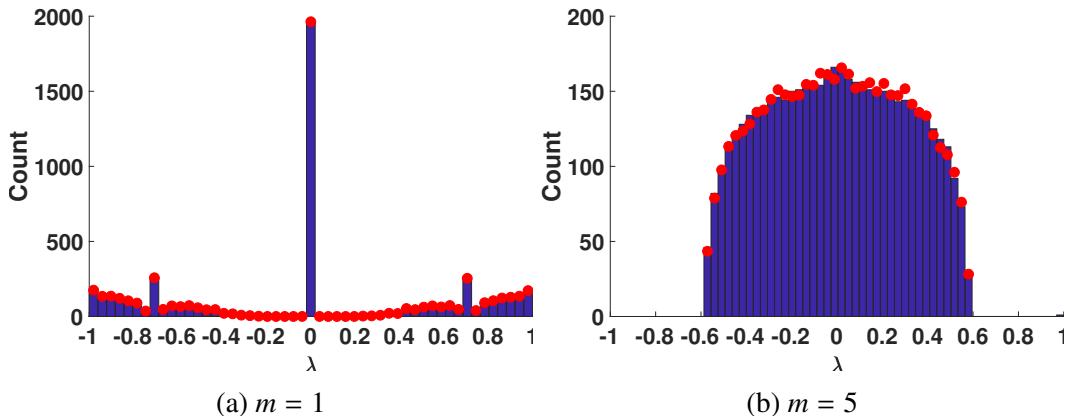


Figure 2.5: Spectral histogram for scale-free model with 5000 nodes and different m . Blue bars are the real spectrum, red points are from KPM (500 moments and 20 probes).

The small-world model generates a random graph by re-wiring edges of a ring lattice with a certain probability p . Here we construct these graphs on 5000 nodes with $p = 0.5$; the pattern in spectrum is insensitive for a wide range of p . In Fig. 2.6, when the graph is sparse with 5000 edges, the spectrum has spikes at 0 and ± 1 , indicating local symmetries, bipartite structure, and disconnected components. With 50000 edges, localized structures disappear and the spectrum has narrower support.

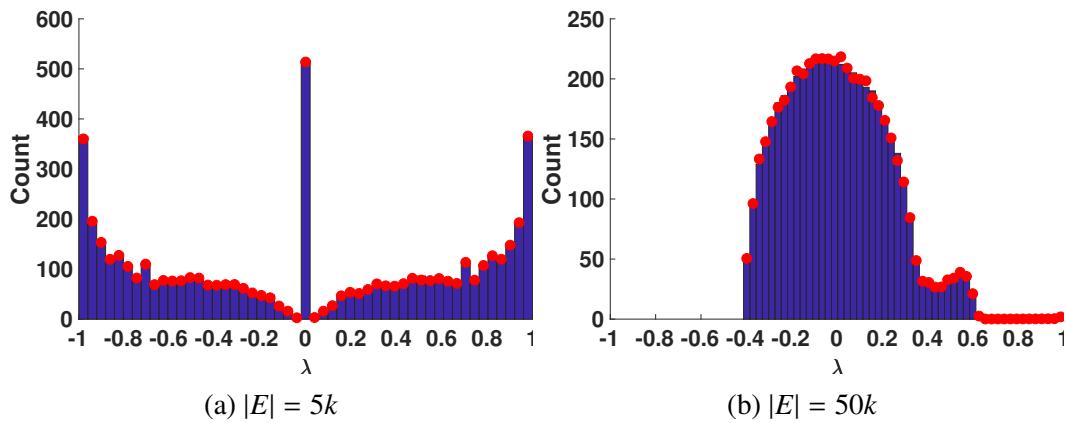


Figure 2.6: Spectral histograms for small-world model with 5000 nodes and re-wiring probability $p = 0.5$, starting with 5000 (2.6a) and 50000 (2.6b) edges. Blue bars are the real spectrum, red points are from KPM (5000 moments and 20 probes).

Finally, we investigate the Block Two-Level Erdős-Rényi (BTER) model [154],

which directly fits an input graph. BTER constructs a similar graph by a two-step process: first create a collection of Erdős-Rényi subgraphs, then interconnect those using a Chung-Lu model [35]. Seshadhri et al. showed their model accurately captures the observable properties of the given graph, including the eigenvalues of the adjacency matrix. Fig. 2.7 compares the DOS/PDOS of the Erdős collaboration network and its BTER counterpart. Unlike the original graph, most 0 eigenvalues in BTER graph come from isolated nodes. The BTER graph also has many more isolated edges ($\lambda = \pm 1$), singly-attached chains ($\lambda = \pm 1/\sqrt{2}$), and singly-attached triangles ($\lambda = -1/2$). We locate these motifs by inspecting nodes with high weights at respective part of the spectrum.

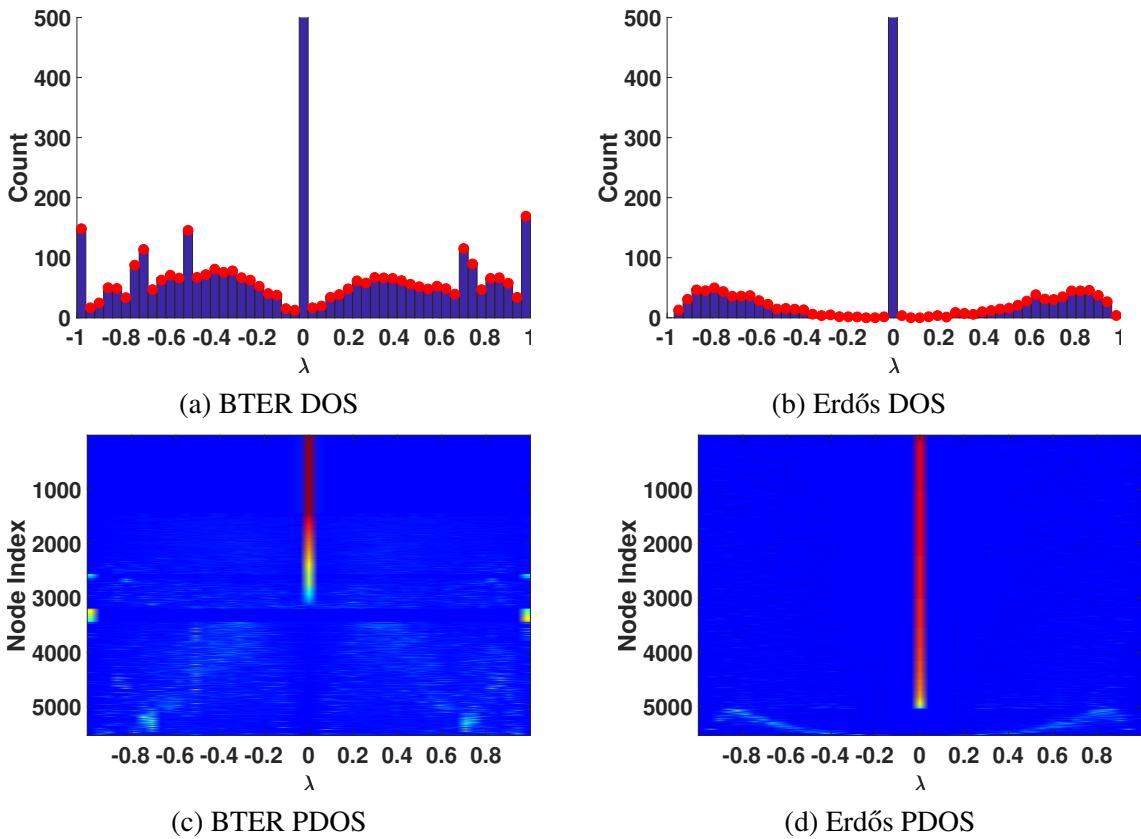


Figure 2.7: Comparison of spectral histogram between Erdős Collaboration Network and the BTER model. Both DOS and PDOS are computed with 500 moments and 20 probe vectors.

2.7 Conclusion

In this paper, we make the computation of spectral densities a practical tool for the analysis of large real-world network. Our approach borrows from methods in solid state physics, but with adaptations that improve performance in the network analysis setting by special handling of graph motifs that leave distinctive spectral fingerprints. We show that the spectral densities are stable to small changes in the graph, as well as providing an analysis of the approximation error in our methods. We illustrate the efficiency of our approach by treating graphs with tens of millions of nodes and billions of edges using only a single compute node. The method provides a compelling visual fingerprint of a graph, and we show how this fingerprint can be used for tasks such as model verification.

Our approach opens the door for the use of complete spectral information in large-scale network analysis. It provides a framework for scalable computation of quantities already used in network science, such as common centrality measures and graph connectivity indices (such as the Estrada index) that can be expressed in terms of the diagonals and traces of matrix functions. But we expect it to serve more generally to define new families of features that describe graphs and the roles nodes play within those graphs. We have shown that graphs from different backgrounds demonstrate distinct spectral characteristics, and thus can be clustered based on those. Looking at LDOS across nodes for role discovery, we can identify the ones with high similarity in their local structures. Moreover, extracting nodes with large weights at various points of the spectrum uncovers motifs and symmetries. In the future, we expect to use DOS/LDOS as graph features for applications in graph clustering, graph matching, role classification, and other tasks.

Acknowledgments. We thank NSF DMS-1620038 for supporting this work.

CHAPTER 3

SCALABLE GAUSSIAN PROCESSES

3.1 Abstract

For applications as varied as Bayesian neural networks, determinantal point processes, elliptical graphical models, and kernel learning for Gaussian processes, one must compute a log determinant of an $n \times n$ positive definite matrix, and its derivatives – leading to prohibitive $O(n^3)$ computations. We propose novel $O(n)$ approaches to estimating these quantities from only fast matrix vector multiplications (MVMs). These stochastic approximations are based on Chebyshev, Lanczos, and surrogate models, and converge quickly even for kernel matrices that have challenging spectra. We leverage these approximations to develop a scalable Gaussian process approach to kernel learning. We find that Lanczos is generally superior to Chebyshev for kernel learning, and that a surrogate approach can be highly efficient and accurate with popular kernels.

3.2 Introduction

There is a pressing need for scalable machine learning approaches to extract rich statistical structure from large datasets. A common bottleneck — arising in determinantal point processes [98], Bayesian neural networks [116], model comparison [117], graphical models [151], and Gaussian process kernel learning [145] — is computing a log determinant over a large positive definite matrix. While we can approximate log determinants by existing stochastic expansions relying on matrix vector multiplications (MVMs), these approaches make assumptions, such as near-uniform eigenspectra [28],

which are unsuitable in machine learning contexts. For example, the popular RBF kernel gives rise to rapidly decaying eigenvalues. Moreover, while standard approaches, such as stochastic power series, have reasonable asymptotic complexity in the rank of the matrix, they require too many terms (MVMs) for the precision necessary in machine learning applications.

Gaussian processes (GPs) provide a principled probabilistic kernel learning framework, for which a log determinant is of foundational importance. Specifically, the *marginal likelihood* of a Gaussian process is the probability of data given only kernel hyper-parameters. This utility function for kernel learning compartmentalizes into automatically calibrated model fit and complexity terms — called *automatic Occam’s razor* — such that the simplest models which explain the data are automatically favored [146, 145], without the need for approaches such as cross-validation, or regularization, which can be costly, heuristic, and involve substantial hand-tuning and human intervention. The automatic complexity penalty, called the *Occam’s factor* [117], is a log determinant of a kernel (covariance) matrix, related to the volume of solutions that can be expressed by the Gaussian process.

Many current approaches to scalable Gaussian processes [e.g., 143, 101, 75] focus on inference assuming a fixed kernel, or use approximations that do not allow for very flexible kernel learning [175], due to poor scaling with number of basis functions or inducing points. Alternatively, approaches which exploit algebraic structure in kernel matrices can provide highly expressive kernel learning [177], but are essentially limited to grid structured data.

Recently, Wilson and Nickisch [176] proposed the *structured kernel interpolation* (SKI) framework, which generalizes structuring exploiting methods to arbitrarily located data. SKI works by providing accurate and fast matrix vector multiplies (MVMs)

with kernel matrices, which can then be used in iterative solvers such as linear conjugate gradients for scalable GP inference. However, evaluating the marginal likelihood and its derivatives, for kernel learning, has followed a scaled eigenvalue approach [177, 176] instead of iterative MVM approaches. This approach can be inaccurate, and relies on a fast eigendecomposition of a structured matrix, which is not available in many consequential situations where fast MVMs are available, including: (i) additive covariance functions, (ii) multi-task learning, (iii) change-points [76], and (iv) diagonal corrections to kernel approximations [157]. Fiedler [55] and Weyl [172] bounds have been used to extend the scaled eigenvalue approach [56, 76], but are similarly limited. These extensions are often very approximate, and do not apply beyond sums of two and three matrices, where each matrix in the sum must have a fast eigendecomposition.

In machine learning there has recently been renewed interest in MVM based approaches to approximating log determinants, such as the Chebyshev [72] and Lanczos [164] based methods, although these approaches go back at least two decades in quantum chemistry computations [14]. Independently, several authors have proposed various methods to compute derivatives of log determinants [115, 158]. But *both* the log determinant *and* the derivatives are needed for efficient GP marginal likelihood learning: the derivatives are required for gradient-based optimization, while the log determinant itself is needed for model comparison, comparisons between the likelihoods at local maximizers, and fast and effective choices of starting points and step sizes in a gradient-based optimization algorithm.

In this paper, we develop novel scalable and general purpose Chebyshev, Lanczos, and surrogate approaches for efficiently and accurately computing both the log determinant and its derivatives simultaneously. Our methods use only fast MVMs, and re-use the same MVMs for both computations. In particular:

- We derive fast methods for simultaneously computing the log determinant and its derivatives by stochastic Chebyshev, stochastic Lanczos, and surrogate models, from MVMs alone. We also perform an error analysis and extend these approaches to higher order derivatives.
- These methods enable fast GP kernel learning whenever fast MVMs are possible, including applications where alternatives such as scaled eigenvalue methods (which rely on fast eigendecompositions) are not, such as for (i) diagonal corrections for better kernel approximations, (ii) additive covariances, (iii) multi-task approaches, and (iv) non-Gaussian likelihoods.
- We illustrate the performance of our approach on several large, multi-dimensional datasets, including a consequential crime prediction problem, and a precipitation problem with $n = 528,474$ training points. We consider a variety of kernels, including deep kernels [178], diagonal corrections, and both Gaussian and non-Gaussian likelihoods.
- We have released code and tutorials as an extension to the GPML library [147] at https://github.com/kd383/GPML_SLD. A Python implementation of our approach is also available through the GPyTorch library: <https://github.com/jrg365/gpytorch>.

When using our approach in conjunction with SKI [176] for fast MVMs, GP kernel learning is $O(n + g(m))$, for m inducing points and n training points, where $g(m) \leq m \log m$. With algebraic approaches such as SKI we also do not need to worry about quadratic storage in inducing points, since symmetric Toeplitz and Kronecker matrices can be stored with at most linear cost, without needing to explicitly construct a matrix.

Although we here use SKI for fast MVMs, we emphasize that the proposed iterative approaches are generally applicable, and can easily be used in conjunction with *any*

method that admits fast MVMs, including classical inducing point methods [143], finite basis expansions [101], and the popular stochastic variational approaches [75]. Moreover, stochastic variational approaches can naturally be combined with SKI to further accelerate MVMs [174].

We start in §3.3 with an introduction to GPs and kernel approximations. In §3.4 we introduce stochastic trace estimation and Chebyshev (§3.4.1) and Lanczos (§3.4.2) approximations. In §3.5, we describe the different sources of error in our approximations. In §3.6 we consider experiments on several large real-world data sets. We conclude in §3.7. The supplementary materials also contain several additional experiments and details.

3.3 Background

A Gaussian process (GP) is a collection of random variables, any finite number of which have a joint Gaussian distribution [e.g., 145]. A GP can be used to define a distribution over functions $f(x) \sim \mathcal{GP}(\mu(x), k(x, x'))$, where each function value is a random variable indexed by $x \in \mathbb{R}^d$, and $\mu : \mathbb{R}^d \rightarrow \mathbb{R}$ and $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ are the mean and covariance functions of the process.

The covariance function is often chosen to be an RBF or Matérn kernel (see the supplementary material for more details). We denote any kernel hyperparameters by the vector θ . To be concise we will generally not explicitly denote the dependence of k and associated matrices on θ .

For any locations $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$, $f_X \sim \mathcal{N}(\mu_X, K_{XX})$ where f_X and μ_X represent the vectors of function values for f and μ evaluated at each of the $x_i \in X$, and K_{XX}

is the matrix whose (i, j) entry is $k(x_i, x_j)$. Suppose we have a vector of corresponding function values $y \in \mathbb{R}^n$, where each entry is contaminated by independent Gaussian noise with variance σ^2 . Under a Gaussian process prior depending on the covariance hyperparameters θ , the log marginal likelihood is given by

$$\mathcal{L}(\theta|y) = -\frac{1}{2} \left[(y - \mu_X)^T \alpha + \log |\tilde{K}_{XX}| + n \log 2\pi \right], \quad (3.1)$$

where $\alpha = \tilde{K}_{XX}^{-1}(y - \mu_X)$ and $\tilde{K}_{XX} = K_{XX} + \sigma^2 I$. Optimization of (3.1) is expensive, since the cheapest way of evaluating $\log |\tilde{K}_{XX}|$ and its derivatives without taking advantage of the structure of \tilde{K}_{XX} involves computing the $O(n^3)$ Cholesky factorization of \tilde{K}_{XX} . $O(n^3)$ computations is too expensive for inference and learning beyond even just a few thousand points.

A popular approach to GP scalability is to replace the exact kernel $k(x, z)$ by an approximate kernel that admits fast computations [143]. Several methods approximate $k(x, z)$ via *inducing points* $U = \{u_j\}_{j=1}^m \subset \mathbb{R}^d$. An example is the subset of regressor (SoR) kernel:

$$k^{\text{SoR}}(x, z) = K_{xU} K_{UU}^{-1} K_{Uz}, \quad (3.2)$$

which is a low-rank approximation [155]. The SoR matrix $\tilde{K}_{XX}^{\text{SoR}} \in \mathbb{R}^{n \times n}$ has rank at most m , allowing us to solve linear systems involving $\tilde{K}_{XX}^{\text{SoR}} = K_{XX}^{\text{SoR}} + \sigma^2 I$ and to compute $\log |\tilde{K}_{XX}^{\text{SoR}}|$ in $O(m^2 n + m^3)$ time. Another popular kernel approximation is the fully independent training conditional (FITC), which is a diagonal correction of SoR so that the diagonal is the same as for the original kernel [157]. Thus kernel matrices from FITC have low-rank plus diagonal structure. This modification has had exceptional practical significance, leading to improved point predictions and much more realistic predictive uncertainty [143, 144], making FITC arguably the most popular approach for scalable Gaussian processes.

Wilson and Nickisch [176] provides a mechanism for fast MVMs through proposing

the structured kernel interpolation (SKI) approximation,

$$K_{XX} \approx WK_{UU}W^T, \quad (3.3)$$

where W is an n -by- m matrix of interpolation weights; the authors of [176] use local cubic interpolation so that W is sparse. The sparsity in W makes it possible to naturally exploit algebraic structure (such as Kronecker or Toeplitz structure) in K_{UU} when the inducing points U are on a grid, for extremely fast matrix vector multiplications with the approximate K_{XX} even if the data inputs X are arbitrarily located. For instance, if K_{UU} is Toeplitz, then each MVM with the approximate K_{XX} costs only $O(n + m \log m)$. By contrast, placing the inducing points U on a grid for classical inducing point methods, such as SoR or FITC, does not result in substantial performance gains, due to the costly cross-covariance matrices K_{xU} and K_{Uz} .

3.4 Methods

Our goal is to estimate, for a symmetric positive definite matrix \tilde{K} ,

$$\log |\tilde{K}| = \text{tr}(\log(\tilde{K})) \quad \text{and} \quad \frac{\partial}{\partial \theta_i} [\log |\tilde{K}|] = \text{tr}\left(\tilde{K}^{-1} \left(\frac{\partial \tilde{K}}{\partial \theta_i}\right)\right), \quad (3.4)$$

where \log is the matrix logarithm [78]. We compute the traces involved in both the log determinant and its derivative via *stochastic trace estimators* [88], which approximate the trace of a matrix using only matrix vector products.

The key idea is that for a given matrix A and a random probe vector z with independent entries with mean zero and variance one, then $\text{tr}(A) = \mathbb{E}[z^T A z]$; a common choice is to let the entries of the probe vectors be Rademacher random variables. In practice, we estimate the trace by the sample mean over n_z independent probe vectors. Often surprisingly few probe vectors suffice.

To estimate $\text{tr}(\log(\tilde{K}))$, we need to multiply $\log(\tilde{K})$ by probe vectors. We consider two ways to estimate $\log(\tilde{K})z$: by a polynomial approximation of \log or by using the connection between the Gaussian quadrature rule and the Lanczos method [72, 164]. In both cases, we show how to re-use the same probe vectors for an inexpensive coupled estimator of the derivatives. In addition, we may use standard radial basis function interpolation of the log determinant evaluated at a few systematically chosen points in the hyperparameter space as an inexpensive surrogate for the log determinant.

3.4.1 Chebyshev Expansion

Chebyshev polynomials are defined by the recursion

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_{j+1}(x) = 2xT_j(x) - T_{j-1}(x) \text{ for } j \geq 1. \quad (3.5)$$

For $f : [-1, 1] \rightarrow \mathbb{R}$ the Chebyshev interpolant of degree m is

$$f(x) \approx p_m(x) := \sum_{j=0}^m c_j T_j(x), \quad \text{where } c_j = \frac{2 - \delta_{j0}}{m+1} \sum_{k=0}^m f(x_k) T_j(x_k), \quad (3.6)$$

and δ_{j0} is the Kronecker delta and $x_k = \cos(\pi(k + 1/2)/(m + 1))$ for $k = 0, 1, 2, \dots, m$; see [60]. Using the Chebyshev interpolant of $\log(1 + \alpha x)$, we approximate $\log|\tilde{K}|$ by

$$\log|\tilde{K}| - n \log \beta = \log|I + \alpha B| \approx \sum_{j=0}^m c_j \text{tr}(T_j(B)), \quad (3.7)$$

when $B = (\tilde{K}/\beta - 1)/\alpha$ has eigenvalues $\lambda_i \in (-1, 1)$.

For stochastic estimation of $\text{tr}(T_j(B))$, we only need to compute $z^T T_j(B) z$ for each given probe vector z . We compute vectors $w_j = T_j(B)z$ and $\partial w_j / \partial \theta_i$ via the coupled recurrences

$$w_0 = z, \quad w_1 = Bz, \quad w_{j+1} = 2Bw_j - w_{j-1} \text{ for } j \geq 1, \quad (3.8)$$

$$\frac{\partial w_0}{\partial \theta_i} = 0, \quad \frac{\partial w_1}{\partial \theta_i} = \frac{\partial B}{\partial \theta_i} z, \quad \frac{\partial w_{j+1}}{\partial \theta_i} = 2 \left(\frac{\partial B}{\partial \theta_i} w_j + B \frac{\partial w_j}{\partial \theta_i} \right) - \frac{\partial w_{j-1}}{\partial \theta_i} \text{ for } j \geq 1. \quad (3.9)$$

This gives the estimators

$$\log |\tilde{K}| \approx \mathbb{E} \left[\sum_{j=0}^m c_j z^T w_j \right] \quad \text{and} \quad \frac{\partial}{\partial \theta_i} \log |\tilde{K}| \approx \mathbb{E} \left[\sum_{j=0}^m c_j z^T \frac{\partial w_j}{\partial \theta_i} \right]. \quad (3.10)$$

Thus, each derivative of the approximation costs two extra MVMs per term.

3.4.2 Gauss Quadrature via Lanczos

We can also approximate $z^T \log(\tilde{K})z$ via a Lanczos decomposition; see [65] for discussion of a Lanczos-based computation of $z^T f(\tilde{K})z$ and [164, 14] for stochastic Lanczos estimation of log determinants. We run m steps of the Lanczos algorithm, which computes the decomposition

$$\tilde{K}Q_m = Q_m T + \beta_m q_{m+1} e_m^T, \quad (3.11)$$

where $Q_m = [q_1, q_2, \dots, q_m] \in \mathbb{R}^{n \times m}$ is a matrix with orthonormal columns such that $q_1 = z / \|z\|$, $T \in \mathbb{R}^{m \times m}$ is tridiagonal, β_m is the residual, and e_m is the m -th Cartesian unit vector. We estimate

$$z^T f(\tilde{K})z \approx e_1^T f(\|z\|^2 T) e_1, \quad (3.12)$$

where e_1 is the first column of the identity. Because the Lanczos algorithm is numerically unstable, there exist several practical implementations to resolve this issue [41, 152]. The approximation (3.12) corresponds to a Gauss quadrature rule for the Riemann-Stieltjes integral of the measure associated with the eigenvalue distribution of \tilde{K} . It is exact when f is a polynomial of degree up to $2m - 1$. This approximation is also exact when \tilde{K} has at most m distinct eigenvalues, which is particularly relevant to Gaussian process regression, since frequently the kernel matrices only have a small number of eigenvalues that are not close to zero.

The Lanczos decomposition also allows us to estimate derivatives of the log deter-

minant at minimal cost. Via the Lanczos decomposition, we have

$$\hat{g} = \mathcal{Q}_m(T^{-1}e_1 \|z\|) \approx \tilde{K}^{-1}z. \quad (3.13)$$

This approximation requires no additional matrix vector multiplications beyond those used to compute the Lanczos decomposition, which we already used to estimate $\log(\tilde{K})z$; in exact arithmetic, this is equivalent to m steps of conjugate gradient (CG). Computing \hat{g} in this way takes $O(mn)$ additional time; subsequently, we only need one matrix-vector multiply by $\partial\tilde{K}/\partial\theta_i$ for each probe vector to estimate $\text{tr}(\tilde{K}^{-1}(\partial\tilde{K}/\partial\theta_i)) = \mathbb{E}[(\tilde{K}^{-1}z)^T(\partial\tilde{K}/\partial\theta_i)z]$.

3.4.3 Diagonal Correction to SKI

The SKI approximation may provide a poor estimate of the diagonal entries of the original kernel matrix for kernels with limited smoothness, such as the Matérn kernel. In general, diagonal corrections to scalable kernel approximations can lead to great performance gains. Indeed, the popular FITC method [157] is exactly a diagonal correction of SoR.

We thus modify the SKI approximation to add a diagonal matrix D ,

$$K_{XX} \approx WK_{UU}W^T + D, \quad (3.14)$$

such that the diagonal of the approximated K_{XX} is exact. In other words, D subtracts the diagonal of $WK_{UU}W^T$ and adds the true diagonal of K_{XX} . This modification is not possible for the scaled eigenvalue method for approximating log determinants in [176], since adding a diagonal matrix makes it impossible to approximate the eigenvalues of K_{XX} from the eigenvalues of K_{UU} .

However, Eq. 3.14 still admits fast MVMs and thus works with our approach for estimating the log determinant and its derivatives. Computing D with SKI costs only $O(n)$ flops since W is sparse for local cubic interpolation. We can therefore compute $(W^T e_i)^T K_{UU} (W^T e_i)$ in $O(1)$ flops.

3.4.4 Estimating higher derivatives

We have already described how to use stochastic estimators to compute the log marginal likelihood and its first derivatives. The same approach applies to computing higher-order derivatives for a Newton-like iteration, to understand the sensitivity of the maximum likelihood parameters, or for similar tasks. The first derivatives of the full log marginal likelihood are

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = -\frac{1}{2} \left[\text{tr} \left(\tilde{K}^{-1} \frac{\partial \tilde{K}}{\partial \theta_i} \right) - \alpha^T \frac{\partial \tilde{K}}{\partial \theta_i} \alpha \right], \quad (3.15)$$

and the second derivatives of the two terms are

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} [\log |\tilde{K}|] = \text{tr} \left(\tilde{K}^{-1} \frac{\partial^2 \tilde{K}}{\partial \theta_i \partial \theta_j} - \tilde{K}^{-1} \frac{\partial \tilde{K}}{\partial \theta_i} \tilde{K}^{-1} \frac{\partial \tilde{K}}{\partial \theta_j} \right), \quad (3.16)$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} [(y - \mu_X)^T \alpha] = 2\alpha^T \frac{\partial \tilde{K}}{\partial \theta_i} \tilde{K}^{-1} \frac{\partial \tilde{K}}{\partial \theta_j} \alpha - \alpha^T \frac{\partial^2 \tilde{K}}{\partial \theta_i \partial \theta_j} \alpha. \quad (3.17)$$

Superficially, evaluating the second derivatives would appear to require several additional solves above and beyond those used to estimate the first derivatives of the log determinant. In fact, we can get an unbiased estimator for the second derivatives with no additional solves, but only fast products with the derivatives of the kernel matrices.

Let z and w be independent probe vectors, and define $g = \tilde{K}^{-1} z$ and $h = \tilde{K}^{-1} w$. Then

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} [\log |\tilde{K}|] = \mathbb{E} \left[g^T \frac{\partial^2 \tilde{K}}{\partial \theta_i \partial \theta_j} z - \left(g^T \frac{\partial \tilde{K}}{\partial \theta_i} w \right) \left(h^T \frac{\partial \tilde{K}}{\partial \theta_j} z \right) \right], \quad (3.18)$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} [(y - \mu_X)^T \alpha] = 2 \mathbb{E} \left[\left(z^T \frac{\partial \tilde{K}}{\partial \theta_i} \alpha \right) \left(g^T \frac{\partial \tilde{K}}{\partial \theta_j} \alpha \right) \right] - \alpha^T \frac{\partial^2 \tilde{K}}{\partial \theta_i \partial \theta_j} \alpha. \quad (3.19)$$

Hence, if we use the stochastic Lanczos method to compute the log determinant and its derivatives, the additional work required to obtain a second derivative estimate is one MVM by each second partial of the kernel for each probe vector and for α , one MVM of each first partial of the kernel with α , and a few dot products.

3.4.5 Radial Basis Functions

Another way to deal with the log determinant and its derivatives is to evaluate the log determinant term at a few systematically chosen points in the space of hyperparameters and fit an interpolation approximation to these values. This is particularly useful when the kernel depends on a modest number of hyperparameters (e.g., half a dozen), and thus the number of points we need to precompute is relatively small. We refer to this method as a surrogate, since it provides an inexpensive substitute for the log determinant and its derivatives. For our surrogate approach, we use radial basis function (RBF) interpolation with a cubic kernel and a linear tail. See e.g. [30, 54, 153, 170] and the supplementary material for more details on RBF interpolation.

3.5 Error Properties

In addition to the usual errors from sources such as solver termination criteria and floating point arithmetic, our approach to kernel learning involves several additional sources of error: we approximate the true kernel with one that enables fast MVMs, we approximate traces using stochastic estimation, and we approximate the actions of $\log(\tilde{K})$ and \tilde{K}^{-1} on probe vectors.

We can compute first-order estimates of the sensitivity of the log likelihood to per-

turbations in the kernel using the same stochastic estimators we use for the derivatives with respect to hyperparameters. For example, if \mathcal{L}^{ref} is the likelihood for a reference kernel $\tilde{K}^{\text{ref}} = \tilde{K} + E$, then

$$\mathcal{L}^{\text{ref}}(\theta|y) = \mathcal{L}(\theta|y) - \frac{1}{2} \left(\mathbb{E} \left[g^T E z \right] - \alpha^T E \alpha \right) + O(\|E\|^2), \quad (3.20)$$

and we can bound the change in likelihood at first order by $\|E\| (\|g\| \|z\| + \|\alpha\|^2)$. Given bounds on the norms of $\partial E / \partial \theta_i$, we can similarly estimate changes in the gradient of the likelihood, allowing us to bound how the marginal likelihood hyperparameter estimates depend on kernel approximations.

If $\tilde{K} = U \Lambda U^T + \sigma^2 I$, the Hutchinson trace estimator has known variance [12]

$$\text{Var} \left[z^T \log(\tilde{K}) z \right] = \sum_{i \neq j} \left[\log(\tilde{K}) \right]_{ij}^2 \leq \sum_{i=1}^n \log(1 + \lambda_j / \sigma^2)^2. \quad (3.21)$$

If the eigenvalues of the kernel matrix without noise decay rapidly enough compared to σ , the variance will be small compared to the magnitude of $\text{tr}(\log \tilde{K}) = 2n \log \sigma + \sum_{i=1}^n \log(1 + \lambda_j / \sigma^2)$. Hence, we need fewer probe vectors to obtain reasonable accuracy than one would expect from bounds that are blind to the matrix structure. In our experiments, we typically only use 5–10 probes — and we use the sample variance across these probes to estimate *a posteriori* the stochastic component of the error in the log likelihood computation. If we are willing to estimate the Hessian of the log likelihood, we can increase rates of convergence for finding kernel hyperparameters.

The Chebyshev approximation scheme requires $O(\sqrt{\kappa} \log(\kappa/\epsilon))$ steps to obtain an $O(\epsilon)$ approximation error in computing $z^T \log(\tilde{K}) z$, where $\kappa = \lambda_{\max} / \lambda_{\min}$ is the condition number of \tilde{K} [72]. This behavior is independent of the distribution of eigenvalues within the interval $[\lambda_{\min}, \lambda_{\max}]$, and is close to optimal when eigenvalues are spread quasi-uniformly across the interval. Nonetheless, when the condition number is large, convergence may be quite slow. The Lanczos approach converges at least twice as fast

as Chebyshev in general [164, Remark 1], and converges much more rapidly when the eigenvalues are *not* uniform within the interval, as is the case with log determinants of many kernel matrices. Hence, we recommend the Lanczos approach over the Chebyshev approach in general. In all of our experiments, the error associated with approximating $z^T \log(\tilde{K})z$ by Lanczos was dominated by other sources of error.

3.6 Experiments

We test our stochastic trace estimator with both Chebyshev and Lanczos approximation schemes on: (1) a sound time series with missing data, using a GP with an RBF kernel; (2) a three-dimensional space-time precipitation data set with over half a million training points, using a GP with an RBF kernel; (3) a two-dimensional tree growth data set using a log-Gaussian Cox process model with an RBF kernel; (4) a three-dimensional space-time crime datasets with a log-Gaussian Cox model with Matérn 3/2 and spectral mixture kernels; and (5) a high-dimensional feature space using the deep kernel learning framework [178]. In the supplementary material we also include several additional experiments to illustrate particular aspects of our approach, including kernel hyperparameter recovery, diagonal corrections, and surrogate methods. Throughout we use the SKI method [176] of Eq. 3.3 for fast MVMs. We find that the Lanczos and surrogate methods are able to do kernel recovery and inference significantly faster and more accurately than competing methods.

3.6.1 Natural sound modeling

Here we consider the natural sound benchmark in [176], shown in Figure 3.1a. Our goal is to recover contiguous missing regions in a waveform with $n = 59,306$ training points. We exploit Toeplitz structure in the K_{UU} matrix of our SKI approximate kernel for accelerated MVMs.

The experiment in [176] only considered scalable inference and prediction, but not hyperparameter learning, since the scaled eigenvalue approach requires all the eigenvalues for an $m \times m$ Toeplitz matrix, which can be computationally prohibitive with cost $O(m^2)$. However, evaluating the marginal likelihood on this training set is not an obstacle for Lanczos and Chebyshev since we can use fast MVMs with the SKI approximation at a cost of $O(n + m \log m)$.

In Figure 3.1b, we show how Lanczos, Chebyshev and surrogate approaches scale with the number of inducing points m compared to the scaled eigenvalue method and FITC. We use 5 probe vectors and 25 iterations for Lanczos, both when building the surrogate and for hyperparameter learning with Lanczos. We also use 5 probe vectors for Chebyshev and 100 moments. Figure 3.1b shows the runtime of the hyperparameter learning phase for different numbers of inducing points m , where Lanczos and the surrogate are clearly more efficient than scaled eigenvalues and Chebyshev. For hyperparameter learning, FITC took several hours to run, compared to minutes for the alternatives; we therefore exclude FITC from Figure 3.1b. Figure 3.1c shows the time to do inference on the 691 test points, while 3.1d shows the standardized mean absolute error (SMAE) on the same test points. As expected, Lanczos and surrogate make accurate predictions much faster than Chebyshev, scaled eigenvalues, and FITC. In short, Lanczos and the surrogate approach are much faster than alternatives for hyperparameter learning with a large number of inducing points and training points.

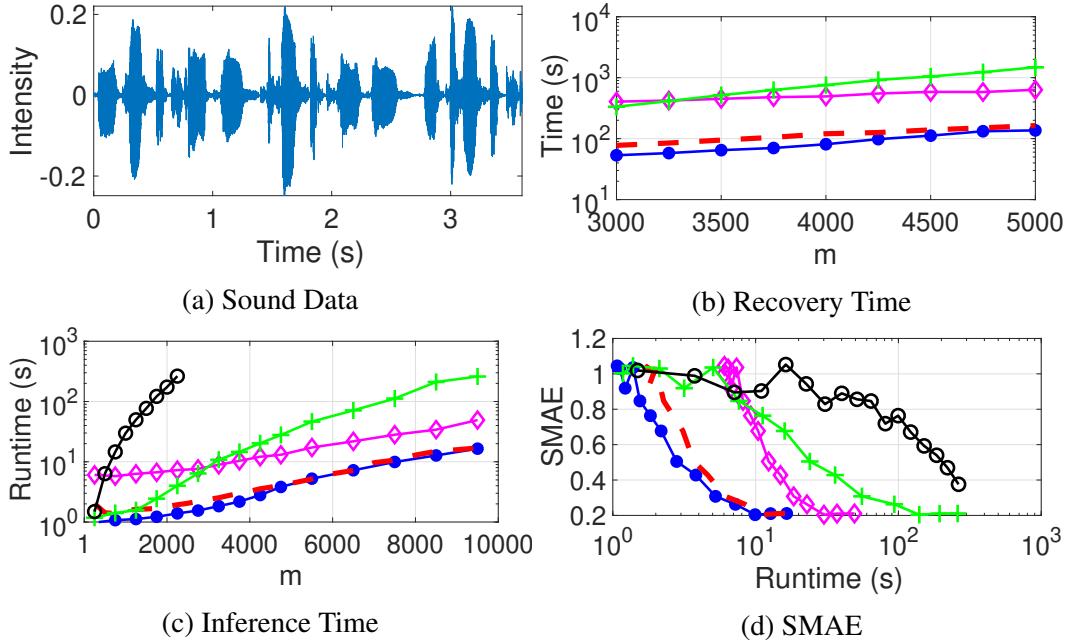


Figure 3.1: Sound modeling using 59,306 training points and 691 test points. The intensity of the time series can be seen in (a). Train time for RBF kernel hyperparameters is in (b) and the time for inference is in (c). The standardized mean absolute error (SMAE) as a function of time for an evaluation of the marginal likelihood and all derivatives is shown in (d). Surrogate is (—), Lanczos is (- - -), Chebyshev is (— ◊ —), scaled eigenvalues is (— + —), and FITC is (— o —).

3.6.2 Daily precipitation prediction

This experiment involves precipitation data from the year of 2010 collected from around 5500 weather stations in the US¹. The hourly precipitation data is preprocessed into daily data if full information of the day is available. The dataset has 628,474 entries in terms of precipitation per day given the date, longitude and latitude. We randomly select 100,000 data points as test points and use the remaining points for training. We then perform hyperparameter learning and prediction with the RBF kernel, using Lanczos, scaled eigenvalues, and exact methods.

¹<https://catalog.data.gov/dataset/u-s-hourly-precipitation-data>

For Lanczos and scaled eigenvalues, we optimize the hyperparameters on the subset of data for January 2010, with an induced grid of 100 points per spatial dimension and 300 in the temporal dimension. Due to memory constraints we only use a subset of 12,000 entries for training with the exact method. While scaled eigenvalues can perform well when fast eigendecompositions are possible, as in this experiment, Lanczos nonetheless still runs faster and with slightly lower mean square error (MSE).

Table 3.1: Prediction Comparison for the Daily Precipitation Data ^a.

Method	#Training Pts	#Induced Pts	MSE	Time [min]
Lanczos	528k	3M	0.613	14.3
Scaled-Eig	528k	3M	0.621	15.9
Exact	12k	-	0.903	11.8

^a The columns are the number of training points, number of induced grid points, mean squared error, and inference time.

Incidentally, we are able to use *3 million* inducing points in Lanczos and scaled eigenvalues, which is enabled by the SKI representation [176] of covariance matrices, for a very accurate approximation. This number of inducing points m is unprecedented for typical alternatives which scale as $\mathcal{O}(m^3)$.

3.6.3 Hickory data

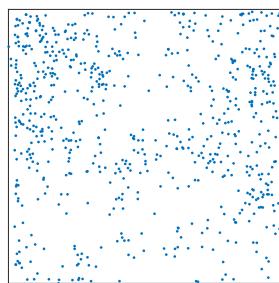
In this experiment, we apply Lanczos to the log-Gaussian Cox process model with a Laplace approximation for the posterior distribution. We use the RBF kernel and the Poisson likelihood in our model. The scaled eigenvalue method does not apply directly to non-Gaussian likelihoods; we thus applied the scaled eigenvalue method in [176] in conjunction with the Fiedler bound in [56] for the scaled eigenvalue comparison. Indeed, a key advantage of the Lanczos approach is that it can be applied whenever fast MVMs are available, which means no additional approximations such as the Fiedler bound are

required for non-Gaussian likelihoods.

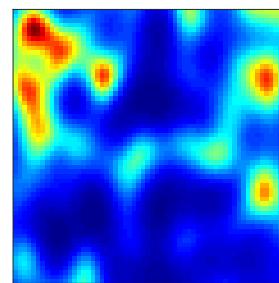
This dataset, which comes from the R package `spatstat`, is a point pattern of 703 hickory trees in a forest in Michigan. We discretize the area into a 60×60 grid and fit our model with exact, scaled eigenvalues, and Lanczos. We see in Table 3.2 that Lanczos recovers hyperparameters that are much closer to the exact values than the scaled eigenvalue approach. Figure 3.2 shows that the predictions by Lanczos are also indistinguishable from the exact computation.

Table 3.2: Hyperparameters Recovered on the Hickory Dataset.

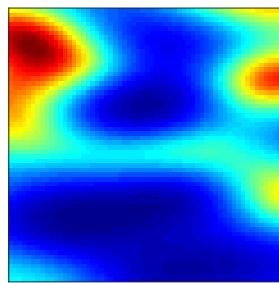
Method	s_f	ℓ_1	ℓ_2	$-\log p(y \theta)$	Time [s]
Exact	0.696	0.063	0.085	1827.56	465.9
Lanczos	0.693	0.066	0.096	1828.07	21.4
Scaled-Eig	0.543	0.237	0.112	1851.69	2.5



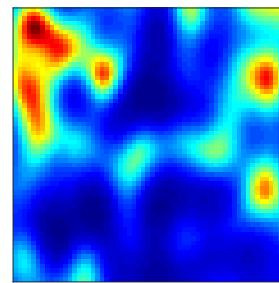
(a) Point Pattern Data



(b) Prediction by Exact



(c) Scaled-Eig



(d) Lanczos

Figure 3.2: Predictions by exact, scaled eigenvalues, and Lanczos on the Hickory dataset.

3.6.4 Crime prediction

In this experiment, we apply Lanczos with the spectral mixture kernel to the crime forecasting problem considered in [56]. This dataset consists of 233,088 incidents of assault in Chicago from January 1, 2004 to December 31, 2013. We use the first 8 years for training and attempt to predict the crime rate for the last 2 years. For the spatial dimensions, we use the log-Gaussian Cox process model, with the Matérn-5/2 kernel, the negative binomial likelihood, and the Laplace approximation for the posterior. We use a spectral mixture kernel with 20 components and an extra constant component for the temporal dimension. We discretize the data into a 17×26 spatial grid corresponding to 1 mile \times 1 mile grid cells. In the temporal dimension we sum our data by weeks for a total of 522 weeks. After removing the cells that are outside Chicago, we have a total of 157,644 observations.

The results for Lanczos and scaled eigenvalues (in conjunction with the Fiedler bound due to the non-Gaussian likelihood) can be seen in Table 3.3. The Lanczos method used 5 Hutchinson probe vectors and 30 Lanczos steps. For both methods we allow 100 iterations of LBFGS to recover hyperparameters and we often observe early convergence. While the root mean square error (RMSE) for Lanczos and scaled eigenvalues happen to be close on this example, the recovered hyperparameters using scaled eigenvalues are very different than for Lanczos. For example, the scaled eigenvalue method learns much larger σ^2 than Lanczos, indicating model misspecification. In general, as the data become increasingly non-Gaussian the Fiedler bound (used for fast scaled eigenvalues on non-Gaussian likelihoods) will become increasingly misspecified, while Lanczos will be unaffected.

Table 3.3: Hyperparameters Recovered, Recovery Time and RMSE for Lanczos and Scaled Eigenvalues on the Chicago Assault Data^a.

Method	ℓ_1	ℓ_2	σ^2	T _{recovery} [s]	T _{prediction} [s]	RMSE _{train}	RMSE _{test}
Lanczos	0.65	0.67	69.72	264	10.30	1.17	1.33
Scaled-Eig	0.32	0.10	191.17	67	3.75	1.19	1.36

^a ℓ_1 and ℓ_2 are the length scales in spatial dimensions. σ^2 is the noise level. T_{recovery} is the time for recovering hyperparameters. T_{prediction} is the time for prediction at all 157,644 observations, including training and testing.

3.6.5 Deep kernel learning

To handle high-dimensional datasets, we bring our methods into the deep kernel learning framework [178] by replacing the final layer of a pre-trained deep neural network (DNN) with a GP. This experiment uses the gas sensor dataset from the UCI machine learning repository. It has 2565 instances with 128 dimensions. We pre-train a DNN, then attach a Gaussian process with RBF kernels to the two-dimensional output of the second-to-last layer. We then further train all parameters of the resulting kernel, *including* the weights of the DNN, through the GP marginal likelihood. In this example, Lanczos and the scaled eigenvalue approach perform similarly well. Nonetheless, we see that Lanczos can effectively be used with SKI on a high dimensional problem to train hundreds of thousands of kernel parameters.

Table 3.4: Prediction RMSE and Per Training Iteration Runtime.

Method	DNN	Lanczos	Scaled-Eig
RMSE	0.1366 ± 0.0387	0.1053 ± 0.0248	0.1045 ± 0.0228
Time [s]	0.4438	2.0680	1.6320

3.7 Conclusion

There are many cases in which fast MVMs can be achieved, but it is difficult or impossible to efficiently compute a log determinant. We have developed a framework for scalable and accurate estimates of a log determinant and its derivatives relying only on MVMs. We particularly consider scalable kernel learning, showing the promise of stochastic Lanczos estimation combined with a pre-computed surrogate model. We have shown the scalability and flexibility of our approach through experiments with kernel learning for several real-world data sets using both Gaussian and non-Gaussian likelihoods, and highly parametrized deep kernels.

Iterative MVM approaches have great promise for future exploration. We have only begun to explore their significant generality. In addition to log determinants, the methods presented here could be adapted to fast posterior sampling, diagonal estimation, matrix square roots, and many other standard operations. The proposed methods only depend on fast MVMs—and the structure necessary for fast MVMs often exists, or can be readily created. We have here made use of SKI [176] to create such structure. But other approaches, such as stochastic variational methods [75], could be used or combined with SKI for fast MVMs, as in [174]. Moreover, iterative MVM methods naturally harmonize with GPU acceleration, and are therefore likely to increase in their future applicability and popularity. Finally, one could explore the ideas presented here for scalable higher order derivatives, making use of Hessian methods for greater convergence rates.

CHAPTER 4

SCALABLE GAUSSIAN PROCESSES WITH DERIVATIVES

4.1 Abstract

Gaussian processes (GPs) with derivatives are useful in many applications, including Bayesian optimization, implicit surface reconstruction, and terrain reconstruction. Fitting a GP to function values and derivatives at n points in d dimensions requires linear solves and log determinants with an $n(d + 1) \times n(d + 1)$ positive definite matrix – leading to prohibitive $O(n^3 d^3)$ computations for standard direct methods. We propose iterative solvers using fast $O(nd)$ matrix-vector multiplications (MVMs), together with pivoted Cholesky preconditioning that cuts the iterations to convergence by several orders of magnitude, allowing for fast kernel learning and prediction. Our approaches, together with dimensionality reduction, allows us to scale Bayesian optimization with derivatives to high-dimensional problems and large evaluation budgets.

4.2 Introduction

Gaussian processes (GPs) provide a powerful learning framework where the *marginal likelihood* is the probability of data given only the kernel hyper-parameters. The marginal likelihood automatically calibrates the model fit and complexity terms to favor the simplest models that explain the data [146, 145]. Computing the model fit term requires solving linear systems with the kernel matrix while the complexity term, or *Oc-cam's factor* [117], is the log determinant of the kernel matrix. The exact kernel learning costs of $O(n^3)$ flops and the prediction cost of $O(n)$ flops per test point are clearly computationally infeasible for large datasets. The situation becomes more challenging if we

consider GPs with both function value and derivative information, in which case training and prediction become $O(n^3 d^3)$ and $O(nd)$ respectively [145, §9.4].

Derivative information is important in many applications, including Bayesian Optimization (BO) [179], implicit surface reconstruction [114], and terrain reconstruction. For many simulation models, derivatives may be computed at little extra cost via finite differences, complex step approximation, an adjoint method, or algorithmic differentiation [57]. But while many scalable approximation methods for Gaussian process regression have been proposed, scalable methods incorporating derivatives have received little attention. In this paper, we propose scalable methods for GPs with derivative information built on the *structured kernel interpolation* (SKI) framework [176], which uses local interpolation to map scattered data onto a large grid of inducing points, enabling fast MVMs using FFTs. As the uniform grids in SKI scale poorly to high-dimensional spaces, we also extend the structured kernel interpolation for products (SKIP) method, which approximates a high-dimensional product kernel as a Hadamard product of low rank Lanczos decompositions [58]. Both SKI and SKIP provide fast approximate kernel MVMs, which are a building block to solve linear systems with the kernel matrix and to approximate log determinants [49].

The specific contributions of this paper are:

- We extend SKI to incorporate derivative information, enabling $O(nd)$ complexity learning and $O(1)$ prediction per test points, relying only on fast MVM with the kernel matrix.
- We also extend SKIP, which enables scalable Gaussian process regression with derivatives in high-dimensional spaces without grids. Our approach allows for $O(nd)$ MVMs.
- We illustrate that preconditioning is critical for fast convergence of iterations for

kernel matrices with derivatives. A pivoted Cholesky preconditioner cuts the iterations to convergence by several orders of magnitude when applied to both SKI and SKIP with derivatives.

- We illustrate the scalability of our approach on several examples including implicit surface fitting of the Stanford bunny, rough terrain reconstruction, and Bayesian optimization.
- We show how our methods, together with active subspace techniques, can be used to extend Bayesian optimization to high-dimensional problems with large evaluation budgets.
- Code, experiments, and figures may be reproduced at https://github.com/ericlee0803/GP_Derivatives.

We start in Section 4.3 by introducing GPs with derivatives and kernel approximations. In Section 4.4, we extend SKI and SKIP to handle derivative information. In Section 4.5, we show representative experiments; and we conclude in Section 4.6.

4.3 Background

A Gaussian process (GP) is a collection of random variables, any finite number of which are jointly Gaussian [145]; it also defines a distribution over functions on \mathbb{R}^d , $f \sim \mathcal{GP}(\mu, k)$, where $\mu : \mathbb{R}^d \rightarrow \mathbb{R}$ is a mean field and $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a symmetric and positive (semi)-definite covariance kernel. For any set of locations $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$, $f_X \sim \mathcal{N}(\mu_X, K_{XX})$ where f_X and μ_X represent the vectors of function values for f and μ evaluated at each of the $x_i \in X$, and $(K_{XX})_{ij} = k(x_i, x_j)$. We assume the observed function value vector $y \in \mathbb{R}^n$ is contaminated by independent Gaussian noise with variance σ^2 .

We denote any kernel hyper-parameters by the vector θ . To be concise, we suppress the dependence of k and associated matrices on θ in our notation. Under a Gaussian process prior depending on the covariance hyper-parameters θ , the log marginal likelihood is given by

$$\mathcal{L}(\theta|y) = -\frac{1}{2} \left[(y - \mu_X)^T \alpha + \log |\tilde{K}_{XX}| + n \log 2\pi \right] \quad (4.1)$$

where $\alpha = \tilde{K}_{XX}^{-1}(y - \mu_X)$ and $\tilde{K}_{XX} = K_{XX} + \sigma^2 I$. The standard direct method to evaluate (4.1) and its derivatives with respect to the hyperparameters uses the Cholesky factorization of \tilde{K}_{XX} , leading to $O(n^3)$ kernel learning that does not scale beyond a few thousand points.

A popular approach to scalable GPs is to approximate the exact kernel with a structured kernel that enables fast MVMs [143]. Several methods approximate the kernel via *inducing points* $U = \{u_j\}_{j=1}^m \subset \mathbb{R}^d$; see, e.g.[143, 101, 75]. Common examples are the subset of regressors (SoR), which exploits low-rank structure, and fully independent training conditional (FITC), which introduces an additional diagonal correction [157]. For most inducing point methods, the cost of kernel learning with n data points and m inducing points scales as $O(m^2n + m^3)$, which becomes expensive as m grows. As an alternative, Wilson proposed the structured kernel interpolation (SKI) approximation,

$$K_{XX} \approx WK_{UU}W^T, \quad (4.2)$$

where U is a uniform grid of inducing points and W is an n -by- m matrix of interpolation weights; the authors of [176] use local cubic interpolation so that W is sparse. If the original kernel is stationary, each MVM with the SKI kernel may be computed in $O(n + m \log(m))$ time via FFTs, leading to substantial performance over FITC and SoR. A limitation of SKI is that the number of grid points increases exponentially with the dimension. This exponential scaling has been addressed by structured kernel interpolation for products (SKIP) [58], which decomposes the kernel matrix for a product

kernel in d -dimensions as a Hadamard (elementwise) product of one-dimensional kernel matrices.

We use fast MVMs to solve linear systems involving \tilde{K} by the method of conjugate gradients. To estimate $\log |\tilde{K}| = \text{tr}(\log(\tilde{K}))$, we apply stochastic trace estimators that require only products of $\log(\tilde{K})$ with random probe vectors. Given a probe vector z , several ideas have been explored to compute $\log(\tilde{K})z$ via MVMs with \tilde{K} , such as using a polynomial approximation of \log or using the connection between the Gaussian quadrature rule and the Lanczos method [72, 164]. It was shown in [49] that using Lanczos is superior to the polynomial approximations and that only a few probe vectors are necessary even for large kernel matrices.

Differentiation is a linear operator, and (assuming a twice-differentiable kernel) we may define a multi-output GP for the function and (scaled) gradient values with mean and kernel functions

$$\mu^\nabla(x) = \begin{bmatrix} \mu(x) \\ \partial_x \mu(x) \end{bmatrix}, \quad k^\nabla(x, x') = \begin{bmatrix} k(x, x') & (\partial_{x'} k(x, x'))^T \\ \partial_x k(x, x') & \partial^2 k(x, x') \end{bmatrix}, \quad (4.3)$$

where $\partial_x k(x, x')$ and $\partial^2 k(x, x')$ represent the column vector of (scaled) partial derivatives in x and the matrix of (scaled) second partials in x and x' , respectively. Scaling derivatives by a natural length scale gives the multi-output GP consistent units, and lets us understand approximation error without weighted norms. As in the scalar GP case, we model measurements of the function as contaminated by independent Gaussian noise.

Because the kernel matrix for the GP on function values alone is a submatrix of the kernel matrix for function values and derivatives together, the predictive variance in the presence of derivative information will be strictly less than the predictive variance without derivatives. Hence, convergence of regression with derivatives is always super-

prior to convergence of regression without, which is well-studied in, e.g. [145, Chapter 7]. Fig. 4.1 illustrates the value of derivative information; fitting with derivatives is evidently much more accurate than fitting function values alone. In higher-dimensional problems, derivative information is even more valuable, but it comes at a cost: the kernel matrix K_{XX}^∇ is of size $n(d + 1)$ -by- $n(d + 1)$. Scalable approximate solvers are therefore vital in order to use GPs for large datasets with derivative data, particularly in high-dimensional spaces.

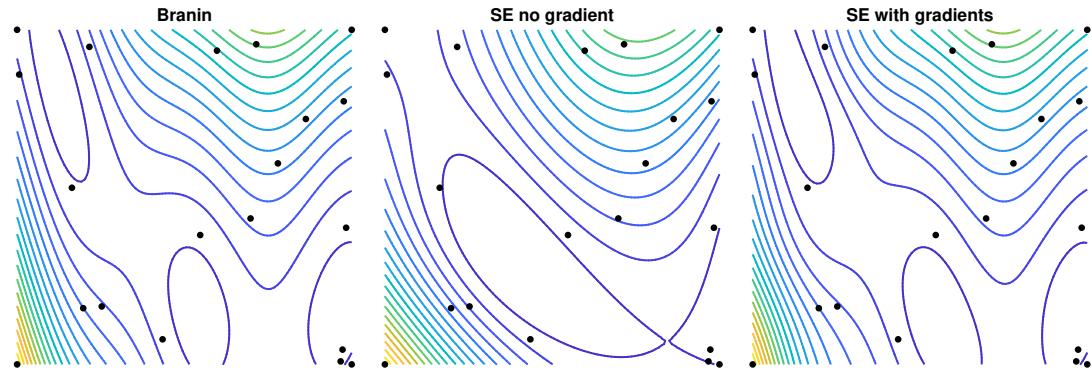


Figure 4.1: An example where gradient information pays off; the true function is on the left. Compare the regular GP without derivatives (middle) to the GP with derivatives (right). Unlike the former, the latter is able to accurately capture critical points of the function.

4.4 Methods

One standard approach to scaling GPs substitutes the exact kernel with an approximate kernel. When the GP fits values and gradients, one may attempt to separately approximate the kernel and the kernel derivatives. Unfortunately, this may lead to indefiniteness, as the resulting approximation is no longer a valid kernel. Instead, we differentiate the approximate kernel, which preserves positive definiteness. We do this for the SKI and SKIP kernels below, but our general approach applies to any approximate MVM.

4.4.1 D-SKI

D-SKI (SKI with derivatives) is the standard kernel matrix for GPs with derivatives, but applied to the SKI kernel. Equivalently, we differentiate the interpolation scheme:

$$k(x, x') \approx \sum_i w_i(x) k(x_i, x') \rightarrow \nabla k(x, x') \approx \sum_i \nabla w_i(x) k(x_i, x'). \quad (4.4)$$

One can use cubic convolutional interpolation [95], but higher order methods lead to greater accuracy, and we therefore use quintic interpolation [125]. The resulting D-SKI kernel matrix has the form

$$\begin{bmatrix} K & (\partial K)^T \\ \partial K & \partial^2 K \end{bmatrix} \approx \begin{bmatrix} W \\ \partial W \end{bmatrix} K_{UU} \begin{bmatrix} W \\ \partial W \end{bmatrix}^T = \begin{bmatrix} WK_{UU}W^T & WK_{UU}(\partial W)^T \\ (\partial W)K_{UU}W^T & (\partial W)K_{UU}(\partial W)^T \end{bmatrix}, \quad (4.5)$$

where the elements of sparse matrices W and ∂W are determined by $w_i(x)$ and $\nabla w_i(x)$ — assuming quintic interpolation, W and ∂W will each have 6^d elements per row. As with SKI, we use FFTs to obtain $O(m \log m)$ MVMs with K_{UU} . Because W and ∂W have $O(n6^d)$ and $O(nd6^d)$ nonzero elements, respectively, our MVM complexity is $O(nd6^d + m \log m)$.

4.4.2 D-SKIP

Several common kernels are *separable*, i.e., they can be expressed as products of one-dimensional kernels. Assuming a compatible approximation scheme, this structure is inherited by the SKI approximation for the kernel matrix without derivatives,

$$K \approx (W_1 K_1 W_1^T) \odot (W_2 K_2 W_2^T) \odot \dots \odot (W_d K_d W_d^T), \quad (4.6)$$

where $A \odot B$ denotes the Hadamard product of matrices A and B with the same dimensions, and W_j and K_j denote the SKI interpolation and inducing point grid matrices in

the j -th coordinate direction. The same Hadamard product structure applies to the kernel matrix with derivatives; for example, for $d = 2$,

$$K^\nabla \approx \begin{bmatrix} W_1 K_1 W_1^T & W_1 K_1 \partial W_1^T & W_1 K_1 W_1^T \\ \partial W_1 K_1 W_1^T & \partial W_1 K_1 \partial W_1^T & \partial W_1 K_1 W_1^T \\ W_1 K_1 W_1^T & W_1 K_1 \partial W_1^T & W_1 K_1 W_1^T \end{bmatrix} \odot \begin{bmatrix} W_2 K_2 W_2^T & W_2 K_2 W_2^T & W_2 K_2 \partial W_2^T \\ W_2 K_2 W_2^T & W_2 K_2 W_2^T & W_2 K_2 \partial W_2^T \\ \partial W_2 K_2 W_2^T & \partial W_2 K_2 W_2^T & \partial W_2 K_2 \partial W_2^T \end{bmatrix}. \quad (4.7)$$

Eq. (4.7) expresses K^∇ as a Hadamard product of one dimensional kernel matrices. Following this approximation, we apply the SKIP reduction [58] and use Lanczos to further approximate Eq. (4.7) as $(Q_1 T_1 Q_1^T) \odot (Q_2 T_2 Q_2^T)$. This can be used for fast MVMs with the kernel matrix. Applied to kernel matrices with derivatives, we call this approach D-SKIP.

Constructing the D-SKIP kernel costs $O(d^2(n + m \log m + r^3 n \log d))$, and each MVM costs $O(dr^2 n)$ flops where r is the effective rank of the kernel at each step (rank of the Lanczos decomposition). We achieve high accuracy with $r \ll n$.

4.4.3 Preconditioning

Recent work [42] has explored several preconditioners for exact kernel matrices. We have had success with preconditioners of the form $M = \sigma^2 I + FF^T$ where $K^\nabla \approx FF^T$ with $F \in \mathbb{R}^{n \times m}$. Solving with the Sherman-Morrison-Woodbury formula (*a.k.a* the matrix inversion lemma) is inaccurate for small σ ; we use the more stable formula $M^{-1}b = \sigma^{-2}(f - Q_1(Q_1^T f))$ where Q_1 is computed in $O(m^2 n)$ time by the economy QR factorization

$$\begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix} = \begin{bmatrix} F \\ \sigma I \end{bmatrix} R. \quad (4.8)$$

In our experiments with solvers for D-SKI and D-SKIP, we have found that a truncated pivoted Cholesky factorization, $K^\nabla \approx (\Pi L)(\Pi L)^T$ works well for the low-rank factoriza-

tion. Computing the pivoted Cholesky factorization is cheaper than MVM-based preconditioners such as Lanczos or truncated eigendecompositions as it only requires the diagonal and the ability to form the rows where pivots are selected. Pivoted Cholesky is a natural choice when inducing point methods are applied as the pivoting can itself be viewed as an inducing point method where the most important information is selected to construct a low-rank preconditioner [73]. The D-SKI diagonal can be formed in $O(nd6^d)$ flops while rows cost $O(nd6^d + m)$ flops; for D-SKIP both the diagonal and the rows can be formed in $O(nd)$ flops.

4.4.4 Dimensionality reduction

In many high-dimensional function approximation problems, only a few directions are relevant. That is, if $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is a function to be approximated, there is often a matrix P with $d < D$ orthonormal columns spanning an *active subspace* of \mathbb{R}^D such that $f(x) \approx f(PP^T x)$ for all x in some domain Ω of interest [40]. The optimal subspace is given by the dominant eigenvectors of the covariance matrix $C = \int_{\Omega} \nabla f(x) \nabla f(x)^T dx$, generally estimated by Monte Carlo integration. Once the subspace is determined, the function can be approximated through a GP on the reduced space, i.e., we replace the original kernel $k(x, x')$ with a new kernel $\check{k}(x, x') = k(P^T x, P^T x')$. Because we assume gradient information, dimensionality reduction based on active subspaces is a natural pre-processing phase before applying D-SKI and D-SKIP.

4.5 Experiments

Our experiments use the squared exponential (SE) kernel, which has product structure and can be used with D-SKIP; and the spline kernel, to which D-SKIP does not directly apply. We use these kernels in tandem with D-SKI and D-SKIP to achieve the fast MVMs derived in Section 4.4. We write D-SE to denote the exact SE kernel with derivatives.

4.5.1 Approximation Benchmark

D-SKI and D-SKIP with the SE kernel approximate the original kernel well, both in terms of MVM accuracy and spectral profile. Comparing D-SKI and D-SKIP to their exact counterparts in Figure 4.2, we see their matrix entries are very close (leading to MVM accuracy near 10^{-5}), and their spectral profiles are indistinguishable. The same is true with the spline kernel.

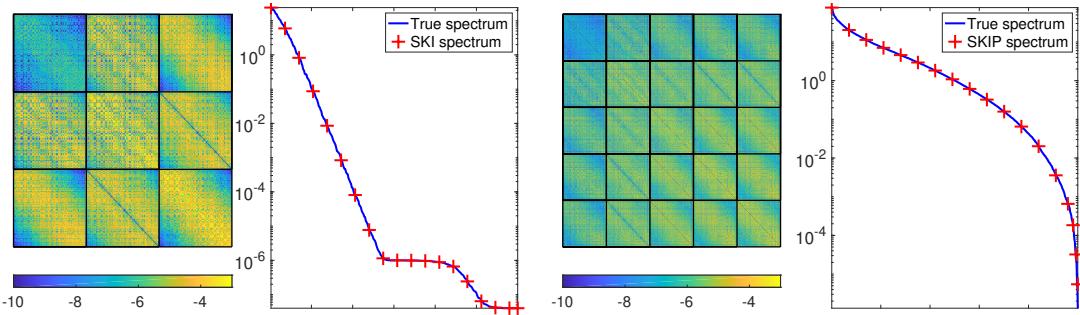


Figure 4.2: (Left two images) \log_{10} error in SKI approximation and comparison to the exact spectrum. (Right two images) \log_{10} error in SKIP approximation and comparison to the exact spectrum.

Additionally, scaling tests in Fig. 4.3 verify the predicted complexity of D-SKI and D-SKIP. We show the relative fitting accuracy of SE, SKI, D-SE, and D-SKI on some

standard test functions in Table 4.1.

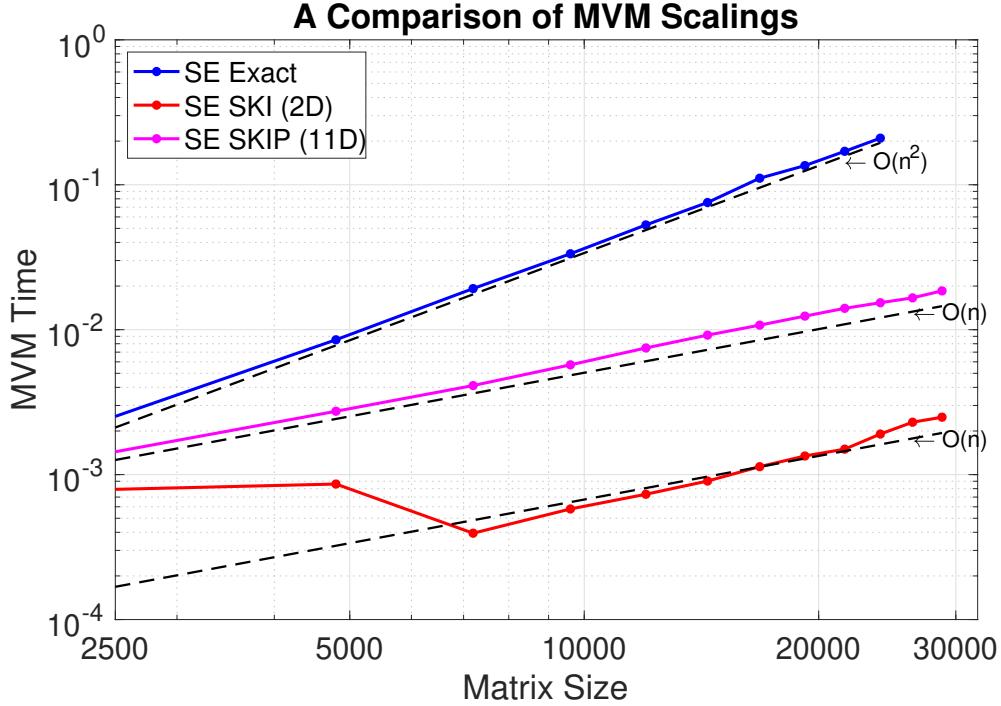


Figure 4.3: Scaling tests for D-SKI in two dimensions and D-SKIP in 11 dimensions. D-SKIP uses fewer data points for identical matrix sizes.

Table 4.1: Relative RMSE on Test Functions Using SKI and Derivatives^a.

	Branin	Franke	Sine Norm	Sixhump	StyTang	Hart3
SE	6.02e-3	8.73e-3	8.64e-3	6.44e-3	4.49e-3	1.30e-2
SKI	3.97e-3	5.51e-3	5.37e-3	5.11e-3	2.25e-3	8.59e-3
D-SE	1.83e-3	1.59e-3	3.33e-3	1.05e-3	1.00e-3	3.17e-3
D-SKI	1.03e-3	4.06e-4	1.32e-3	5.66e-4	5.22e-4	1.67e-3

^a Relative RMSE error measured on 10000 testing points. Test functions from [159] includes five 2D functions (Branin, Franke, Sine Norm, Six-hump, and Styblinski-Tang) and one 3D function (Hartman). We train the SE kernel on 4000 points, the D-SE kernel on $4000/(d + 1)$ points, and SKI and D-SKI with SE kernel on 10000 points to achieve comparable runtimes between methods.

4.5.2 Dimensionality reduction

We apply active subspace pre-processing to the 20 dimensional Welsh test function in [24]. The top six eigenvalues of its gradient covariance matrix are well separated from the rest as seen in Fig. 4.4a. However, the function is far from smooth when projected onto the leading 1D or 2D active subspace, as Figs. 4.4b to 4.4d indicates, where the color shows the function value.

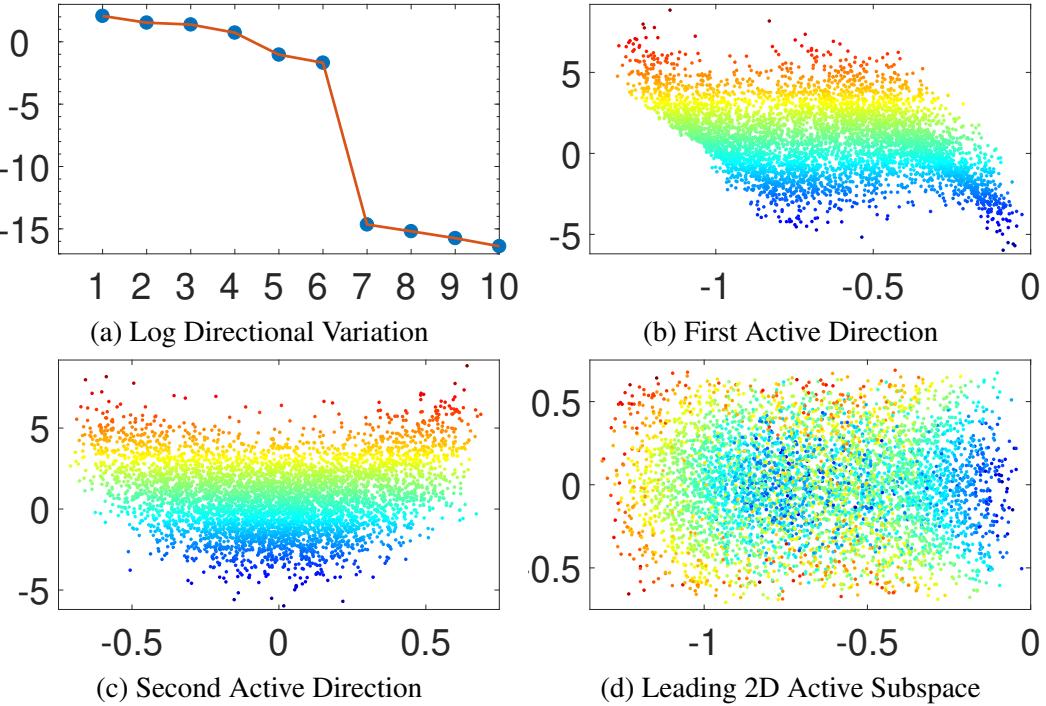


Figure 4.4: Fig. 4.4a shows the top 10 eigenvalues of the gradient covariance. Welsh is projected onto the first and second active direction in 4.4b and 4.4c. After joining them together, we see in 4.4d that points of different color are highly mixed, indicating a very spiky surface.

We therefore apply D-SKI and D-SKIP on the 3D and 6D active subspace, respectively, using 5000 training points, and compare the prediction error against D-SE with 190 training points because of our scaling advantage. Table 4.2 reveals that while the 3D active subspace fails to capture all the variation of the function, the 6D active subspace

is able to do so. These properties are demonstrated by the poor prediction of D-SKI in 3D and the excellent prediction of D-SKIP in 6D.

Table 4.2: Relative RMSE and SMAE prediction error for Welsh ^a.

	D-SE	D-SKI (3D)	D-SKIP (6D)
RMSE	4.900e-02	2.267e-01	3.366e-03
SMAE	4.624e-02	2.073e-01	2.590e-03

^a The D-SE kernel is trained on $4000/(d + 1)$ points, with D-SKI and D-SKIP trained on 5000 points. The 6D active subspace is sufficient to capture the variation of the test function.

4.5.3 Rough terrain reconstruction

Rough terrain reconstruction is a key application in robotics [62, 97], autonomous navigation [70], and geostatistics. Through a set of terrain measurements, the problem is to predict the underlying topography of some region. In the following experiment, we consider roughly 23 million non-uniformly sampled elevation measurements of Mount St. Helens obtained via LiDAR [39]. We bin the measurements into a 970×950 grid, and downsample to a 120×117 grid. Derivatives are approximated using a finite difference scheme.

We randomly select 90% of the grid for training and the remainder for testing. We do not include results for D-SE, as its kernel matrix has dimension roughly 4×10^4 . We plot contour maps predicted by SKI and D-SKI in Fig. 4.5 — the latter looks far closer to the ground truth than the former. This is quantified in the following table:

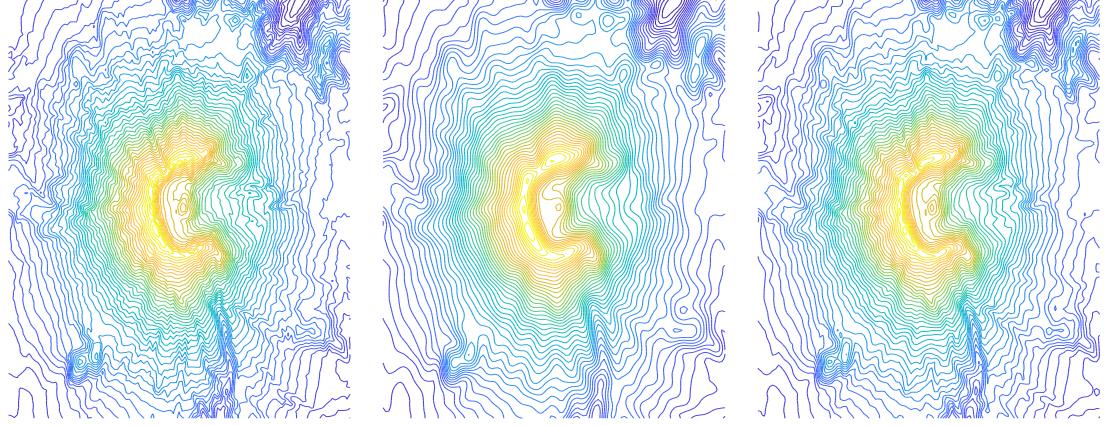


Figure 4.5: On the left is the true elevation map of Mount St. Helens. In the middle is the elevation map calculated with the SKI. On the right is the elevation map calculated with D-SKI.

Table 4.3: Hyperparameters Recovered, Recovery SMAE, and Recovery Time for SKI and D-SKI on Mountain St. Helens Data^a.

	ℓ	s	σ_1	σ_2	SMAE _{test}	SMAE _{all}	Time[s]
SKI	35.196	207.689	12.865	n.a.	0.0308	0.0357	37.67
D-SKI	12.630	317.825	6.446	2.799	0.0165	0.0254	131.70

^a σ_1 and σ_2 are the noise parameters for value and gradient, respectively.

4.5.4 Implicit surface reconstruction

Reconstructing surfaces from point cloud data and surface normals is a standard problem in computer vision and graphics. One popular approach is to fit an implicit function that is zero on the surface with gradients equal to the surface normal. Local Hermite RBF interpolation has been considered in prior work [114], but this approach is sensitive to noise. In our experiments, using a GP instead of splining reproduces implicit surfaces with very high accuracy. In this case, a GP with derivative information is required, as the function values are all zero.

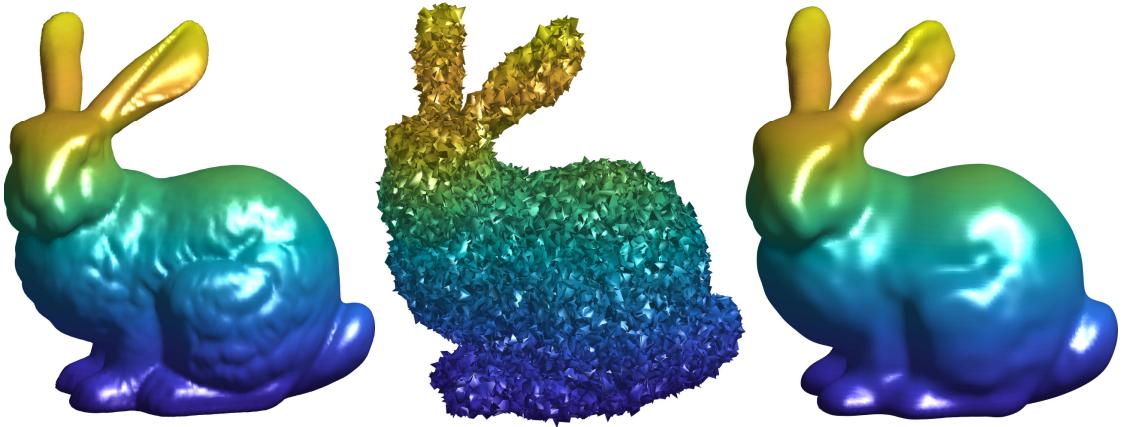


Figure 4.6: (Left) Original surface (Middle) Noisy surface (Right) SKI reconstruction from noisy surface ($s = 0.4$, $\sigma = 0.12$).

In Figure 4.6, we fit the Stanford bunny using 25000 points and associated normals, leading to a K^∇ matrix of dimension 10^5 , clearly far too large for exact training. We therefore use SKI with the thin-plate spline kernel, with a total of 30 grid points in each dimension. The left image is a ground truth mesh of the underlying point cloud and normals. The middle image shows the same mesh, but with heavily noised points and normals. Using this noisy data, we fit a GP and reconstruct a surface shown in the right image, which looks very close to the original.

4.5.5 Bayesian optimization with derivatives

Prior work examines Bayesian optimization (BO) with derivative information in low-dimensional spaces to optimize model hyperparameters [179]. Wang et al. consider high-dimensional BO (without gradients) with random projections uncovering low-dimensional structure [166]. We propose BO with derivatives and dimensionality reduction via active subspaces, detailed in Algorithm 1.

Algorithm 2: BO with derivatives and active subspace learning

while *Budget not exhausted* **do**

Calculate active subspace projection $P \in \mathbb{R}^{D \times d}$ using sampled gradients
 Optimize acquisition function, $u_{n+1} = \arg \max \mathcal{A}(u)$ with $x_{n+1} = P u_{n+1}$
 Sample point x_{n+1} , value f_{n+1} , and gradient ∇f_{n+1}
 Update data $\mathcal{D}_{i+1} = \mathcal{D}_i \cup \{x_{n+1}, f_{n+1}, \nabla f_{n+1}\}$
 Update hyperparameters of GP with gradient defined by kernel $k(P^T x, P^T x')$

end

Algorithm 1 estimates the active subspace and fits a GP with derivatives in the reduced space. Kernel learning, fitting, and optimization of the acquisition function all occur in this low-dimensional subspace. In our tests, we use the expected improvement (EI) acquisition function, which involves both the mean and predictive variance. We consider two approaches to rapidly evaluate the predictive variance $v(x) = k(x, x) - K_{xx} \tilde{K}^{-1} K_{xx}$ at a test point x . In the first approach, which provides a biased estimate of the predictive variance, we replace K^{-1} with the preconditioner solve computed by pivoted Cholesky; using the stable QR-based evaluation algorithm, we have

$$v(x) \approx \hat{v}(x) \equiv k(x, x) - \sigma^{-2} (\|K_{Xx}\|^2 - \|Q_1^T K_{Xx}\|^2). \quad (4.9)$$

We note that the approximation $\hat{v}(x)$ is always a (small) overestimate of the true predictive variance $v(x)$. In the second approach, we use a randomized estimator as in [22] to compute the predictive variance at many points X' simultaneously, and use the pivoted Cholesky approximation as a control variate to reduce the estimator variance:

$$v_{X'} = \text{diag}(K_{X'X'}) - \mathbb{E}_z [z \odot (K_{X'X} \tilde{K}^{-1} K_{XX'} z - K_{X'X} M^{-1} K_{XX'} z)] - \hat{v}_{X'}. \quad (4.10)$$

The latter approach is unbiased, but gives very noisy estimates unless many probe vectors z are used. Both the pivoted Cholesky approximation to the predictive variance and

the randomized estimator resulted in similar optimizer performance in our experiments.

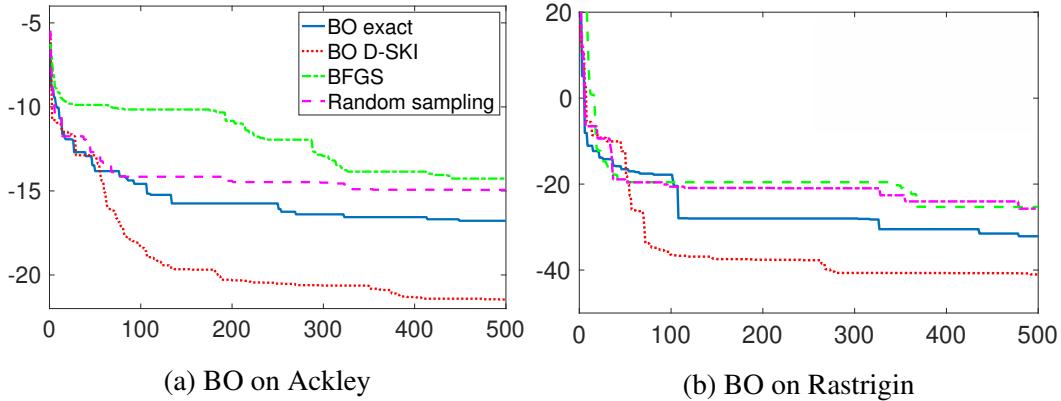


Figure 4.7: In the following experiments, 5D Ackley and 5D Rastrigin are embedded into a 50 dimensional space. We run Algorithm 1, comparing it with BO exact, multi-start BFGS, and random sampling. D-SKI with active subspace learning clearly outperforms the other methods.

To test Algorithm 1, we mimic the experimental set up in [166]: we minimize the 5D Ackley and 5D Rastrigin test functions [159], randomly embedded respectively in $[-10, 15]^{50}$ and $[-4, 5]^{50}$. We fix $d = 2$, and at each iteration pick two directions in the estimated active subspace at random to be our active subspace projection P . We use D-SKI as the kernel and EI as the acquisition function. The results of these experiments are shown in Figure 4.7a and Fig. 4.7b, in which we compare Algorithm 1 to three other baseline methods: BO with EI and no gradients in the original space; multi-start BFGS with full gradients; and random search. In both experiments, the BO variants perform better than the alternatives, and our method outperforms standard BO.

4.6 Conclusion

When gradients are available, they are a valuable source of information for Gaussian process regression; but inclusion of d extra pieces of information per point naturally

leads to new scaling issues. We introduce two methods to deal with these scaling issues: D-SKI and D-SKIP. Both are structured interpolation methods, and the latter also uses kernel product structure. We have also discussed practical details —preconditioning is necessary to guarantee convergence of iterative methods and active subspace calculation reveals low-dimensional structure when gradients are available. We present several experiments with kernel learning, dimensionality reduction, terrain reconstruction, implicit surface fitting, and scalable Bayesian optimization with gradients. For simplicity, these examples all possessed full gradient information; however, our methods trivially extend if only partial gradient information is available.

There are several possible avenues for future work. D-SKIP shows promising scalability, but it also has large overheads, and is expensive for Bayesian optimization as it must be recomputed from scratch with each new data point. We believe kernel function approximation via Chebyshev interpolation and tensor approximation will likely provide similar accuracy with greater efficiency. Extracting low-dimensional structure is highly effective in our experiments and deserves an independent, more thorough treatment. Finally, our work in scalable Bayesian optimization with gradients represents a step towards the unified view of global optimization methods (i.e. Bayesian optimization) and gradient-based local optimization methods (i.e. BFGS).

CHAPTER 5

ROBUST LARGE-VOCABULARY TOPIC MODELING

5.1 Abstract

Across many data domains, co-occurrence statistics about the joint appearance of objects are powerfully informative. By transforming unsupervised learning problems into decompositions of co-occurrence statistics, spectral inference provides transparent algorithms and optimality guarantees for non-linear dimensionality reduction or latent topic analysis. As object vocabularies grow, however, it becomes rapidly more expensive to store and run inference algorithms on co-occurrence statistics. Current rectification techniques, which preprocess real data in order to overcome model-data mismatch, are even more expensive, as they require iterative projections that destroy sparsity of the co-occurrence data. In this paper, we propose novel approaches that can simultaneously compress and rectify the co-occurrence statistics, scaling gracefully with the size of vocabulary and the dimension of latent space. We also present new algorithms that are capable of learning latent variables from the compressed statistics without losing visible precision, and verify that they perform comparably to previous approaches on both textual and non-textual data.

5.2 Introduction

Understanding the low-dimensional geometry of noisy and complex data is a fundamental problem of unsupervised learning. Probabilistic models explain data generation processes in terms of low-dimensional latent variables. Inferring a posterior distribution for

these latent variables provides us with a compact representation for various exploratory analyses and downstream tasks. However, exact inference is often intractable due to entangled interactions between the latent variables [25, 2, 1, 142]. Variational inference transforms the posterior approximation into an optimization problem over simpler distributions with independent parameters [91, 165, 26], while Markov Chain Monte Carlo enables users to sample from the desired posterior distribution [130, 131, 150]. However, these likelihood-based methods require numerous iterations without any guarantee beyond local improvement at each step [99].

When the data consists of collections of discrete objects, co-occurrence statistics summarize interactions between objects. Collaborative filtering learns low-dimensional representations of individual items, which are useful for recommendation systems, by explicitly decomposing the co-occurrence of items that are jointly consumed by certain users [102, 107]. Word-vector models learn low-dimensional embeddings of individual words, which encode useful linguistic biases for neural networks, by implicitly decomposing the co-occurrence of words that appear together [140, 106]. If co-occurrence provides a rich enough set of unbiased moments about an underlying generative model, spectral methods can provably learn posterior configurations from co-occurrence information alone, without iterating through individual training examples [8, 5, 81, 3].

However, two major limitations hinder users from taking advantage of spectral inference based on co-occurrence. First, the second-order co-occurrence matrix grows quadratically in the size of the objects (e.g. items, words, products). Pruning these objects is an option, but for a retailer selling millions of products, a low-dimensional representation of a small subset of the products is inadequate. Second, inference quality is poor in real data that does not necessarily follow our computational model. Whereas likelihood-based methods have an intrinsic capability to fit the data to the model despite

their mismatch, sample noise can destroy the performance of spectral methods even if the data is synthesized from the model [99]. *Rectification*, a process of projecting empirical co-occurrence onto a set consistent with the geometry of the model, can improve the performance of spectral inference in the face of model mismatch [102]. But running multiple projections dominates overall inference cost even when the vocabulary is small. In addition, the rectification process makes the co-occurrence dense, exacerbating storage costs when dealing with large vocabularies.

In this paper, we propose the Epsilon Non-Negative rectification (ENN) and the Low-rank Anchor Word algorithm (LAW). Given a vocabulary of N objects and the user-specified latent dimension K , ENN simultaneously compresses and rectifies the co-occurrence matrix $\mathbf{C} \in \mathbb{R}^{N \times N}$ into $\mathbf{Y}\mathbf{Y}^T$ with $\mathbf{Y} \in \mathbb{R}^{N \times K}$. Each entry of the decomposition $(\mathbf{Y}\mathbf{Y}^T)_{ij}$ tightly approximates the rectified co-occurrence \mathbf{C}_{ij}^* but is allowed to be a tiny negative value above $-\epsilon$. Then LAW learns the latent clusters (e.g., topics of documents or genres of items) and their correlations provided only with \mathbf{Y} , guaranteeing the same performance as running the original Anchor Word algorithm on \mathbf{C}^* if $\mathbf{Y}\mathbf{Y}^T \geq 0$. In contrast, we also formulate the Proximal Alternating Linearized Minimization rectification algorithm (PALM) that approximates the rectified co-occurrence \mathbf{C}^* by $\mathbf{Y}\mathbf{Y}^T$ with $\mathbf{Y} \geq 0$. While the non-negativity of \mathbf{Y} in the PALM approach allows LAW to perform exact inference, the fact that PALM enforces more constraints than ENN means that PALM also provides a less faithful approximation to the original co-occurrence data. Our experiments on various textual and non-textual datasets show that ENN learns a high-quality factor \mathbf{Y} for which LAW provides results of quality comparable to those based on the full co-occurrence \mathbf{C}^* ; in contrast, while PALM works comparably in some settings, in others there is a visible loss of accuracy.

We also adopt a randomized algorithm that constructs a low-rank approximation of

the full co-occurrence C directly from the raw data. While PALM requires the full co-occurrence, ENN can work directly with the low-rank initialization, eliminating the need to ever store a full co-occurrence matrix. Note also that the second-order methods rely on the *separability assumption*,¹ which has been another criticism in theory despite their superior performance in practice [103]. Our analyses show that models based on large vocabularies find more separable anchor objects, learning stable latent clusters without much sensitivity to sample noise. Overall, we complete a robust and scalable pipeline: that efficiently performs quality posterior inference for unsupervised learning from co-occurrence information within time and space complexity linear in N .

The major contributions of our paper are:

- We introduce two efficient rectifications, ENN and PALM, that compress quadratic and noisy co-occurrence information on the fly with a linear space rectified representation.
- We develop a low-rank algorithm (LAW) for anchor-word-based topic modeling that works directly on the compressed rectified representations and provides near-exact performance.
- We propose a robust and scalable pipeline, LR-JSMF, that learns topic models with a small number of passes directly over the data. This new pipeline scales to large vocabularies that were previously intractable for spectral inference, and offers a $\sim 100 \times$ speedup over previous methods for general data sets.

¹Each cluster has at least one anchor object which dedicates exclusively to that cluster, nothing else.

5.3 Background and Related Work

Our new algorithms build on the the **Joint-Stochastic Matrix Factorization (JSMF)** framework [102], which we now describe. Let $\mathbf{H} \in \mathbb{R}^{N \times M}$ be the object-example matrix whose m -th column vector \mathbf{h}_m counts the occurrences of each of the N objects in the vocabulary in example m . We denote the total number of objects in example m by n_m . Given a user-specified number of clusters K , we seek to learn an object-cluster matrix $\mathbf{B} \in \mathbb{R}^{N \times K}$ where B_{ik} is the conditional probability of observing object i given latent cluster k . Instead of learning \mathbf{B} directly from the sparse and noisy observations \mathbf{H} , JSMF begins with constructing the joint-stochastic co-occurrence $\mathbf{C} \in \mathbb{R}^{N \times N}$ by

$$\mathbf{C} = \hat{\mathbf{H}}\hat{\mathbf{H}}^T - \hat{\mathbf{H}}_{diag}, \quad \hat{\mathbf{h}}_m = \frac{\mathbf{h}_m}{\sqrt{n_m(n_m - 1)M}}, \quad \hat{\mathbf{H}}_{diag} = \text{diag}\left(\sum_{m=1}^M \frac{\mathbf{h}_m}{n_m(n_m - 1)M}\right). \quad (5.1)$$

Then the original Anchor Word algorithm decomposes \mathbf{C} into $\mathbf{B}\mathbf{A}\mathbf{B}^T$ by Algorithm 3, where $\mathbf{A} \in \mathbb{R}^{K \times K}$ is the cluster correlation matrix, whose entry A_{kl} captures the joint probability between two latent clusters k and l .² In the limit, using data generated according to the correct probabilistic models, \mathbf{A} must agree with the second-moment of the cluster proportions, which is given as a Bayesian prior in the models.

As with other spectral algorithms for latent variable models [81, 4], the decomposition as described so far may fail to learn high-quality clusters due to *model-data mismatch* [99]. Under the probabilistic model assumed to generate the data, the expected value of the co-occurrence should not only be normalized to sum to one (*NOR*) and be entry-wise non-negative (*NN*), but it should also be positive semi-definite with rank equal to the number of clusters K (\mathcal{PSD}_K) [102]. However, the empirical \mathbf{C} from real data is often indefinite and full-rank due to sample noise³ and the unbiased construction

² \mathbf{C} is proven to be a by far more robust estimator than \mathbf{H} in [7]. But actual construction of \mathbf{C} in [8] is slightly misleading without dividing by M . We report the full equations in this paper.

³Rectification still improves the quality of clusters on the synthetic data that is generated from our models.

Algorithm 3: Anchor Word algorithm (AW)

Input: Object co-occurrence $\mathbf{C} \in \mathbb{R}^{N \times N}$
 The number of clusters K

Output: Anchor objects $\mathbf{S} = \{s_1, \dots, s_K\}$
 Latent clusters $\mathbf{B} \in \mathbb{R}^{N \times K}$
 Cluster correlations $\mathbf{A} \in \mathbb{R}^{K \times K}$

begin

- L_1 -normalize the rows of \mathbf{C} to form $\bar{\mathbf{C}}$.
- Find \mathbf{S} via the column pivoted QR on $\bar{\mathbf{C}}^T$.
- Find $\check{\mathbf{B}}$ with $\check{\mathbf{B}}_{ki} = p(\text{cluster } k \mid \text{object } i)$ by solving N simplex-constrained least squares in parallel to minimize $\|\bar{\mathbf{C}} - \check{\mathbf{B}}^T \bar{\mathbf{C}}_{S^*}\|_F$.
- Recover \mathbf{B} from $\check{\mathbf{B}}$ by the Bayes' rule.
- Recover \mathbf{A} by $\mathbf{B}_{S^*}^{-1} \mathbf{C}_{SS} \mathbf{B}_{S^*}^{-1}$.

end

Algorithm 4: Rectified AW algorithm (RAW)

Input/Output: Same as Algorithm 3

begin

- $\mathbf{C}_0 \leftarrow \mathbf{C}$
- repeat** with $t = 0, 1, 2, \dots$

 - $(\mathbf{U}, \Lambda_K) \leftarrow \text{Truncated-Eig}(\mathbf{C}_t, K)$
 - $\Lambda_K^+ \leftarrow \max(\Lambda_K, 0)$
 - $\mathbf{C}^{\mathcal{PSD}_K} \leftarrow \mathbf{U} \Lambda_K^+ \mathbf{U}^T$
 - $\mathbf{C}^{\mathcal{NOR}} \leftarrow \mathbf{C}^{\mathcal{PSD}_K} + \frac{1 - \sum_{ij} C_{ij}^{\mathcal{PSD}_K}}{N^2} \mathbf{e} \mathbf{e}^T$
 - $\mathbf{C}^{\mathcal{NN}} \leftarrow \max(\mathbf{C}^{\mathcal{NOR}}, 0)$
 - $\mathbf{C}_{t+1} \leftarrow \mathbf{C}^{\mathcal{NN}}$

- until** converging to a certain \mathbf{C}^*
- $(\mathbf{S}, \mathbf{B}, \mathbf{A}) \leftarrow \text{AW}(\mathbf{C}^*, K)$ (Algorithm 3)

end

of \mathbf{C} in Equation (5.1), which penalizes all diagonal entries. The **Rectified Anchor Word (RAW)** algorithm has an additional rectification step that enforces that \mathbf{C} should enjoy the expected structures before running the main algorithm. In [102], an Alternating Projection (AP) rectification as given in Algorithm 4 is used to overcome the gap between the underlying assumptions of our models and the actual data.

Rectification is also important for addressing the issue of *outlier bias*. Real data often exhibits rare objects that are only present in a few examples. The corresponding co-occurrence of these objects are inevitably sparse with large variance, but the greedy anchor selection favors choosing these outliers due to L_1 normalization of \mathbf{C} . Previous work tried to bypass this problem by oversampling clusters by the number of outliers under some additional identifiability assumptions [61]. This approach is not always feasible, especially for a large vocabulary that introduces many low frequency objects. When synonyms and short documents cause undesirable sparsity to Latent Semantic Analysis [100], projection onto the leading eigen-subspace blurs sparse co-occurrences. Similarly, \mathcal{PSD}_K -projection turns out to significantly reduce outlier bias, and the remaining projections are useful for maintaining the probabilistic structures of \mathbf{C} , which then allow users to recover \mathbf{B} and \mathbf{A} in Algorithm 3.

Handling a *large vocabulary* is another major challenge for spectral methods. Even if we limit our focus only to second-order models, the space complexity of RAW is already $O(N^2)$, growing rapidly with increasing vocabulary. We are unable to exploit the high sparsity of \mathbf{C} as a single iteration of the AP-rectification makes \mathbf{C} significantly denser. The three projections in AP-rectification and the rest of the anchor word algorithm in Algorithm 3 have time complexities of $O(N^2K)$, $O(N^2)$, $O(N^2)$ and $O(N^2K)$, respectively, and so pose a difficulty when scaling to a large vocabulary size N . On the other hand, the *separability assumption* is crucial for the second-order models, and while there has been a line of research that tries to relax this assumption [17, 86], it has been formally shown that most topic models are indeed separable if their vocabulary sizes are sufficiently larger than the number of clusters [46], again emphasizing the urgency of an approach with better time and space scaling in the vocabulary size.

5.4 Low-rank Rectification and Compression

The rectified co-occurrence C^* in Section 5.3 must be of rank K and positive semidefinite, hinting at an opportunity to represent it in terms of an outer product $\mathbf{Y}\mathbf{Y}^T$ for some $\mathbf{Y} \in \mathbb{R}^{N \times K}$. One idea for achieving this structure is to use a low-rank representation $C_t = \mathbf{Y}_t\mathbf{Y}_t^T$ throughout the rectification in Algorithm 4. Another way to obtain this structure is to directly minimize $\|C - \mathbf{Y}\mathbf{Y}^T\|_F$ with the necessary constraints. In this section, we propose two new algorithms to simultaneously compress and rectify the input by representing $C \in \mathbb{R}^{N \times N}$ by a low-rank outer product $\mathbf{Y}\mathbf{Y}^T$.

Algorithm 5: ENN-rectification (ENN)

Input: Object co-occurrence $C \in \mathbb{R}^{N \times N}$
 The number of clusters K
Output: Rectified compression $\mathbf{Y} \in \mathbb{R}^{N \times K}$

begin

- $\mathbf{E} \leftarrow \mathbf{0} \in \mathbb{R}^{N \times N}$ (sparse format)
- $C^{op} : \mathbf{x} \rightarrow C\mathbf{x}$ (Implicit operator)
- repeat** with $t = 0, 1, 2, \dots$

 - $(\mathbf{U}, \Lambda_K) \leftarrow \text{Truncated-Eig}(C^{op}, K)$
 - $\Lambda_K^+ \leftarrow \max(\Lambda_K, 0)$
 - $\mathbf{Y} \leftarrow \mathbf{U}(\Lambda_K^+)^{1/2}$
 - $E_{ij} \leftarrow \max(-\mathbf{Y}_{i*}\mathbf{Y}_{j*}^T, 0)$
 - $r \leftarrow (1 - \|\mathbf{Y}^T \mathbf{e}\|_2^2 - \sum_{ij} E_{ij})/N^2$
 - $C^{op} : \mathbf{x} \rightarrow \mathbf{Y}(\mathbf{Y}^T \mathbf{x}) + \mathbf{E}\mathbf{x} + r(\mathbf{e}^T \mathbf{x})\mathbf{e}$

- until** \mathbf{E} converges

end

5.4.1 Epsilon Non-Negative Rectification (ENN)

The alternating projection iteration in Algorithm 4 produces low-rank semi-definite intermediate matrices in factored form at each iteration. By construction, the positive

Algorithm 6: PALM-rectification (PALM)

Input: Object co-occurrence $C \in \mathbb{R}^{N \times N}$
 The number of clusters K
Output: Rectified compression $Y \in \mathbb{R}^{N \times K}$

```

begin
     $(V, D) \leftarrow \text{Truncated-Eig}(C, K)$ 
     $(X_0, Y_0) \leftarrow (V \sqrt{D}, V \sqrt{D})$ 
    repeat
         $c_t \leftarrow \gamma L_1(Y_t)$ 
         $X'_{t+1} \leftarrow X_t - (1/c_t) \nabla_X J(X_t, Y_t)$ 
         $X_{t+1} \leftarrow \max(X'_{t+1}, 0)$ 
         $d_t \leftarrow \gamma L_2(X_{t+1})$ 
         $Y'_{t+1} \leftarrow Y_t - (1/d_t) \nabla_Y J(X_{t+1}, Y_t)$ 
         $Y_{t+1} \leftarrow \max(Y'_{t+1}, 0)$ 
    until  $Y$  converges
end

```

semi-definite projection (\mathcal{PSD}_K) and the normalization projection (\mathcal{NOR}) produce positive semi-definite matrices of rank K and $K + 1$, respectively. Unfortunately, the final projection to enforce elementwise non-negativity (\mathcal{NN}) destroys this low rank structure. However, the \mathcal{NN} projection only makes significant changes to a few elements; that is, the output of the \mathcal{NN} projection at step t is nearly rank $K + 1$ plus a sparse correction E_t . The Epsilon Non-Negative Rectification algorithm (Algorithm 5) has the same structure as Algorithm 4, but with the key difference that it returns a sparse-plus-low-rank representation of the \mathcal{NN} projection rather than materializing a dense representation. Matrix-vector products with this sparse-plus-low-rank representation require $O(NK + \text{nnz}(E_t))$ time, and $O(K)$ such matrix-vector products can be used in a Lanczos eigen-solver to compute the truncated eigen-decomposition at the start of the next iteration.

Maintaining a sparse correction matrix E_t at each step lets the ENN approach avoid the storage overheads of the original alternating projection algorithm. To overcome the quadratic time cost at each iteration, though, we need to avoid explicitly computing ev-

every element of the intermediate $\mathbf{Y}\mathbf{Y}^T$ in the course of the NN projection. However, we can bound the magnitude of many elements of $\mathbf{Y}\mathbf{Y}^T$ by the Cauchy-Schwartz inequality: $|C_{ij}| \leq \|\mathbf{y}_i\|_2 \|\mathbf{y}_j\|_2$ where \mathbf{y}_i and \mathbf{y}_j denote columns of \mathbf{Y}^T . Let I denote the index set indices $\{i : \|\mathbf{y}_i\|_2^2 > \epsilon\} \subseteq [N]$ for given ϵ ; then every large entry of C belongs to either $\mathbf{Y}_{I*}\mathbf{Y}^T$ or $\mathbf{Y}(\mathbf{Y}^T)_{*I}$. As C is symmetric, checking the negative entries in $\mathbf{Y}_{I*}\mathbf{Y}^T$ is sufficient to find a symmetric correction \mathbf{E} that guarantees $\mathbf{Y}\mathbf{Y}^T + \mathbf{E} \geq -\epsilon$. We refer to this property as *Epsilon Non-Negativity*: ϵ balances the trade-off between the effect of leaving small negative entries versus increasing the size of I to look up. We limit $|I|$ to be $O(K)$ based on the common sampling complexity of a suitable set of rows for a near-optimal rank-K approximation⁴.

5.4.2 Proximal Alternating Linearized Minimization Rectification (PALM)

To avoid small negative entries, we investigate another rectified compression algorithm that directly minimizes $\|\mathbf{C} - \mathbf{Y}\mathbf{Y}^T\|_F$ subject to the stronger NN -constraint $\mathbf{Y} \geq 0$ and the usual NOR -constraint $\|\mathbf{Y}^T \mathbf{e}\|_2 = 1$. Concretely, we try to

$$\text{minimize } J(\mathbf{X}, \mathbf{Y}) := \frac{1}{2} \|\mathbf{C} - \mathbf{X}\mathbf{Y}^T\|_F^2 + \frac{s}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 \quad \text{subject to } \mathbf{X} \geq 0, \mathbf{Y} \geq 0. \quad (5.2)$$

\mathcal{PSD}_K - and NOR -constraints are implicitly satisfied by jointly minimizing the two terms in the objective function J , whereas NN -constraint is explicit in the formulation. Thus we can apply the Proximal Alternating Linearized Minimization [27] for learning \mathbf{Y} given \mathbf{C} ; the relevant proximal operator is NN projection of \mathbf{Y} , which takes $O(NK)$ at most.

⁴This choice is standard in literature on low-rank approximation via column subset selection.

Note that J is semi-algebraic (as it is a real polynomial) with two partial derivatives: $\nabla_{\mathbf{X}} J = (\mathbf{X}\mathbf{Y}^T - \mathbf{C})\mathbf{Y} + s(\mathbf{X} - \mathbf{Y})$ and $\nabla_{\mathbf{Y}} J = (\mathbf{Y}\mathbf{X}^T - \mathbf{C})\mathbf{X} + s(\mathbf{Y} - \mathbf{X})$. So, the following lemma guarantees the global convergence.

Lemma 3. *For any fixed \mathbf{Y} , $\nabla_{\mathbf{X}} J(\mathbf{X}, \mathbf{Y})$ is globally Lipschitz continuous with the moduli $L_1(\mathbf{Y}) = \|\mathbf{Y}^T \mathbf{Y} + s\mathbf{I}_K\|_2$. So is $\nabla_{\mathbf{Y}} J(\mathbf{X}, \mathbf{Y})$ given any fixed \mathbf{X} with $L_2(\mathbf{X}) = \|\mathbf{X}^T \mathbf{X} + s\mathbf{I}_K\|_2$.*

Proof.

$$\begin{aligned} & \|\nabla_{\mathbf{X}} J(\mathbf{X}, \mathbf{Y}) - \nabla_{\mathbf{X}} J(\mathbf{X}', \mathbf{Y})\|_F \\ &= \|(\mathbf{Y}^T \mathbf{Y} + s\mathbf{I}_K)(\mathbf{X} - \mathbf{X}')\|_F \\ &\leq \|\mathbf{Y}^T \mathbf{Y} + s\mathbf{I}_K\|_2 \cdot \|\mathbf{X} - \mathbf{X}'\|_F \end{aligned}$$

The proof is symmetric for the other case with $L_2(\mathbf{X}) = \|\mathbf{X}^T \mathbf{X} + s\mathbf{I}_K\|_2$. \square

Algorithm 6 shows the PALM-rectification with the adaptive control of the learning rates based on the tight 2-norm Lipschitz modulis at each step t .

5.5 Low-rank Anchor Word Algorithm

The output for both methods in §5.4 is a compressed co-occurrence matrix $\mathbf{C} = \mathbf{Y}\mathbf{Y}^T$. In this section, we present the **Low-rank Anchor Word algorithm (LAW)** that reduces the time complexity of finding anchor objects from $O(N^2K)$ to $O(NK^2)$ by taking advantage of this form. We note that LAW applies whenever \mathbf{C} is in a low-rank representation, which does not have to be derived from our methods. Moreover, it is exact for non-negative \mathbf{C} , but it is robust in practice when we allow small negative entries in \mathbf{C} , as in the case with ENN.

Algorithm 7: Low-rank AW (LAW)

Input: Object co-occurrence $\mathbf{C} = \mathbf{Y}\mathbf{Y}^T$
Output: Anchor objects $S = \{s_1, \dots, s_K\}$
Latent clusters $\mathbf{B} \in \mathbb{R}^{N \times K}$
Cluster correlations $\mathbf{A} \in \mathbb{R}^{K \times K}$

begin

- Calculate row sums $\mathbf{d} = \mathbf{Y}(\mathbf{Y}^T \mathbf{e})$.
- Normalize $\bar{\mathbf{Y}} \leftarrow \text{diag}(\mathbf{d})^{-1} \mathbf{Y}$.
- Compute QR decomposition of $\mathbf{Y} = \mathbf{Q}\mathbf{R}$.
- Form $\mathbf{X} = \bar{\mathbf{Y}}\mathbf{R}^T$.
- Select S using column pivoted QR on \mathbf{X}^T .
- Solve n simplex-constrained least square problems to minimize $\|\mathbf{X} - \check{\mathbf{B}}\mathbf{X}_{S*}\|_F$.
- Recover \mathbf{B} from $\check{\mathbf{B}}$ using Bayes' rule.
- Recover $\mathbf{A} = \mathbf{B}_{S*}^{-1} \mathbf{Y}_{S*} \mathbf{Y}_{S*}^T \mathbf{B}_{S*}^{-1}$.

end

Algorithm 8: Low-rank JSMF (LR-JSMF)

Input: Raw object-example $\mathbf{H} \in \mathbb{R}^{N \times M}$
Output: Anchor objects $S = \{s_1, \dots, s_K\}$
Latent clusters $\mathbf{B} \in \mathbb{R}^{N \times K}$
Cluster correlations $\mathbf{A} \in \mathbb{R}^{K \times K}$

begin

- Get $\hat{\mathbf{H}}, \hat{\mathbf{H}}_{\text{diag}}$ from \mathbf{H} by (5.1).
- $C_{op} : \mathbf{x} \rightarrow \hat{\mathbf{H}}(\hat{\mathbf{H}}^T \mathbf{x}) - \hat{\mathbf{H}}_{\text{diag}} \mathbf{x}$
- $(\mathbf{U}, \Lambda_K) \leftarrow \text{Randomized-Eig}(C_{op}, K)$
- Initialize ENN with \mathbf{U}, Λ_K .
- $\mathbf{Y} \leftarrow \text{ENN-rectification}$
- $(S, \mathbf{B}, \mathbf{A}) \leftarrow \text{LAW}(\mathbf{Y})$ (Algorithm 7)

end

The first step is to L_1 -normalize the rows of \mathbf{C} . Given $\mathbf{C} \geq 0$, the L_1 -norm of each row is simply the sum of all its entries, so we can calculate the row norms by $\mathbf{d} = \mathbf{Y}(\mathbf{Y}^T \mathbf{e})$. To obtain the normalized \mathbf{C} , we simply scale the rows of \mathbf{Y} , and $\overline{\mathbf{C}} = (\text{diag}(\mathbf{d})^{-1} \mathbf{Y}) \mathbf{Y}^T = \overline{\mathbf{Y}} \mathbf{Y}^T$. These steps cost $O(NK)$.

Next, we need to apply column pivoted QR to $\overline{\mathbf{C}}^T$ in order to identify the pivots as our anchor objects \mathbf{S} . By taking the QR decomposition $\mathbf{Y} = \mathbf{Q}\mathbf{R}$, $\overline{\mathbf{C}}^T$ can be further transformed into $\mathbf{Q}\mathbf{R}\overline{\mathbf{Y}}^T$. Notice that $\overline{\mathbf{C}}^T$ is an orthogonal embedding of $\mathbf{X}^T = \mathbf{R}\overline{\mathbf{Y}}^T$ onto a higher-dimensional space, which preserves the column L_2 -norms. Lemma 4 shows that column pivoted QR on $\overline{\mathbf{C}}^T$ and on $\mathbf{R}\overline{\mathbf{Y}}^T$ are equivalent, which allows us to lower the computation cost from $O(N^2K)$ to $O(NK^2)$.

Lemma 4. *Let \mathbf{S} be the set of pivots that have been selected by column pivoted QR on $\overline{\mathbf{C}}^T = \mathbf{Q}\mathbf{X}^T$. Given the QR decomposition, $\overline{\mathbf{X}}_{\mathbf{S}*}^T = \mathbf{P}\mathbf{T}$, then $\overline{\mathbf{C}}_{\mathbf{S}*}^T = (\mathbf{Q}\mathbf{P})\mathbf{T}$ is the corresponding QR decomposition for the columns of $\overline{\mathbf{C}}$. For any remaining row $i \in [N] \setminus \mathbf{S}$,*

$$\|(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\overline{\mathbf{X}}_{i*}^T\|_2 = \|(\mathbf{I} - (\mathbf{Q}\mathbf{P})(\mathbf{Q}\mathbf{P})^T)\overline{\mathbf{C}}_{i*}^T\|_2 \quad (5.3)$$

Therefore, the next column pivot is identical for $\overline{\mathbf{C}}^T$ and $\overline{\mathbf{X}}^T$. By induction, column pivoted QR on $\overline{\mathbf{C}}^T$ and $\overline{\mathbf{X}}^T$ return the same pivots.

Proof. Because both \mathbf{Q} and \mathbf{P} have orthonormal columns,

$$(\mathbf{Q}\mathbf{P})^T(\mathbf{Q}\mathbf{P}) = \mathbf{P}^T(\mathbf{Q}^T\mathbf{Q})\mathbf{P} = \mathbf{P}^T\mathbf{P} = \mathbf{I} \quad (5.4)$$

Thus, $\mathbf{Q}\mathbf{P}$ and \mathbf{T} forms the QR decomposition of $\overline{\mathbf{C}}_{\mathbf{S}*}^T$. The residual of a remaining column $i \in [N] \setminus \mathbf{S}$ is $(\mathbf{I} - (\mathbf{Q}\mathbf{P})(\mathbf{Q}\mathbf{P})^T)\overline{\mathbf{C}}_{i*}^T$ and $(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\overline{\mathbf{X}}_{i*}^T$ for $\overline{\mathbf{C}}^T$ and $\overline{\mathbf{X}}^T$, respectively. Simplify the former gives us

$$(\mathbf{I} - (\mathbf{Q}\mathbf{P})(\mathbf{Q}\mathbf{P})^T)\overline{\mathbf{C}}_{i*}^T$$

$$\begin{aligned}
&= (\mathbf{I} - (\mathbf{Q}\mathbf{P})(\mathbf{Q}\mathbf{P})^T)\mathbf{Q}\overline{\mathbf{X}}_{i*}^T \\
&= \mathbf{Q}\overline{\mathbf{X}}_{i*}^T - \mathbf{Q}\mathbf{P}\mathbf{P}^T\mathbf{Q}^T\mathbf{Q}\overline{\mathbf{X}}_{i*}^T \\
&= \mathbf{Q}(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\overline{\mathbf{X}}_{i*}^T
\end{aligned}$$

Finally,

$$\|(\mathbf{I} - (\mathbf{Q}\mathbf{P})(\mathbf{Q}\mathbf{P})^T)\overline{\mathbf{C}}_{i*}^T\|_2^2 = \overline{\mathbf{X}}_{i*}(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\mathbf{Q}^T\mathbf{Q}(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\overline{\mathbf{X}}_{i*}^T = \|(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\overline{\mathbf{X}}_{i*}^T\|_2^2$$
(5.5)

Because the next pivot is selected as the column whose residual has the largest L_2 -norm, Eq. 5.5 indicates that the same pivot will be selected for $\overline{\mathbf{C}}^T$ and $\overline{\mathbf{X}}^T$. Inductively, the anchors \mathbf{S} recovered by column pivoted QR on those matrices are equivalent. \square

Following the recovery of \mathbf{S} , AW solves N independent simplex-constrained least square problems $\|\overline{\mathbf{C}}_{i*} - \check{\mathbf{B}}_{i*}^T \overline{\mathbf{C}}_{S*}\|_2$. Again we can leverage the L_2 -norm preserving property,

$$\|\overline{\mathbf{C}}_{i*} - \check{\mathbf{B}}_{i*}^T \overline{\mathbf{C}}_{S*}\|_2 = \|\mathbf{X}_{i*}\mathbf{Q}^T - \check{\mathbf{B}}_{i*}^T \mathbf{X}_{S*}\mathbf{Q}^T\|_2 = \|\mathbf{X}_{i*} - \check{\mathbf{B}}_{i*}^T \mathbf{X}_{S*}\|_2$$
(5.6)

and reduce the dimension of the least-square problems from N to K , thereby the complexity from $O(N^2K)$ to $O(NK^2)$. The remaining part of the algorithm follows exactly as AW.

Low-rank Joint Stochastic Matrix Factorization (LR-JSMF) We complete our scalable framework of processing co-occurrence statistic by introducing a direct initialization method from the raw object-example data for ENN. This allows us to avoid creating and storing \mathbf{C} , which is a burden of memory when N becomes sufficiently large. In Algorithm 5, \mathbf{C} only appears in the initial truncated eigen-decomposition, after which we maintain the compressed operator \mathbf{C}_{op} independent of it. On the other hand, we just need the matrix-vector multiplication by \mathbf{C} for the iterative methods in

initialization. Using the generative formula in Eq. 5.1, we are able to implicitly apply \mathbf{C} to vectors as an outer-product plus diagonal operator, in terms of \mathbf{H} , at $O(NMK)$ computation cost. To further reduce the number of times the operator is applied, we adopt the one-pass randomized eigen-decomposition by Halko et al. [71]. This technique enables us to initialize with a single pass over the dataset, without concurrently storing the entire \mathbf{H} in memory. A limitation is when the number of clusters is large and the gap between the K -th eigenvalue and the ones below is small, we will have to incorporate a few power iterations for refinement, as suggested by the original paper. This will result in a multi-pass method, but still far more efficient on large object size and parallelization-friendly.

5.6 Experimental Results

A good factorization should be accurate, meaningful, and fast. In two series of experiments, we show that LR-JSMF maintains model quality while running in a fraction of the space and time needed for the original JSMF method. Previous methods have required truncation of the vocabulary even to run on consumer-grade computers. We show not only that we are able to handle increasingly large vocabularies without loss of speed, but that using larger vocabularies measurably improves model quality relative to truncated vocabularies.

For the first series of experiments, we measure the accuracy of each rectification component as well as the entire pipeline of LR-JSMF. To produce a strong baseline, we begin with constructing the full co-occurrence \mathbf{C} from each of our datasets \mathbf{H} by (5.1), and produce the rectified \mathbf{C}_{AP} by running Alternating Projection (AP) on \mathbf{C} . Next we compress \mathbf{C} into \mathbf{Y}_{ENN} and \mathbf{Y}_{PALM} by running ENN (50 iterations, $|I| = 10K + 1000$) and

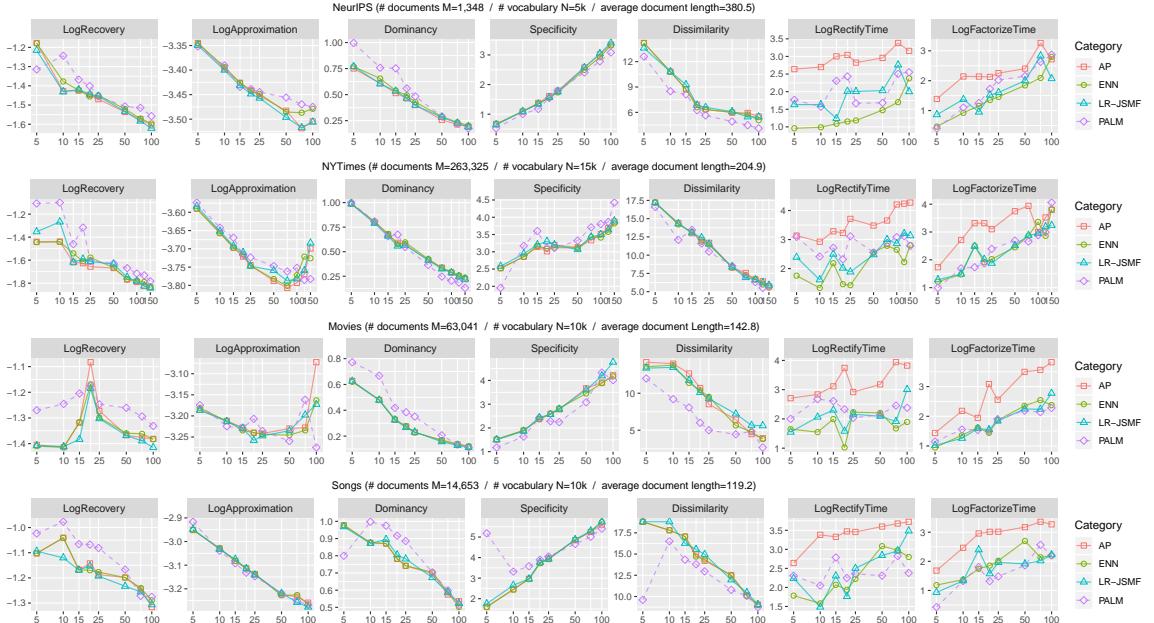


Figure 5.1: Experiment on four datasets. ENN and LR-JSMF mostly agree with AP, whereas PALM has slight inconsistency. The general information of each dataset is above the corresponding row. Recovery, approximation, and runtimes are in \log_{10} scale. Note that ENN and LR-JSMF are almost two orders of magnitude faster than AP. The x -axis indicates the number of clusters K . Lower numbers in y -axis are better except Specificity and Dissimilarity.

PALM (100 iterations, $s = 1e^{-4}$). For testing our complete low-rank pipeline, we also construct (\mathbf{V}, \mathbf{D}) directly from the raw data \mathbf{H} by the randomized eigen-decomposition in Algorithm 8, learning the compressed statistics $\mathbf{Y}_{LR-JSMF}$ again by running ENN initialized with $\mathbf{V} \sqrt{\mathbf{D}}$. Then we run the Anchor Word algorithm (AW) on \mathbf{C}_{AP} and the Low-rank Anchor Word algorithm (LAW) on each of \mathbf{Y}_{ENN} , \mathbf{Y}_{PALM} , and $\mathbf{Y}_{LR-JSMF}$.

The goal of rectification is to apply spectral inference to data that does not follow our modeling assumptions, so we evaluate on real data. In addition to two standard datasets from the UCI Machine Learning repository (NeurIPS papers and New York Times articles), we also use two non-textual datasets (Movies and Songs) previously used to demonstrate the performance of full algorithm with AP-rectification in [102]. Although our ultimate goal is to extend JSMF to large vocabularies, we use the same restricted vocabulary as [102] for a fair comparison in the first series of experiment.

Figure 5.1 shows the overall performance of the learned topic clusters from these four datasets with increasing number of clusters K . Low **Recovery** error $\frac{1}{N} \sum_i \|\overline{\mathbf{C}}_{i*} - \check{\mathbf{B}}_{i*} \overline{\mathbf{C}}_{S*}\|_2$ implies that the learned anchor objects successfully reconstruct the co-occurrence space of the entire objects. Low **Approximation** error $\|\mathbf{C} - \mathbf{BAB}^T\|_2$ means that our factorization captures most of information given in the unbiased co-occurrence statistics. In real data, low **Dominancy** $\frac{1}{K} \sum_k \mathbf{A}_{kk}$ implies that our models learn more correlations between clusters. High **Specificity** $\frac{1}{K} \sum_k \text{KL}(\mathbf{B}_{*k} \| \sum_i \mathbf{C}_{*i})$ indicates that the learned clusters are different enough from the corpus distribution, whereas high **Dissimilarity** counts the average number of objects in each cluster that do not occur among top 20 in other clusters, showing the interpretable difference across the learned clusters. We do not report the Cluster Coherence because it often measures deceptively [32]. The first five columns show that ENN and LR-JSMF learn approximately same clusters as JSMF with the full AP, showing no visible loss in accuracy across all settings. More importantly, the randomness we introduced into LR-JSMF results in a very low variance over a number of runs. This is important as the stability of spectral inference is a major advantage relative to MCMC or Variational Inference. Although PALM deviates a small amount from the other three methods in a few cases, it mostly achieves the same level of accuracy and follows the overall trend closely. In terms of runtimes, all of our methods have clear advantage over AP, gaining $1 \sim 2$ orders of magnitude speedup in most situations. Even for applications on relatively small vocabulary sizes, our algorithms shows a notable improvement in efficiency.

For the second series of experiments, we create eight corpora $\{\mathbf{H}_N, \mathbf{H}_{2N}, \dots, \mathbf{H}_{8N}\}$ for each dataset \mathbf{H} by tailoring their vocabulary sizes as multiples of a base vocabulary of N objects. In this case we are not able to compare LR-JSMF to previous methods because we cannot store the full co-occurrence matrices for the larger vocabulary: these models would be impossible. Figure 5.2 illustrates the overall performance of

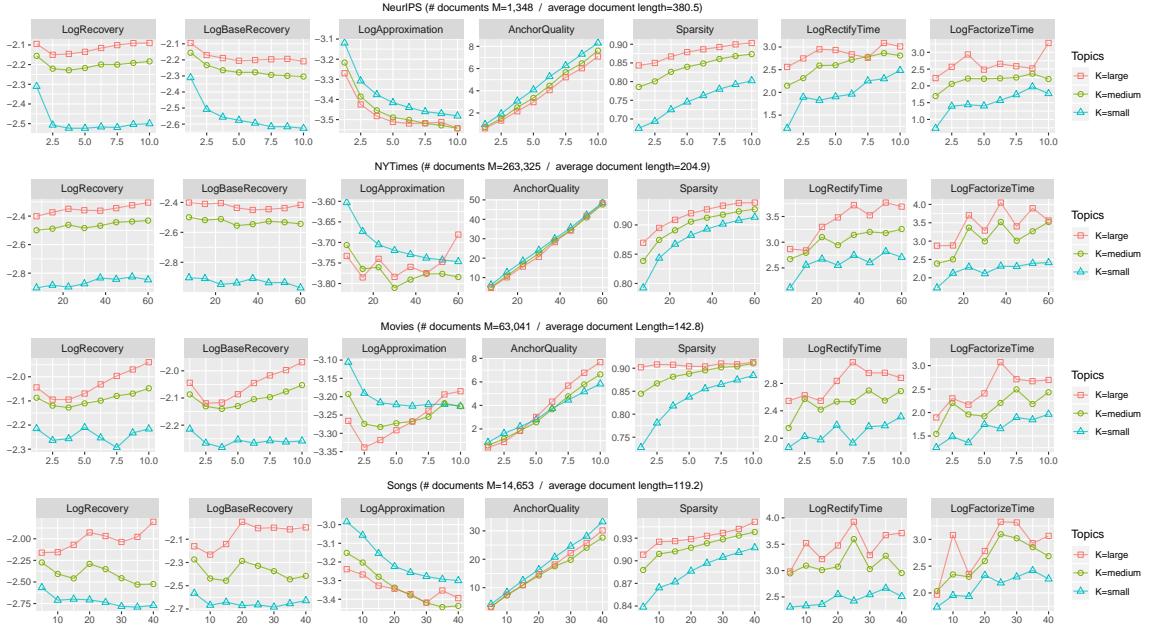


Figure 5.2: As we increase the vocabulary size for four collections, anchor quality and sparsity improve, but running time is stable. The x -axis indicates the vocabulary size N in thousands. Values above 15k will not fit in memory on standard hardware with previous algorithms.

the learned small/medium/large size clusters with increasing vocabulary size N . Low **BaseRecovery** means that the anchor objects from the models with larger vocabulary better reconstruct the objects in our base vocabulary (H_N). High **AnchorQuality** indicates that the average rank of the anchor objects s_k in every other topic clusters than k is high, implying the anchor objects rarely contribute to other clusters than their own. High **Sparsity** ($\frac{1}{K} \sum_k \frac{\sqrt{N} - (\|b_k\|_1 / \|b_k\|_2)}{\sqrt{N} - 1}$) [80] says that our topic clusters are more concentrated on specific objects.

We observe the quality of anchors increases with increasing vocabulary size, verifying that using larger vocabularies helps better satisfy the separability assumption. We also verify that a large vocabulary often better approximates the co-occurrence statistics and better reconstructs the co-occurrence space of the *base* vocabulary, but these patterns are not always consistent in non-textual datasets. In contrast, Sparsity consistently improves, increasing the interpretability of the learned clusters. Most excitingly, the run-

ning times of ENN and LAW show the scalability of our new rectification and low-rank algorithm, thereby demonstrating that LR-JSMF is an efficient and robust pipeline.

Finally, we have also inspected the qualitative behavior of the recovered clusters, as we increase the vocabulary size. The topic clusters become significantly more specific, while the clustering of objects is more conspicuous. Figure 5.3 shows how using a larger vocabulary size can lead to more distinguishable topics, especially as it allows us to make use of words that are relatively rare, but used in much narrower contexts. Going from left to right, we can observe that the set of topic words become more and more specific: for instance, the topic corresponding to the third row is slightly vague when observing just the left half of the row, while as we increase the vocabulary size beyond 5000, we gain access to highly topic-specific words such as `hjb` (Hamilton-Jacobi-Bellman equation) or `pid` (Proportional Integral Derivative), which signifies the row’s pertinence to dynamical and control systems. The flip side of the figure shows that words we would normally consider as non-topical can often be assigned high contributions towards certain topics. The strong red shade on the bottom left indicates that words such as “equivalent” or “cambridge” are strongly connected to the machine learning literature.

5.7 Conclusion

Spectral algorithms provide an appealing alternative for identifying interpretable low-rank subspaces by simple factorizations of higher-order moments. But this simplicity is also a weakness: violations of modeling assumptions destroy performance unless they are handled through rectification, and the size of the moment matrices limits us to small vocabularies. In this paper, we developed an efficient and scalable framework:

		Vocabulary Size						
		1250	2500	3750	5000 (base)	6250	7500	8750
refractory interconnection seen detail transmission considered	interspike	bursting	neuron	signalling	ipsp	tst		
	marder	stomatogastric	circuit	meilijson	stg	tam		
	abbott	konishi	synaptic	quiescent	substances	hyperpolarized		
	acad	axonal	cell	ryckebusch	inactivation	memorized		
	pyloric	modulatory	layer	leech	depolarized	transposed		
	bird	ionic	signal	silent	shapiro	tsividis		
san additional amount considered developed significant	male	kanji	recognition	radical	sicl	subword		
	henderson	phonemic	layer	npm	joe	phoneroes		
	jackel	subsystem	hidden	demi	chinese	otherfilter		
	recog	dtw	word	shikano	hanazawa	sdnn		
	dictionary	strokes	speech	letterform	lexicon	perplexity		
	ocr	gender	net	tebelskis	preprocessed	males		
cambridge pendulum requires con bellman plan	discounted	tutor	control	hjb	jacobi	ovi		
	bradtke	lqr	action	rein	forcement	pid		
	discount	disturbances	dynamic	biped	viscosity	idm		
	eligibility	disturbance	optimal	trol	sel	umass		
	indirect	hamilton	reinforcement	handicapped	bizzi	missile		
	amherst	smdp	controller	gullapalli	swinging	queueing		
directional seen dark neuroscience soc supported	transparent	luminance	cell	unoriented	aftereffect	moc		
	adelson	ruderman	field	heeger	blast	ori		
	geniculate	andersen	visual	thalamus	mae	taube		
	amacrine	bergen	motion	directionally	knierim	muller		
	deg	selectively	direction	hayashi	mexican	swindle		
	mennaughton	lond	image	werblin	specimen	skagg		
equivalent computing cambridge considered simply detail	fix	solvable	gaussian	birmingham	boxplot	pbr		
	satisfying	distribu	noise	kolmogorov	dependences	owi		
	royal	barber	approximation	gmm	cb2	danish		
	opper	parametrized	hidden	winther	trigonometric	ylz		
	leibler	const	bound	imation	colt	diabetes		
	treatment	eter	matrix	statist	minimizer	devroye		

Figure 5.3: Losses or gains in topic words depending on the vocabulary size. Each row represents a topic from the NeurIPS dataset, with the top 6 topical words shown in the middle column. The red and green cells denote topic words that are lost or gained by shifting the vocabulary size from the default size 5000, respectively. The intensities of the colors indicate the words’ contributions towards the specific topic.

Low-Rank Joint Stochastic Matrix Factorization. We provide theoretical advances in compressed matrix factorization, leading to high-quality low-rank non-negative approximations without quadratic blowup. The method provides orders of magnitude speedups for rectification even on small vocabularies. Perhaps most importantly, we can now apply reliable, high-quality factorizations of high-dimensional data sets on laptop-grade hardware, massively increasing the applicability and potential use of these algorithms.

CHAPTER 6

WEIGHTED K-MEANS FOR ELECTRONIC STRUCTURE CALCULATION

6.1 Abstract

The recently developed interpolative separable density fitting (ISDF) decomposition is a powerful way for compressing the redundant information in the set of orbital pairs, and has been used to accelerate quantum chemistry calculations in a number of contexts. The key ingredient of the ISDF decomposition is to select a set of non-uniform grid points, so that the values of the orbital pairs evaluated at such grid points can be used to accurately interpolate those evaluated at all grid points. The set of non-uniform grid points, called the interpolation points, can be automatically selected by a QR factorization with column pivoting (QRCP) procedure. This is the computationally most expensive step in the construction of the ISDF decomposition. In this work, we propose a new approach to find the interpolation points based on the centroidal Voronoi tessellation (CVT) method, which offers a much less expensive alternative to the QRCP procedure when ISDF is used in the context of hybrid functional electronic structure calculations. The CVT method only uses information from the electron density, and can be efficiently implemented using a K-Means algorithm. We find that this new method achieves comparable accuracy to the ISDF-QRCP method, at a cost that is negligible in the overall hybrid functional calculations. For instance, for a system containing 1000 silicon atoms simulated using the HSE06 hybrid functional on 2000 computational cores, the cost of QRCP-based method for finding the interpolation points costs 38.1 s, while the CVT procedure only takes 0.7 s. We also find that the ISDF-CVT method also enhances the smoothness of the potential energy surface in the context of *ab initio* molecular dynamics (AIMD) simulations with hybrid functionals.

6.2 Introduction

Orbital pairs of the form $\{\varphi_i(\mathbf{r})\psi_j(\mathbf{r})\}_{i,j=1}^N$, where φ_i, ψ_j are single particle orbitals, appear ubiquitously in quantum chemistry. A few examples include the Fock exchange operator, the MP2 amplitude, and the polarizability operator [160, 121]. When N is proportional to the number of electrons N_e in the system, the total number of orbital pairs is $N^2 \sim O(N_e^2)$. On the other hand, the number of degrees of freedom needed to resolve all orbital pairs on a dense grid is only $O(N_e)$. Hence as N_e becomes large, the set of all orbital pairs contains apparent redundant information. In order to compress the redundant information and to design more efficient numerical algorithms, many algorithms in the past few decades have been developed. Pseudospectral decomposition [129, 149], Cholesky decomposition [20, 96, 6, 120], density fitting (DF) or resolution of identity (RI) [148, 168], and tensor hypercontraction (THC) [138, 139] are only a few examples towards this goal. When the single particle orbitals φ_i, ψ_j are already localized functions, “local methods” or “linear scaling methods” [63, 29, 69, 132] can be applied to construct such decomposition with cost that scales linearly with respect to N_e . Otherwise, the storage cost of the matrix to represent all orbital pairs on a grid is already $O(N_e^3)$, and the computational cost of compressing the orbital pairs is then typically $O(N_e^4)$.

Recently, Lu and Ying developed a new decomposition called the interpolative separable density fitting (ISDF) [113], which takes the following form

$$\varphi_i(\mathbf{r})\psi_j(\mathbf{r}) \approx \sum_{\mu=1}^{N_\mu} \zeta_\mu(\mathbf{r}) (\varphi_i(\hat{\mathbf{r}}_\mu)\psi_j(\hat{\mathbf{r}}_\mu)). \quad (6.1)$$

For a given \mathbf{r} , if we view $[\psi_i(\mathbf{r})\psi_j(\mathbf{r})]$ as a row of the matrix $\{\psi_i\psi_j\}$ discretized on a dense grid, then the ISDF decomposition states that all such matrix rows can be approximately expanded using a linear combination of matrix rows with respect to a selected

set of *interpolation points* $\{\hat{r}_\mu\}_{\mu=1}^{N_\mu}$. The coefficients of such linear combination, or *interpolating vectors*, are denoted by $\{\zeta_\mu(\mathbf{r})\}_{\mu=1}^{N_\mu}$. Here N_μ can be interpreted as the numerical rank of the ISDF decomposition. Compared to the standard density fitting method, the three-tensor $\{\varphi_i(\hat{r}_\mu)\psi_j(\hat{r}_\mu)\}$ with three indices i, j, μ takes a separable form. This reduces the storage cost of the decomposed tensor from $O(N_e^3)$ to $O(N_e^2)$, and the computational cost from $O(N_e^4)$ to $O(N_e^3)$. Note that if the interpolation points $\{\hat{r}_\mu\}_{\mu=1}^{N_\mu}$ are chosen to be on a uniform grid, then the ISDF decomposition reduces to the pseudospectral decomposition, where $N_\mu \sim O(N_e)$ but with a large preconstant. For instance, the pseudospectral decomposition can be highly inefficient for molecular systems, where the grid points in the vacuum contribute nearly negligibly to the orbital pairs. On the other hand, by selecting the interpolation points carefully, e.g. through a randomized QR factorization with column pivoting (QRCP) procedure [64], the number of interpolation points can be significantly reduced. The QRCP based ISDF decomposition has been applied to accelerate a number of applications, at least in the context of pseudopotential approximation where the wavefunctions are smooth, including two-electron integral computation [113], correlation energy in the random phase approximation [112], density functional perturbation theory [109], and hybrid density functional calculations [84]. For example, when iterative solvers are used for hybrid density functional calculations, the Fock exchange operator V_X defined in terms of a set of orbitals $\{\varphi_i\}$ needs to be repeatedly applied to another set of Kohn-Sham orbitals $\{\varphi_j\}$

$$(V_X[\{\varphi_i\}]\psi_j)(\mathbf{r}) = - \sum_{i=1}^{N_e} \varphi_i(\mathbf{r}) \int K(\mathbf{r}, \mathbf{r}') \varphi_i(\mathbf{r}') \psi_j(\mathbf{r}') d\mathbf{r}'. \quad (6.2)$$

where $K(\mathbf{r}, \mathbf{r}')$ is the kernel for the Coulomb or the screened Coulomb operator. The integration in Eq. 6.2 is often carried out by solving Poisson-like equations, using e.g. a fast Fourier transform (FFT) method, and the computational cost is $O(N_e^3)$ with a large preconstant. This is typically the most time consuming component in hybrid functional calculations, and can be accelerated by the ISDF decomposition for the orbital pairs

$$\{\varphi_i \psi_j\}.$$

In Ref. [113], the interpolation points and the interpolation vectors are determined simultaneously through randomized QRCP applied to $\{\psi_i(\mathbf{r})\psi_j(\mathbf{r})\}$ directly. We recently found that the randomized QRCP procedure has $O(N_e^3)$ complexity but with a relatively large preconstant, and may not be competitive enough when used repeatedly. In order to overcome such difficulty, we proposed a different approach in Ref. [84] that determines the two parts separately and reduces the computational cost. We use the relatively expensive randomized QRCP procedure to find the interpolation points in advance, and only recompute the interpolation vectors whenever $\{\psi_i(\mathbf{r})\psi_j(\mathbf{r})\}$ has been updated using an efficient least squares procedure that exploits the separable nature of the matrix to be approximated. As a result, we can significantly accelerate hybrid functional calculations using the ISDF decomposition in all but the first SCF iteration.

In this work, we further remove the need of performing the QRCP decomposition completely and, hence, significantly reduce the computational cost. Note that an effective choice of the set of interpolation points should satisfy the following two conditions.

(1) The distribution of the interpolation points should roughly follow the distribution of the electron density. In particular, there should be more points when the electron density is high, and less or even zero points if the electron density is very low. (2) The interpolation points should not be very close to each other. Otherwise, matrix rows represented by the interpolation points are nearly linearly dependent, and the matrix formed by the interpolation vectors will be highly ill-conditioned. The QRCP procedure satisfies both (1) and (2) simultaneously, and thus is an effective way for selecting the interpolation points. Here we demonstrate that (1) and (2) can also be satisfied through a much simpler centroidal Voronoi tessellation (CVT) procedure applied to a weight vector such as the electron density.

The Voronoi tessellation technique has been widely used in computer science [11], and scientific and engineering applications such as image processing[50], pattern recognition [135], and numerical integration [19]. The concept of Voronoi tessellation can be simply understood as follows. Given a discrete set of weighted points, the CVT procedure divides a domain into a number of regions, each consisting of a collection of points that are closest to its weighted centroid. Here we choose the electron density as the weight, and the centroids as the interpolation points. The centroids must be located where the electron density is significant, and hence satisfy the requirement (1). The centroids are also mutually separated from each other by a finite distance due to the nearest neighbor principle, and hence satisfy the requirement (2). Although detailed analysis of the error stemming from such a choice of interpolation points is very difficult for general nonlinear functions, we find that the CVT procedure approximately minimizes the residual of the ISDF decomposition in Eq. 6.1. In practice, the CVT procedure only applies to one vector (the electron density) instead of $O(N_e^2)$ vectors and hence is very efficient.

We apply the ISDF-CVT method to accelerate hybrid functional calculations in a planewave basis set. We perform such calculations for different systems with insulating (liquid water), semiconducting (bulk silicon), and metallic (disordered silicon aluminum alloy) characters, as well as ab initio molecular dynamics (AIMD) simulations. We find that the ISDF-CVT method achieves similar accuracy to that obtained from the ISDF-QRCP method, with significantly improved efficiency. For instance, for a bulk silicon system containing 1000 silicon atoms computed on 2000 computational cores with kinetic energy cutoff being 10 Ha, the QRCP procedure finds the interpolation points with 38.1 s, while the CVT procedure only takes 0.7 s. Since the solution of the CVT procedure is continuous with respect to changes in the electron density, we also find that the CVT procedure produces a smoother potential energy surface than that by the QRCP

procedure in the context of AIMD simulations.

The remainder of the paper is organized as follows. We briefly introduce the ISDF decomposition in §6.3. In §6.4 we describe the ISDF-CVT procedure and its implementation for hybrid functional calculations. We present numerical results of the ISDF-CVT method in §6.5, and conclude in §6.6.

6.3 Interpolative Separable Density Fitting (ISDF) decomposition

In this section, we briefly introduce the ISDF decomposition [113] evaluated using the method developed in Ref. [84], which employs a separate treatment of the interpolation points and interpolation vectors.

First, assume the interpolation points $\{\hat{\mathbf{r}}_\mu\}_{\mu=1}^{N_\mu}$ are known, then the interpolation vectors can be efficiently evaluated using a least squares method as follows. Using a linear algebra notation, Eq. 6.1 can be written as

$$\mathbf{Z} \approx \Theta \mathbf{C}, \quad (6.3)$$

where each column of \mathbf{Z} is given by $Z_{ij}(\mathbf{r}) = \varphi_i(\mathbf{r})\psi_j(\mathbf{r})$ sampled on a dense real space grids $\{\mathbf{r}_i\}_{i=1}^{N_g}$, and $\Theta = [\zeta_1, \zeta_2, \dots, \zeta_{N_\mu}]$ contains the interpolating vectors. Each column of \mathbf{C} indexed by (i, j) is given by

$$[\varphi_i(\hat{\mathbf{r}}_1)\psi_j(\hat{\mathbf{r}}_1), \dots, \varphi_i(\hat{\mathbf{r}}_\mu)\psi_j(\hat{\mathbf{r}}_\mu), \dots, \varphi_i(\hat{\mathbf{r}}_{N_\mu})\psi_j(\hat{\mathbf{r}}_{N_\mu})]^T. \quad (6.4)$$

Eq. 6.3 is an overdetermined linear system with respect to the interpolation vectors Θ . The least squares approximation to the solution is given by

$$\Theta = \mathbf{Z} \mathbf{C}^T (\mathbf{C} \mathbf{C}^T)^{-1}. \quad (6.5)$$

It may appear that the matrix-matrix multiplications $\mathbf{Z}\mathbf{C}^T$ and $\mathbf{C}\mathbf{C}^T$ take $\mathcal{O}(N_e^4)$ operations because the size of \mathbf{Z} is $N_g \times N^2$ and the size of \mathbf{C} is $N_\mu \times N^2$. However, both multiplications can be carried out with fewer operations due to the separable structure of \mathbf{Z} and \mathbf{C} . The computational complexity for computing the interpolation vectors is $\mathcal{O}(N_e^3)$, and numerical results indicate that the preconstant is also much smaller than that involved in hybrid functional calculations [84]. Hence the interpolation vectors can be obtained efficiently using the least squares procedure.

The problem for finding a suitable set of interpolation points $\{\hat{\mathbf{r}}_\mu\}_{\mu=1}^{N_\mu}$ can be formulated as the following linear algebra problem. Consider the discretized matrix \mathbf{Z} of size $N_g \times N^2$, and find N_μ rows of \mathbf{Z} so that the rest of the rows of \mathbf{Z} can be approximated by the linear combination of the selected N_μ rows. This is called an interpolative decomposition [31], and a standard method to achieve such a decomposition is the QR factorization with column pivoting (QRCP) procedure [31] as

$$\mathbf{Z}^T \boldsymbol{\Pi} = \mathbf{Q} \mathbf{R}. \quad (6.6)$$

Here \mathbf{Z}^T is the transpose of \mathbf{Z} , \mathbf{Q} is an $N^2 \times N_g$ matrix that has orthonormal columns, \mathbf{R} is an upper triangular matrix, and $\boldsymbol{\Pi}$ is a permutation matrix chosen so that the magnitude of the diagonal elements of \mathbf{R} form a non-increasing sequence. The magnitude of each diagonal element of \mathbf{R} indicates how important the corresponding column of the permuted \mathbf{Z}^T is, and whether the corresponding grid point should be chosen as an interpolation point. The QRCP factorization can be terminated when the $(N_\mu + 1)$ -th diagonal element of \mathbf{R} becomes less than a predetermined threshold. The leading N_μ columns of the permuted \mathbf{Z}^T are considered to be linearly independent numerically. The corresponding grid points are chosen as the interpolation points. The indices for the chosen interpolation points $\{\hat{\mathbf{r}}_\mu\}$ can be obtained from indices of the nonzero entries of the first N_μ columns of the permutation matrix $\boldsymbol{\Pi}$.

The QRCP decomposition satisfies the requirements (1) and (2) discussed in §6.2. First, QRCP permutes matrix columns of \mathbf{Z}^T with large norms to the front, and pushes matrix columns of \mathbf{Z}^T with small norms to the back. Note that the square of the vector 2-norm of the column of \mathbf{Z}^T labeled by \mathbf{r} is just

$$\sum_{i,j=1}^N \varphi_i^2(\mathbf{r})\varphi_j^2(\mathbf{r}) = \left(\sum_{i=1}^N \varphi_i^2(\mathbf{r}) \right) \left(\sum_{j=1}^N \psi_j^2(\mathbf{r}) \right). \quad (6.7)$$

In the case when φ_i, ψ_j are the set of occupied orbitals, the norm of each column of \mathbf{Z}^T is simply the electron density. Hence the interpolation points chosen by QRCP will occur where the electron density is significant. Second, once a column is selected, all other columns are immediately orthogonalized with respect to the chosen column. Hence nearly linearly dependent matrix columns will not be selected repeatedly. As a result, the interpolation points chosen by QRCP are well separated spatially.

It turns out that the direct application of the QRCP procedure (Eq. 6.6) still requires $O(N_e^4)$ computational complexity. The key idea used in Ref. [113] to lower the cost is to randomly subsample columns of the matrix \mathbf{Z} to form a smaller matrix $\tilde{\mathbf{Z}}$ of size $N_g \times \tilde{N}_\mu$, where \tilde{N}_μ is only slightly larger than N_μ . Applying the QRCP procedure to this subsampled matrix $\tilde{\mathbf{Z}}$ approximately yields the choice of interpolation points, but the computational complexity is reduced to $O(N_e^3)$. In the context of hybrid density functional calculations, the cost of the randomized QRCP method can be comparable to that of applying the exchange operator in the planewave basis set [84]. However, the ISDF decomposition can still significantly reduce the computational cost, since the interpolation points only need to be performed once for a fixed geometric configuration.

6.4 Centroidal Voronoi Tessellation based ISDF decomposition

In this section, we demonstrate that the interpolation points can also be selected from a Voronoi tessellation procedure. For a d -dimensional space, the Voronoi tessellation partitions a set of points $\{\mathbf{r}_i\}_{i=1}^{N_g}$ in \mathbb{R}^d into a number of disjoint cells. The partition is based on the distance of each point to a finite set of points, called its generators. In our context, let $\{\hat{\mathbf{r}}_\mu\}_{\mu=1}^{N_\mu}$ denote such a set of generators, and the corresponding cell, C_μ , of a given generator $\hat{\mathbf{r}}_\mu$ is defined as a cluster of points

$$C_\mu = \{\mathbf{r}_i \mid \text{dist}(\mathbf{r}_i, \hat{\mathbf{r}}_\mu) < \text{dist}(\mathbf{r}_i, \hat{\mathbf{r}}_\nu) \text{ for all } \nu \neq \mu\}. \quad (6.8)$$

The distance can be chosen to be any metric, e.g. the L_2 distance as $\text{dist}(\mathbf{r}, \mathbf{r}') = \|\mathbf{r} - \mathbf{r}'\|_2$. In the case when the distances of a point \mathbf{r} to $\hat{\mathbf{r}}_\mu, \hat{\mathbf{r}}_\nu$ are exactly the same, we may arbitrarily assign \mathbf{r} to one of the clusters.

The Centroidal Voronoi tessellation (CVT) is a specific type of Voronoi tessellation in which the generator $\hat{\mathbf{r}}_\mu$ is chosen to be the centroid of its cell. Given a weight function $\rho(\mathbf{r})$ (such as the electron density), the centroid of a cluster C_μ is defined as

$$\mathbf{c}(C_\mu) = \frac{\sum_{\mathbf{r}_j \in C_\mu} \mathbf{r}_j \rho(\mathbf{r}_j)}{\sum_{\mathbf{r}_j \in C_\mu} \rho(\mathbf{r}_j)}. \quad (6.9)$$

Combined with the L_2 distance, CVT can be viewed as a minimization problem over both all possible partition of the cells and the centroids as [118]

$$\{C_\mu^*, \mathbf{c}_\mu^*\} = \arg \min_{\{C_\mu, \mathbf{c}_\mu\}} \sum_{\mu=1}^{N_\mu} \sum_{\mathbf{r}_k \in C_\mu} \rho(\mathbf{r}_k) \|\mathbf{r}_k - \mathbf{c}_\mu\|^2, \quad (6.10)$$

and the interpolation points are then chosen to be the minimizers $\hat{\mathbf{r}}_\mu = \mathbf{c}_\mu(C_\mu^*) = \mathbf{c}_\mu^*$. Following the discussion in §6.2, the electron density as the weight function (6.10) enforces that the interpolation points should locate at points where the electron density is significant and hence satisfies the requirement (1). Since the cells C_μ^* are disjoint, the

centroids c_μ^* are also separated by a finite distance away from each other and hence satisfies the requirement (2). Because the ISDF decomposition is a highly nonlinear process, in general we cannot expect the choice of interpolation points from CVT decomposition to maximally reduce the error of the decomposition. Instead, we demonstrate that the choice of the interpolation points from CVT approximately minimizes the residual for the ISDF decomposition, and hence provides a heuristic solution to the problem of finding interpolation points.

Theorem 5. *When the set of electron orbitals $\{\varphi_i\}$ are Lipschitz continuous, CVT method approximately minimizes the residual error of the ISDF decomposition.*

Proof. For simplicity we assume the limiting case where $\varphi_i = \psi_i$, and hence each row of Z is $Z(\mathbf{r}) = [\varphi_i(\mathbf{r})\varphi_j(\mathbf{r})]_{i,j=1}^{N_\mu}$.

Now suppose we cluster all matrix rows of Z into sub-collections $\{C_\mu\}_{\mu=1}^{N_\mu}$, and for each C_μ we choose a representative matrix row $Z(r_\mu)$. Then the error of the ISDF can be approximately characterized as

$$R = \sum_{\mu=1}^{N_\mu} \sum_{r_k \in C_\mu} \|Z(r_k) - \text{Proj}_{\text{Span}\{Z(r_\mu)\}} Z(r_k)\|^2, \quad (6.11)$$

where the projection is defined according to the L_2 inner product as

$$\text{Proj}_{\text{Span}\{Z(r_\mu)\}} Z(r_k) = \frac{Z(r_k) \cdot Z(r_\mu)}{Z(r_\mu) \cdot Z(r_\mu)} Z(r_\mu). \quad (6.12)$$

Let Φ be the $N_g \times N$ matrix with each row $\Phi(\mathbf{r}) = [\varphi_i(\mathbf{r})]_{i=1}^N$, then the electron density $\rho(\mathbf{r})$ is equal to $\Phi(\mathbf{r}) \cdot \Phi(\mathbf{r})$. Using the relation

$$Z(r_\mu) \cdot Z(r_\mu) = (\Phi(r_\mu) \cdot \Phi(r_\mu))^2 = \rho(r_\mu)^2, \quad (6.13)$$

we have

$$R = \sum_{\mu=1}^{N_\mu} \sum_{r_k \in C_\mu} \rho(r_k)^2 \left(1 - \frac{(\Phi(r_k) \cdot \Phi(r_\mu))^4}{\rho(r_k)^2 \rho(r_\mu)^2} \right) \quad (6.14)$$

$$= \sum_{\mu=1}^{N_\mu} \sum_{r_k \in C_\mu} \rho(r_k)^2 [1 - \cos^4(\theta(r_k, r_\mu))]. \quad (6.15)$$

Here $\theta(r_k, r_\mu)$ is the angle between the vectors $\Phi(r_k)$ and $\Phi(r_\mu)$. Since

$$\rho(r_k)[1 - \cos^4(\theta(r_k, r_\mu))] \leq 2\Phi(r_k) \cdot \Phi(r_k) \sin^2(\theta(r_k, r_\mu)) \quad (6.16)$$

$$\leq 2 \|\Phi(r_k) - \Phi(r_\mu)\|^2, \quad (6.17)$$

we have

$$R \leq 2 \sum_{\mu=1}^{N_\mu} \sum_{r_k \in C_\mu} \rho(r_k) \|\Phi(r_k) - \Phi(r_\mu)\|^2 \quad (6.18)$$

$$\approx 2 \sum_{\mu=1}^{N_\mu} \sum_{r_k \in C_\mu} \rho(r_k) \|\nabla_r \Phi(r_\mu)\|^2 \|r_k - r_\mu\|^2. \quad (6.19)$$

If we bound the gradient of $\Phi(r)$ by its Lipschitz constant, or simply neglect the spatial inhomogeneity in the electron orbitals, we arrive at the minimization criterion for the centroidal Voronoi tessellation decomposition. \square

Many algorithms have been developed to efficiently compute the Voronoi tessellation [124]. One most widely used method is the Lloyd's algorithm [111], which in discrete case is equivalent to the K-Means algorithm [118]. The K-Means algorithm is an iterative method that greedily minimizes the objective by taking alternating steps between $\{C_\mu\}$ and $\{\mathbf{c}_\mu\}$. In this work, we adopt a weighted version of the K-Means algorithm, which is demonstrated in Algorithm 9. Note that the K-Means algorithm can be straightforwardly parallelized. We distribute the grid points evenly at the beginning. The classification step is the most time consuming step, and can be locally computed for each group of grid points. After this step, the weighted sum and total weight of all clusters can be reduced from and broadcast to all processors for the next iteration.

Algorithm 9: Weighted K-Means Algorithm to Find Interpolation Points for Density Fitting

Input : Grid points $\{\mathbf{r}_i\}_{i=1}^{N_g}$, Weight function $\rho(\mathbf{r})$, Initial centroids $\{\mathbf{c}_\mu^{(0)}\}$

Output: Interpolation points $\{\hat{\mathbf{r}}_\mu\}_{\mu=1}^{N_\mu}$

Set $t \leftarrow 0$

do

Classification step:

for $i = 1$ **to** N_g

| Assign point \mathbf{r}_i to the cluster $C_\mu^{(t)}$ **if** $\mathbf{c}_\mu^{(t)}$ is the closest centroid to \mathbf{r}_i

end

Update step:

for $\mu = 1$ **to** N_μ

| $\mathbf{c}_\mu^{(t+1)} \leftarrow \sum_{\mathbf{r}_j \in C_\mu^{(t)}} \mathbf{r}_j \rho(\mathbf{r}_j) / \sum_{\mathbf{r}_j \in C_\mu^{(t)}} \rho(\mathbf{r}_j)$

end

Set $t \leftarrow t + 1$

while $\{\mathbf{c}_\mu^{(t)}\}$ not converged and maximum steps not reached

for $\mu = 1$ **to** N_μ

| **Set** $\hat{\mathbf{r}}_\mu \leftarrow \mathbf{c}_\mu^{(t)}$

end

In order to demonstrate the CVT procedure, we consider the weight function $\rho(\mathbf{r})$ given by the summation of four Gaussian functions in a 2D domain. The initial choice of centroids, given by 40 uniformly distributed random points, together with its associated Voronoi tessellation are plotted in Figure 6.1 (a). Figure 6.1 (b) demonstrates the converged centroids and the associated Voronoi tessellation using the weighted K-Means algorithm. We observe that the centroids concentrate on where the weight function is significant, and are well-separated.

We also show how the interpolation points are placed and moved in real chemical systems, i.e. the ammonia-borane (BH_3NH_3) decomposition reaction process. Figure 6.2 (a) shows the electron density of the molecule at the compressed, equilibrium, and dissociated configurations, respectively, according to the energy landscape in Fig. 6.2 (c). We plot the interpolation points found by the weighted K-Means algo-

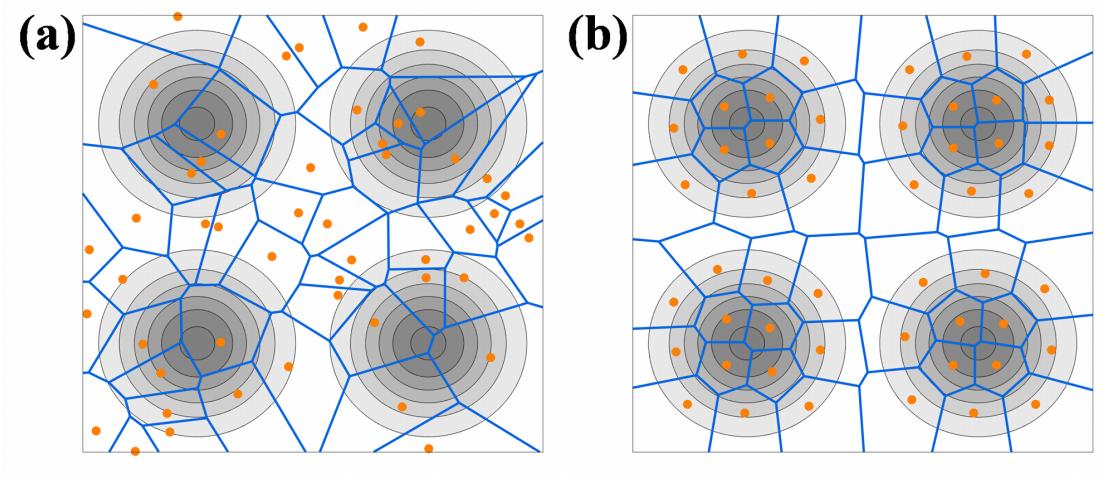


Figure 6.1: Schematic illustration of the CVT procedure in a 2D domain, including (a) initial random choice of centroids and Voronoi tessellation and centroidal Voronoi tessellation generated by the weighted K-Means algorithm. The weight function is given by the linear superposition of 4 Gaussian functions.

rithm in Fig. 6.2 (b). At the compressed configuration, all the interpolation points are distributed evenly around the molecule. As the bond length increases, some interpolation points are transferred from BH_3 to NH_3 . Finally at the dissociated configuration, NH_3 has more interpolation points around the molecule, since there are more electrons in NH_3 than BH_3 . Along the decomposition reaction process, both the transfer of the interpolation points and the potential energy landscape are smooth with respect to the change of the bond length.

6.5 Numerical results

We demonstrate the accuracy and efficiency of the ISDF-CVT method for hybrid functional calculations by using the DGDFT (Discontinuous Galerkin Density Functional Theory) software package [108, 82, 83, 15, 180]. DGDFT is a massively parallel electronic structure software package designed for large scale DFT calculations involving

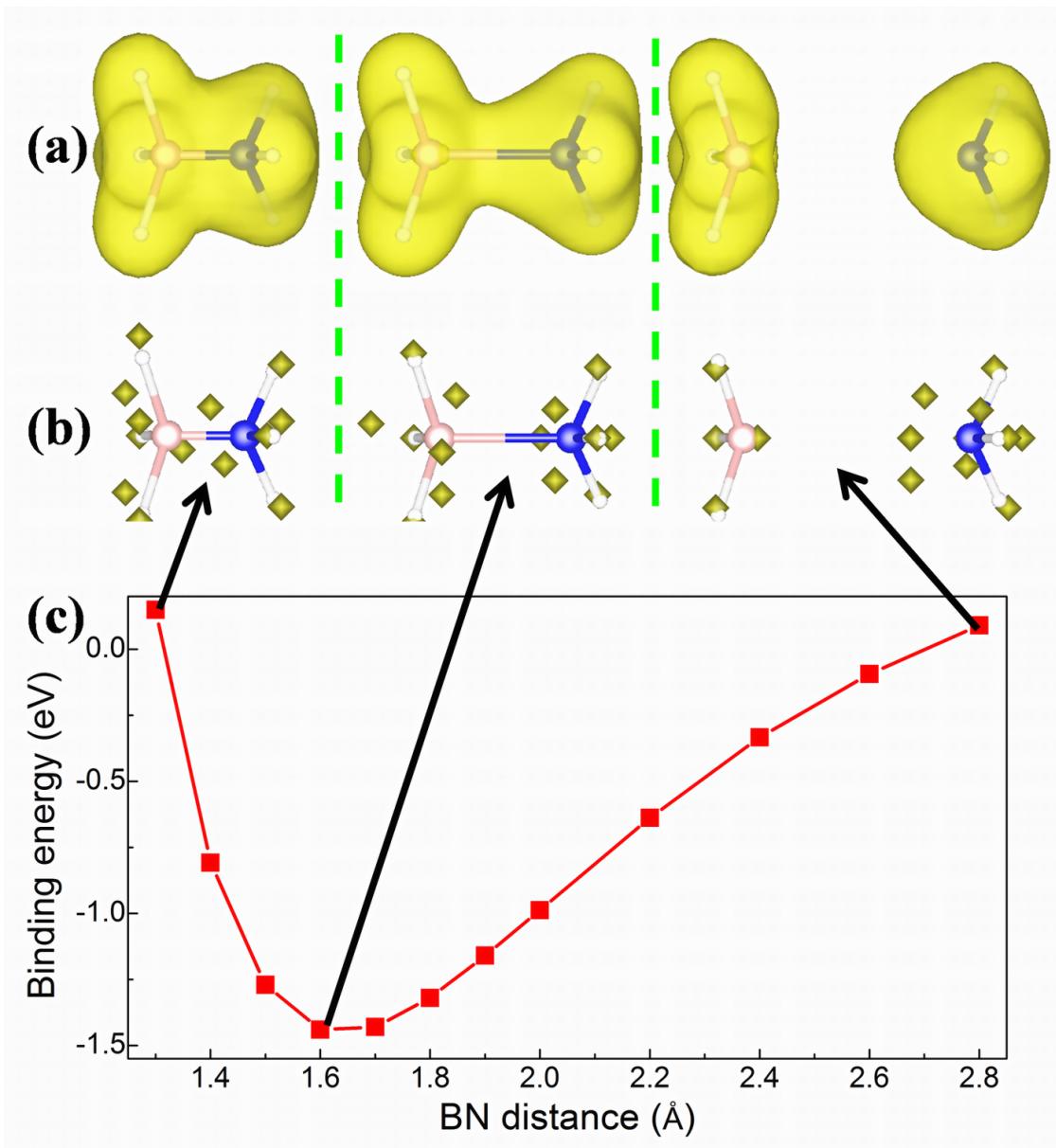


Figure 6.2: The decomposition reaction process of BH_3NH_3 computed with hybrid functional (HSE06) calculations by using the CVT procedure to select interpolation points, including (a) the electron density (yellow isosurfaces), (b) the interpolation points (yellow squares) $\{\hat{\mathbf{r}}_\mu\}_{\mu=1}^{N_\mu}$ ($N_\mu = 8$) selected from the real space grid points $\{\mathbf{r}_i\}_{i=1}^{N_g}$ ($N_g = 100^3$ and $E_{\text{cut}} = 60 \text{ Ha}$) when the BN distance respectively is 1.3, 1.7 and 2.8 \AA and (c) the binding energy as a function of BN distance for BH_3NH_3 in a $10 \text{ \AA} \times 10 \text{ \AA} \times 10 \text{ \AA}$ box. The white, pink and blue pink balls denote hydrogen, boron and nitrogen atoms, respectively.

up to tens of thousands of atoms. It includes a self-contained module called PWDFT for performing planewave based electronic structure calculations (mostly for benchmark and validation purposes). We implemented the ISDF-CVT method in PWDFT. We use the Message Passing Interface (MPI) to handle data communication. We use the Hartwigsen-Goedecker-Hutter (HGH) norm-conserving pseudopotential [74]. The atomic valence electron configuration is $1s^1$ for the H atom, $2s^22p^1$ for the B atom, $2s^22p^3$ for the N atom, $2s^22p^4$ for the O atom, $3s^23p^2$ for the Si atom in our DFT calculations, respectively. All calculations use the HSE06 functional [77], carried out on the Edison systems at the National Energy Research Scientific Computing Center (NERSC). Each node consists of two Intel “Ivy Bridge” processors with 24 cores in total and 64 gigabyte (GB) of memory. Our implementation only uses MPI. The number of cores is equal to the number of MPI ranks used in the simulation.

In this section, we demonstrate the performance of the ISDF-CVT method for accelerating hybrid functional calculations by using three types of systems [85]. They consist of bulk silicon systems (Si_{64} , Si_{216} , and Si_{1000}), a bulk water system with 64 molecules ($(H_2O)_{64}$), and a disordered silicon aluminum alloy system ($Al_{176}Si_{24}$). Bulk silicon systems (Si_{64} , Si_{216} and Si_{1000}) and bulk water system ($(H_2O)_{64}$) are semiconducting with a relatively large energy gap $E_{gap} > 1.0$ eV, and the $Al_{176}Si_{24}$ system is metallic with a small energy gap $E_{gap} < 0.1$ eV. All systems are closed shell systems, and the number of occupied bands is $N_{band} = N_e/2$, where N_e is the number of valence electrons. In order to compute the energy gap in the systems, we also include two unoccupied bands in all calculations.

6.5.1 Accuracy: Si_{216} and $\text{Al}_{176}\text{Si}_{24}$

We demonstrate the accuracy of the CVT-based ISDF decomposition in the hybrid functional calculation for semiconducting Si_{216} and metallic $\text{Al}_{176}\text{Si}_{24}$ systems, respectively. Although there is no general theoretical guarantee for the convergence of the K-Means algorithm and the convergence can depend sensitively on the initialization [9, 10], we find that, in the current context, initialization to have little impact on the final accuracy of the approximation. Hence we use random initialization for the K-Means algorithm. In all calculations, the adaptively compressed exchange (ACE) technique is used to accelerate hybrid functional calculations without loss of accuracy [110]. The results obtained in this work are labeled as ACE-ISDF (CVT), which are compared against those obtained from the previous work based on the QRCP decomposition [84] labeled as ACE-ISDF (QRCP). In both cases, we introduce a rank parameter c to control the trade-off between efficiency and accuracy, by setting the number of interpolation points $N_\mu = cN_e$. We measure the error using the valence band maximum (VBM) energy level, the conduction band minimum (CBM) energy level, the energy gap, the Hartree-Fock exchange energy, the total energy, and the atomic forces, respectively. We remark that, in ISDF-CVT and ISDF-QRCP, the atomic force is computed directly using the Hellmann-Feynman formula, thereby neglects the Pulay force contribution from the change of the interpolation points. On the other hand, there is no Pulay contribution in the ACE formulation, and the Hellmann-Feynman force F_I^{ACE} can be used as the reference solution.

The last three quantities are defined as

$$\begin{aligned}\Delta E_{\text{HF}} &= \left| E_{\text{HF}}^{\text{ACE-ISDF (CVT)}} - E_{\text{HF}}^{\text{ACE}} \right| / N_A \\ \Delta E &= \left| E^{\text{ACE-ISDF (CVT)}} - E^{\text{ACE}} \right| / N_A \\ \Delta F &= \max_I \| F_I^{\text{ACE-ISDF (CVT)}} - F_I^{\text{ACE}} \| \end{aligned}$$

where N_A is the number of atoms and I is the atom index.

Table 6.1 shows that the accuracy of the ACE-ISDF (CVT) method can systematically improve as the rank parameter c increases. When the rank parameter is large enough (≥ 20.0), the results from ACE-ISDF (CVT) are fully comparable (the energy error is below 10^{-6} Ha/atom and the force error is below 10^{-5} Ha/Bohr) to those obtained from the benchmark calculations. Furthermore, for a moderate choice of the rank parameter $c = 6.0$, the error of the energy per atom reaches below the chemical accuracy of 1 kcal/mol (1.6×10^{-3} Ha/atom), and the error of the force is around 10^{-3} Ha/Bohr. This is comparable to the accuracy obtained from ACE-ISDF (QRCP), and to e.g. linear scaling methods for insulating systems with reasonable amount of truncation needed to achieve significant speedup [45]. In fact, when compared with ACE-ISDF (QRCP) in Figure 6.3, we find that the CVT based ISDF decomposition achieves slightly higher accuracy, though there is no theoretical guarantee for this to hold in general. The last column of Table 6.1 shows the runtime of the K-Means algorithm. As c increases, the number of interpolation points as well as the number of cells increases proportionally. Hence we observe that the runtime of K-Means scales linearly with respect to c .

6.5.2 Efficiency: Si₁₀₀₀

We report the efficiency of the ISDF-CVT method by performing hybrid DFT calculations for a bulk silicon system with 1000 atoms ($N_{\text{band}} = 2000$) on 2000 computational cores as shown in Table 6.2, with respect to various choices of the kinetic energy cutoff (E_{cut}). With the number of interpolation points fixed at $N_\mu = 12000$, both QRCP and K-Means scales linearly with the number of grid points N_g . Yet the runtime of K-Means is around two orders of magnitude faster than QRCP. The determination of interpolation

Table 6.1: Accuracy of ACE-ISDF Based Hybrid Functional Calculations (HSE06) Obtained by Using the CVT method To Select Interpolation Points, with Varying Rank Parameter c for Semiconducting Si_{216} and Metallic $\text{Al}_{176}\text{Si}_{24}$ Systems^a.

c	E_{VBM}	E_{CBM}	E_{gap}	ΔE_{HF}	ΔE	ΔF	T_{KMEANS}
ACE-ISDF: Semiconducting Si_{216} ($N_{\text{band}} = 432$)							
4.0	6.7467	8.3433	-1.5967	2.69E-03	3.08E-03	5.04E-03	0.228
5.0	6.6852	8.2231	-1.5379	9.46E-04	1.12E-03	2.29E-03	0.248
6.0	6.6640	8.1522	-1.4882	3.76E-04	4.62E-04	1.05E-03	0.301
7.0	6.6550	8.1163	-1.4613	1.55E-04	1.98E-04	6.49E-04	0.312
8.0	6.6510	8.1030	-1.4520	7.33E-05	9.55E-05	3.07E-04	0.349
9.0	6.6490	8.0980	-1.4490	3.60E-05	4.96E-05	2.30E-04	0.398
10.0	6.6479	8.0959	-1.4480	1.78E-05	2.64E-05	1.30E-04	0.477
12.0	6.6472	8.0945	-1.4473	4.46E-06	8.91E-06	8.37E-05	0.530
16.0	6.6469	8.0937	-1.4468	1.51E-07	1.41E-06	3.20E-05	0.773
20.0	6.6468	8.0935	-1.4467	4.06E-07	3.33E-07	1.20E-05	0.830
24.0	6.6468	8.0935	-1.4467	2.99E-07	1.06E-07	5.18E-06	0.931
ACE	6.6468	8.0934	-1.4466	0.00E+00	0.00E+00	0.00E+00	-
ACE-ISDF: Metallic $\text{Al}_{176}\text{Si}_{24}$ ($N_{\text{band}} = 312$)							
4.0	7.9258	8.0335	-0.1076	3.80E-03	4.03E-03	8.01E-03	0.430
5.0	7.8537	7.9596	-0.1059	1.60E-03	1.69E-03	3.18E-03	0.535
6.0	7.8071	7.9127	-0.1056	6.07E-04	6.39E-04	1.48E-03	0.611
7.0	7.7843	7.8860	-0.1017	2.07E-04	2.17E-04	1.03E-03	0.731
8.0	7.7749	7.8749	-0.1000	7.43E-05	7.77E-05	4.40E-04	0.948
9.0	7.7718	7.8710	-0.0992	3.02E-05	3.20E-05	1.98E-04	0.947
10.0	7.7709	7.8697	-0.0989	1.48E-05	1.60E-05	1.80E-04	1.096
12.0	7.7703	7.8690	-0.0987	4.64E-06	5.60E-06	8.51E-05	1.305
16.0	7.7702	7.8688	-0.0986	6.35E-07	1.41E-06	3.24E-05	1.646
20.0	7.7701	7.8687	-0.0986	1.70E-08	5.30E-07	1.91E-05	2.037
ACE	7.7701	7.8687	-0.0986	0.00E+00	0.00E+00	0.00E+00	-

^a The unit for VBM (E_{VBM}), CBM (E_{CBM}) and the energy gap E_{gap} is eV. The unit for the error in the Hartree-Fock exchange energy ΔE_{HF} and the total energy ΔE is Ha/atom, and the unit for the error in atomic forces ΔF is Ha/Bohr. We use the results from the ACE-enabled hybrid functional calculations as the reference. The last column shows the time (in seconds) for K-Means with different c values, with 434 cores for Si_{216} and 314 cores for $\text{Al}_{176}\text{Si}_{24}$ on Edison.

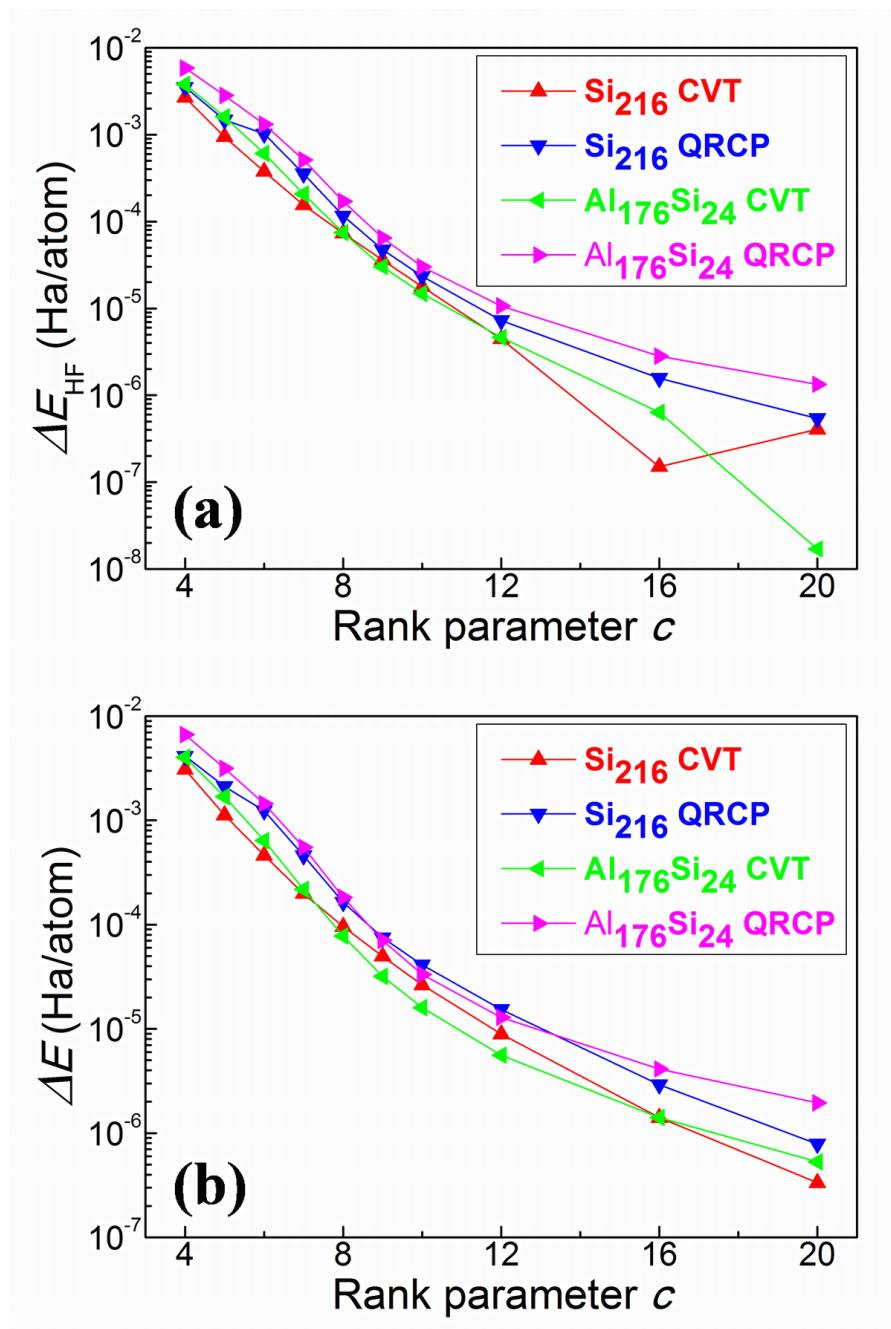


Figure 6.3: The accuracy of ACE-ISDF based hybrid functional calculations (HSE06) obtained by using the CVT and QRCP procedures to select the interpolation points, with varying rank parameter c from 4 to 20 for Si_{216} and $Al_{176}Si_{24}$, including the error of (a) Hartree-Fock exchange energy ΔE_{HF} (Ha/atom) and (b) total energy ΔE (Ha/atom).

vectors, which consists of solving a least-square problem, previously costs a fifth of the ISDF runtime but now becomes the dominating component in CVT-based ISDF decomposition. Notice that the ISDF method allows us to reduce the number of Poisson-like equations from $N_e^2 = 4 \times 10^6$ to $N_\mu = 12000$, which results in a significant speedup in terms of the cost of the FFT operations.

Table 6.2: Wall Clock Time (in seconds) Spent in the Components of the ACE-ISDF and ACE Enabled Hybrid DFT Calculations Related to the Exchange Operator, for Si_{1000} on 2002 Edison cores at Different E_{cut} Levels^a.

Si ₁₀₀₀		ACE-ISDF			ACE
E_{cut}	N_g	IP _{QRCP}	IP _{KMEANS}	IV (FFT)	FFT
10	74 ³	38.06	0.70	12.48 (0.33)	85.15
20	104 ³	126.39	1.24	36.48 (0.71)	143.54
30	128 ³	240.87	2.03	68.50 (1.43)	268.88
40	148 ³	434.16	3.26	108.18 (3.10)	783.27

^a Interpolation points are selected via either the QRCP or CVT procedure with the same rank parameter $c = 6.0$. N_g is the number of grid points in real space.

6.5.3 Ab Initio Molecular Dynamics: Si₆₄ and (H₂O)₆₄

In this section, we demonstrate the accuracy of the ACE-ISDF (CVT) method in the context of AIMD simulations for a bulk silicon system Si₆₄ under the NVE ensemble [59], and a liquid water system (H₂O)₆₄ under the NVT ensemble [59], respectively. For the Si₆₄ system, the initial MD structure (initial temperature $T = 300$ K) is optimized by hybrid DFT calculations, and we perform the simulation ($E_{\text{cut}} = 20$ Ha) for 1.0 ps with a MD time step of 1.0 fs. For the (H₂O)₆₄ system, we perform the simulation ($E_{\text{cut}} = 60$ Ha) for 2.0 ps with a MD time step of 0.5 fs to sample the radial distribution function after equilibrating the system starting from a prepared initial guess [47]. In this case, the Van der Waals (VdW) interaction is modeled at the level of the DFT-D2

method [68]. We use a single level Nose-Hoover thermostat [134, 79] at $T = 295$ K, and the choice of mass of the Nose-Hoover thermostat is 85000 au.

In the AIMD simulation, the interpolation points need to be recomputed for each atomic configuration. At the initial MD step, although the initialization strategy does not impact the accuracy of the physical observable, it can affect the convergence rate of the K-Means algorithm. We measure the convergence in terms of the fraction of points that switch clusters during two consecutive iterations. Figure 6.4 (a) shows the convergence of the K-Means algorithm with interpolation points initially chosen from a random distribution and from the QRCP solution, respectively. We find that the K-Means algorithm spends around half the number of iterations to wait for 0.1% of the points to settle on the respective clusters. However, these points often belong to the boundary of the clusters and have little effect on the positions of the centroids (interpolation points). Therefore, we decide to terminate K-Means algorithm whenever the fraction of points that switch clusters falls below the 0.1% threshold. It is evident that QRCP initialization leads to faster convergence than random sampling. However, in the AIMD simulation, a very good initial guess of the interpolation points can be simply obtained from those from the previous MD step. Figure 6.4 (b) shows that the number of K-Means iterations in the MD simulation can be very small, which demonstrates the effectiveness of this initialization strategy.

Figure 6.5 (a) shows that both the CVT-based and QRCP-based ISDF decomposition lead to controlled energy drift, defined as $E_{\text{drift}}(t) = (E_{\text{tot}}(t) - E_{\text{tot}}(0))/E_{\text{tot}}(0)$. In the NVE simulation on bulk silicon system Si_{64} , the energy drift per atom is 6.6×10^{-5} , 7.5×10^{-5} and 2.5×10^{-5} Ha/ps respectively for the ISDF-CVT, ISDF-QRCP, and the conventional nested two-level SCF iteration procedure, indicating that ISDF is a promising method for reducing the cost of hybrid functional calculations with controllable loss of accuracy.

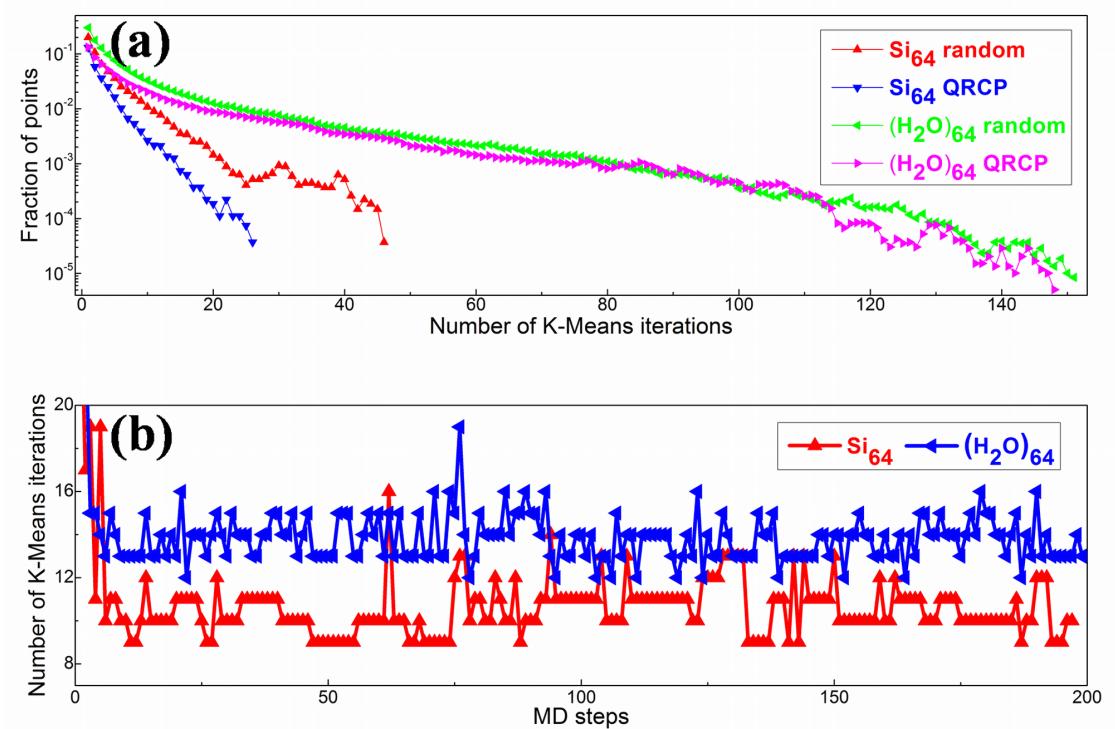


Figure 6.4: Comparison of the ISDF-CVT method by using either random or QRCP initialization for hybrid DFT AIMD simulations on bulk silicon system Si_{64} and liquid water system $(\text{H}_2\text{O})_{64}$, including (a) the fraction of points what switch cluster in each K-Means iteration and (b) the number of K-Means iterations during each MD step.

Figure 6.5 (b) shows the total potential energy obtained by the three methods along the MD trajectory, and the difference among the three methods is more noticeable. This is due to the fact that ISDF decomposition is a low rank decomposition for the pair product of orbitals, which leads to error in the Fock exchange energy and hence the total potential energy. Nonetheless, we find that such difference mainly results in a shift of the potential energy surface along the MD trajectory, and hence has little affect on physical observables defined via relative potential energy differences. Furthermore, the CVT method yields a potential energy trajectory that is much smoother compared to that obtained from QRCP. This is because the interpolation points obtained from CVT are driven by the electron density, which varies smoothly along the MD trajectory. Such

properties do not hold for the QRCP method. This means that the CVT method can be more effective when a smooth potential energy surface is desirable, such as in the case of geometry optimization. The absolute error of the potential energy from the CVT method is coincidentally smaller than that from QRCP, but again we are not aware of any reason for this behavior to hold in general.

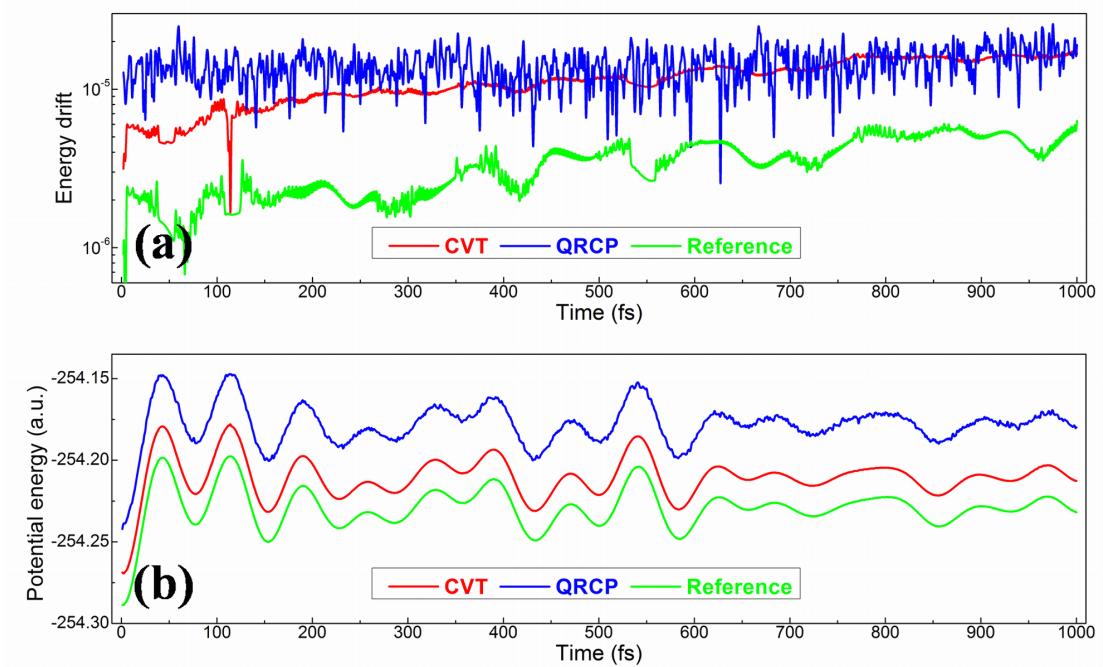


Figure 6.5: Comparison of hybrid HSE06 DFT AIMD simulations by using the ISDF-CVT and ISDF-QRCP methods as well as exact nested two-level SCF iteration procedure as the reference on the bulk silicon Si_{64} , including (a) relatively energy drift and (b) potential energy during MD steps.

We also apply the ACE-ISDF (CVT) and ACE-ISDF (QRCP) methods for hybrid DFT AIMD simulations on liquid water system $(\text{H}_2\text{O})_{64}$ under the NVT ensemble to sample the radial distribution function in Figure 6.6. We find that the results from all three methods agree very well, and our result is in quantitative agreement with previous hybrid functional calculations [47], which uses a different exchange-correlation functional (PBE0) and Van der Waals functional (TS-vdW) [161].

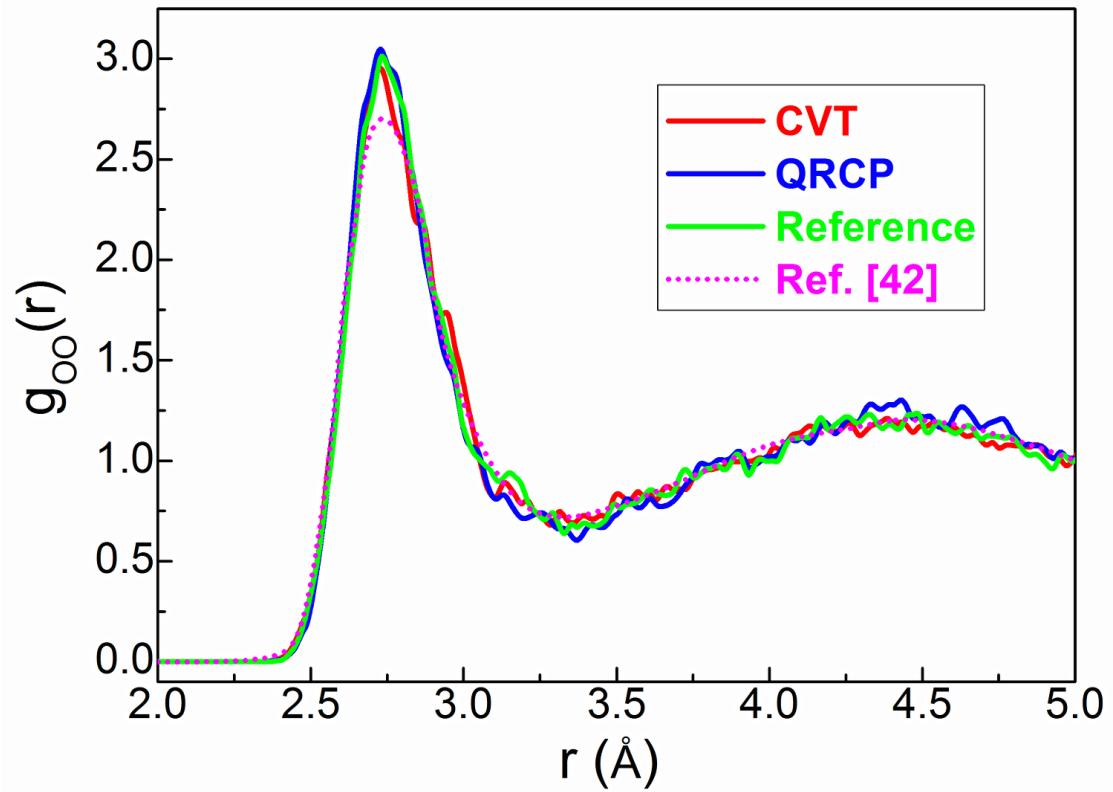


Figure 6.6: The oxygen-oxygen radial distribution functions $g_{OO}(r)$ of liquid water system (H_2O)₆₄ at $T = 295$ K obtained from hybrid HSE06 + DFT-D2 AIMD simulations with the ISDF-CVT and ISDF-QRCP methods, exact nested two-level SCF iteration procedure (as the reference) as well as previous hybrid PBE0 + TS-vdW calculation [47].

6.6 Conclusion

In this work, we demonstrate that the interpolative separable density fitting decomposition (ISDF) can be efficiently performed through a separated treatment of interpolation points and interpolation vectors. We find that the centroidal Voronoi tessellation method (CVT) provides an effective choice of interpolation points using only the electron density as the input information. The resulting interpolation points are by design inhomogeneous in the real space, concentrated at regions where the electron density is significant, and are well separated from each other. These are all key ingredients for obtaining a low

rank decomposition that is accurate and a well conditioned set of interpolation vectors. We demonstrate that the CVT-based ISDF decomposition can be an effective strategy for reducing the cost hybrid functional calculations for large systems. The CVT-based method achieves similar accuracy when compared with that obtained from QRCP, with significantly improved efficiency. Since the solution of the CVT method depends continuously with respect to the electron density, we also find that the CVT method produces a smoother potential energy surface than that by the QRCP method in the context of ab initio molecular dynamics simulation. Our analysis indicates that it might be possible to further improve the quality of the interpolation points by taking into account the gradient information in the weight vector. We also expect that the CVT-based strategy can also be useful in other contexts where the ISDF decomposition is applicable, such as ground state calculations with rung-5 exchange-correlation functionals, and excited state calculations. These will be explored in the future work.

6.7 Acknowledgments

This work was partly supported by the National Science Foundation under grant No. DMS-1652330, the DOE under grant No. DE-SC0017867, the DOE CAMERA project (L. L.), and by the DOE Scientific Discovery through Advanced Computing (SciDAC) program (K. D., W. H. and L. L.). The authors thank the National Energy Research Scientific Computing (NERSC) center and the Berkeley Research Computing (BRC) program at the University of California, Berkeley for making computational resources available. We thank Anil Damle and Robert Saye for useful discussions.

CHAPTER 7
CONCLUSION

BIBLIOGRAPHY

- [1] E A. Erosheva. Bayesian estimation of the grade of membership model. *Bayesian Stat.*, 7, 01 2003.
- [2] Edoardo M Airoldi, David M Blei, Stephen E Fienberg, and Eric P Xing. Mixed membership stochastic blockmodels. *Journal of machine learning research*, 9 (Sep):1981–2014, 2008.
- [3] A. Anandkumar, D. Hsu, and S. Kakade. A method of moments for mixture models and hidden markov models. In *COLT*, 2012.
- [4] Anima Anandkumar, Dean P. Foster, Daniel Hsu, Sham Kakade, and Yi-Kai Liu. A spectral algorithm for latent Dirichlet allocation. In *NIPS*, 2012.
- [5] Animashree Anandkumar, Daniel Hsu, and Sham M Kakade. A method of moments for mixture models and hidden markov models. In *Conference on Learning Theory*, pages 33–1, 2012.
- [6] F. Aquilante, T. B. Pedersen, and R. Lindh. Low-cost evaluation of the exchange Fock matrix from Cholesky and density fitting representations of the electron repulsion integrals. *J. Chem. Phys.*, 126:194106, 2007.
- [7] S. Arora, R. Ge, and A. Moitra. Learning topic models – going beyond SVD. In *FOCS*, 2012.
- [8] Sanjeev Arora, Rong Ge, Yonatan Halpern, David Mimno, Ankur Moitra, David Sontag, Yichen Wu, and Michael Zhu. A practical algorithm for topic modeling with provable guarantees. In *ICML*, 2013.
- [9] David Arthur and Sergei Vassilvitskii. How slow is the k-means method? In *Proceedings of the twenty-second annual symposium on Computational geometry*, pages 144–153. ACM, 2006.

- [10] David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035. Society for Industrial and Applied Mathematics, 2007.
- [11] Franz Aurenhammer. Voronoi diagrams—a survey of a fundamental geometric data structure. *ACM Computing Surveys (CSUR)*, 23(3):345–405, 1991.
- [12] Haim Avron and Sivan Toledo. Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix. *J. ACM*, 58(2):8:1–8:34, 2011. doi: 10.1145/1944345.1944349. URL <http://dx.doi.org/10.1145/1944345.1944349>.
- [13] Haim Avron and Sivan Toledo. Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix. *Journal of the ACM (JACM)*, 58(2):8, 2011.
- [14] Zhaojun Bai, Mark Fahey, Gene H Golub, M Menon, and E Richter. Computing partial eigenvalue sums in electronic structure calculations. Technical report, Tech. Report SCCM-98-03, Stanford University, 1998.
- [15] A. S. Banerjee, L. Lin, W. Hu, C. Yang, and J. E. Pask. Chebyshev polynomial filtered subspace iteration in the discontinuous galerkin method for large-scale electronic structure calculations. *J. Chem. Phys.*, 145:154101, 2016.
- [16] Anirban Banerjee. *The spectrum of the graph Laplacian as a tool for analyzing structure and evolution of networks*. PhD thesis, 2008.
- [17] Trupti Bansal, Chiranjib Bhattacharyya, and Ravindran Kannan. A provable SVD-based algorithm for learning topics in dominant admixture corpus. In *NIPS*, 2014.

- [18] Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- [19] Axel D Becke. A multicenter numerical integration scheme for polyatomic molecules. *The Journal of chemical physics*, 88(4):2547–2553, 1988.
- [20] N. H. F. Beebe and J. Linderberg. Simplifications in the generation and transformation of two-electron integrals in molecular calculations. *Int. J. Quantum Chem.*, 12:683, 1977.
- [21] Costas Bekas, Effrosyni Kokiopoulou, and Yousef Saad. An estimator for the diagonal of a matrix. *Applied Numerical Mathematics*, 57(11-12):1214–1229, 2007.
- [22] Costas Bekas, Efi Kokiopoulou, and Yousef Saad. An estimator for the diagonal of a matrix. *Applied Numerical Mathematics*, 57(11-12):1214–1229, November 2007. ISSN 0168-9274. doi: 10.1016/j.apnum.2007.01.003. URL <http://dx.doi.org/10.1016/j.apnum.2007.01.003>.
- [23] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, NIPS’01, pages 585–591, Cambridge, MA, USA, 2001. MIT Press.
- [24] Einat Neumann Ben-Ari and David M Steinberg. Modeling data from computer experiments: an empirical comparison of kriging with MARS and projection pursuit regression. *Quality Engineering*, 19(4):327–338, 2007.
- [25] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. *JMLR*, 2003.
- [26] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference:

- A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [27] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2):459–494, 2014.
- [28] Christos Boutsidis, Petros Drineas, Prabhanjan Kambadur, Eugenia-Maria Kon-topoulou, and Anastasios Zouzias. A randomized algorithm for approximating the log determinant of a symmetric positive definite matrix. *arXiv preprint arXiv:1503.00374*, 2015.
- [29] D. R. Bowler and T. Miyazaki. O(n) methods in electronic structure calculations. *Rep. Prog. Phys.*, 75:036503, 2012. doi: doi:10.1088/0034-4885/75/3/036503.
- [30] Martin Dietrich Buhmann. Radial basis functions. *Acta Numerica 2000*, 9:1–38, 2000.
- [31] T. F. Chan and P. C. Hansen. Some applications of the rank revealing QR factorization. *SIAM J. Sci. Statist. Comput.*, 13:727–741, 1992.
- [32] Jonathan Chang, Jordan Boyd-Graber, Chong Wang, Sean Gerrish, and David M. Blei. Reading tea leaves: How humans interpret topic models. In *NIPS*, 2009.
- [33] Isaac Chavel. *Eigenvalues in Riemannian geometry*, volume 115. Academic press, 1984.
- [34] Jeff Cheeger. A lower bound for the smallest eigenvalue of the laplacian. In *Proceedings of the Princeton conference in honor of Professor S. Bochner*, 1969.
- [35] Fan Chung and Linyuan Lu. Connected components in random graphs with given expected degree sequences. *Annals of combinatorics*, 6(2):125–145, 2002.

- [36] Fan Chung and Linyuan Lu. *Complex graphs and networks*. Number 107 in CBMS Regional Conference Series in Mathematics. American Mathematical Soc., 2006.
- [37] Fan RK Chung and Fan Chung Graham. *Spectral graph theory*. Number 92. American Mathematical Soc., 1997.
- [38] David Cohen-Steiner, Weihao Kong, Christian Sohler, and Gregory Valiant. Approximating the spectrum of a graph. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1263–1271. ACM, 2018.
- [39] Puget Sound LiDAR Consortium. Mount Saint Helens LiDAR data. University of Washington, 2002.
- [40] Paul G. Constantine. *Active subspaces: Emerging ideas for dimension reduction in parameter studies*. SIAM, 2015.
- [41] Jane K Cullum and Ralph A Willoughby. *Lanczos algorithms for large symmetric eigenvalue computations: Vol. I: Theory*. SIAM, 2002.
- [42] Kurt Cutajar, Michael Osborne, John Cunningham, and Maurizio Filippone. Preconditioning kernel matrices. pages 2529–2538, 2016.
- [43] D. M. Cvetković, M. Doob, and H. Sachs. *Spectra of Graphs: Theory and Applications*. Wiley, third edition, 1998.
- [44] Dragoš Cvetkovic, Slobodan Simic, and Peter Rowlinson. *An introduction to the theory of graph spectra*. Cambridge University Press, 2009.
- [45] William Dawson and François Gygi. Performance and accuracy of recursive sub-

- space bisection for hybrid dft calculations in inhomogeneous systems. *J. Chem. Theory Comput.*, 11(10):4655–4663, 2013. doi: 10.1021/acs.jctc.5b00826.
- [46] Weicong Ding, Prakash Ishwar, and Venkatesh Saligrama. Most large topic models are approximately separable. In *ITA, 2015*, pages 199–203. IEEE, 2015.
- [47] R. A. DiStasio, B. Santra, Z. Li, X. Wu, and R. Car. The individual and collective effects of exact exchange and dispersion interactions on the ab initio structure of liquid water. *J. Chem. Phys.*, 141:084502, 2014.
- [48] William E Donath and Alan J Hoffman. Lower bounds for the partitioning of graphs. In *Selected Papers Of Alan J Hoffman: With Commentary*, pages 437–442. World Scientific, 2003.
- [49] Kun Dong, David Eriksson, Hannes Nickisch, David Bindel, and Andrew G. Wilson. Scalable log determinants for Gaussian process kernel learning. pages 6330–6340, 2017.
- [50] Qiang Du, Max Gunzburger, Lili Ju, and Xiaoqiang Wang. Centroidal voronoi tessellation algorithms for image compression, segmentation, and multichannel restoration. *Journal of Mathematical Imaging and Vision*, 24(2):177–194, 2006.
- [51] Francois Ducastelle and Françoise Cyrot-Lackmann. Moments developments and their application to the electronic charge distribution of d bands. *Journal of Physics and Chemistry of Solids*, 31(6):1295–1306, 1970.
- [52] Nicole Eikmeier and David F Gleich. Revisiting power-law distributions in spectra of real world networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 817–826, 2017.

- [53] Illés J Farkas, Imre Derényi, Albert-László Barabási, and Tamas Vicsek. Spectra of “real-world” graphs: Beyond the semicircle law. *Physical Review E*, 64(2):026704, 2001.
- [54] Gregory E Fasshauer. *Meshfree approximation methods with MATLAB*, volume 6. World Scientific, 2007.
- [55] M. Fiedler. Hankel and Loewner matrices. *Linear Algebra and Its Applications*, 58:75–95, 1984.
- [56] Seth Flaxman, Andrew Wilson, Daniel Neill, Hannes Nickisch, and Alex Smola. Fast kronecker inference in gaussian processes with non-gaussian likelihoods. In *International Conference on Machine Learning*, pages 607–616, 2015.
- [57] Alexander Forrester, Andy Keane, et al. *Engineering design via surrogate modelling: a practical guide*. John Wiley & Sons, 2008.
- [58] Jacob R. Gardner, Geoff Pleiss, Ruihan Wu, Kilian Q. Weinberger, and Andrew G. Wilson. Product kernel interpolation for scalable Gaussian processes. *arXiv preprint arXiv:1802.08903*, 2018.
- [59] J Willard Gibbs. *Elementary principles in statistical mechanics*. Courier Corporation, 2014.
- [60] Amparo Gil, Javier Segura, and Nico Temme. *Numerical Methods for Special Functions*. SIAM, 2007.
- [61] Nicolas Gillis and Stephen A Vavasis. Fast and robust recursive algorithms for separable nonnegative matrix factorization. *IEEE transactions on pattern analysis and machine intelligence*, 36(4):698–714, 2014.

- [62] David Gingras, Tom Lamarche, Jean-Luc Bedwani, and Érick Dupuis. Rough terrain reconstruction for rover motion planning. In *Proceedings of the Canadian Conference on Computer and Robot Vision (CRV)*, pages 191–198. IEEE, 2010.
- [63] S. Goedecker. Linear scaling electronic structure methods. *Rev. Mod. Phys.*, 71: 1085–1123, 1999.
- [64] G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins Univ. Press, Baltimore, fourth edition, 2013.
- [65] Gene Golub and Gérard Meurant. *Matrices, Moments and Quadrature with Applications*. Princeton University Press, 2010.
- [66] Gene H Golub and Gérard Meurant. Matrices, moments and quadrature ii; how to compute the norm of the error in iterative methods. *BIT Numerical Mathematics*, 37(3):687–705, 1997.
- [67] Carolyn Gordon, David L. Webb, and Scott Wolpert. One cannot hear the shape of a drum. *Bull. Amer. Math. Soc.*, 27:134–138, 1992.
- [68] S. Grimme. Semiempirical GGA-type density functional constructed with a long-range dispersion correction. *J. Comput. Chem.*, 27(15):1787–1799, 2006. doi: 10.1002/jcc.20495.
- [69] M. Guidon, J. Hutter, and J. Vandevondele. Auxiliary density matrix methods for Hartree-Fock exchange calculations. *J. Chem. Theory Comput.*, 6:2348–2364, 2010.
- [70] Raia Hadsell, J. Andrew Bagnell, Daniel F. Huber, and Martial Hebert. Space-carving kernels for accurate rough terrain estimation. *International Journal of Robotics Research*, 29:981–996, July 2010.

- [71] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011.
- [72] Insu Han, Dmitry Malioutov, and Jinwoo Shin. Large-scale log-determinant computation through stochastic Chebyshev expansions. In *ICML*, pages 908–917, 2015.
- [73] Helmut Harbrecht, Michael Peters, and Reinhold Schneider. On the low-rank approximation by the pivoted Cholesky decomposition. *Applied Numerical Mathematics*, 62(4):428–440, 2012.
- [74] C. Hartwigsen, S. Goedecker, and J. Hutter. Relativistic separable dual-space gaussian pseudopotentials from H to Rn. *Phys. Rev. B*, 58:3641, 1998. doi: 10.1103/PhysRevB.58.3641.
- [75] J Hensman, N Fusi, and N.D. Lawrence. Gaussian processes for big data. In *Uncertainty in Artificial Intelligence (UAI)*. AUAI Press, 2013.
- [76] William Herlands, Andrew Wilson, Hannes Nickisch, Seth Flaxman, Daniel Neill, Wilbert Van Panhuis, and Eric Xing. Scalable Gaussian processes for characterizing multidimensional change surfaces. *Artificial Intelligence and Statistics*, 2016.
- [77] J. Heyd, G. E. Scuseria, and M. Ernzerhof. Erratum: "hybrid functionals based on a screened coulomb potential" [J. Chem. Phys. 118, 8207 (2003)]. *J. Chem. Phys.*, 124:219906, 2006. doi: 10.1063/1.2204597.
- [78] Nicholas J Higham. *Functions of matrices: theory and computation*. SIAM, 2008.

- [79] W. G. Hoover. Canonical dynamics: Equilibrium phase-space distributions. *Phys. Rev. A*, 31:1695, 1985. doi: 10.1103/PhysRevA.31.1695.
- [80] Patrik O. Hoyer. Non-negative matrix factorization with sparseness constraints. *JMLR*, 2004.
- [81] Daniel Hsu, Sham M Kakade, and Tong Zhang. A spectral algorithm for learning hidden markov models. *Journal of Computer and System Sciences*, 78(5):1460–1480, 2012.
- [82] W. Hu, L. Lin, and C. Yang. DGDFT: A massively parallel method for large scale density functional theory calculations. *J. Chem. Phys.*, 143(12):124110, 2015. doi: 10.1063/1.4931732.
- [83] W. Hu, L. Lin, and C. Yang. Edge reconstruction in armchair phosphorene nanoribbons revealed by discontinuous galerkin density functional theory. *Phys. Chem. Chem. Phys.*, 17(47):31397–31404, 2015. doi: 10.1039/C5CP00333D.
- [84] Wei Hu, Lin Lin, and Chao Yang. Interpolative separable density fitting decomposition for accelerating hybrid density functional calculations with applications to defects in silicon. *J. Chem. Theory Comput.*, accepted, 2017. doi: 10.1021/acs.jctc.7b00807.
- [85] Wei Hu, Lin Lin, and Chao Yang. Projected commutator diis method for accelerating hybrid functional electronic structure calculations. *J. Chem. Theory Comput.*, accepted, 2017. doi: 10.1021/acs.jctc.7b00892.
- [86] Kejun Huang, Xiao Fu, and Nikolaos D. Sidiropoulos. Anchor-free correlated topic modeling: Identifiability and algorithm. In *NIPS*, 2016.
- [87] B. Huffaker, M. Fomenkov, and k. claffy. Internet Topology Data Comparison.

Technical report, Cooperative Association for Internet Data Analysis (CAIDA), May 2012.

- [88] Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 19(2):433–450, 1990.
- [89] Dunham Jackson. *Über die Genauigkeit der Annäherung stetiger Funktionen durch ganze rationale Funktionen gegebenen Grades und trigonometrische Summen gegebener Ordnung*. Dieterich’schen Universität Buchdruckerei, 1911.
- [90] John David Jackson. *Mathematics for Quantum Mechanics: An Introductory Survey of Operators, Eigenvalues, and Linear Vector Spaces*. Dover Publications, 2006.
- [91] M. Jordan, Z. Ghahramani, T. Jaakola, and L. Saul. Introduction to variational methods for graphical models. *Machine Learning*, pages 183–233, 1999.
- [92] Mark Kac. Can one hear the shape of a drum? *The American Mathematical Monthly*, 73(4):1–23, 1966.
- [93] Leonid Vasilevich Kantorovich and Gennady S Rubinstein. On a space of completely additive functions. *Vestnik Leningrad. Univ*, 13(7):52–59, 1958.
- [94] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. on Scientific Computing*, 20(1), 1998.
- [95] Robert Keys. Cubic convolution interpolation for digital image processing. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 29(6):1153–1160, 1981.

- [96] H. Koch, A. Sánchez de Merás, and T. B. Pedersen. Reduced scaling in electronic structure calculations using cholesky decompositions. *J. Chem. Phys.*, 118:9481–9484, 2003.
- [97] Kurt Konolige, Motilal Agrawal, and Joan Sola. Large-scale visual odometry for rough terrain. In *Robotics Research*, pages 201–212. Springer, 2010.
- [98] Alex Kulesza, Ben Taskar, et al. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2–3):123–286, 2012.
- [99] Alex Kulesza, N Raj Rao, and Satinder Singh. Low-rank spectral learning. In *Artificial Intelligence and Statistics*, pages 522–530, 2014.
- [100] Thomas K Landauer, Peter W Foltz, and Darrell Laham. An introduction to latent semantic analysis. *Discourse processes*, 25(2-3):259–284, 1998.
- [101] Q. Le, T. Sarlos, and A. Smola. Fastfood-computing Hilbert space expansions in loglinear time. In *Proceedings of the 30th International Conference on Machine Learning*, pages 244–252, 2013.
- [102] Moontae Lee, David Bindel, and David Mimno. Robust spectral inference for joint stochastic matrix factorization. In *NIPS*, 2015.
- [103] Moontae Lee, David Bindel, and David Mimno. From correlation to hierarchy: Practical topic modeling via spectral inference. In *12th INFORMS Workshop on Data Mining and Decision Analytics*, 2017.
- [104] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- [105] Bruno Lévy. Laplace-Beltrami eigenfunctions towards an algorithm that “under-

- stands” geometry. In *Shape Modeling and Applications, 2006. SMI 2006. IEEE International Conference on*, pages 13–13. IEEE, 2006.
- [106] Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *NIPS*, 2014.
- [107] Dawen Liang, Jaan Altosaar, Laurent Charlin, and David M Blei. Factorization meets the item embedding: Regularizing matrix factorization with item co-occurrence. In *Proceedings of the 10th ACM conference on recommender systems*, pages 59–66. ACM, 2016.
- [108] L. Lin, J. Lu, L. Ying, and W. E. Adaptive local basis set for kohn-sham density functional theory in a discontinuous galerkin framework i: Total energy calculation. *J. Comput. Phys.*, 231(4):2140–2154, 2012. doi: 10.1016/j.jcp.2011.11.032.
- [109] L. Lin, Z. Xu, and L. Ying. Adaptively compressed polarizability operator for accelerating large scale ab initio phonon calculations. *Multiscale Model. Simul.*, 15:29–55, 2017.
- [110] Lin Lin. Adaptively compressed exchange operator. *J. Chem. Theory Comput.*, 12(5):2242–2249, 2016. doi: 10.1021/acs.jctc.6b00092.
- [111] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982.
- [112] J. Lu and K. Thicke. Cubic scaling algorithms for RPA correlation using interpolative separable density fitting. *J. Comput. Phys.*, 351:187 – 202, 2017.
- [113] J. Lu and L. Ying. Compression of the electron repulsion integral tensor in tensor hypercontraction format with cubic scaling cost. *J. Comput. Phys.*, 302:329–335, 2015. doi: 10.1016/j.jcp.2015.09.014.

- [114] Ives Macedo, Joao Paulo Gois, and Luiz Velho. Hermite radial basis functions implicits. *Computer Graphics Forum*, 30(1):27–42, 2011.
- [115] D MacKay and MN Gibbs. Efficient implementation of gaussian processes. *Neural Computation*, 1997.
- [116] David JC MacKay. *Bayesian methods for adaptive models*. PhD thesis, California Institute of Technology, 1992.
- [117] David JC MacKay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.
- [118] J. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proc. of the Fifth Berkeley Symp. On Math. Stat. and Prob.*, volume 1, pages 281–297, 1967.
- [119] Jan R Magnus and Heinz Neudecker. Matrix differential calculus with applications in statistics and econometrics. *Wiley series in probability and mathematical statistics*, 1988.
- [120] S. Manzer, P. R. Horn, N. Mardirossian, and M. Head-Gordon. Fast, accurate evaluation of exact exchange: The occ-RI-K algorithm. *J. Chem. Phys.*, 143(2):024113, 2015.
- [121] R. Martin. *Electronic Structure – Basic Theory and Practical Methods*. Cambridge Univ. Pr., West Nyack, NY, 2004.
- [122] H. P. McKean. Selberg’s trace formula as applied to a compact riemann surface. *Communications on Pure and Applied Mathematics*, 25(3):225–246, 1972.
- [123] Frank McSherry. Spectral partitioning of random graphs. In *focs*, page 529. IEEE, 2001.

- [124] NN Medvedev. The algorithm for three-dimensional voronoi polyhedra. *Journal of computational physics*, 67(1):223–229, 1986.
- [125] Erik H. W. Meijering, Karel J. Zuiderveld, and Max A. Viergever. Image reconstruction by convolution with symmetrical piecewise n th-order polynomial kernels. *IEEE Transactions on Image Processing*, 8(2):192–201, 1999.
- [126] M. Mihail. Conductance and convergence of markov chains-a combinatorial treatment of expanders. In *30th Annual Symposium on Foundations of Computer Science*. IEEE, 1989. doi: 10.1109/sfcs.1989.63529. URL <https://doi.org/10.1109/sfcs.1989.63529>.
- [127] Bojan Mohar. Isoperimetric numbers of graphs. *Journal of Combinatorial Theory, Series B*, 47(3):274–291, 1989.
- [128] Ravi Montenegro, Prasad Tetali, et al. Mathematical aspects of mixing times in markov chains. *Foundations and Trends® in Theoretical Computer Science*, 1(3):237–354, 2006.
- [129] R. B. Murphy, M. D. Beachy, R. A. Friesner, and M. N. Ringnalda. Pseudospectral localized møller–plesset methods: Theory and calculation of conformational energies. *J. Chem. Phys.*, 103:1481, 1995.
- [130] Radford M Neal. *Probabilistic inference using Markov chain Monte Carlo methods*, 1993.
- [131] Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*, 2(11):2, 2011.
- [132] Frank Neese, Frank Wennmohs, Andreas Hansen, and Ute Becker. Efficient, approximate and parallel hartree–fock and hybrid dft calculations. a “chain-of-

- spheres” algorithm for the hartree–fock exchange. *Chem. Phys.*, 356:98–109, 2009.
- [133] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*, pages 849–856, 2002.
- [134] S. Nosé. A unified formulation of the constant temperature molecular dynamics methods. *J. Chem. Phys.*, 81:511, 1984. doi: 10.1063/1.447334.
- [135] Robert L Ogniewicz and Olaf Kübler. Hierarchic voronoi skeletons. *Pattern recognition*, 28(3):343–359, 1995.
- [136] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The PageRank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [137] B. N. Parlett. The software scene in the extraction of eigenvalues from sparse matrices. *SIAM Journal on Scientific and Statistical Computing*, 5(3):590–604, sep 1984. doi: 10.1137/0905042. URL <https://doi.org/10.1137/0905042>.
- [138] R. M. Parrish, E. G. Hohenstein, T. J. Martínez, and C. D. Sherrill. Tensor hypercontraction. II. Least-squares renormalization. *J. Chem. Phys.*, 137:224106, 2012.
- [139] R. M. Parrish, E. G. Hohenstein, T. J. Martínez, and C. D. Sherrill. Discrete variable representation in electronic structure theory: Quadrature grids for least-squares tensor hypercontraction. *J. Chem. Phys.*, 138:194107, 2013.
- [140] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global vectors for word representation. In *EMNLP*, 2014.

- [141] Alex Pothen, Horst D. Simon, and Kan-Pu Liou. Partitioning sparse matrices with eigenvectors of graphs. *SIAM J. Matrix Anal. Appl.*, 11(3):430–452, 1990.
- [142] Jonathan K Pritchard, Matthew Stephens, and Peter Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959, 2000.
- [143] Joaquin Quiñonero-Candela and Carl Edward Rasmussen. A unifying view of sparse approximate gaussian process regression. *Journal of Machine Learning Research*, 6(Dec):1939–1959, 2005.
- [144] Joaquin Quinonero-Candela, Carl Edward Rasmussen, and Christopher KI Williams. Approximation methods for Gaussian process regression. *Large-scale kernel machines*, pages 203–223, 2007.
- [145] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for Machine Learning*. The MIT Press, 2006.
- [146] Carl Edward Rasmussen and Zoubin Ghahramani. Occam’s razor. In *Neural Information Processing Systems (NIPS)*, 2001.
- [147] Carl Edward Rasmussen and Hannes Nickisch. Gaussian processes for machine learning (GPML) toolbox. *Journal of Machine Learning Research (JMLR)*, 11: 3011–3015, Nov 2010.
- [148] X. Ren, P. Rinke, V. Blum, J. Wieferink, A. Tkatchenko, A. Sanfilippo, K. Reuter, and M. Scheffler. Resolution-of-identity approach to Hartree-Fock, hybrid density functionals, RPA, MP2 and GW with numeric atom-centered orbital basis functions. *New J. Phys.*, 14:053020, 2012.
- [149] G. Reynolds, T. J. Martinez, and E. A. Carter. Local weak pairs spectral and pseudospectral singles and doubles configuration interaction. *J. Chem. Phys.*, 105:6455, 1996.

- [150] Christian Robert and George Casella. *Monte Carlo statistical methods*. Springer Science & Business Media, 2013.
- [151] Havard Rue and Leonhard Held. *Gaussian Markov random fields: theory and applications*. CRC Press, 2005.
- [152] Youcef Saad. *Numerical methods for large eigenvalue problems*. Manchester University Press, 1992.
- [153] Robert Schaback and Holger Wendland. Kernel techniques: from machine learning to meshless methods. *Acta Numerica*, 15:543–639, 2006.
- [154] Comandur Seshadhri, Tamara G Kolda, and Ali Pinar. Community structure and scale-free collections of erdős-rényi graphs. *Physical Review E*, 85(5):056109, 2012.
- [155] Bernhard W Silverman. Some aspects of the spline smoothing approach to non-parametric regression curve fitting. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–52, 1985.
- [156] Alistair Sinclair and Mark Jerrum. Approximate counting, uniform generation and rapidly mixing markov chains. *Information and Computation*, 82(1):93–133, jul 1989. doi: 10.1016/0890-5401(89)90067-9. URL [https://doi.org/10.1016/0890-5401\(89\)90067-9](https://doi.org/10.1016/0890-5401(89)90067-9).
- [157] Edward Snelson and Zoubin Ghahramani. Sparse Gaussian processes using pseudo-inputs. In *Advances in neural information processing systems (NIPS)*, volume 18, page 1257. MIT Press, 2006.
- [158] Michael L Stein, Jie Chen, Mihai Anitescu, et al. Stochastic approximation of score functions for gaussian processes. *The Annals of Applied Statistics*, 7(2):1162–1191, 2013.

- [159] S. Surjanovic and D. Bingham. Virtual library of simulation experiments: Test functions and datasets. <http://www.sfu.ca/ ssurjano>, 2018.
- [160] A. Szabo and N.S. Ostlund. *Modern Quantum Chemistry: Introduction to Advanced Electronic Structure Theory*. McGraw-Hill, New York, 1989.
- [161] A. Tkatchenko and M. Scheffler. Accurate molecular van der waals interactions from ground-state electron density and free-atom reference data. *Phys. Rev. Lett.*, 102:073005, 2009.
- [162] Lloyd N Trefethen. *Approximation theory and approximation practice*, volume 128. Siam, 2013.
- [163] Luca Trevisan. Max cut and the smallest eigenvalue. *SIAM Journal on Computing*, 41(6):1769–1786, 2012.
- [164] Shashanka Ubaru, Jie Chen, and Yousef Saad. Fast estimation of $\text{tr}(F(A))$ via stochastic Lanczos quadrature.
- [165] M. Wainwright and M. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, pages 1–305, 2008.
- [166] Ziyu Wang, Masrour Zoghi, Frank Hutter, David Matheson, Nando De Freitas, et al. Bayesian optimization in high dimensions via random embeddings. In *Proceedings of the International Joint Conferences on Artificial Intelligence*, pages 1778–1784, 2013.
- [167] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*, 393(6684):440, 1998.

- [168] F. Weigend. A fully direct RI-HF algorithm: Implementation, optimised auxiliary basis sets, demonstration of accuracy and efficiency. *Phys. Chem. Chem. Phys.*, 4:4285–4291, 2002.
- [169] Alexander Weiße, Gerhard Wellein, Andreas Alvermann, and Holger Fehske. The kernel polynomial method. *Reviews of modern physics*, 78(1), 2006.
- [170] Holger Wendland. *Scattered data approximation*, volume 17. Cambridge university press, 2004.
- [171] Hermann Weyl. Über die asymptotische verteilung der eigenwerte. *Nachrichten von der Gesellschaft der Wissenschaften zu Göttingen, Mathematisch-Physikalische Klasse*, 1911:110–117, 1911.
- [172] Hermann Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). *Mathematische Annalen*, 71(4):441–479, 1912.
- [173] Eugene P. Wigner. On the distribution of the roots of certain symmetric matrices. *Annals of Mathematics*, pages 325–327, 1958.
- [174] Andrew G Wilson, Zhiting Hu, Ruslan R Salakhutdinov, and Eric P Xing. Stochastic variational deep kernel learning. In *Advances in Neural Information Processing Systems*, pages 2586–2594, 2016.
- [175] Andrew Gordon Wilson. *Covariance kernels for fast automatic pattern discovery and extrapolation with Gaussian processes*. PhD thesis, University of Cambridge, 2014.
- [176] Andrew Gordon Wilson and Hannes Nickisch. Kernel interpolation for scalable structured Gaussian processes (KISS-GP). *International Conference on Machine Learning (ICML)*, 2015.

- [177] Andrew Gordon Wilson, Elad Gilboa, Nehorai Arye, and John P Cunningham.
Fast kernel learning for multidimensional pattern extrapolation. In *Advances in Neural Information Processing Systems*, pages 3626–3634, 2014.
- [178] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing.
Deep kernel learning. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pages 370–378, 2016.
- [179] Jian Wu, Matthias Poloczek, Andrew G Wilson, and Peter Frazier.
Bayesian optimization with gradients. pages 5273–5284, 2017.
- [180] G. Zhang, L. Lin, W. Hu, C. Yang, and J. E. Pask.
Adaptive local basis set for Kohn–Sham density functional theory in a discontinuous Galerkin framework ii: Force, vibration, and molecular dynamics calculations. *J. Comput. Phys.*, 335: 426–443, 2017.