



THE UNIVERSITY OF QUEENSLAND
A U S T R A L I A

OPTIMIZING PERFORMANCE
IN GAUSSIAN PROCESSES

MICHAEL CICCOTOSTO-CAMP

SUPERVISOR: FRED (FARBOD) ROOSTA

Co-SUPERVISORS: ANDRIES POTGIETER, YAN ZHAO, CHRIS VAN DER HEIDE

BACHELOR OF MATHEMATICS (HONOURS)
JUNE 2022

THE UNIVERSITY OF QUEENSLAND
SCHOOL OF MATHEMATICS AND PHYSICS

CONTENTS

ACKNOWLEDGEMENTS	iv
SYMBOLS AND NOTATION	v
INTRODUCTION	1
1. GAUSSIAN PROCESSES	4
KERNELS	4
1.1. REPRODUCING KERNEL HILBERT SPACES	6
1.2. GAUSSIAN RADIAL BASIS KERNEL	9
1.3. KERNEL MACHINES	10
1.3.1. INTRODUCTION TO SUPPORT VECTOR MACHINES FOR BINARY CLASSIFICATION	10
1.3.2. INTRODUCTION TO GAUSSIAN PROCESSES FOR REGRESSION	12
1.4. GAUSSIAN PROCESSES FOR REGRESSION	16
1.4.1. NOISE-FREE OBSERVATIONS	17
1.4.2. PREDICTION WITH NOISY OBSERVATIONS	20
1.5. GAUSSIAN PROCESSES FOR CLASSIFICATION	21
1.5.1. LINEAR MODELS FOR CLASSIFICATION	21
1.5.2. LAPACE APPROXIMATION FOR POSTERIOR	24
1.5.3. PREDICTIONS	26
2. THE NYSTROM METHOD	28
2.1. THE NYSTROM METHOD	28
2.2. COLUMN PROBABILITIES	32
2.3. LEVERAGE SCORES	34
2.3.1. STATISTICAL LEVERAGE SCORES	35
2.3.2. RANK- k STATISTICAL LEVERAGE SCORES	35
2.3.3. RIDGE LEVERAGE SCORES	36
2.4. DATA PROBABILITIES	37
3. RANDOM FOURIER FEATURES	39
3.1. THEORY AND COMPUTATION	39
3.2. ORTHOGONAL RANDOM FEATURES	43
3.3. RANDOM ORTHO-MATRICES AND STRUCTURED ORTHOGONAL RANDOM MATRICES	44
4. KRYLOV SUBSPACE METHODS	47
4.1. KRYLOV SUBSPACES	47
4.2. GRAM-SCHMIDT PROCESS AND QR FACTORISATIONS	50
4.3. ARNOLDI AND LANCZOS ALGORITHM	52
4.4. OPTIMALITY CONDITIONS	56
4.5. CONJUGATE GRADIENT ALGORITHM	58
4.6. MINIMUM RESIDUAL	63
5. APPLICATIONS AND RESULTS	67
5.1. GAUSSIAN PROCESSES PREDICTION REVIEWED	67
5.2. EXPERIMENTAL SETUP	67
5.2.1. KERNEL MATRIX APPROXIMATION TESTING	68

5.2.2. KRYLOV SUBSPACE METHODS APPROXIMATION TESTING	69
5.3. DISCUSSION	69
5.3.1. KERNEL MATRIX APPROXIMATION	70
5.3.2. KRYLOV SUBSPACE METHODS	72
5.4. CONCLUSION	75
REFERENCES	77
APPENDIX A. SUPPLEMENTARY RESULTS	82
A.1. GRAM MATRIX SPECTRAL VALUES	82
A.2. NYSTROM SCORES	85
A.3. RIDGE LEVERAGE SCORES	91
A.4. NYSTROM ERRORS	97
A.5. NYSTROM SAMPLING DISTRIBUTION CONSTRUCTION TIMES	109
A.6. NYSTROM SKETCH TIMES	112
A.7. NYSTROM MEMORY USAGE	118
A.8. RFF ERRORS	124
A.9. RFF TIMES	128
A.10. RFF ADDITIONAL MEMORY USAGE	130
A.11. KERNEL COMPARISONS	132
A.12. PREDICTION ERRORS	143
A.13. COMPARISON OF COMBINING TECHNIQUES	146

ACKNOWLEDGEMENTS

The completion of this thesis is the product of, not just my own work, but also the constant support and encouragement of many friends and colleges. First and foremost I would like to express my deepest gratitude to my supervisor Dr Fred Roosta for the incredible project opportunity and the many hours spent helping me navigate the challenges of my honors year. Additionally, this endeavor could not have been possible without the generous support of A/Prof Andries Potgieter, Dr Yan Zhao, Dr Chris van der Heide and Dr Alan Huang, all of who were instrumental in provided insightful feedback and advice. I also wish to thank Bailey Whitbread, Oscar Smee and Georgianna Berthaly-Martyn for supplying copies of their own honors thesis to use as a guide.

I am incredible grateful for my family, without their unwavering love and encouragement I would not be in the position I am today. In particular, I would like to thank my father, Dr David Camp, for meticulously reading and providing corrections to my entire thesis. Finally, I wish to thank my many friends, classmates and room 67-561 fellow-honors-doers, for their patience and support throughout the year.

SYMBOLS AND NOTATION

Matrices are capitalized bold face letters while vectors are lowercase bold face letters.

<i>Syntax</i>	<i>Meaning</i>
\triangleq	An equality which acts as a statement.
$ A $	The determinate of a matrix.
$\langle \cdot, \cdot \rangle_{\mathcal{H}}$	The inner product with respect to the Hilbert space \mathcal{H} , sometimes abbreviated as $\langle \cdot, \cdot \rangle$ if the Hilbert space is clear from context.
$\ \cdot\ _{\mathcal{V}}$	The norm of a vector with respect to the vector space \mathcal{V} , sometimes abbreviated as $\ \cdot\ $ if the vector space is clear from context.
x^T, X^T	The transpose operator.
x^*, X^*	The hermitian operator.
$a.*b$ or $A.*B$	Element-wise vector (matrix) multiplication, similar to Matlab.
\propto	Proportional to.
∇ or ∇_f	The partial derivative (with respect to f).
$\nabla\nabla$	The Hessian.
\sim	Distributed according to, example $x \sim \mathcal{N}(0, 1)$
0 or 0_n or $0_{n \times m}$	The zero vector (matrix) of appropriate length (size) or the zero vector of length n or the zero matrix with dimensions $n \times m$.
1 or 1_n or $1_{n \times m}$	The one vector (matrix) of appropriate length (size) or the one vector of length n or the one matrix with dimensions $n \times m$.
$\mathbb{1}_{n \times m}$	The matrix with ones along the diagonal and zeros on off diagonal elements. The identity matrix in the case $n = m$.

$\mathbf{A}_{(.,.)}$	Index slicing to extract a submatrix from the elements of $\mathbf{A} \in \mathbb{R}^{n \times m}$, similar to indexing slicing from the python and Matlab programming languages. Each parameter can receive a single value or a 'slice' consisting of a start and an end value separated by a semicolon. The first and second parameter describe what row and columns should be selected, respectively. A single value means that only values from the single specified row/column should be selected. A slice tells us that all rows/columns between the provided range should be selected. Additionally if now start and end values are specified in the slice then all rows/columns should be selected. For example, the slice $\mathbf{A}_{(1:3,j:j')}$ is the $\mathbb{R}^{3 \times (j'-j+1)}$ submatrix containing the first three rows of \mathbf{A} and columns j to j' . As another example, $\mathbf{A}_{(:,j)}$ is the j^{th} column of \mathbf{A} .
\mathbf{A}^\dagger	Denotes the unique psuedo inverse or Moore-Penore inverse of \mathbf{A} .
\mathbb{C}	The complex numbers.
C	The classes in a classification problem.
cholesky (\mathbf{A})	A function to compute the Cholesky decomposition, \mathbf{L} , of the matrix \mathbf{A} , where $\mathbf{L}\mathbf{L}^\top = \mathbf{A}$.
$\text{cov}(\mathbf{f})$	Gaussian process posterior covariance.
d	The number of features in the data set.
D	The dimension of the feature space of the feature mapping constructed in the Random Fourier Feature method.
\mathcal{D}	The dataset, $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
$\text{diag}(\mathbf{w})$	Vector argument, a diagonal matrix containing the elements of vector \mathbf{w} .
$\text{diag}(\mathbf{W})$	Matrix argument, a vector containing the diagonal elements of the matrix \mathbf{W} .
\mathbb{E} or $\mathbb{E}_{q(x)}[z(x)]$	Expectation, or expectation of $z(x)$ where $x \sim q(x)$.
\mathcal{GP}	Gaussian process $f \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$, the function f is distributed as a Guassian process with mean function $m(\mathbf{x})$ and covariance function $k(\mathbf{x}, \mathbf{x}')$.

$k(\cdot, \cdot)$	A covariance or kernel matrix.
$\text{l.s}\{\mathbf{x}_i\}_{i=1}^n$	The linear-span of $\{\mathbf{x}_i\}_{i=1}^n$, that is, $\{\sum_{k=1}^n \lambda_k \mathbf{x}_k \mid \lambda_k \in \mathbb{R}\}$.
$\mathbf{K}_{WW'}$	For two data sets $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n]^\top \in \mathbb{R}^{n \times d}$ and $\mathbf{W}' = [\mathbf{w}'_1, \mathbf{w}'_2, \dots, \mathbf{w}'_m]^\top \in \mathbb{R}^{n' \times d}$ the matrix $\mathbf{K}_{WW'} \in \mathbb{R}^{n \times n'}$ has elements $(\mathbf{K}_{WW'})_{i,j} = k(\mathbf{w}_i, \mathbf{w}'_j)$.
$\text{lin-solve}(\mathbf{A}, \mathbf{B})$	A function used to solve $\mathbf{X} = \mathbf{A}^{-1} \mathbf{B}$ for the linear system $\mathbf{AX} = \mathbf{B}$.
$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ or $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$	(the variable \mathbf{x} has a) Multivariate Gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$.
n and n_*	The number of training (and tests) cases.
N	The dimension of the feature space.
\mathbb{N}	The natural numbers, $\mathbb{N} = \{1, 2, 3, \dots\}$.
$\mathcal{O}(\cdot)$	Big-O notation. If a function f is a member of $\mathcal{O}(g)$ then the absolute value of $f(x)$ is at most a positive multiple of $g(x)$ for all sufficiently large values of x .
$y \mid x$ and $p(x \mid y)$	A conditonal random variable y given x and its probability density.
\mathbf{Q}, \mathbf{V}	Typically used to denote a matrix with orthonormal structure.
\mathbb{R}	The real numbers.
$\text{tr}(\mathbf{A})$	The trace of a matrix.
\mathbb{V} or $\mathbb{V}_{q(x)}[z(x)]$	Variance, the variance of $z(x)$ when $x \sim q(x)$.
\mathcal{X}	Input space.
\mathbf{X}	The $n \times d$ matrix of training inputs.
\mathbf{X}_*	The $n_* \times d$ matrix of test inputs.
\mathbf{x}_i	The i^{th} training input.

\mathbb{Z}

The integers, $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$.

INTRODUCTION

Time series prediction (and related regressional tasks) is a subject of high interest across many disciplines of science and mathematics. The history of time series can be traced back to the birth of science in ancient Greece where Aristotle devised a systematic approach to weather forecasting in 350 BC in his famous treatise *Meteorologica*. This method was later used to help predict when certain meteorological induced events, such as the flooding of the Nile river [HHF73]. Statistical modelling for time series prediction would not come until the 20th century where development of AutoRegressive Moving Average (ARMA) models which were first mentioned by Yule [Yul27] in 1927 and later popularized by Box and Jenkins in their book *Time Series Analysis* published in 1970 [Box08].

Given a data set of n observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, where each input $x_i \in \mathbb{R}_{++}$ is a time value and $y_i \in \mathbb{R}$ is a output or experimental observation that acts a function of time, the goal of time series prediction is to try and best predict a value y_* at a novel time x_* . With computing power becoming ever more advanced and affordable, many have turned to Machine Learning (ML) to develop sophisticated models to address the problem of creating accurate yet computationally inexpensive time series predictors. Broadly speaking, ML is any class of heuristic algorithm that attempts to refine and develop some model to perform a "simple" task by learning through user provided input. ML is founded on the idea that any form of task learning is done through sensory input taken from the surrounding environment. More formally speaking, ML attempts to generate a function $f : X \rightarrow Y$, for some input set X and observation or output set Y , where the outputs given by f closely aligns to actual observations. It is tacitly assumed that the phenomena we are studying follows laws which admit mathematical formulation and that experimental results can be reproduced to some degree of accuracy. Typically, experiments will never produce exact values of the underpinning law, g . Instead experimental observations, y_i , will include a small amount of random error so that $y_i = g(x_i) + \varepsilon_i$ where $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$.

A ML model will attempt to make accurate predictions using some simplified formulation of the world. The distribution corresponding to the probability of a prediction within the context of the "state of the world" is referred to as the *likelihood*. The uncertainty within the likelihood stems from the predictive limits of the model. These limitation usually arise as a consequence of selecting a model which is either too simple or complex. The "state of the world" is sometimes internally captured by the model as a set of mutable parameters θ . The process of taking observations and using them to form predictions is called *inference* which, in some sense, is synonymous with learning [VdW19].

ML can be applied to time series prediction in a fairly straight forward manner by simply teaching a ML algorithm the time series data set, \mathcal{D} , to hopefully produce a function f that serves as a good approximant for event prediction.

In this thesis we shall focus on a particular class of ML algorithms called Bayesian models which, unsurprisingly, makes use of Bayesian statistics to drive inference. In Bayesian models a *prior* distribution is used to quantify the uncertainty of the current state of the model before any observations are made. The model can then be updated once data is observed by using the likelihood to give a *posterior* distribution which represents the reduced uncertainty after "teaching" the model new observations. Methods of teaching a model how to change its behavior using a new set of observations often involves the use of a

loss function L . The loss function is used as an aid in deciding what action, a , should be taken in to best minimize uncertainty. The best action, roughly speaking, can be evaluated as

$$a_{\text{opt}} = \arg \min_a \int L(y_*, a) p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) dy_*.$$

Interestingly, the best action does not rely so much on the model's internalized parameters but rather on the predictive distribution $p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})$ [VdW19]. This key insight has spawned a class of ML algorithms that focuses on inferring the function f directly by computing $p(f | \mathcal{D})$ instead of finding optimal internal parameters using $p(\theta | \mathcal{D})$ [Mur12]. Models that perform inference in this manner are called *non-parameteric* models. Within the *non-parameteric* model paradigm, the predictive distribution can be represented as

$$p(y_* | x_*, \mathbf{X}, \mathbf{y}) = \int p(y_* | f, x_*) p(f | \mathbf{X}, \mathbf{y}) df$$

and once new data is observed the posterior can be updated using Baye's rule

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(f, f_* | \mathbf{y}) = \frac{p(\mathbf{y} | f) p(f, f_*)}{p(\mathbf{y})}$$

[Ras06]. This thesis will focus on a particular non-parameteric Bayesian ML model called Gaussian processes (GPs). The over arching idea of GPs is to assign a prior probability to every possible function mapping from X to Y . While this does not appear to be computationally tractable due to the seemingly uncountable infinite number of mappings that would require checking, it turns out, these computations can infact be carried out given we are only seeking predictions at a finite number of points using a finite number of observations. GPs occupy a special place within the realm of ML since they account for uncertainty in a principled way, are relatively simple to implement and are highly modular allowing them to easily be incorporated into a larger systems. It is no surprise then that while other kernel methods (such as kernelized k^{th} nearest neighbors and ridge regression) are still overshadowed by their neural network cousins, GPs have made a quiet comeback in the ML community [Cao18].

The following example highlights a particular GP success story: a team of researchers led by Associate Professor Andries Potgieter at QAAFI (UQ) are currently investigating new digital approaches to accurately derive crop phenology stages (i.e. mid green, peak, flowering, grain filling and harvest) measured at field scale across large regions. Such methods can be used to better inform farmers and industry on the optimised time to plant various crops to minimize crop loss from environmental stresses such as frost and fungal disease. This involves analysing crop growth from previous seasons (i.e. 2018-2021) to forecast when certain phenological stages will take place in the current harvest. Outputs form this tool will allow producers to accurately map the temporal and spatial extend of phenology at a field and farm levels across different regions and seasons. This problem is readily converted into a time series problem. Originally, Potgieter's team surveyed a number of different parameteric models to carry out forecasting. However, the parameteric models we serverely limited in their ability to inform when key phenological stages would take place. After seeing the success of applying GPs to other remote sensing tasks [SD22] the team investigated the use of GPs in their own research to find that they could produce much higher resolution predictions from which they could infer a far richer phenological timeline [Pot13]. A comparision of using GPs over other parameteric models is shown in Figure 1. Potgieter's team found that the only draw back to using GPs was the lengthy run time required to create predictions and fears that

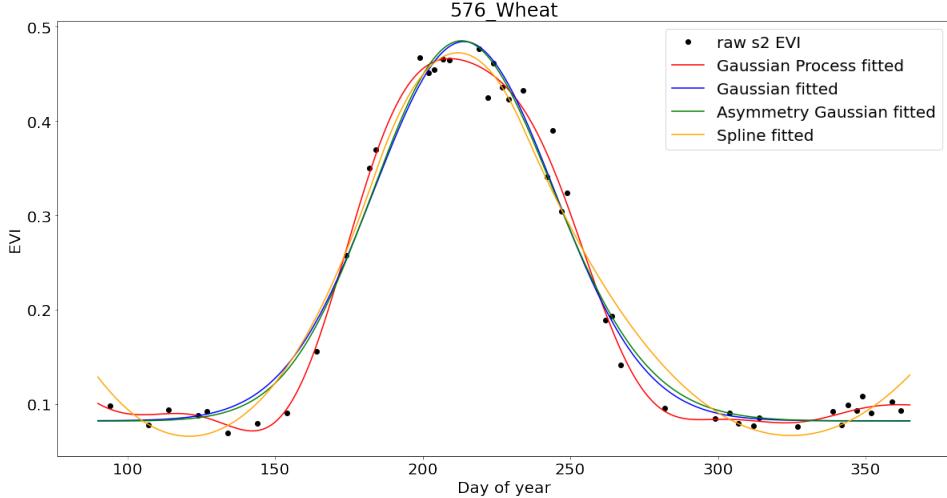


FIGURE 1. Potgieter’s team found that GPs were superior in terms of predicting a phenological timeline for a number of common seasonal crops over other parametric models.

collecting new data each season will only exacerbate the issue. This is a common problem shared by anyone wanting to use GPs. Due to their unwieldy $\mathcal{O}(n^3)$ runtime, where n is the number of observations, GPs become impractical to apply on datasets with $n > 10^5$ samples. As such, the goal of this thesis is to explore various avenues one can take to replace some of the more intense calculations of GPs with computationally more efficient approximations without overly sacrificing accuracy.

Chapter 1 will give a more mathematical treatment of GPs starting from the ground through a review of some fundamental material from functional analysis also the theory behind the motivation of GPs before finally concluding with concrete algorithms for GP regression and classification. Chapters 2 and 3 will cover techniques for approximating a large matrix used with GPs that provides information on how similar each observation is to one other. Chapter 4 then gives alternative methods for solving linear systems, an essential component required for the GP algorithm to work.

1. GAUSSIAN PROCESSES

The aim of this chapter is to explore the theory behind GPs. First, some essential theory from functional analysis on kernels and reproducing kernel Hilbert spaces will be reviewed which are not only used in GPs but are found in a vast array of machine learning models, aptly named kernel machines. Afterward, we shall go through the underlying statistics that drive GP prediction and use it to form algorithms for both regression and classification tasks. Note that most of the theory presented here is only for real-values data sets although most the time complex-valued generalizations do exist.

Kernels. Often in machine learning we are often met with the challenge of how to best represent data instances as fixed size feature vectors $\mathbf{x}_i \in X$. For certain objects it might not be obvious at all how to represent the data as a fixed length vector. Good examples of variable length data include textual documents and genomic data. For these data types we can define a method of measuring similarity between objects which requires them to first be converted to a fixed length feature vector first [Mur12]. To do this we begin by mapping the feature vectors into a Hilbert space \mathcal{H} which enriches the vector space with an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ and a norm $\| \cdot \|_{\mathcal{H}} : \mathcal{H} \rightarrow \mathbb{R}$. Input data is transformed into feature space vectors via a non-linear feature mapping $\Phi : X \rightarrow \mathcal{H}$. The benefit of using feature maps in this way is that a non-linear decision boundary can be constructed using linear models. In some instances a similarity measure can be computed directly using a function $k : X \times X \rightarrow \mathbb{R}$, instead of needing to construct a Φ and then computing the inner product of the transformed instances. Functions that act directly on our data instances are known as kernel functions and using them to avoid computation associated with the underlying feature space is known as the kernel trick [Ste08]. These ideas are stated more formally in Definition 1.

Definition 1 (Kernel). Let X be a non-empty set. Then a function $k : X \times X \rightarrow \mathbb{R}$ is called a kernel on X if there exists a Hilbert space and a map $\Phi : X \rightarrow \mathcal{H}$ such that for all $\mathbf{x}, \mathbf{x}' \in X$ we have $k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{H}}$. We call the Φ the feature map and \mathcal{H} the feature space of k .

It is worth noting that almost no conditions are placed on the set X , allowing it to accommodate virtually any form of data. It is not surprising then that neither the feature map nor the feature space are uniquely determined by the kernel. As shown by the example from Steinwart and Christmann [Ste08], when $X = \mathbb{R}$ and $k(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$ where $\mathbf{x}, \mathbf{x}' \in X$, we can see that k is a kernel using the feature map $\Phi(\mathbf{x}) = \mathbf{x}$ and $\mathcal{H} = \mathbb{R}$. However, another suitable feature map for this particular kernel is $\Phi'(\mathbf{x}) = (\mathbf{x}/\sqrt{2}, \mathbf{x}/\sqrt{2})$ with a corresponding feature space of $\mathcal{H} = \mathbb{R}^2$ since

$$\langle \Phi'(\mathbf{x}), \Phi'(\mathbf{x}') \rangle_{\mathbb{R}^2} = \frac{\mathbf{x}'}{\sqrt{2}} \cdot \frac{\mathbf{x}}{\sqrt{2}} + \frac{\mathbf{x}'}{\sqrt{2}} \cdot \frac{\mathbf{x}'}{\sqrt{2}} = \mathbf{x} \cdot \mathbf{x}'$$

for $\mathbf{x}, \mathbf{x}' \in X$. While there might be numerous functions that provide some notion of similarity between data entries, these functions might not be valid kernels. Instead of needing to construct a feature map and feature space to verify that a chosen function is a valid kernel using Definition 1, we can make use of a much simpler set of criteria. Before embarking on this train of thought, we need to define the following.

Definition 2 (Positive Definite and Positive Semidefinite). A function $k : X \times X \rightarrow \mathbb{R}$ is positive semidefinite if for all $n \in \mathbb{N}$ and $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ and all $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ we have

$$(1) \quad \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_j, \mathbf{x}_i) \geq 0.$$

Furthermore, k is said to be positive definite if for mutually distinct $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ equality (1) only holds for $\alpha_1 = \dots = \alpha_n = 0$ [Ste08].

Definition 3 (Symmetric). A function $k : X \times X \rightarrow \mathbb{R}$ is called symmetric if $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$ for any inputs $\mathbf{x}', \mathbf{x} \in X$ [Ste08].

Definition 4 (Gram Matrix). For fixed $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ the matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ where $\mathbf{K}_{i,j} \triangleq k(\mathbf{x}_j, \mathbf{x}_i)$ is the Gram matrix [Ste08].

Note that checking if a function is positive (semi) definite is equivalent to checking that any Gram matrix produced by a function is positive (semi) definite. If k is a real valued kernel corresponding to the feature map Φ , then k is symmetric by virtue of the fact that the inner product of a real Hilbert space is symmetric. Moreover k is positive definite since for $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ and $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ we have

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_j, \mathbf{x}_i) \\ &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle_{\mathcal{H}} \\ &= \left\| \sum_i^n \alpha_i \Phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \\ &\geq 0. \end{aligned}$$

The following theorems tell us that it is not only necessary for a kernel to be positive semi definite but it is also a sufficient condition.

Theorem 5. A function $k : X \times X \rightarrow \mathbb{R}$ is a kernel if and only if it is symmetric and positive semidefinite [Ste08, page 118].

Proof. In view of the above discussion, it suffices to show that a symmetric and positive definite function is a kernel. To this end, we write

$$\mathcal{H}_{\text{pre}} \triangleq \left\{ \sum_{i=1}^n \alpha_i k(\cdot, \mathbf{x}_i) \mid n \in \mathbb{N}, \alpha_1, \dots, \alpha_n \in \mathbb{R}, \mathbf{x}_1, \dots, \mathbf{x}_n \in X \right\}$$

and for $f \triangleq \sum_{i=1}^n \alpha_i k(\cdot, \mathbf{x}_i) \in \mathcal{H}_{\text{pre}}$ and $g \triangleq \sum_{j=1}^m \beta_j k(\cdot, \mathbf{x}'_j) \in \mathcal{H}_{\text{pre}}$, we define

$$\langle f, g \rangle \triangleq \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(\mathbf{x}'_j, \mathbf{x}_i).$$

Note that this definition is independent of the representation of f since we have $\langle f, g \rangle = \sum_{j=1}^m \beta_j f(\mathbf{x}'_j)$. Furthermore, by the symmetry of k , we find $\langle f, g \rangle = \sum_{i=1}^m \alpha_i g(\mathbf{x}_i)$, that is, the definition is also independent of the representation of g . Moreover, the definition immediately shows that $\langle \cdot, \cdot \rangle$ is bilinear and symmetric, and since k is positive definite, $\langle \cdot, \cdot \rangle$ is also positive, that is, $\langle f, f \rangle \geq 0$ for all $f \in \mathcal{H}_{\text{pre}}$. Consequently, $\langle \cdot, \cdot \rangle$ satisfies the Cauchy-Schwarz inequality, that is,

$$|\langle f, g \rangle|^2 \leq \langle f, f \rangle \cdot \langle g, g \rangle, \quad f, g \in \mathcal{H}_{\text{pre}}.$$

Now let $f \triangleq \sum_{i=1}^n \alpha_i k(\cdot, \mathbf{x}_i)$ with $\langle f, f \rangle = 0$. Then for all $\mathbf{x} \in X$, we have

$$|f(\mathbf{x})|^2 = \left| \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) \right|^2 = |\langle f, k(\cdot, \mathbf{x}) \rangle|^2 \leq \langle k(\cdot, \mathbf{x}), k(\cdot, \mathbf{x}) \rangle \cdot \langle f, f \rangle$$

and hence we find that $f = \mathbf{0}$. Therefore we have shown that $\langle \cdot, \cdot \rangle$ is an inner product on \mathcal{H}_{pre} . Let \mathcal{H} be a completion of \mathcal{H}_{pre} and $I : \mathcal{H}_{\text{pre}} \rightarrow \mathcal{H}$ be the corresponding isometric embedding. Then \mathcal{H} is a Hilbert space and we have

$$\langle I k(\cdot, \mathbf{x}'), I k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = \langle k(\cdot, \mathbf{x}'), k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_{\text{pre}}} = k(\mathbf{x}, \mathbf{x}')$$

for all \mathbf{x}, \mathbf{x}' , that is, $\mathbf{x} \mapsto I k(\cdot, \mathbf{x})$, $\mathbf{x} \in X$, defines a feature map of k . \square

1.1. Reproducing Kernel Hilbert Spaces. We shall now shift our attention towards reproducing kernel Hilbert spaces (RKHS) and describe their relation to kernels, and see that in some sense, the RKHS of a kernel k is the smallest feature space for a kernel. The formal definition of a RKHS is stated in Definition 6.

Definition 6 (RKHS). Let $X \neq \emptyset$ and \mathcal{H} be a real Hilbert space over X

- (1) A function $k : X \times X \rightarrow \mathbb{R}$ is called a *reproducing kernel* if we have $k(\cdot, \mathbf{x}) \in \mathcal{H}$ for all $\mathbf{x} \in X$ and the *reproducing property*

$$f(\mathbf{x}) = \langle f, k(\cdot, \mathbf{x}) \rangle$$

holds for all $f \in \mathcal{H}$ and $\mathbf{x} \in X$.

- (2) The space \mathcal{H} is called a *reproducing kernel Hilbert space* over X if for all $\mathbf{x} \in X$ the Dirac functional $\delta_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}$ defined by $\delta_{\mathbf{x}}(f) \triangleq f(\mathbf{x})$, $f \in \mathcal{H}$ is continuous.

[Ste08]

An important property of the RKHS is that the convergence in the norm implies pointwise convergence. Specifically, in a RKHS for any sequence of functions $\{f_n\} \subset \mathcal{H}$ where $\|f_n - f\| \rightarrow 0$ we have $|\delta_{\mathbf{x}}(f_n) - \delta_{\mathbf{x}}(f)| = |f_n(\mathbf{x}) - f(\mathbf{x})| \rightarrow 0$. Note that because the evaluation function is both linear and continuous, then it is also bounded in the sense that there is a $c \in \mathbb{R}$, $c > 0$ such that for every $f \in \mathcal{H}$ and a fixed $\mathbf{x} \in X$ we have $|\delta_{\mathbf{x}}(f)| \leq c \|f\|_{\mathcal{H}}$ [Ber96]. This property of uniform convergence implying pointwise convergence is important since it tells us that if functions $f, g \in \mathcal{H}$ are close in norm then their evaluation at any point is also similar. The following lemma ties together the definition of an RKHS, reproducing kernel and a kernel.

Lemma 7. Let \mathcal{H} be a Hilbert function space over X that has a reproducing kernel k . Then \mathcal{H} is a RKHS and \mathcal{H} is also a feature space of k where the feature map $\Phi : X \rightarrow \mathcal{H}$ is given by

$$\Phi(\mathbf{x}) = k(\cdot, \mathbf{x})$$

for some $\mathbf{x} \in X$. We call Φ the canonical feature map.

Proof. Since the reproducing property tells us that any Dirac functional can be represented by the reproducing kernel this means

$$|\delta_{\mathbf{x}}(f)| = |f(\mathbf{x})| = |\langle f, k(\cdot, \mathbf{x}) \rangle| \leq \|k(\cdot, \mathbf{x})\|_{\mathcal{H}} \cdot \|f\|_{\mathcal{H}}$$

for all $\mathbf{x} \in X$, $f \in \mathcal{H}$. This shows continuity of $\delta_{\mathbf{x}}$ for $\mathbf{x} \in X$. In order to show that Φ is a feature map, fix an $\mathbf{x}' \in X$ and set $f = k(\cdot, \mathbf{x}')$. Then for $\mathbf{x} \in X$, the reproducing property yields

$$\langle \Phi(\mathbf{x}'), \Phi(\mathbf{x}) \rangle_{\mathcal{H}} = \langle k(\cdot, \mathbf{x}'), k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = \langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = f(\mathbf{x}) = k(\mathbf{x}', \mathbf{x}).$$

□

This tells us that every Hilbert space with a reproducing kernel is a RKHS. We can also show the converse, that is, every RKHS has a unique reproducing kernel seen in Theorem 8.

Theorem 8. Let \mathcal{H} be a RKHS over X . Then $k : X \times X \rightarrow \mathbb{R}$ defined by $k(\mathbf{x}', \mathbf{x}) = \langle \delta_{\mathbf{x}}, \delta_{\mathbf{x}'} \rangle_{\mathcal{H}}$, $\mathbf{x}, \mathbf{x}' \in X$ is the only reproducing kernel of \mathcal{H} [Ste08, page 120].

Proof. We first show that k is a reproducing kernel. To this end, we know from the Frechet-Riesz Representation theorem (see [Ste08, page 504]) that there exists a anti-linear isomorphism, $I : \mathcal{H}' \rightarrow \mathcal{H}$ that assigns to every functional in \mathcal{H}' the representing element in \mathcal{H} , that is, $g'(f) = \langle f, Ig' \rangle$ for all $f \in \mathcal{H}$, $g' \in \mathcal{H}'$. Then, for all $\mathbf{x}, \mathbf{x}' \in X$, we have

$$k(\mathbf{x}, \mathbf{x}') = \langle \delta_{\mathbf{x}}, \delta_{\mathbf{x}'} \rangle_{\mathcal{H}'} = \langle Ig(\mathbf{x}), Ig(\mathbf{x}') \rangle_{\mathcal{H}} = \delta_{\mathbf{x}}(Ig(\mathbf{x}')) = Ig(\mathbf{x}'),$$

which shows $k(\cdot, \mathbf{x}') = Ig(\mathbf{x}')$ for all $\mathbf{x}' \in X$. From this we immediately obtain

$$f(\mathbf{x}') = \delta_{\mathbf{x}'}(f) = \langle f, Ig(\mathbf{x}') \rangle = \langle f, k(\cdot, \mathbf{x}') \rangle,$$

that is, k has the reproducing property. Now let \tilde{k} be an arbitrary reproducing kernel of \mathcal{H} . For $\mathbf{x}' \in X$, the basis representation of $\tilde{k}(\cdot, \mathbf{x}')$ then yields

$$\tilde{k}(\cdot, \mathbf{x}') = \sum_{i \in I} \langle \tilde{k}(\cdot, \mathbf{x}'), e_i \rangle e_i = \sum_{i \in I} \overline{e_i(\mathbf{x}')} e_i,$$

where the convergence is with respect to $\|\cdot\|_{\mathcal{H}}$. Since convergence in the norm implies pointwise convergence, we obtain $\tilde{k}(\mathbf{x}, \mathbf{x}') = \sum_{i \in I} e_i(\mathbf{x}) \overline{e_i(\mathbf{x}')}$. Finally, since \tilde{k} and $\{e_i\}_{i \in I}$ were arbitrarily chosen, we find that $\tilde{k} = k$, that is, k is the only reproducing kernel of \mathcal{H} . □

Theorem 8 shows that a RKHS is uniquely determined by its kernel. In fact the other direction can also be shown to afford a one-to-one correspondence between kernels and RKHS. This is known as the Moore-Aronszajn theorem. However, before we prove it, we shall need the following result from [Ber03, page 15].

Theorem 9. Let \mathcal{H}_0 be any subspace of \mathbb{C}^X , the space of complex functions on X , on which an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ is defined, with associated norm $\|\cdot\|_{\mathcal{H}_0}$. In order that there exists a Hilbert space \mathcal{H} such that

- a) $\mathcal{H}_0 \subset \mathcal{H} \subset \mathbb{C}^X$ and the topology defined on \mathcal{H}_0 by the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ coincides with the topology induced on \mathcal{H}_0 by \mathcal{H}

- b) \mathcal{H} has a reproducing kernel k
- c) it is necessary and sufficient that the evaluation functionals $(e_t)_{t \in X}$ are continuous on \mathcal{H}_0
- d) it is also necessary and sufficient any Cauchy sequence (f_n) in \mathcal{H}_0 converging pointwise to 0 converges also to 0 in the norm sense.

[Ber03, page 15].

Theorem 10 (Moore-Aronszajn). Suppose k is a symmetric positive definite kernel on a set X . Then there is a unique Hilbert space \mathcal{H} of functions for which k is the reproducing kernel. The subspace \mathcal{H}_0 of \mathcal{H} spanned by the functions $(k(\cdot, \mathbf{x})_{\mathbf{x} \in X})$ is dense in \mathcal{H} and \mathcal{H} is the set of functions on X which are pointwise limits of Cauchy sequences in \mathcal{H}_0 with the inner product

$$(2) \quad \langle f, g \rangle_{\mathcal{H}_0} = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \overline{\beta_j} k(\mathbf{y}_j, \mathbf{x}_i)$$

where $f = \sum_{i=1}^n \alpha_i k(\cdot, \mathbf{x}_i)$ and $g = \sum_{j=1}^m \beta_j k(\cdot, \mathbf{y}_j)$ [Ber03, page 20].

Proof. First remark that the complex number $\langle f, g \rangle$ defined by (2) does not depend on the representations not necessarily unique of f and g :

$$\langle f, g \rangle_{\mathcal{H}_0} \sum_{i=1}^n \alpha_i \overline{g(\mathbf{x}_i)} = \sum_{j=1}^m \overline{\beta_j} f(\mathbf{y}_j),$$

this shows that $\langle f, g \rangle_{\mathcal{H}_0}$ depends on f and g only through their values. Then, taking

$$f = \sum_{i=1}^n k(\cdot, \mathbf{x}_i) \quad \text{and} \quad g = k(\cdot, \mathbf{x})$$

we get

$$\langle f, k(\cdot, \mathbf{x}) \rangle = \sum_{i=1}^n \alpha_i \overline{g(\mathbf{x}_i)} = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) = f(\mathbf{x}).$$

Thus the inner product with $k(\cdot, \mathbf{x})$ "reproduces" the values of functions in \mathcal{H}_0 . In particular

$$\|k(\cdot, \mathbf{x})\|_{\mathcal{H}_0}^2 = \langle k(\cdot, \mathbf{x}), k(\cdot, \mathbf{x}) \rangle = k(\mathbf{x}, \mathbf{x}).$$

As k is a positive type function, $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ is a semi-positive hermitian form on $\mathcal{H}_0 \times \mathcal{H}_0$. Now, suppose that $\langle f, f \rangle_{\mathcal{H}_0} = 0$. From the Cauchy-Schwarz inequality we have

$$\forall \mathbf{x} \in X \quad |f(\mathbf{x})| = |\langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_0}| \leq \langle f, f \rangle_{\mathcal{H}_0}^{\frac{1}{2}} [k(\mathbf{x}, \mathbf{x})]^{\frac{1}{2}} = 0$$

and $f \equiv 0$. Let us consider \mathcal{H}_0 endowed with the topology associated with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$ and check conditions c) and d) in Theorem 9. Let f and g in \mathcal{H}_0

$$\begin{aligned} \forall \mathbf{x} \in X \quad |e_{\mathbf{x}}(f) - e_{\mathbf{x}}(g)| &= |\langle f - g, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_0}| \\ &\geq \|f - g\|_{\mathcal{H}_0} [k(\mathbf{x}, \mathbf{x})]^{\frac{1}{2}}. \end{aligned}$$

Therefore the evaluation functionals are continuous on \mathcal{H}_0 and c) is satisfied.

Let us now check d). Let (f_n) be a Cauchy sequence (hence bounded) in \mathcal{H}_0 converging pointwise to $\mathbf{0}$ and let $A > 0$ be an upper bound for $(\|f_n\|_{\mathcal{H}_0})$. Let $\varepsilon > 0$ and $N(\varepsilon)$ such that

$$n > N(\varepsilon) \Rightarrow \|f_{N(\varepsilon)} - f_n\|_{\mathcal{H}_0} < \frac{\varepsilon}{A}.$$

Fix $k, \alpha_1, \dots, \alpha_k$ and $\mathbf{x}_1, \dots, \mathbf{x}_k$ such that

$$f_{N(\varepsilon)} = \sum_{i=1}^k \alpha_i k(\cdot, \mathbf{x}_i).$$

As

$$\|f_n\|_{\mathcal{H}_0}^2 = \langle f_n - f_{N(\varepsilon)}, f_n \rangle_{\mathcal{H}_0} + \langle f_{N(\varepsilon)}, f_n \rangle_{\mathcal{H}_0},$$

we have, for $n > N(\varepsilon)$,

$$\|f_n\|_{\mathcal{H}_0}^2 < \varepsilon + \sum_{i=1}^k \alpha_i f_n(\mathbf{x}_i),$$

hence $\limsup_{n \rightarrow \infty} \|f_n\|^2 \leq \varepsilon$. As ε is arbitrary this entails that (f_n) converges to 0 in the norm sense. We are now in a position to apply Theorem 9 to \mathcal{H}_0 : there exists a Hilbert space \mathcal{H} of functions on X satisfying a) and b). \mathcal{H} is the set of functions f for which there exists a Cauchy sequence (f_n) in \mathcal{H}_0 converging pointwise to f . Such a sequence (f_n) is also converging to f in the norm sense: \mathcal{H}_0 is dense in \mathcal{H} . Therefore \mathcal{H} is unique and

$$\begin{aligned} \forall \mathbf{x} \in X \quad f(\mathbf{x}) &= \lim_{n \rightarrow \infty} f_n(\mathbf{x}) = \lim_{n \rightarrow \infty} \langle f_n, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_0} \\ &= \langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} \end{aligned}$$

thus k is the reproducing kernel of \mathcal{H} . □

The elements of a RKHS will often inherit the analytical properties of its corresponding kernel. This means that kernels provide a mechanism for generating spaces of functions with useful analytical properties.

1.2. Gaussian Radial Basis Kernel. Although there are many kernels to use at our disposal, we turn our attention to a specific class of kernel that shall be used extensively in the upcoming theory and experimentation.

Definition 11 (Gaussian Radial Basis Kernel). *For $d \in \mathbb{N}$, $\sigma \in \mathbb{R}_{++}$ and $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^d$ we define*

$$k_\sigma(\mathbf{z}, \mathbf{z}') \triangleq \exp \left(-\sigma^{-2} \sum_{j=1}^d (z_j - z'_j)^2 \right).$$

Then k_σ is a real valued kernel called the Gaussian Radial Basis Kernel (RBF) kernel with bandwidth σ . Moreover k_σ can be computed as

$$\exp \left(\frac{-\|\mathbf{z} - \mathbf{z}'\|_2^2}{\sigma^2} \right)$$

[Ste08].

The Gaussian RBF kernel makes for a very simple and intuitive measurement of similarity between its inputs. One geometric interpretation of the Gaussian RBF kernel is that as the radius of the smallest d -sphere containing $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^d$ grows the corresponding measurement of similarity decays exponentially. A visual representation of this decay is shown in Figure 2.

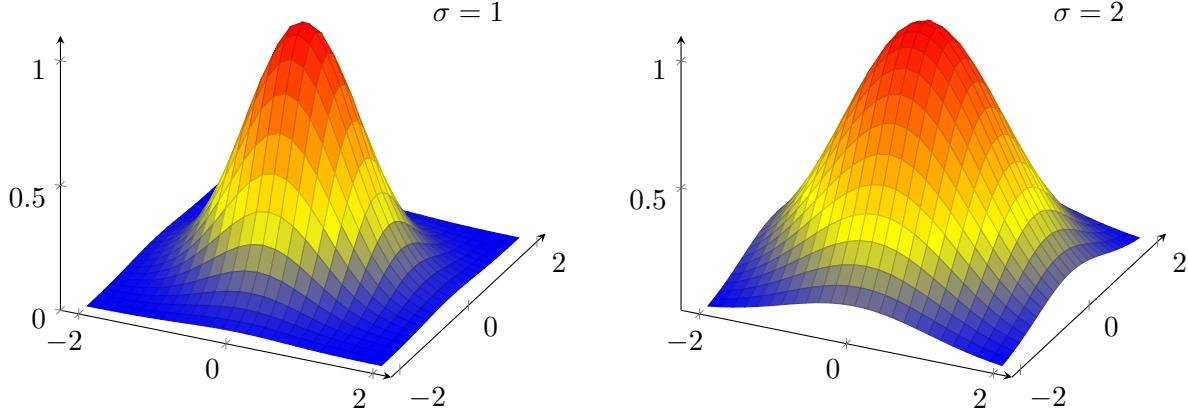


FIGURE 2. A graph of the Gaussian RBF from Definition 11 for $d = 2$. Evidently, a larger value of σ slows the rate of decay increasing the similarity between the same pair of samples.

This kernel is infinitely differentiable meaning it has mean square derivatives of all orders and is therefore very smooth. In fact, some argue that such strong smoothness makes it unrealistic for modelling natural phenomena [Ras06, Ste99]. Nonetheless, Gaussian RBF kernels remain the one of the most widely used kernels in literature.

1.3. Kernel Machines.

1.3.1. Introduction to Support Vector Machines for Binary Classification. In this section, we will be investigating at two different machine learning models that make use of kernels to perform classification and regression. The first class of kernel machines to be evaluated are support vector machines (SVM). SVMs were originally designed for binary classification and as such only a model for binary classification is presented, although extensions exist that allow regression and multi-class classification.

For the binary classification problem we are tasked with labelling new samples with either one of two classes, -1 or 1 . We shall assume our input space consists of vectors from \mathbb{R}^d and that we provided with a labelled training set $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$. One simple method to classify samples is by creating an affine linear hyperplane satisfying

$$(3) \quad \begin{aligned} \langle \mathbf{w}, \mathbf{x}_i \rangle + b &> 0, & y_i = +1 \\ \langle \mathbf{w}, \mathbf{x}_i \rangle + b &< 0, & y_i = -1 \end{aligned}$$

for some $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ where $\|\mathbf{w}\|_2 = 1$ [Ste08]. Moreover we would like \mathbf{w} and b to maximise the margin, that is the maximal distance between the separating hyperplane and the points in D . The specific

\mathbf{w} and b obtained through the training set is denoted \mathbf{w}_D and b_D and the resulting descision function is defined as

$$f_D(\mathbf{x}) \triangleq \text{sign}(\langle \mathbf{w}_D, \mathbf{x} \rangle + b_D).$$

There are, however, a number of shortcomings to this model. The most obvious is that our training data may not be linearly separable in \mathbb{R}^d meaning no such \mathbf{w}_D and b_D exist. Moreover, when noise is introduced to the data set this model will prioritize finding a hyperplane that perfectly separates the two classes, making no compromises in misclassifying points, subjecting it to overfitting. SVMs were introduced by Boser *et al.* [Bos92] to address the first issue of separability. Their approach was to lift the input vector into a more malleable Hilbert space H_0 using a feature map. The inputs are then classified within the new space. Unfortunately this method does nothing to address the second issue of over fitting and, if anything, actually worsens it. Cortes and Vapnik [Cor95] attempted to address this second issue by introducing slack variables to (3) so that now we need to satisfy $y_i (\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b) \geq 1 - \xi_i$ for some $\xi_i \in \mathbb{R}_{++}$. These constraints can be re-written as

$$\xi_i \geq 1 - y_i (\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b)$$

and combining this with our slack constraints (that is $\xi_i \geq 0$) yields

$$\xi_i \geq \max\{0, 1 - y_i (\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b)\} = L_{\text{hinge}}(y_i, \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b)$$

where L_{hinge} is the hinge loss defined as

$$L_{\text{hinge}}(y, \eta) \triangleq \max\{0, 1 - y\eta\}.$$

This optimization problem can be re-written in the form

$$\min_{(\mathbf{w}, b) \in H_0 \times \mathbb{R}} \lambda \|\mathbf{w}\|_{H_0} + \frac{1}{n} \sum_{i=1}^n L_{\text{hinge}}(y_i, f_{(\mathbf{w}, b)})$$

where $f_{(\mathbf{w}, b)} : X \rightarrow \mathbb{R}$ is defined as

$$f_{(\mathbf{w}, b)} \triangleq \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b$$

[Ste08]. Unfortunately, this new embedding requires us to solve for optimal parameters in a very high, or even infinite, dimension vector space. To get around this, the Lagrange approach is typically used to solve the corresponding dual problem. When the hinge loss is used the dual problem becomes

$$(4) \quad \begin{aligned} & \max_{\alpha \in [0, C]^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle \\ & \text{subject to } \sum_{i=1}^n y_i \alpha_i = 0 \end{aligned}$$

Notice that in the dual problem, we find that inner products are only taken with vectors that have the feature map applied to them allowing us to employ the kernel trick described in Section 1 so that (4) becomes

$$\begin{aligned} & \max_{\alpha \in [0, C]^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \\ & \text{subject to } \sum_{i=1}^n y_i \alpha_i = 0. \end{aligned}$$

1.3.2. Introduction to Gaussian Processes for Regression. The next machine learning model of interest that uses kernels are gaussian processes. To motivate this model, consider the time series data in Figure 3 (A).

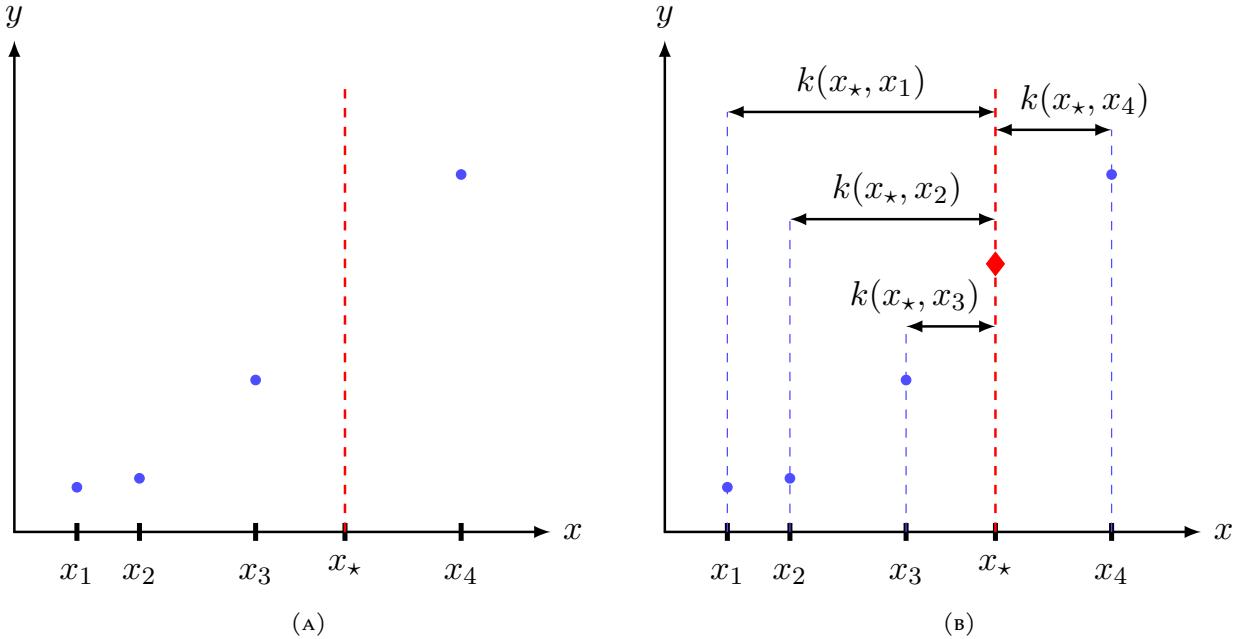


FIGURE 3. Panel (A) shows depicts the classical problem of time series prediction, guessing a value for x_* given values for surrounding times. Panel (B) shows a suitable choice for the value at time x_* with the reasoning that closely surrounding values should have greater influence over inference. Example taken from the 2013 Machine Learning course CPSC 540 instructed by Nando de Freitas [NdF13].

In this diagram there is a number of observed values $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)\}$ (blue labels) as well as a missing observation at time x_* . This is a classic problem of time series prediction, that is what is a good prediction for the missing value at time x_* ? Perhaps something close to the red diamond seen in Figure 3 (B)? Why does this red diamond seem like a good choice? Because for known data for which the measurement of similarity is small, we expect the corresponding outputs should also be similar since most natural phenomena are continuous by nature. This reasoning is used to underpin GPs, that is, training inputs that are more similar to the input value should have greater influence over the prediction.

Like SVMs, we can motivate the mathematical model of a GP through a linear model. Since GPs are designed for regression tasks, we shall only focus on GP regression although we will see later that GPs can be extended to perform binary classification and even multi-class classification. To begin, consider the following linear regression model

$$(5) \quad f(x) \triangleq \langle w, x \rangle$$

where we again assume that \mathbf{x} belongs to \mathbb{R}^d and that $\mathbf{w} \in \mathbb{R}^d$ is a weight vector. Notice the striking resemblances to the linear classifier used in our derivation for the SVM model, although this time we are using the value computed by the inner product directly to infer instead of fitting it over a sign function to force it into a binary class. Suppose we have independently sampled observations $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ to a noisy version of f

$$y_i = f(\mathbf{x}_i) + \varepsilon_i$$

where $\varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_n^2)$. Together the assumption of noise and the base linear model give rise to a likelihood, or more specifically, a probability density over the observations given the inputs and weight parameters. Due to the assumption of independence in our observations

$$\begin{aligned} (6) \quad p(y | \mathbf{X}, \mathbf{w}) &= \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left(-\frac{(y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2}{2\sigma_n^2}\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma_n^2} \left(\sum_{i=1}^n (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2\right)\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma_n^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2\right) \\ &= \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma_n^2 \mathbf{1}_{n \times n}) \end{aligned}$$

where $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top \in \mathbb{R}^n$ and $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ [Ras06, page 9]. Within the Bayesian paradigm, a prior is required to represent our beliefs about the parameters in the absence of any information. Typically, the following prior is used for the weight vector

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$$

where Σ_p is an appropriate covariance matrix. Ideally, we would like to know the posterior pdf $p(\mathbf{w} | \mathbf{y}, \mathbf{X})$ which refines our choices of \mathbf{w} by taking into account our observations. The posterior can be computed using Bayes rule

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}) \propto p(\mathbf{y} | \mathbf{w}, \mathbf{X}) p(\mathbf{w}).$$

(6) gives us a probability for $p(\mathbf{y} | \mathbf{w}, \mathbf{X})$ and since $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$ then

$$p(\mathbf{w}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} \mathbf{w}^\top \Sigma_p^{-1} \mathbf{w}\right)$$

[Kro14]. This means, up to proportionality

$$\begin{aligned} p(\mathbf{w} | \mathbf{X}, \mathbf{y}) &\propto \exp\left(-\frac{1}{2\sigma_n^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})\right) \exp\left(-\frac{1}{2} \mathbf{w}^\top \Sigma_p^{-1} \mathbf{w}\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{w} - \bar{\mathbf{w}})^\top \left(\frac{1}{\sigma_n^2} \mathbf{X}^\top \mathbf{X} + \Sigma_p^{-1}\right) (\mathbf{w} - \bar{\mathbf{w}})\right) \end{aligned}$$

where $\bar{\mathbf{w}} \triangleq \sigma_n^{-2} (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \Sigma_p^{-1})^{-1} \mathbf{X}^\top \mathbf{y}$. Notice that this again is the kernel of a multivariate Gaussian distribution with mean $\bar{\mathbf{w}}$ and covariance \mathbf{A}^{-1} where $\mathbf{A} \triangleq \sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \Sigma_p^{-1}$ so that

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\mathbf{w}, \mathbf{A}^{-1}).$$

To make a prediction of our target function for an input, \mathbf{x}_* , outside our observed values we can take the average over all possible parameter values weighted by the posterior to predict $f_* \triangleq f(\mathbf{x}_*)$ which yields

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int_{\mathbb{R}^d} p(f_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} = \mathcal{N}(\mathbf{x}_*^\top \bar{\mathbf{w}}, \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*).$$

This gives another Gaussian distribution whose mean is the mean of the posterior distribution of the weight vectors multiplied by the input vector, and whose covariance in the quadratic form of the covariance of the weight vectors again with the input vectors. This makes sense since it tells us that the uncertainty of the model grows quadratically with the magnitude of the input.

We can now employ the kernel trick in the exact same manner in the derivation of the SVM model, that is, by using a feature mapping Φ to lift the inputs of our linear regression model from (5) into a higher dimension and more workable Hilbert space so that our model now becomes

$$f(\mathbf{x}) \triangleq \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle.$$

The derivation for the new model is identical with the only difference being that \mathbf{x}_* is replaced with $\Phi(\mathbf{x}_*)$ and \mathbf{X} is replaced with $\Phi(\mathbf{X}) \triangleq [\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_n)]^\top \in \mathbb{R}^{n \times N}$ where N is the dimension of the feature space. The new predictive distribution can be expressed as

$$(7) \quad f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}\left(\frac{1}{\sigma_n^2} \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{X})^\top \mathbf{y}, \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{x}_*)\right)$$

where \mathbf{A} is now $\mathbf{A} \triangleq \frac{1}{\sigma_n^2} \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) + \Sigma_p^{-1} \in \mathbb{R}^{N \times N}$. From this, it becomes evident that the inverse of \mathbf{A} is required to compute both the mean and the covariance. This is not favourable since this would require knowledge of the feature space into which the feature map sends inputs. Moreover, computing \mathbf{A}^{-1} may become impractical if the dimension of the feature space is incredibly large. Remember, the whole point of the kernel trick is to avoid any computation that involves direct knowledge of \mathcal{H} but rather to use a kernel k to bypass these obstacles and indirectly produce inner products of the data applied to the feature map. With this in mind, let us try and find different expressions for the mean and the covariance of (7) that will enable us to apply the kernel trick. Before starting, we need to find a suitable expression for the mean. First define the notation

$$\mathbf{K}_{WW'} \triangleq \Phi(\mathbf{W}) \Sigma_p \Phi(\mathbf{W}')^\top \in \mathbb{R}^{n \times n'}$$

where $\mathbf{W} \in \mathbb{R}^{n \times d}$ and $\mathbf{W}' \in \mathbb{R}^{n' \times d}$ are two data matrices. Consider the following

$$\begin{aligned} & \mathbf{A} \Sigma_p \Phi(\mathbf{X})^\top \\ &= (\sigma_n^{-2} \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) + \Sigma_p^{-1}) \Sigma_p \Phi(\mathbf{X})^\top \\ &= \sigma_n^{-2} \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) \Sigma_p \Phi(\mathbf{X})^\top + \Phi(\mathbf{X})^\top \\ &= \sigma_n^{-2} \Phi(\mathbf{X})^\top (\Phi(\mathbf{X}) \Sigma_p \Phi(\mathbf{X})^\top + \sigma_n^2 \mathbb{1}_{n \times n}) \\ &= \sigma_n^{-2} \Phi(\mathbf{X})^\top (\mathbf{K}_{XX} + \sigma_n^2 \mathbb{1}_{n \times n}) \end{aligned}$$

meaning

$$\begin{aligned} \sigma_n^{-2} \Phi(\mathbf{X})^\top (\mathbf{K}_{XX} + \sigma_n^2 \mathbb{1}_{n \times n}) &= \mathbf{A} \Sigma_p \Phi(\mathbf{X})^\top \\ \sigma_n^{-2} \mathbf{A}^{-1} \Phi(\mathbf{X})^\top (\mathbf{K}_{XX} + \sigma_n^2 \mathbb{1}_{n \times n}) &= \Sigma_p \Phi(\mathbf{X})^\top \\ \sigma_n^{-2} \mathbf{A}^{-1} \Phi(\mathbf{X})^\top &= \Sigma_p \Phi(\mathbf{X})^\top (\mathbf{K}_{XX} + \sigma_n^2 \mathbb{1}_{n \times n})^{-1} \end{aligned}$$

so that the current mean of

$$\frac{1}{\sigma_n^2} \Phi(\mathbf{x}_*)^\top \mathbf{A}^{-1} \Phi(\mathbf{X})^\top \mathbf{y}$$

in (7) can be replaced with

$$\Phi(\mathbf{x}_*)^\top \Sigma_p \Phi(\mathbf{X})^\top (\mathbf{K}_{\mathbf{XX}} + \sigma_n^2 \mathbb{1}_{n \times n})^{-1} \mathbf{y}.$$

To find a more suitable expression for the covariance matrix, we will need the assistance of the matrix inversion lemma stated without proof in Lemma 12.

Lemma 12 (Matrix Inversion Lemma). *For $\mathbf{Z} \in \mathbb{K}^{n \times m}$, $\mathbf{W} \in \mathbb{K}^{m \times m}$ and $\mathbf{U}, \mathbf{V}^\top \in \mathbb{K}^{n \times m}$ then*

$$(\mathbf{Z} + \mathbf{UWV})^{-1} = \mathbf{Z}^{-1} - \mathbf{Z}^{-1} \mathbf{U} (\mathbf{W}^{-1} + \mathbf{VZ}^{-1} \mathbf{U})^{-1} \mathbf{VZ}^{-1}$$

assuming the relevant inverses exist [Pre92, page 75].

Proof. First note that

$$\mathbf{U} + \mathbf{UWVZ}^{-1} = \mathbf{UW} (\mathbf{W}^{-1} + \mathbf{VZU}) = (\mathbf{Z} + \mathbf{UWV}) \mathbf{Z}^{-1} \mathbf{U}$$

and

$$(\mathbf{Z} + \mathbf{UWV})^{-1} \mathbf{UW} = \mathbf{ZU} (\mathbf{W}^{-1} + \mathbf{VZ}^{-1} \mathbf{U})^{-1}.$$

Now

$$\begin{aligned} \mathbf{Z}^{-1} &= (\mathbf{Z} + \mathbf{UWV})^{-1} (\mathbf{Z} + \mathbf{UWV}) \mathbf{Z}^{-1} \\ &= (\mathbf{Z} + \mathbf{UWV})^{-1} (\mathbb{1}_{n \times n} + \mathbf{UWVZ}^{-1}) \\ &= (\mathbf{Z} + \mathbf{UWV})^{-1} + (\mathbf{Z} + \mathbf{UWV})^{-1} \mathbf{UWVZ}^{-1} \\ &= (\mathbf{Z} + \mathbf{UWV})^{-1} + \mathbf{Z}^{-1} \mathbf{U} (\mathbf{W}^{-1} + \mathbf{VZ}^{-1} \mathbf{U})^{-1} \mathbf{VZ}^{-1}. \end{aligned}$$

□

Consider

$$(8) \quad \mathbf{A} = \Sigma_p^{-1} + \Phi(\mathbf{X})^\top (\sigma_n^{-2} \mathbb{1}_{n \times n}) \Phi(\mathbf{X})$$

then applying the matrix inversion lemma by setting $\mathbf{Z}^{-1} = \Sigma_p$, $\mathbf{W}^{-1} = \sigma_n^2 \mathbb{1}_{n \times n}$ and $\mathbf{V} = \mathbf{U} = \Phi(\mathbf{X})$ (8) then becomes

$$\Sigma_p - \Sigma_p \Phi(\mathbf{X})^\top (\sigma_n^2 \mathbb{1}_{n \times n} + \Phi(\mathbf{X}) \Sigma_p \Phi(\mathbf{X})^\top)^{-1} \Phi(\mathbf{X}) \Sigma_p.$$

Thus (7) can be equivalently formulated as

$$(9) \quad f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\Phi(\mathbf{x}_*)^\top \Sigma_p \Phi(\mathbf{X})^\top (\mathbf{K}_{\mathbf{XX}} + \sigma_n^2 \mathbb{1}_{n \times n})^{-1} \mathbf{y}, \Phi(\mathbf{x}_*)^\top \Sigma_p \Phi(\mathbf{x}_*) - \Phi(\mathbf{x}_*)^\top \Sigma_p \Phi(\mathbf{X})^\top (\sigma_n^2 \mathbb{1}_{n \times n} + \Phi(\mathbf{X}) \Sigma_p \Phi(\mathbf{X})^\top)^{-1} \Phi(\mathbf{X}) \Sigma_p \Phi(\mathbf{x}_*)).$$

The astute reader may have noticed the very suggestive notation of labelling matrices of the form $\Phi(\mathbf{W}) \Sigma_p \Phi(\mathbf{W}')^\top$ as $\mathbf{K}_{\mathbf{WW}'}$ as though it may have some sort of connection to a kernel. To make this even more obvious, notice that each occurrence of the feature map in both expressions for the mean and covariance in (9) can be replaced with a $\mathbf{K}_{\mathbf{WW}'}$ for some appropriate choice of \mathbf{W} and \mathbf{W}' yielding the more notationally cleaner expression

$$(10) \quad f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mathbf{K}_{\mathbf{x}_* \mathbf{X}} (\mathbf{K}_{\mathbf{XX}} + \sigma_n^2 \mathbb{1}_{n \times n})^{-1} \mathbf{y}, k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{K}_{\mathbf{x}_* \mathbf{X}} (\sigma_n^2 \mathbb{1}_{n \times n} + \mathbf{K}_{\mathbf{XX}})^{-1} \mathbf{K}_{\mathbf{x}_* \mathbf{X}}^\top).$$

To get a better idea of the connection to kernels, since Σ_p is a symmetric positive semi definite matrix, it defines an inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\Sigma_p} = \mathbf{y}^\top \Sigma_p \mathbf{x}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{K}^N$$

[Wan, page 34] so that

$$(11) \quad (\mathbf{K}_{WW'})_{ij} = \langle \Phi(\mathbf{w}_i), \Phi(\mathbf{w}_j) \rangle_{\Sigma_p} = k(\mathbf{w}_i, \mathbf{w}_j)$$

where k is the kernel with feature map Φ and inner product $\langle \cdot, \cdot \rangle_{\Sigma_p}$. In fact, when $\mathbf{W} = \mathbf{W}'$, (11) is exactly the Gram matrix with said kernel. Thus GPs are another great example of models that take advantage of the kernel trick. We shall see in the coming chapters on how exactly we can compute predictions at novel points.

1.4. Gaussian Processes for Regression. We saw in Section 1.3 that, unlike most other machine learning models, GPs infer over a distribution of functions $p(f | \mathcal{D})$ instead of a vector of parameteric values $p(\boldsymbol{\theta} | \mathcal{D})$. Naively, one may attempt to find a suitable f by fixing a class of functions \mathcal{F} and then search over this class to find a function that best represents the data. However, this may not work well if there is not enough richness in \mathcal{F} to represent the data. Instead, a suitable f is selected by first assigning a prior probability to every possible function using the training data and then to select the function with the highest probability. To keep this computation tractable we only evalute our predicted function at a finite number of points. The prediction itself is found by taking the mean over all functions with respect to the posterior conditioned on the observed data which is assumed to be jointly Gaussian with the input value. This gives rise to Gaussian process more formally stated in Definition 13.

Definition 13 (Gaussian Process). *A Gaussian process (GP) is a collection of random variables with index set I , such that every finite subset of random variables has a joint Gaussian distribution [Ras06, Mur12].*

A GP is completely characterized by a mean function $m(\mathbf{x})$ and a kernel, which in the context of GPs is sometimes called a covariance function, $k(\mathbf{x}, \mathbf{x}')$ on a real process as

$$\begin{aligned} m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})] \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]. \end{aligned}$$

GPs define a prior over all possible functions which can be used to create a posterior once enough data has been observed. The prior is used to represent the functions we expect to see before any observations are made. Although defining a prior over all possible functions may seem computationally intractable, we actually only need to define a distribution over a finite number of points. Before any observations are made, we typically assume that the mean function is the constant zero function, that is $m(\mathbf{x}) = 0$. A function $f(\mathbf{x})$ sampled from a GP with mean $m(\mathbf{x})$ and covariance $k(\mathbf{x}, \mathbf{x}')$ is written as

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

Since a GP is a collection of random variables it must satisfy the consistency requirement, that is, observing some of the values should not change the distribution of any small subset of unobserved values. More specifically if

$$(\mathbf{y}_1, \mathbf{y}_2) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

then

$$\begin{aligned} \mathbf{y}_1 &\sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{1,1}) \\ \mathbf{y}_2 &\sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{2,2}) \end{aligned}$$

where $\boldsymbol{\Sigma}_{1,1}$ and $\boldsymbol{\Sigma}_{2,2}$ are the relevant sub matrices. Again, we shall us the notation that for set of data $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n]^\top \in \mathbb{R}^{n \times d}$ and $\mathbf{W}' = [\mathbf{w}'_1, \mathbf{w}'_2, \dots, \mathbf{w}'_m]^\top \in \mathbb{R}^{n' \times d}$ we use the notation

$$(K_{WW'})_{i,j} \triangleq k(\mathbf{w}_i, \mathbf{w}'_j)$$

where $K_{WW'} \in \mathbb{R}^{n \times n'}$. The covariance function completely characterized by its kernel. Unless otherwise stated, the kernel or covariance function used in examples and experimentation is the Gaussian RBF kernel, Definition 11.

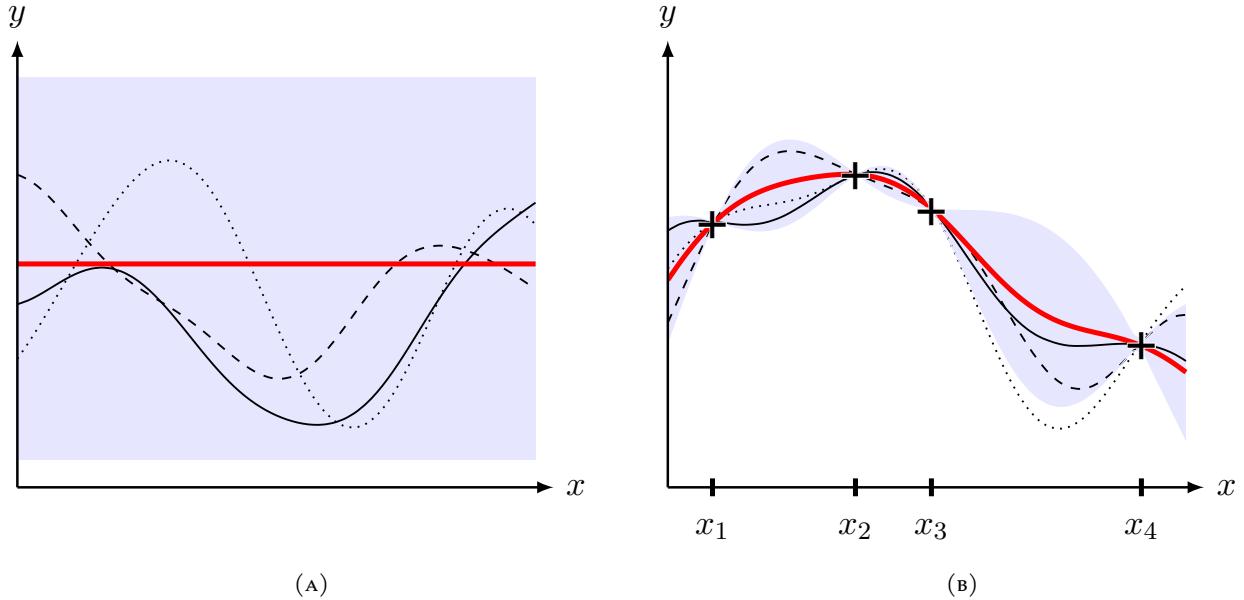


FIGURE 4. Panel (A) shows three function drawn from the prior distribution. Panel (B) shows three function drawn from the prior distribution after four observations have been made. In both panels the mean function is drawn in red, sampled functions in black and twice the standard deviation shaded in light blue.

Figure 4 (A) shows three samples drawn from the prior before any observations are made. GPs also allow us to compute the pointwise variance which can provide some measure of variability for predicted values. The blue shaded area of Figure 4 (A) represents twice the standard deviation about the mean.

1.4.1. Noise-free observations. Typically when using GP we would like to incorporate data from observations, or training data, into our predictions on unobserved values. Let us suppose there is some obsevered data $\mathcal{D} = \{(\mathbf{x}_i, f_i) \mid i \in \{1, 2, \dots, n\}\}$ which is (unrealistically) noise-free that we would like to model as a GP. In other words, for any sample in our dataset we can be certain that the observed value is the true

value of the underlying function we wish to model. Then for the observed data

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{XX}}).$$

We would then like to make a prediction for unobserved values say $\mathbf{X}_* = [\mathbf{x}_{1*}, \mathbf{x}_{2*}, \dots, \mathbf{x}_{n*}]$ with value f_* as has a distribution of

$$\mathbf{f}_* \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{X}_*\mathbf{X}_*}).$$

Here \mathbf{f} and \mathbf{f}_* are independent but we would like to give them some sort of correlation. We can do this by having them originate from the same joint distribution. According to the prior, we can write the joint distribution of the training points \mathbf{f} and the test points \mathbf{f}_* as

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{XX}} & \mathbf{K}_{\mathbf{X}_*\mathbf{X}}^\top \\ \mathbf{K}_{\mathbf{X}_*\mathbf{X}} & \mathbf{K}_{\mathbf{X}_*\mathbf{X}_*} \end{bmatrix}\right).$$

While the above does give us some information that \mathbf{f}_* is related to the observed data and the test inputs, it does not provide any method to evaluate \mathbf{f}_* . To do this we shall need the assistance of the following theorem.

Theorem 14. (*Marginals and conditionals of an multivariate normal distributions* [Mur12, page 142]) Suppose $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2]$ is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$$

then the posterior conditional is given by

$$\begin{aligned} \mathbf{x}_2 | \mathbf{x}_1 &\sim \mathcal{N}(\mathbf{x}_2 | \boldsymbol{\mu}_{2|1}, \boldsymbol{\Sigma}_{2|1}) \\ \boldsymbol{\mu}_{2|1} &= \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \\ \boldsymbol{\Sigma}_{2|1} &= \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12}. \end{aligned}$$

Proof. First define $\mathbf{z} = \mathbf{x}_2 - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-2} \mathbf{x}_1 \triangleq \mathbf{x}_2 + \mathbf{A} \mathbf{x}_1$, so that

$$\begin{aligned} \text{cov}(\mathbf{z}, \mathbf{x}_1) &= \text{cov}(\mathbf{x}_2, \mathbf{x}_1) + \text{cov}(\mathbf{A} \mathbf{x}_1, \mathbf{x}_1) \\ &= \boldsymbol{\Sigma}_{21} + \mathbf{A} \mathbb{V}[\mathbf{x}_1] \\ &= \boldsymbol{\Sigma}_{21} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-2} \boldsymbol{\Sigma}_{11} \\ &= 0. \end{aligned}$$

This shows that \mathbf{z} and \mathbf{x}_1 are uncorrelated and, since they are also jointly normal, also independent. Now, since $\mathbb{E}[\mathbf{z}] = \boldsymbol{\mu}_2 + \mathbf{A} \boldsymbol{\mu}_1$, it follows that

$$\begin{aligned} \boldsymbol{\mu}_{2|1} &= \mathbb{E}[\mathbf{x}_2 | \mathbf{x}_1] \\ &= \mathbb{E}[\mathbf{z} - \mathbf{A} \mathbf{x}_1 | \mathbf{x}_1] \\ &= \mathbb{E}[\mathbf{z} | \mathbf{x}_1] - \mathbb{E}[\mathbf{A} \mathbf{x}_1 | \mathbf{x}_1] \\ &= \boldsymbol{\mu}_2 + \mathbf{A}(\boldsymbol{\mu}_1 - \mathbf{x}) \\ &= \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-2} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \end{aligned}$$

Which shows the first part. Next

$$\begin{aligned}
\Sigma_{2|1} &= \mathbb{V}[\mathbf{x}_2 \mid \mathbf{x}_1] \\
&= \mathbb{V}[\mathbf{z} - \mathbf{A}\mathbf{x}_1 \mid \mathbf{x}_1] \\
&= \mathbb{V}[\mathbf{z} \mid \mathbf{x}_1] + \mathbb{V}[\mathbf{A}\mathbf{x}_1 \mid \mathbf{x}_1] - \mathbf{A} \operatorname{cov}(\mathbf{z}, -\mathbf{x}_1) - \operatorname{cov}(\mathbf{z}, -\mathbf{x}_1) \mathbf{A}^\top \\
&= \mathbb{V}[\mathbf{z} \mid \mathbf{x}_1] \\
&= \mathbb{V}[\mathbf{z}].
\end{aligned}$$

Finally

$$\begin{aligned}
\Sigma_{2|1} &= \mathbb{V}[\mathbf{z}] \\
&= \mathbb{V}[\mathbf{x}_2 + \mathbf{A}\mathbf{x}_1] \\
&= \mathbb{V}[\mathbf{x}_2] + \mathbf{A}\mathbb{V}[\mathbf{x}_1]\mathbf{A}^\top + \mathbf{A} \operatorname{cov}(\mathbf{x}_2, \mathbf{x}_1) + \operatorname{cov}(\mathbf{x}_1, \mathbf{x}_2)\mathbf{A}^\top \\
&= \Sigma_{22} + \Sigma_{21}\Sigma_{11}^{-2}\Sigma_{11}\Sigma_{11}^{-2}\Sigma_{12} - 2\Sigma_{21}\Sigma_{11}^{-2}\Sigma_{12} \\
&= \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-2}\Sigma_{12}
\end{aligned}$$

as wanted. \square

Thus, once data has been observed, finding a mean and covariance for f_* involves a direct application of Theorem 14 which gives

$$f_* \mid \mathbf{K}_{X_* X}^\top, \mathbf{K}_{XX}, f \sim \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

where

$$\begin{aligned}
\boldsymbol{\mu}_* &= \mathbf{0} + \mathbf{K}_{X_* X} \mathbf{K}_{XX}^{-1} (\mathbf{f} - \mathbf{0}) \\
&= \mathbf{K}_{X_* X} \mathbf{K}_{XX}^{-1} \mathbf{f}
\end{aligned}$$

and

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{X_* X_*} - \mathbf{K}_{X_* X} \mathbf{K}_{XX}^{-1} \mathbf{K}_{X_* X}$$

meaning we can write a distribution for f_* as

$$(12) \quad f_* \mid \mathbf{K}_{X_* X}^\top, \mathbf{K}_{XX}, f \sim \mathcal{N}\left(\mathbf{K}_{X_* X} \mathbf{K}_{XX}^{-1} \mathbf{f}, \mathbf{K}_{X_* X_*} - \mathbf{K}_{X_* X} \mathbf{K}_{XX}^{-1} \mathbf{K}_{X_* X}^\top\right).$$

Function values from the unobserved inputs X_* , that is f_* , can be estimated using the joint posterior distribution by evaluating the mean of 12. Figure 4 (B) shows these computations given a data set $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)\}$. Notice that the variance tightens around the observed values since (assuming no noise in our data is present) we now know for certain this is how our target function should behave at x_1, x_2, x_3 and x_4 . Clearly, specifying the properties of the prior is important since it fixes the properties of the functions considered during inference.

1.4.2. *Prediction with Noisy observations.* When attempting to model our value function we usually do not have access to the value function itself but a noisy version thereof, $y = f(\mathbf{x}) + \varepsilon$ where $\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$ meaning the prior on the noisy values becomes

$$\text{cov}(\mathbf{y}) = \mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{I}.$$

The reason why noise is only added along the diagonal follows from the assumption of independence of noise in our data. We can write out the new distribution of the observed noisy values along the points at which we wish to test the underlying function as

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n} & \mathbf{K}_{\mathbf{X}_*\mathbf{X}}^\top \\ \mathbf{K}_{\mathbf{X}_*\mathbf{X}} & \mathbf{K}_{\mathbf{X}_*\mathbf{X}_*} \end{bmatrix} \right).$$

Using a similar we arrive at a similar condition distribution of $\mathbf{f}_* | \mathbf{K}_{\mathbf{X}_*\mathbf{X}}^\top, \mathbf{K}_{\mathbf{X}\mathbf{X}}, \mathbf{f}$ we arrive at one of the most fundamental equations for GP regression tasks

$$(13) \quad \mathbf{f}_* | \mathbf{K}_{\mathbf{X}_*\mathbf{X}}^\top, \mathbf{K}_{\mathbf{X}\mathbf{X}}, \mathbf{y} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*))$$

where

$$(14) \quad \bar{\mathbf{f}}_* \triangleq \mathbf{K}_{\mathbf{X}_*\mathbf{X}} [\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n}]^{-1} \mathbf{y}$$

$$(15) \quad \text{cov}(\mathbf{f}_*) = \mathbf{K}_{\mathbf{X}_*\mathbf{X}_*} - \mathbf{K}_{\mathbf{X}_*\mathbf{X}} [\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n}]^{-1} \mathbf{K}_{\mathbf{X}_*\mathbf{X}}^\top$$

Remarkably, this gives us the the exact same posterior distribution ascertained from the weight space derivation in (10). Notice that the prediction of the mean in (14) is a linear combination of the observations, sometimes referred to as a *linear predictor*. Another way of looking at the prediction is seeing it as a linear combination of n kernel evaluations centered at the input \mathbf{x}_*

$$\mathbf{f}_* = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}_*)$$

where $\boldsymbol{\alpha} = [\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n}]^{-1} \mathbf{y}$. Intuitively, this expression can be understood by realising that, despite defining the GP using a joint Gaussian distribution over the observations, when making predictions GPs only care about the $(n+1)$ -dimension distribution defined by the n observations and the single test point. When the GP is marginalized by taking the relevant submatrix block of the covariance matrix, this conditioning gives us our desired 1-dimensional prediction.

Also notice that the covariance does not depend on observations but scales quadratically to the norm of the testing inputs. This is a key feature of GPs. The variance is comprised of the difference between the prior covariance, $\mathbf{K}_{\mathbf{X}_*\mathbf{X}_*}$, and positive term $\mathbf{K}_{\mathbf{X}_*\mathbf{X}} [\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n}]^{-1} \mathbf{K}_{\mathbf{X}_*\mathbf{X}}^\top$ which represents knowledge given by the observations about the underlying function.

Algorithm 1 shows one possible implementation for computing the mean and covariance of a single test input.

Algorithm 1: Unoptimized GPR

input : Observations \mathbf{X}, \mathbf{y} and a test input \mathbf{x}_* .
output: A prediction \bar{f}_* with its corresponding variance $\mathbb{V}[f_*]$.

- 1 $\mathbf{L} = \text{cholesky}(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n})$
- 2 $\boldsymbol{\alpha} = \text{lin-solve}(\mathbf{L}^\top, \text{lin-solve}(\mathbf{L}, \mathbf{y}))$
- 3 $\bar{f}_* = \mathbf{K}_{\mathbf{x}_*\mathbf{X}} \boldsymbol{\alpha}$
- 4 $\mathbf{v} = \text{lin-solve}(\mathbf{L}, \mathbf{K}_{\mathbf{x}_*\mathbf{X}})$
- 5 $\mathbb{V}[f_*] = \mathbf{K}_{\mathbf{x}_*\mathbf{x}_*} - \mathbf{v}^\top \mathbf{v}$
- 6 **return** $\bar{f}_*, \mathbb{V}[f_*]$

A Cholesky decomposition is typically used since \mathbf{L} can be used twice to assist in solving both the linear systems in the mean and covariance. Unfortunately, a Cholesky decomposition incurs a runtime of $\mathcal{O}(n^3)$ where n is the number of samples [Tre97, page 176] making it impractical for large data sets. In the later chapters we shall consider other methods for solving these linear systems.

1.5. Gaussian Processes for Classification. As with most classification models, the Gaussian processes classifier (GPC) seeks an estimate for the joint probability $p(y, \mathbf{x})$ where $\mathbf{x} \in \mathbb{R}^d$ is an input, as in the regression case, but y is now a class taking on a discrete and finite number of values $\{\mathcal{C}_i\}_{i=1}^C$. Using Baye's theorem the joint probability density can be decomposed into either $p(y)p(\mathbf{x} | y)$ or $p(\mathbf{x})p(\mathbf{y} | \mathbf{x})$ giving rise to the *generative* and *discriminative* approaches respectively [Ras06, page 34]. The generative approach models the prior probabilities of each class, $p(\mathcal{C}_i)$, as well as the class conditional probabilities for each input $p(\mathbf{x} | \mathcal{C}_i)$ and computes the posterior as

$$p(y | \mathbf{x}) = \frac{p(y)p(\mathbf{x} | y)}{\sum_{i=1}^C p(\mathcal{C}_i)p(\mathbf{x} | \mathcal{C}_i)}.$$

On the other hand, the discriminative method focuses on modelling $p(y | \mathbf{x})$ directly. With both these paradigms at our disposal, which one would be preferred for our GPC? While there are strengths and weaknesses associated with both models, the discriminative approach is usually chosen as it has a rather attractive property of directly modeling what we require, that is $p(y | \mathbf{x})$. Additionally, the density estimation of $p(\mathbf{x} | \mathcal{C}_i)$ using in the generative model presents a number of difficulties, especially for larger values of d . If we are only focused on classifying inputs, the generative approach could mean trying to solve a harder problem than what is necessary. For this reason we focus on GPCs that adopt the discriminative approach.

1.5.1. Linear Models for Classification. We can start by reviewing linear models for the simplest form of classification, that is binary classification. Adopting the notation from SVM (see Section 1.3.1) literature, the binary classification problem involves assigning an input \mathbf{x} to a class of either -1 or $+1$. For a linear model likelihood can be formulated as

$$(16) \quad p(y = +1 | \mathbf{x}, \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$$

given a weight vector \mathbf{w} and where $\sigma(z)$ is chosen to be any sigmoid function, see Definition 15.

Definition 15 (Sigmoid Function). *A sigmoid function is a monotonically increasing function mapping from \mathbb{R} to $[0, 1]$* [Ras06, page 35].

In this text, the commonly used logistic function

$$(17) \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

will take the role of the sigmoid function in (16), graphed in Figure 5. This type of model is aptly named the logistic regression. Unlike GPR, the likelihood is no longer a Gaussian distribution. Instead it follows

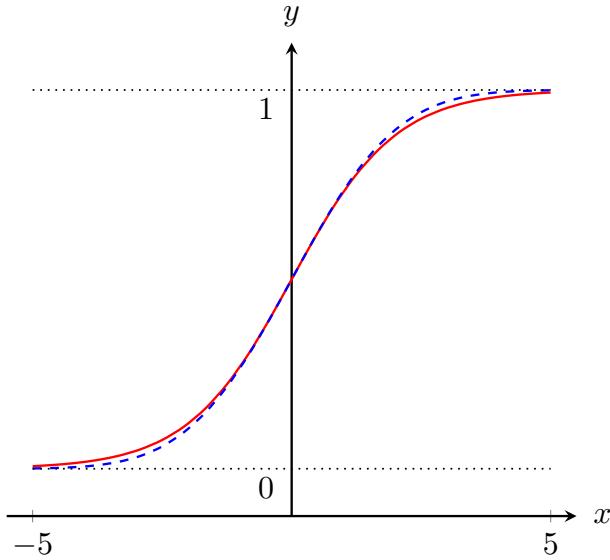


FIGURE 5. The logistic function from (17) (solid red) juxtaposed with a close approximation, the scaled probit function (dashed blue).

the Bernoulli distribution

$$p(y | \mathbf{x}, \mathbf{w}) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)^{\frac{1+y}{2}} (1 - \sigma(\langle \mathbf{x}, \mathbf{w} \rangle))^{\frac{1-y}{2}}$$

which for symmetric likelihood functions can be written more concisely as

$$p(y_i | \mathbf{x}_i, \mathbf{w}) = \sigma(y_i f_i)$$

where

$$(18) \quad f_i \triangleq f(\mathbf{x}_i) = \langle \mathbf{x}, \mathbf{w} \rangle.$$

Thus, the logistic regression model can be written as the log ratio of the likelihoods of the input belonging to either class, that is

$$\text{logit}(\mathbf{x}) \triangleq \langle \mathbf{x}, \mathbf{w} \rangle = \log \left(\frac{p(y = +1)}{p(y = -1)} \right)$$

where logit is commonly referred to as the logit transformation [Ras06, page 37]. For a given dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ we assume each observation is independently generated conditioned over $f(\mathbf{x})$. Similar

to GPR, a Gaussian prior is used for the weights so that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_p)$ giving an un-normalised log posterior of

$$\log p(\mathbf{w} | \mathbf{X}, \mathbf{y}) \propto -\frac{1}{2}\mathbf{w}^\top \Sigma_p^{-1} \mathbf{w} + \sum_{i=1}^n \log \sigma(y_i f_i).$$

However, unlike GPR an analytic form for the mean and variance for the posterior is not available due to the non-Gaussian nature of the likelihood, although, when using the logistic function it is easy enough to show that the log likelihood is concave as a function of \mathbf{w} for a fixed dataset. This means a number of numerical optimization techniques, such as Newton's method or the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm [Fle00] can be used to solve these values.

The idea behind Gaussian process classification for binary classes is that a Gaussian process prior is placed over a latent function $f(x)$ where the output is then "squashed" through a sigmoid function to obtain a prior on

$$\pi(\mathbf{x}) \triangleq p(y = +1 | \mathbf{x}) = \sigma(f(\mathbf{x})).$$

This construction is illustrated in Figure 6 and provides a natural extension to the linear logistic regression model.

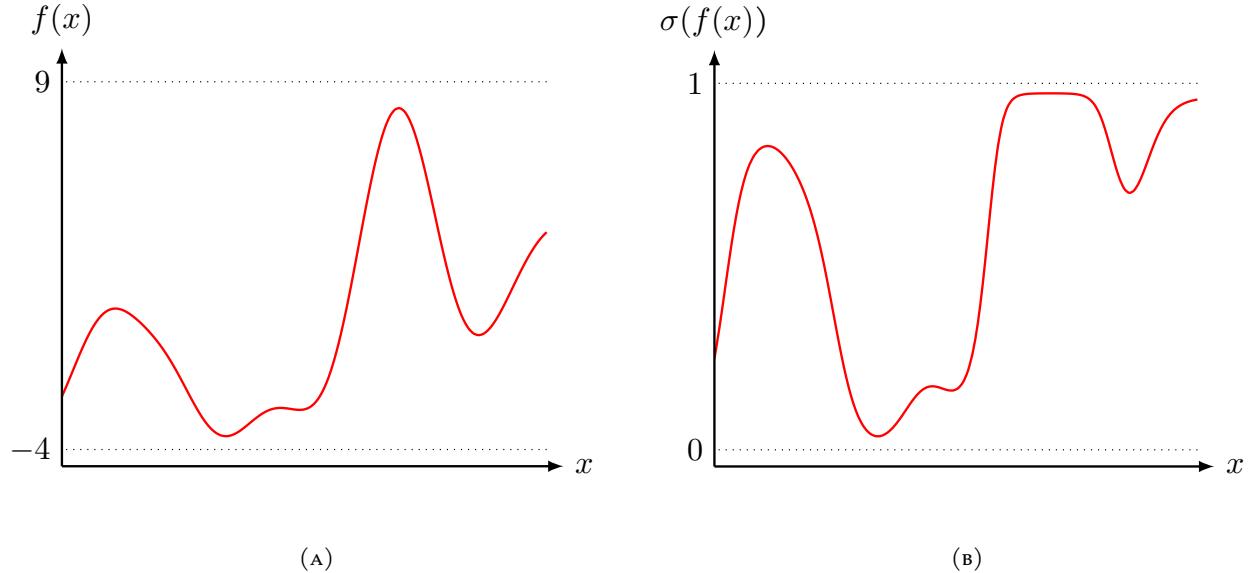


FIGURE 6. The latent function f , panel (A), is transformed using a sigmoid function, panel (B), to provide a probabilistic interpretation of x belonging to the class +1.

Specifically, the linear model from (18) is replaced with a GPR model and the Gaussian prior on the weights with a GPR weight prior with

$$p\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix}\right) = \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_{XX} & \mathbf{K}_{x^* X}^\top \\ \mathbf{K}_{x^* X} & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix}\right)$$

where $f_* = f(\mathbf{x}_*)$ and $\mathbf{f} = f(\mathbf{X})$. For classification tasks, we assume that each observation has received the correct label which is why no noise is added to the covariance matrix.

Note that values of f are also never observed within the phenomena we are modelling, nor are we particularly interested in them. The function f serves the role of a *nuisance function* and acts solely as a convenience tool within our formulations. The ultimate goal is to make predictions for π , not f , and the goal of the coming sections will be to eventually integrate out f .

Subsequently, predictions for $\pi_\star = \pi(\mathbf{x}_\star)$ are made by averaging over all possible latent functions weighted by the posterior giving the prediction

$$(19) \quad \bar{\pi}_\star \triangleq p(y_\star = +1 \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_\star) = \int \sigma(f_\star) p(f_\star \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_\star) df_\star$$

While this is a sound model, computing predictions is not so straight forward since the integral in (19) is not analytically tractable as applying σ to f_\star no longer produces a Gaussian distribution. Later on we will see how we can make use of our numerical toolbox to derive a good approximation for $\bar{\pi}_\star$.

1.5.2. Lapace Approximation for Posterior. We saw that the integral in (19) could not be used to make predictions for $\bar{\pi}_\star$ analytically. In this section we shall address how the distribution for the latent process, $p(f_\star \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_\star)$, can be approximated to provide a numerically tractable succedaneum. Using Baye's theorem

$$\begin{aligned} p(f_\star \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_\star) &= \int p(f_\star, \mathbf{f} \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_\star) d\mathbf{f} \\ &= \frac{1}{p(\mathbf{y} \mid \mathbf{X}, \mathbf{x}_\star)} \int p(f_\star \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f}) p(\mathbf{f} \mid \mathbf{X}) p(\mathbf{y} \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f}) d\mathbf{f} \\ &= \int p(f_\star \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f}) p(\mathbf{f} \mid \mathbf{X}, \mathbf{y}) d\mathbf{f} \end{aligned}$$

using the fact that $p(\mathbf{y} \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f}, f_\star) = p(\mathbf{y} \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f})$ [Bis06, Ras06]. The conditional distribution $p(f_\star \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f})$ can be derived as

$$p(f_\star \mid \mathbf{X}, \mathbf{x}_\star, \mathbf{f}) = \mathcal{N}\left(f_\star \mid \mathbf{K}_{\mathbf{x}_\star \mathbf{X}} \mathbf{K}_{\mathbf{X} \mathbf{X}}^{-1} \mathbf{y}, k(\mathbf{x}_\star, \mathbf{x}_\star) - \mathbf{K}_{\mathbf{x}_\star \mathbf{X}} \mathbf{K}_{\mathbf{X} \mathbf{X}}^{-1} \mathbf{K}_{\mathbf{x}_\star \mathbf{X}}^\top\right)$$

through the use of (14) and (15). Unfortunately

$$p(\mathbf{f} \mid \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} \mid \mathbf{f}) p(\mathbf{f} \mid \mathbf{X})}{p(\mathbf{y} \mid \mathbf{X})}$$

does not follow a Gaussian distribution. Instead we can use a Lapace approximation to estimate $p(\mathbf{f} \mid \mathbf{X}, \mathbf{y})$ as a Gaussian distribution. Briefly, the Lapace approximation works by assuming the distribution at hand, $p(\mathbf{z})$, can be modelled as

$$p(\mathbf{z}) = \frac{1}{c} q(\mathbf{z})$$

where $q(\mathbf{z})$ is multivariate Gaussian and c is some normalization constant [Bis06, page 214]. To do this, first the centre of $q(\mathbf{z})$ is placed at the mode of $p(\mathbf{z})$. The mode of $p(\mathbf{z})$ is

$$\mathbf{z}_0 = \arg \min_{\mathbf{z}} p(\mathbf{z})$$

which can be computed by solving

$$(20) \quad \nabla p(\mathbf{z}_0) = \mathbf{0}.$$

To ensure the covariance of the synthesized multivariate Gaussian behaves similar to the original distribution we can make use of an important property of the Gaussian distribution which is its logarithm

being a quadratic function of its inputs. Taking the Taylor series expansion of $\ln q(\mathbf{z})$ centered at \mathbf{z}_0 yields

$$\ln q(\mathbf{z}) \simeq \ln q(\mathbf{z}_0) - \frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0)$$

where

$$\mathbf{A} = -\nabla \nabla \ln q(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0}.$$

Exponentiating both sides gives

$$(21) \quad \begin{aligned} q(\mathbf{z}) &\simeq q(\mathbf{z}_0) \exp\left(-\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0)\right) \\ &\propto \mathcal{N}(\mathbf{z} | \mathbf{z}_0, \mathbf{A}^{-1}). \end{aligned}$$

Returning to our original problem of estimating $p(\mathbf{f} | \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | \mathbf{X})$ as a Gaussian distribution, the prior $p(\mathbf{f} | \mathbf{X})$ follows a Gaussian distribution with zero mean and covariance $\mathbf{K}_{\mathbf{XX}}$ and the distribution of $p(\mathbf{y} | \mathbf{f})$ (assuming independence of samples) can be written as

$$p(\mathbf{y} | \mathbf{f}) = \prod_{i=1}^n \sigma(y_i f_i).$$

To find a Laplace approximation for $p(\mathbf{f} | \mathbf{X}, \mathbf{y})$ we only need to consider an unnormalized posterior when maximizing with respect to \mathbf{f} since $p(\mathbf{y} | \mathbf{f})$ does not depend on \mathbf{f} . Thus, the log of the unnormalized posterior is

$$\begin{aligned} \Psi(\mathbf{f}) &\triangleq \ln p(\mathbf{y} | \mathbf{f}) + \ln p(\mathbf{f} | \mathbf{X}) \\ &= -\sum_{i=1}^n \ln(1 + \exp(y_i f_i)) - \frac{1}{2} \mathbf{f}^T \mathbf{K}_{\mathbf{XX}}^{-1} \mathbf{f} - \frac{1}{2} \ln |\mathbf{K}_{\mathbf{XX}}| - \frac{n}{2} \ln 2\pi. \end{aligned}$$

The gradient and Hessian of the unnormalized posterior then becomes

$$\begin{aligned} \nabla \Psi(\mathbf{f}) &= \nabla \ln p(\mathbf{y} | \mathbf{f}) - \mathbf{K}_{\mathbf{XX}}^{-1} \mathbf{f} = (\mathbf{t} - \boldsymbol{\pi}) - \mathbf{K}_{\mathbf{XX}}^{-1} \mathbf{f} \\ \nabla \nabla \Psi(\mathbf{f}) &= \nabla \nabla \ln p(\mathbf{y} | \mathbf{f}) - \mathbf{K}_{\mathbf{XX}}^{-1} = -\mathbf{W} - \mathbf{K}_{\mathbf{XX}}^{-1} \end{aligned}$$

where $\pi_i = p(y_i = +1 | f_i) = \sigma(f_i)$, $\mathbf{t} = (\mathbf{y} + \mathbf{1}) / 2 \in \mathbb{R}^n$ and $\mathbf{W} \triangleq -\nabla \nabla \ln p(\mathbf{y} | \mathbf{f})$ is a diagonal matrix (since the distribution of y_i only depends on f_i and not $f_{j \neq i}$) with entries $W_{ii} = \sigma(y_i f_i)$ [Bis06, Ras06]. From (20), the mode of $\hat{\mathbf{f}}$ of Ψ can be computed as

$$(22) \quad \begin{aligned} \nabla \Psi(\hat{\mathbf{f}}) &= \mathbf{0} = (\mathbf{t} - \boldsymbol{\pi}) - \mathbf{K}_{\mathbf{XX}}^{-1} \hat{\mathbf{f}} \\ \iff \hat{\mathbf{f}} &= \mathbf{K}_{\mathbf{XX}}(\mathbf{t} - \boldsymbol{\pi}). \end{aligned}$$

Since $\mathbf{t} - \boldsymbol{\pi}$ is a non-linear function, a non-linear optimization technique method is required to solve $\hat{\mathbf{f}}$ in (22). Since the Hessian of $\Psi(\mathbf{f})$ is available, Newton's method is typically employed as fast iterative method to approximate $\hat{\mathbf{f}}$ where $\hat{\mathbf{f}}$ is updated as

$$\hat{\mathbf{f}}^{\text{new}} = \mathbf{K}_{\mathbf{XX}} (\mathbb{1}_{n \times n} + \mathbf{W} \mathbf{K}_{\mathbf{XX}})^{-1} (\mathbf{W} \hat{\mathbf{f}}^{\text{old}} + \nabla \ln(\mathbf{y} | \hat{\mathbf{f}}^{\text{old}})).$$

Once a suitable mode is found, using (21), the Lapacian approximation for $p(\mathbf{f} | \mathbf{X}, \mathbf{y})$ becomes

$$(23) \quad p(\mathbf{f} | \mathbf{X}, \mathbf{y}) \simeq q(\mathbf{f} | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\hat{\mathbf{f}}, (\mathbf{K}_{\mathbf{XX}}^{-1} + \mathbf{W})^{-1}).$$

1.5.3. *Predictions.* With the Laplace approximation for $p(\mathbf{f} | \mathbf{X}, \mathbf{y})$ ((23)) and an exact probability distribution for $p(f_\star | \mathbf{X}, \mathbf{x}_\star, \mathbf{f})$, a mean for the latent process, $p(f_\star | \mathbf{X}, \mathbf{y}, \mathbf{x}_\star)$, can now be computed by invoking (14) to give

$$(24) \quad \begin{aligned} \mu_{f_\star} &= \mathbb{E}[f_\star | \mathbf{X}, \mathbf{y}, \mathbf{x}_\star] = \mathbf{K}_{\mathbf{x}_\star \mathbf{X}} \mathbf{K}_{\mathbf{X} \mathbf{X}}^{-1} \hat{\mathbf{f}} \\ &= \mathbf{K}_{\mathbf{x}_\star \mathbf{X}} \nabla \ln(\mathbf{y} | \hat{\mathbf{f}}) \\ &= \mathbf{K}_{\mathbf{x}_\star \mathbf{X}} (\mathbf{t} - \boldsymbol{\pi}). \end{aligned}$$

Similarly, the variance can be computed using (15) to give

$$(25) \quad \sigma_{f_\star}^2 = \mathbb{V}[f_\star | \mathbf{X}, \mathbf{y}, \mathbf{x}_\star] = k(\mathbf{x}_\star, \mathbf{x}_\star) - \mathbf{K}_{\mathbf{x}_\star \mathbf{X}} (\mathbf{K}_{\mathbf{X} \mathbf{X}} + \mathbf{W}^{-1})^{-1} \mathbf{K}_{\mathbf{x}_\star \mathbf{X}}^\top.$$

Using (19), predictions can now be made as

$$(26) \quad \bar{\pi}_\star \simeq \int \sigma(f_\star) q(f_\star | \mathbf{X}, \mathbf{y}, \mathbf{x}_\star) df_\star$$

where $q(f_\star | \mathbf{X}, \mathbf{y}, \mathbf{x}_\star)$ is a multivariate Gaussian distribution with mean and variance given by equations (24) and (25) respectively. Notice that the prediction given in (26) is a convolution of a Gaussian and logistic function which unfortunately cannot be evaluated analytically. However, Spiegelhalter and Lauritzen [Spi90] show that a good approximation can be found by replacing the sigmoid function with the probit function $\Phi(\lambda a)$ which is simply the cumulative distribution function (CDF) of the standard Gaussian distribution. To get the best approximation using the probit function, the constant factor λ is adjusted to equate their slopes at the origin. The value of λ that gives this equality is $\lambda = \sqrt{\pi}/8$. The similarity between the sigmoid function and probit function rescaled by a factor of $\sqrt{\pi}/8$ is illustrated in Figure 5. The reason for replacing the sigmoid function with a probit function is that the convolution of a Gaussian distribution and probit function can be analytically evaluated as

$$(27) \quad \int \Phi(\lambda a) \mathcal{N}(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{\frac{1}{2}}}\right).$$

Again apply the approximation $\sigma(a) \simeq \Phi(\lambda a)$ to left hand side of (27) gives the following estimate for the convolution of a Gaussian and sigmoid function

$$(28) \quad \int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da \simeq \sigma\left(\frac{\mu}{(1 + \pi\sigma^2/8)^{\frac{1}{2}}}\right)$$

[Bis06, page 219]. The integral used to approximate $\bar{\pi}_\star$ in (26) can now be estimated using (28) to give

$$\bar{\pi}_\star = \sigma\left(\frac{\mu_{f_\star}}{(1 + \pi\sigma_{f_\star}^2/8)^{\frac{1}{2}}}\right).$$

This theory justifies Algorithm 2 which creates predictions based on the GPC method.

Algorithm 2: Unoptimized GPC

input : Observations \mathbf{X}, \mathbf{y} and a test input \mathbf{x}^* .
output: A prediction \bar{f}_* with its corresponding variance $\mathbb{V}[f_*]$.

```

1  $t = (\mathbf{y} + \mathbf{1}) / 2$ 
2  $f = \mathbf{0}$ 
3 repeat
4    $\mathbf{W} = \text{diag}(\sigma(\mathbf{y}.^*\mathbf{f}))$ 
5    $\boldsymbol{\alpha} = \text{lin-solve}(\mathbf{1}_{n \times n} + \mathbf{W}\mathbf{K}_{\mathbf{X}\mathbf{X}}, \mathbf{K}_{\mathbf{X}\mathbf{X}})$ 
6    $\mathbf{f} = \boldsymbol{\alpha}(t - \sigma(\mathbf{f}) + \mathbf{W}\mathbf{f})$ 
7 until convergence
8  $\mu_{f_*} = \mathbf{K}_{\mathbf{x}_*\mathbf{X}}(t - \sigma(\mathbf{f}))$ 
9  $\sigma_{f_*}^2 = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{K}_{\mathbf{x}_*\mathbf{X}}(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \mathbf{W}^{-1})^{-1}\mathbf{K}_{\mathbf{x}_*\mathbf{X}}^\top$ 
10  $\bar{\pi}_* = \sigma\left(\mu_{f_*}/\left(1 + \pi\sigma_{f_*}^2/8\right)^{\frac{1}{2}}\right)$ 
11 return  $\bar{\pi}_*, \mu_{f_*}, \sigma_{f_*}^2$ 
```

2. THE NYSTROM METHOD

In Chapter 1 we saw that GP regression and classification relied on a Gram matrix (see Definition 4) to produce predictions. Unfortunately, from a computational perspective, constructing the Gram matrix for a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ brings about a nasty bottle neck on account of the $\mathcal{O}(n^2)$ kernel evaluations. Even before the rise of ML, a lot of research was devoted to creating numerical methods that quickly construct a low rank approximation of large matrices, \mathbf{A} , which ordinarily are a computational burden to build exactly. These methods are based on the idea of capturing the columns space of the matrix that best describes the the action of \mathbf{A} as an operator. Mahoney provides an enlightened summary as to why the column space is of paramount importance in these approximation techniques

"To understand why sampling columns (or rows) from a matrix is of interest, recall that matrices are "about" their columns and rows that is, linear combinations are taken with respect to them; one all but understands a given matrix if one understands its column space, row space, and null spaces; and understanding the subspace structure of a matrix sheds a great deal of light on the linear transformation that the matrix represents."

[MWM11, page 13]

Moreover, this class of algorithms lend very nice forms when \mathbf{A} possess positive definite structure, which is exactly the case for Gram matrices.

2.1. The Nyström Method. Attempting to compute an entire kernel matrix can be a computational headache prompting an investigation of estimative alternatives. The approximation techniques studied in this chapter have been spurred on by the John-Lindenstrauss lemma stated in Lemma 16.

Lemma 16 (John-Lindenstrauss). *Given $0 < \varepsilon < 0$, any set of n points, X , in a high dimensional Euclidean space can be embedded into a ℓ -dimensional Euclidean space where $\ell = \mathcal{O}(\ln(n))$ via some linear map $\Omega \in \mathbb{R}^{n \times \ell}$ which satisfies*

$$(1 - \varepsilon) \|\mathbf{u} - \mathbf{v}\|^2 \leq \|\Omega\mathbf{u} - \Omega\mathbf{v}\|^2 \leq \varepsilon \|\mathbf{u} - \mathbf{v}\|^2$$

for any $\mathbf{u}, \mathbf{v} \in X$ [MWM11, page 15].

The John-Lindenstrauss lemma tells us that $\mathbf{Q}\mathbf{Q}^*\mathbf{A}$ will serve as a good approximation to some matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ where $\mathbf{Q}\mathbf{Q}^*$, in some sense, projects onto a rank- k subspace of \mathbf{A} 's column space. This is because if $\mathbf{Q}\mathbf{Q}^*$ closely matches the behavior of Ω from the lemma then the pair-wise distances between points before and after applying $\mathbf{Q}\mathbf{Q}^*$ should remain fairly similar. To state this a little more explicitly, for a matrix \mathbf{A} and a positive error tolerance ε we seek a matrix $\mathbf{Q} \in \mathbb{R}^{n \times k_\varepsilon}$ with orthonormal columns such that

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|_F \leq \varepsilon$$

which can be expressed in a more short hand notation as

$$(29) \quad \mathbf{A} \simeq \mathbf{Q}\mathbf{Q}^*\mathbf{A}.$$

This is commonly called the *fixed precision approximation problem*. To simplify algorithmic development, a value of k is specified in advance (instead of ε , thus removing k 's dependence on ε) which is instead given the name *fixed rank problem*. Within the fixed rank problem framework, when \mathbf{A} is hermitian, the

matrix $\mathbf{Q}\mathbf{Q}^*$ acts as a good projection for both the columns and row space of \mathbf{A} so that we have both $\mathbf{A} \simeq \mathbf{Q}\mathbf{Q}^*\mathbf{A}$ and $\mathbf{A} \simeq \mathbf{A}\mathbf{Q}\mathbf{Q}^*$ meaning

$$(30) \quad \mathbf{A} \simeq \mathbf{Q}\mathbf{Q}^*(\mathbf{A}) \simeq \mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q}\mathbf{Q}^*.$$

Furthermore, if \mathbf{A} is positive semi-definite we can improve the quality of our approximation of our approximation at almost no additional cost [Hal11, page 32]. Using the approximation from (30)

$$\begin{aligned} (31) \quad \mathbf{A} &\simeq \mathbf{Q}(\mathbf{Q}^*\mathbf{A}\mathbf{Q})\mathbf{Q}^* \\ &= \mathbf{Q}(\mathbf{Q}^*\mathbf{A}\mathbf{Q})(\mathbf{Q}^*\mathbf{A}\mathbf{Q})^\dagger(\mathbf{Q}^*\mathbf{A}\mathbf{Q})\mathbf{Q}^* \\ &\simeq (\mathbf{A}\mathbf{Q})(\mathbf{Q}^*\mathbf{A}\mathbf{Q})^\dagger(\mathbf{Q}^*\mathbf{A}). \end{aligned}$$

This is known as the Nystrom method. Since any Gram matrix is positive semi-definite, we can always applied the Nystrom method to find an approximation to it. A general Nystrom framework is presented in Algorithm 3.

Algorithm 3: General Nystrom Framework

input : A positive semi-definite matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, a matrix $\mathbf{Q} \in n \times k$ that satisfies (29).

output: A rank k approximation $\bar{\mathbf{A}} \simeq \mathbf{A}$.

- 1 $\mathbf{C} = \mathbf{A}\mathbf{Q}$
 - 2 $\mathbf{W} = \mathbf{Q}^*\mathbf{C}$
 - 3 **return** $\mathbf{C}\mathbf{W}^\dagger\mathbf{C}^*$
-

However, Algorithm 3 assumes that \mathbf{Q} has already been computed. Naturally, the next question is then how to efficiently construct a suitable matrix \mathbf{Q} that satisfies (29)? We can do this through a very popular column sampling technique ubiquitous in numerical linear algebra literature. This technique has been driven by Theorem 17.

Theorem 17. Every $\mathbf{A} \in \mathbb{R}^{n \times m}$ matrix contains a k -column submatrix \mathbf{C} for which

$$\|\mathbf{A} - \mathbf{C}\mathbf{C}^\dagger\mathbf{A}\|_F \leq \sqrt{1 + k(n - k)} \cdot \|\mathbf{A} - \mathbf{A}_k\|$$

where \mathbf{A}_k is the best rank- k approximation of \mathbf{A} [Hal11, page 11].

Before we delve further into this column sampling Nystrom method, we must first cover the random matrix multiplication algorithm which serves as a backbone for this technique. Therefore, let $\mathbf{A} \in \mathbb{R}^{n \times m}$ be a target matrix we would like to approximate and suppose that \mathbf{A} can be represented as the sum of 'simpler' (for example, sparse or low-rank) matrices, \mathbf{A}_i , so that

$$(32) \quad \mathbf{A} = \sum_{i=1}^I \mathbf{A}_i.$$

The basic idea is to consider a Monte-Carlo approximation of (32) that randomly selects \mathbf{A}_i according to the distribution $\{p_i\}_{i=1}^I$ to give an estimate

$$(33) \quad \mathbf{A} \simeq \frac{1}{c} \sum_{t=1}^c p_{i_t}^{-1} \mathbf{A}_{i_t}$$

where c is the number of samples and each summand is rescaled by a factor of $p_{i_t}^{-1}$ to ensure our estimate is unbiased [PGMaJT21, pages 24-27]. The random matrix multiplication algorithm works by attempting to find a Monte-Carlo estimate for \mathbf{AB} , where $\mathbf{A} \in \mathbb{R}^{n \times I}$ and $\mathbf{B} \in \mathbb{R}^{I \times m}$. Recall that any matrix multiplication can be written in its outer product form

$$\mathbf{AB} = \sum_{i=1}^I \mathbf{A}_{(:,i)} \mathbf{B}_{(i,:)}$$

[FR20, Dri06]. A straight forward way to approximate this using the Monte-Carlo estimate is to simply set each \mathbf{A}_i in (32) to the corresponding rank-1 outer product summand $\mathbf{A}_{(:,i)} \mathbf{B}_{(i,:)}$. This justifies the random matrix multiplication algorithm seen in Algorithm 4 [PDaMWM17, page 16].

Algorithm 4: Random Matrix Multiplication

input : $\mathbf{A} \in \mathbb{R}^{n \times I}$ and $\mathbf{B} \in \mathbb{R}^{I \times m}$, the number of samples $1 \leq c \leq I$ and a probability distribution over I , $\{p_i\}_{i=1}^I$.

output: Matrices $\mathbf{C} \in \mathbb{R}^{n \times c}$ and $\mathbf{R} \in \mathbb{R}^{c \times m}$ such that $\mathbf{CR} \simeq \mathbf{AB}$.

```

1 for  $t = 1, \dots, c$  do
2   Pick  $i_t \in \{1, \dots, I\}$  with  $\mathbb{P}[i_t = k] = p_k$ , independently and with replacement.
3    $\mathbf{C}_{(:,t)} = \frac{1}{\sqrt{cp_{i_t}}} \mathbf{A}_{(:,i_t)}$ 
4    $\mathbf{R}_{(t,:)} = \frac{1}{\sqrt{cp_{i_t}}} \mathbf{B}_{(i_t,:)}$ 
5 end
6 return  $\mathbf{CR} = \sum_{t=1}^c \frac{1}{cp_{i_t}} \mathbf{A}_{(:,i_t)} \mathbf{B}_{(i_t,:)}$ 

```

This algorithm makes this idea a little more precise, taking in the two matrices to multiply together as well as a probability distribution over I to provide an estimate for \mathbf{AB} of the form

$$\mathbf{AB} \simeq \sum_{t=1}^c \frac{1}{cp_{i_t}} \mathbf{A}_{(:,i_t)} \mathbf{B}_{(i_t,:)}.$$

Equivalently, the above can be restated as the product of two matrices \mathbf{CR} formed by Algorithm 4, where \mathbf{C} consists of c randomly selected rescaled columns of \mathbf{A} and \mathbf{R} is c randomly selected rescaled rows of \mathbf{B} . Notice that

$$\mathbf{CR} = \sum_{t=1}^c \mathbf{C}_{(:,i_t)} \mathbf{R}_{(i_t,:)} = \sum_{t=1}^c \left(\frac{1}{\sqrt{cp_{i_t}}} \mathbf{A}_{(:,i_t)} \right) \left(\frac{1}{\sqrt{cp_{i_t}}} \mathbf{B}_{(i_t,:)} \right) = \frac{1}{c} \sum_{t=1}^c \frac{1}{p_{i_t}} \mathbf{A}_{(:,i_t)} \mathbf{B}_{(i_t,:)}.$$

To make development easier, let us define a sampling and rescaling matrix, usually referred to as a sketching matrix, $\mathbf{S} \in \mathbb{R}^{n \times c}$ to be the matrix with elements $S_{it,t} = 1/\sqrt{cp_{i_t}}$ if the i_t column of \mathbf{A} is chosen during the t^{th} trial and all other entries of \mathbf{S} are set to 0. Then we have

$$\mathbf{C} = \mathbf{AS} \quad \text{and} \quad \mathbf{R} = \mathbf{S}^\top \mathbf{B}$$

so that

$$(34) \quad \mathbf{C}\mathbf{R} = \mathbf{A}\mathbf{S}\mathbf{S}^\top\mathbf{B} \simeq \mathbf{AB}.$$

Notice that \mathbf{S} is generally a very sparse matrix and therefore is generally not constructed explicitly and instead the matrix products \mathbf{AS} and $\mathbf{S}^\top\mathbf{B}$ are done through row and column rescaling [PDAWMW17, page 17]. Lemma 18 provides some bounds on \mathbf{CR} as an estimate for \mathbf{AB} .

Lemma 18. *Let \mathbf{C} and \mathbf{R} be constructed as described in Algorithm 4, then*

$$\mathbb{E}[(\mathbf{CR})_{ij}] = (\mathbf{AB})_{ij}.$$

That is, \mathbf{CR} is an unbiased estimate of \mathbf{AB} . Furthermore

$$\mathbb{V}[(\mathbf{CR})_{ij}] \leq \frac{1}{c} \sum_{k=1}^n \frac{\mathbf{A}_{ik}^2 \mathbf{B}_{kj}^2}{p_k}.$$

Proof. For some fixed pair i, j for each $t = 1, \dots, c$ define $\mathbf{X}_t = \left(\frac{\mathbf{A}_{(:,i_t)} \mathbf{B}_{(i_t,:)}}{cp_{i_t}} \right)_{ij} = \frac{\mathbf{A}_{(i,i_t)} \mathbf{B}_{(i_t,j)}}{cp_{i_t}}$. Thus, for any t ,

$$\mathbb{E}[\mathbf{X}_t] = \sum_{k=1}^n p_k \frac{\mathbf{A}_{ik} \mathbf{B}_{kj}}{cp_k} = \frac{1}{c} \sum_{k=1}^n \mathbf{A}_{ik} \mathbf{B}_{kj} = \frac{1}{c} (\mathbf{AB})_{ij}.$$

Since we have $(\mathbf{CR})_{ij} = \sum_{t=1}^c \mathbf{X}_t$, it follows that

$$\mathbb{E}[(\mathbf{CR})_{ij}] = \mathbb{E}\left[\sum_{t=1}^c \mathbf{X}_t\right] = \sum_{t=1}^c [\mathbb{E}\mathbf{X}_t] = (\mathbf{AB})_{ij}.$$

Hence, \mathbf{CR} is an unbiased estimator of \mathbf{AB} , regardless of the choice of the sampling probabilities. Using the fact that $(\mathbf{CR})_{ij}$ is the sum of c independent random variables, we get

$$\mathbb{V}[(\mathbf{CR})_{ij}] = \mathbb{V}\left[\sum_{t=1}^c \mathbf{X}_t\right] = \sum_{t=1}^c \mathbb{V}[\mathbf{X}_t].$$

Using the fact $\mathbb{V}[\mathbf{X}_t] \leq \mathbb{E}[\mathbf{X}_t^2] = \sum_{k=1}^n \frac{\mathbf{A}_{ik}^2 \mathbf{B}_{kj}^2}{c^2 p_k}$, we get

$$\mathbb{V}[(\mathbf{CR})_{ij}] = \sum_{t=1}^c \mathbb{V}[\mathbf{X}_t] \leq c \sum_{k=1}^n \frac{\mathbf{A}_{ik}^2 \mathbf{B}_{kj}^2}{c^2 p_k} = \frac{1}{c} \frac{\mathbf{A}_{ik}^2 \mathbf{B}_{kj}^2}{p_k}.$$

□

So how does this help us with the Nystrom method? Consider using the random matrix multiplication algorithm to approximate the matrix multiplication of a Gram matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ and $\mathbb{1}_{n \times n}$. (34) gives

$$\mathbf{K}\mathbf{S}\mathbf{S}^\top\mathbb{1}_{n \times n} = \mathbf{K}\mathbf{S}\mathbf{S}^\top \simeq \mathbf{K}.$$

We see now that the sketching matrix produced by Algorithm 4 provides a sketching matrix \mathbf{S} that satisfies the properties of \mathbf{Q} from (29) meaning that \mathbf{S} can be used in place of \mathbf{Q} within the Nystrom estimate from (31). These ideas are used together in Algorithm 5 that uses the column sampling technique from

Algorithm 4 together with the general Nystrom framework (Algorithm 3) lends a new column sampling Nystrom method to approximate a Gram matrix for a provided dataset and probability distribution [PDaMWM05, AGaMWM13].

Algorithm 5: Nystrom Method via Column Sampling

input : Data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$, the number of samples $1 \leq c \leq n$ and a probability distribution over n , $\{p_i\}_{i=1}^n$.

output: An approximation of the Gram matrix corresponding to \mathbf{X} , that is $\bar{\mathbf{K}} \simeq \mathbf{K}$ where $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$.

- 1 Initialize \mathbf{C} as an empty $n \times c$ matrix.
- 2 Pick c columns with the probability of choosing the k^{th} column ($1 \leq k \leq n$) as $\mathbb{P}[k = i] = p_i$, independently and with replacement and let I a list of indices of the sampled columns.
- 3 **for** $i \in I$ **do**
- 4 $\mathbf{K}_{(:,i)} = [k(\mathbf{x}_1, \mathbf{x}_i), \dots, k(\mathbf{x}_n, \mathbf{x}_i)]^\top$
- 5 $\mathbf{C}_{(:,i)} = \mathbf{K}_{(:,i)} / \sqrt{cp_i}$
- 6 **end**
- 7 $\mathbf{W} = \mathbf{K}_{(I,I)} \in \mathbb{R}^{c \times c}$
- 8 Rescale each entry of \mathbf{W} , \mathbf{W}_{ij} , by $1/c\sqrt{p_ip_j}$.
- 9 Compute \mathbf{W}^\dagger
- 10 **return** $\mathbf{C}\mathbf{W}^\dagger\mathbf{C}^*$

As we can tell from the algorithms inputs, this requires some sort of probability distribution to select the columns. As seen in Lemma 18 any probability distribution will provide an unbiased estimate. However, some probability distributions can be used to lower the variance faster than others. Naively, we could just employ uniform sampling where each column is selected with equal probability although it should be cautioned that this is seldom a good idea since uniform sampling tends to over sample landmarks from one large cluster while under sampling or even missing entire small but important clusters. As a result, the approximation for \mathbf{K} will decline [CMaCM17, page 3]. This is demonstrated in graphic form in Figure 7.

To combat this issue, alternative probabilities density can be constructed to take into account a measure of importance in landmark selection. Indeed there has been a plethora of research that has shown the importance of using data-dependent non-uniform probability distributions to obtain provably better error bounds within the Nystrom framework [PDaMWM05, AGaMWM13, CMaCM17, PDe11, MBCaCMaCM15, Kum09]. A few of the more common distributions will be discussed in the coming sections.

2.2. Column Probabilities. Recall that the Nystrom method from Algorithm 5 is largely dependent on the random matrix multiplication algorithm (see Algorithm 4) to produce a suitable sketching matrix. Moreover, improvements in the sketching matrix produced by the random matrix multiplication algorithm are reflected as smaller errors in the Nystrom approximation. Now, consider using the random matrix multiplication algorithm to approximate $\mathbf{A}\mathbf{A}^\top$ by setting $\mathbf{B} = \mathbf{A}$. The output is an approximation

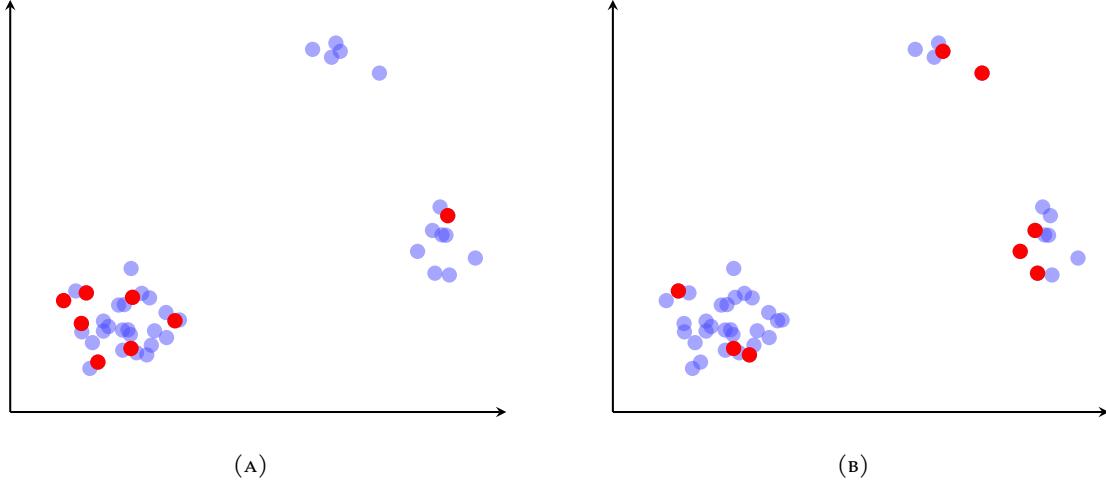


FIGURE 7. Employing uniform sampling in the column sampling Nystrom estimate can lead to oversampling from denser parts of the data set (Panel (A)). Instead data dependent probability densities are commonly used to better cover the relevant data (Panel (B)). Example taken from [CMaCM17, page 4].

of the form

$$\mathbf{A}\mathbf{A}^\top \simeq \mathbf{C}\mathbf{C}^\top = \mathbf{C}\mathbf{R}.$$

The probability distribution

$$p_i = \frac{\|\mathbf{A}_{(:,i)}\|_2^2}{\|\mathbf{A}\|_F}.$$

aims to minimize the error between $\mathbf{A}\mathbf{A}^\top$ and the approximation $\mathbf{C}\mathbf{C}^\top$. As a result, we should expect that \mathbf{C} becomes a better estimate for $\mathbf{A}\mathbf{S}$, implying that the sketching matrix, \mathbf{S} , is using a better sampling and landmark selection criteria. Drineas and Mahoney give a precise bound on this error presented in Theorem 19 [PDaMWM05, page 2158].

Theorem 19. *Given $\mathbf{A} \in \mathbb{R}^{n \times I}$, $1 \leq c \leq I$ and the probability distribution $\{p_i\}_{i=1}^I$ described in Section 2.2. Construct \mathbf{C} using Algorithm 4, then*

$$\mathbb{E} [\|\mathbf{A}\mathbf{A}^\top - \mathbf{C}\mathbf{C}^\top\|_F] \leq \frac{1}{\sqrt{c}} \|\mathbf{A}\|_F^2$$

[PDaMWM05, page 2158].

To show Theorem 19, we can actually prove something a little more general.

Lemma 20. *Given $\mathbf{A} \in \mathbb{R}^{n \times I}$, $\mathbf{B} \in \mathbb{R}^{I \times m}$, $1 \leq c \leq I$ and the probability distribution $\{p_i\}_{i=1}^I$ as follows*

$$p_i = \frac{\|\mathbf{A}_{(:,i)}\|_2 \|\mathbf{B}_{(i,:)}\|_2}{\sum_{j=1}^I \|\mathbf{A}_{(:,j)}\|_2 \|\mathbf{B}_{(j,:)}\|_2}.$$

Construct \mathbf{C} using Algorithm 4, using the probability distribution described above, then

$$\mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F] \leq \frac{1}{\sqrt{c}} \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2.$$

This choice of probability distribution minimises $\mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F]$ among all possible sampling probabilities [Dri06, pages 9-12].

Proof. First note that

$$\mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F^2] = \sum_{k=1}^n \sum_{j=1}^m \mathbb{E} [(AB - CR)_{kj}^2] = \sum_{k=1}^n \sum_{j=1}^m \mathbb{V} [(CR)_{kj}].$$

Thus from Lemma 18, it follows that

$$\begin{aligned} & \mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F^2] \\ &= \frac{1}{c} \sum_{i=1}^I \frac{1}{p_i} \left(\sum_{k=1}^n A_{ki}^2 \right) \left(\sum_{j=1}^m B_{ij}^2 \right) - \frac{1}{c} \|\mathbf{AB}\|_F^2 \\ &= \frac{1}{c} \sum_{i=1}^I \frac{1}{p_i} \|A_{(:,i)}\|_2^2 \|B_{(i,:)}\|_2^2 - \frac{1}{c} \|\mathbf{AB}\|_F^2. \end{aligned}$$

Substituting in a probability of

$$p_i = \frac{\|A_{(:,i)}\|_2 \|B_{(i,:)}\|_2}{\sum_{j=1}^I \|A_{(j,:)}\|_2 \|B_{(:,j)}\|_2},$$

yields

$$\begin{aligned} \mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F^2] &= \frac{1}{c} \left(\sum_{i=1}^I \|A_{(:,i)}\|_2 \|B_{(i,:)}\|_2 \right)^2 - \frac{1}{c} \|\mathbf{AB}\|_F^2 \\ &\leq \frac{1}{c} \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2. \end{aligned}$$

To verify that this choice of probability distribution minimises $\mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F]$ define the function

$$f(p_1, \dots, p_n) = \sum_{i=1}^I \frac{1}{p_i} \|A_{(:,i)}\|_2^2 \cdot \|B_{(i,:)}\|_2^2$$

which characterises the dependence of $\mathbb{E} [\|\mathbf{AB} - \mathbf{CR}\|_F]$ on the probability distribution. To minimise f subject to $\sum_{i=1}^I p_i = 1$, we introduce the Lagrange multiplier λ and define the function

$$g(p_1, \dots, p_n) = f(p_1, \dots, p_n) + \lambda \left(\sum_{i=1}^I p_i - 1 \right).$$

The minimum is then

$$0 = \frac{\partial g}{\partial p_i} = -\frac{1}{p_i^2} \|A_{(:,i)}\|_2^2 \cdot \|B_{(i,:)}\|_2^2 + \lambda.$$

Thus

$$p_i = \frac{\|A_{(:,i)}\|_2 \cdot \|B_{(i,:)}\|_2}{\sqrt{\lambda}} = \frac{\|A_{(:,i)}\|_2 \cdot \|B_{(i,:)}\|_2}{\sum_{j=1}^I \|A_{(j,:)}\|_2 \|B_{(:,j)}\|_2}$$

where the second equality comes from solving for $\sqrt{\lambda}$ in $\sum_{i=1}^{I-1} p_i = 1$. These probabilities are indeed minimizing since $\frac{\partial^2 g}{\partial p_i^2} > 0$ for every i such that $\|A_{(:,i)}\|_2^2 \cdot \|B_{(i,:)}\|_2^2 > 0$. \square

2.3. Leverage Scores.

2.3.1. *Statistical Leverage Scores.* Our next distribution originates from the least-squares problem. Briefly, in an over constrained least-squares problem, where $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{b} \in \mathbb{R}^n$, for $m \ll n$ there usually are not any $\mathbf{x} \in \mathbb{R}^m$ for which $\mathbf{Ax} = \mathbf{b}$. Instead, alternative criteria must be employed to seek a \mathbf{x} which in some way comes "closest" to satisfying this equality. Perhaps one of the more popular criterion is to minimize the ℓ^2 -norm, that is

$$\mathbf{x}_{opt} = \arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|$$

[MWM11, page 19-21]. The optimal value for \mathbf{x} can be solved as $\mathbf{x}_{opt} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$. The least-squares solution is commonly used to find the best weight vector (in this case \mathbf{x}) for a linear model, given a dataset. Fitted or predicted values are usually obtained from $\hat{\mathbf{b}} = \mathbf{H}\mathbf{b}$ where the projector onto the column space of \mathbf{A}

$$\mathbf{H} = \mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$$

is sometimes referred to as the *hat matrix*. The element \mathbf{H}_{ij} has the direct interpretation as the influence or statistical leverage exerted on $\hat{\mathbf{b}}_i$. Thus, examining the hat matrix can reveal to us columns of \mathbf{A} which bear a significant impact on $\hat{\mathbf{b}}$ [Hoa78, page 17]. Relatedly, if the element \mathbf{H}_{ii} is particularly large this is indicative of the i^{th} column of \mathbf{A} having a strong influence over values of $\hat{\mathbf{b}}$, justifying the interpretation of \mathbf{H}_{ii} as *statistical leverage scores*.

The statistical leverage scores are maximised when $\mathbf{A}_{(:,i)}$ is linearly independent from \mathbf{A} 's other columns and decreases when it aligns with many other columns or when the value of $\|\mathbf{A}_{(:,i)}\|$ is small [MBCaC-MaCM15, page 5]. To compute the statistical leverage scores, if $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$ is the SVD of \mathbf{A} , then

$$\begin{aligned} \mathbf{H}_{ii} &= (\mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top)_{ii} \\ &= (\mathbf{U}\Sigma^2(\Sigma^2)^{-1}\mathbf{U})_{ii} \\ &= \|\mathbf{U}_{(i,:)}\|_2^2. \end{aligned}$$

Note that \mathbf{H}_{ii} may not constitute as a probability distribution, as may the other leverage scores which we will soon discuss. This is easily remedied through normalisation, in this case dividing each statistical leverage score by $\text{tr}(\mathbf{H})$. The idea behind using statistical leverage scores as a probability distribution in the Nystrom method is that they help prioritize selection of columns that are more linearly independent from other columns so that the range of our approximate better aligns with the range of our original \mathbf{A} .

2.3.2. *Rank- k Statistical Leverage Scores.* We can generalize this notion of statistical leverage scores to include lower rank approximations. Let $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the compact SVD of a \mathbf{A} real $n \times m$ matrix. Setting $r = \min\{n, m\}$, the compact SVD can be partitioned as

$$\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2] \in \mathbb{R}^{n \times r}, \quad \Sigma = \begin{bmatrix} \Sigma_1 & \\ & \Sigma_2 \end{bmatrix} \in \mathbb{R}^{r \times r}, \quad \mathbf{V} = [\mathbf{V}_1, \mathbf{V}_2] \in \mathbb{R}^{m \times r}.$$

Here \mathbf{U}_1 contains the first $k \leq r$ columns of \mathbf{U} , \mathbf{V}_1 the first k rows of \mathbf{V} and Σ_1 is a $k \times k$ matrix containing the top k singular values across its diagonal. The matrix $\mathbf{A}_k = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1$ serves as the best rank- k approximation to \mathbf{A} . The statistical leverage scores relative to the best rank- k approximation are again \mathbf{H}_{ii} , but this time \mathbf{H} is computed only using the best rank- k approximation of \mathbf{A} , that is \mathbf{A}_k . These low

rank scores can be evaluated as

$$\ell_i^k \triangleq \left(\mathbf{A}_k (\mathbf{A}_k^\top \mathbf{A}_k)^{-1} \mathbf{A}_k^\top \right)_{ii} = \left\| (\mathbf{U}_1)_{(i,:)} \right\|_2^2.$$

What makes low-rank statistical leverage scores particularly appealing is that they can be approximated quickly with a truncated SVD [AGaMWM13, pages 3-4].

2.3.3. Ridge Leverage Scores. The low rank leverage scores we saw in Section 2.3.2 will not always be unique and can be sensitive to perturbations [MBCaCMaCM15, page 6]. Consequently these scores can vary drastically when \mathbf{A} is modified slightly or when we only have access to partial information on the matrix. This undermines the possibility of computing good quality low rank approximations from statistical leverage scores. This shortcoming is addressed in the next class of leverage score, that is, ridge leverage scores. Ridge leverage scores are similar to statistical leverage scores although a ridge regression term (hence the name) is added to the hat matrix with a regularization parameter λ . The λ -ridge leverage score is defined as

$$r_i^\lambda \triangleq \left(\mathbf{A} (\mathbf{A}^\top \mathbf{A} + \lambda \mathbb{1}_{n \times n})^{-1} \mathbf{A}^\top \right)_{ii}.$$

A regularization parameter of

$$\lambda = \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}$$

is typically used since this choice of λ will guarantee that the sum of the ridge leverage scores (keep in mind that the raw ridge leverage scores do not necessarily form a probability distribution) is bounded by $2k$, stated more formally in Lemma 21.

Lemma 21. *When using a regularization parameter of $\lambda = \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}$ we have $\sum_{i=1}^n r_i^\lambda \leq 2k$ [MBCaCMaCM15, pages 6-7].*

Proof. Writing r_i^λ using the SVD of \mathbf{A} where $\lambda = \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}$ gives

$$\begin{aligned} r_i^\lambda &= \mathbf{A}_{(i,:)} \left(\mathbf{U} \Sigma \mathbf{U}^\top + \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k} \mathbf{U} \mathbf{U}^\top \right)^{-1} \mathbf{A}_{(i,:)}^\top \\ &= \mathbf{A}_{(i,:)} \left(\mathbf{U} \bar{\Sigma}^2 \mathbf{U}^\top \right)^{-1} \mathbf{A}_{(i,:)}^\top \\ &= \mathbf{A}_{(i,:)} \left(\mathbf{U} \bar{\Sigma}^{-2} \mathbf{U}^\top \right) \mathbf{A}_{(i,:)}^\top \end{aligned}$$

where $\bar{\Sigma}_{ii}^2 = \sigma_i^2(\mathbf{A}) + \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}$. Then

$$\begin{aligned} \sum_{i=1}^n r_i^\lambda &= \text{tr} \left(\mathbf{A}^\top \mathbf{U} \bar{\Sigma}^{-2} \mathbf{U}^\top \mathbf{A} \right) \\ &= \text{tr} \left(\mathbf{V} \Sigma \bar{\Sigma}^{-2} \Sigma \mathbf{V}^\top \right) \\ &= \text{tr} \left(\Sigma^2 \bar{\Sigma}^{-2} \right). \end{aligned}$$

Here we have

$$\left(\Sigma^2 \bar{\Sigma}^{-2} \right)_{ii} = \frac{\sigma_i^2(\mathbf{A})}{\sigma_i^2(\mathbf{A}) + \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}}.$$

For $i \leq k$ we simply upper bound this by 1, yielding

$$\text{tr}(\Sigma^2 \bar{\Sigma}^{-2}) = k + \sum_{i=k+1}^n \frac{\sigma_i^2(\mathbf{A})}{\sigma_i^2(\mathbf{A}) + \frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}} \leq k + \sum_{i=k+1}^n \frac{\sigma_i^2(\mathbf{A})}{\frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}} = k + \frac{\sum_{i=k+1}^n \sigma_i^2(\mathbf{A})}{\frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k}} \leq k + k.$$

□

From now on (unless otherwise stated) the regularization parameter seen in Section 2.3.3 will always be used for ridge leverage scores where the notation

$$r_i^k \triangleq \left(\mathbf{A} \left(\mathbf{A}^\top \mathbf{A} + \left(\frac{\|\mathbf{A} - \mathbf{A}_k\|_F^2}{k} \right) \mathbb{1}_{n \times n} \right)^{-1} \mathbf{A}^\top \right)_{ii}$$

is employed to show that the best rank- k matrix is utilized in the regularization parameter. Adding regularization to the hat matrix offers a smoother alternative which ‘washes out’ small singular directions meaning they are sampled with proportionally lower probability [MBCaCMaCM15, page 6]. Alaoui and Mahoney [Ala15] prove that ridge leverage scores provide theoretically better bounds over uniform sampling techniques when the number of sampled columns is proportional to $\text{tr}(\mathbf{H}_\lambda) \cdot \ln(n)$ where \mathbf{H}_λ is the hat matrix with added regularization, that is $\mathbf{H}_\lambda = \mathbf{A} (\mathbf{A}^\top \mathbf{A} + \lambda \mathbb{1}_{n \times n})^{-1} \mathbf{A}^\top$. With the rising popularity of ridged leverage scores, a number of iterative methods have been devised (and continue to be developed) that take advantage of parallel computing to provide fast approximations [PGMaJT21, page 90].

2.4. Data Probabilities. Here we present two novel Nystrom probability sampling distributions that make use of the data matrix, \mathbf{X} , directly without requiring any knowledge of the associated kernel matrix. The first of these two methods is the data-column sampling distribution defined simply as

$$p_i = \|\mathbf{X}_{(i,:)}\|_2 / \|\mathbf{X}\|_F^2.$$

This sampling distribution can be thought of as the data-column dual of the kernel column sampling probabilities seen in Section 2.2. The second sampling method is the \mathbf{Q} matrix column sampling distribution defined as

$$p_i = \|\mathbf{Q}_{(i,:)}\|_2.$$

where $\mathbf{X} = \mathbf{Q}\mathbf{R}$ is the QR factorisation (see Theorem 35) of \mathbf{X} . This sampling technique can be thought of as the data matrix equivalent of leverage score sampling. There are numerous benefits computing the sampling probabilities using the data matrix directly without needing any knowledge about the corresponding matrix. To start, the obvious benefit is that no properties or elements (along with their corresponding calculations) of the kernel matrix are required saving large amounts of computation time. Second, if the parameters of the kernel function are changed, or an entirely different kernel function is used altogether, then the data-column sampling probabilities need not be re-evaluated.

The main grounds for the theoretical justification for using the data matrix directly to produce a sampling distribution comes from a result shown in [Kar10], presented below.

Theorem 22 (Inner Product of Kernel Matrices). *Let us assume that we observe n i.i.d. random vectors, $X_i \in \mathbb{R}^p$. let us consider the kernel matrix \mathbf{M} with entries*

$$\mathbf{M}_{ij} = f\left(\frac{\mathbf{X}_i^\top \mathbf{X}_j}{p}\right).$$

We assume that

- (a) We have n/p and p/n remain bounded as $n \rightarrow \infty$.
- (b) Σ_p is a positive semi-definite $p \times p$ matrix, and $\|\Sigma_p\|_2 = \sigma_1(\Sigma_p)$ remains bounded in p , that is, there exists $K > 0$, such that $\sigma_1(\Sigma_p) \leq K$, for all p .
- (c) Σ_p/p has a finite limit, that is, there exists $l \in \mathbb{R}$ such that $\lim_{p \rightarrow \infty} \text{tr}(\Sigma_p)/p = l$.
- (d) $\mathbf{X}_i = \Sigma_p^{1/2} \mathbf{Y}_i$.
- (e) The entries of \mathbf{Y}_i , a p -dimensional random vector, are i.i.d. Also, denoting by $\mathbf{Y}_i(k)$ the k^{th} entry of \mathbf{Y}_i , we assume that $\mathbb{E}(\mathbf{Y}_i(k)) = 0$, $\mathbb{V}(\mathbf{Y}_i(k)) = 1$ and $\mathbb{E}(|\mathbf{Y}_i(k)|^{4+\varepsilon}) < \infty$ for some $\varepsilon > 0$.
- (f) f is a C^1 function in a neighborhood of $l = \lim_{p \rightarrow \infty} \text{tr}(\Sigma_p)/p$ and a C^3 function in a neighborhood of 0.

Under these assumptions, the kernel matrix \mathbf{M} can (in probability) be approximated consistently in the operator norm, when p and n tend to ∞ , by the matrix \mathbf{K} , where

$$\mathbf{K} = \left(f(0) + f''(0) \frac{\text{tr}(\Sigma_p^2)}{2p^2} \right) \mathbf{1}\mathbf{1}' + f'(0) \frac{\mathbf{X}\mathbf{X}^\top}{p} + v_p \mathbf{1},$$

and where

$$v_p = f\left(\frac{\text{tr}(\Sigma_p)}{p}\right) - f(0) - f'(0) \frac{\text{tr}(\Sigma_p)}{p}.$$

In other words,

$$\|\mathbf{M} - \mathbf{K}\|_2 \rightarrow 0$$

in probability, when $p \rightarrow \infty$ [Kar10, page 9].

When using the RBF kernel these assumptions typically hold. This theorem shows that as the size of the kernel matrix becomes large, it begins to behave like the underpinning data matrix, up to a factor of a constant.

3. RANDOM FOURIER FEATURES

As seen in Section 1 GPs rely heavily on the Gram matrix (see Definition 4) to create predictions based on training data $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ where $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ and $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top \in \mathbb{R}^n$. Unfortunately, the size of the Gram matrix scales quadratically with the number of samples making it difficult to train using data sets with more than 10^5 samples. Instead the kernel function itself can be factorized allowing one to convert training and kernel evaluation into the corresponding operations of a linear machine by mapping data into a relatively low-dimensional randomized feature space. This idea was first introduced by Rahimi and Recht [Rah08] where they proposed that, instead of using a kernel function to implicitly lift data into a higher dimensional feature space, an explicit feature map $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ could be used to approximate k as $k(\mathbf{x}, \mathbf{y}) = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle_{\mathbb{R}^N} \simeq \langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle_{\mathbb{R}^D}$ where D is chosen so that $n \gg D$. Thus once $\varphi(\mathbf{x}_i)$ has been computed for each \mathbf{x}_i , every entry of the Gram matrix can be swiftly approximated as

$$K_{ij} = K_{ji} \simeq \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{y}_j) \rangle_{\mathbb{R}^D}.$$

Already there have been numerous applications of this technique in GPs that have seen improved time performance with little loss in prediction accuracy [Pot21].

3.1. Theory and Computation. Contrary to the kernel trick, the Random Fourier Features (RFF) technique approximates $\langle \Phi(\cdot), \Phi(\cdot) \rangle_{\mathbb{R}^N}$ through an explicit feature mapping φ . The RFF technique hinges on Bochners theorem, stated without proof in Theorem 23, which characterises positive definite functions.

Theorem 23 (Bochner's). *A continuous and shift-invariant function $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y}) = k(\Delta)$ is positive definite (see Definition 2) if and only if it can be represented as*

$$k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{C}^d} \exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle) \mu_k(d\boldsymbol{\omega})$$

where μ_k is a positive finite measure on the frequencies of $\boldsymbol{\omega}$ [Hah33, Liu21].

The spectral distribution μ_k can be represented as finite measure induced by the Fourier transformation. Choosing a kernel for which $k(\mathbf{0}) = 1$ normalizes μ_k to a probability distribution $p(\cdot)$. For instance, the spectral distribution of the Gaussian RBF kernel is

$$(35) \quad p(\mathbf{w}) = \frac{1}{\sqrt{(2\pi)^D \left| \frac{\sigma^2}{2} \mathbb{1}_{D \times D} \right|}} \exp \left(-\frac{1}{2} \mathbf{w}^\top \left(\frac{\sigma^2}{2} \mathbb{1}_{D \times D} \right)^{-1} \mathbf{w} \right)$$

[Rah08, page 3]. One caveat in Bochner's theorem is that it requires our kernel to be shift-invariant (sometimes also referred to as stationary) as stated in Definition 24.

Definition 24 (Shift-Invariant). *A kernel $k : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{C}$ is called shift-invariant if $k(\mathbf{x}, \mathbf{y}) = g(\mathbf{x} - \mathbf{y})$ for some positive definite function $g : \mathbb{R}^N \rightarrow \mathbb{C}$ [HAE16, page 3].*

Clearly, the Gaussian RBF kernel is shift-invariant since it only relies on the bounding radius of \mathbf{x} and \mathbf{y} . Thus, from Bochner's theorem, a positive definite shift-invariant kernel with $k(\mathbf{0}) = 1$ can be computed

as

$$(36) \quad k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{C}^d} \exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle) p(\boldsymbol{\omega}) d\boldsymbol{\omega}.$$

The main idea of RFF is to approximate the integral in (36) using the following Monte-Carlo estimate

$$\begin{aligned} k(\mathbf{x} - \mathbf{y}) &= \int_{\mathbb{C}^d} \exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle) p(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &= \mathbb{E}_{\boldsymbol{\omega} \sim p(\cdot)} (\exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle)) \\ &\simeq \frac{1}{D} \sum_{j=1}^D \exp(i\langle \boldsymbol{\omega}_j, \mathbf{x} - \mathbf{y} \rangle) \\ &= \sum_{j=1}^D \left(\frac{1}{\sqrt{D}} \exp(i\langle \boldsymbol{\omega}_j, \mathbf{x} \rangle) \right) \overline{\left(\frac{1}{\sqrt{D}} \exp(i\langle \boldsymbol{\omega}_j, \mathbf{y} \rangle) \right)} \\ &= \langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle_{\mathbb{C}^D} \end{aligned}$$

where $\boldsymbol{\omega}_i \stackrel{\text{iid}}{\sim} p(\cdot)$ using the feature map

$$(37) \quad \varphi(\mathbf{x}) = \frac{1}{\sqrt{D}} [z(\boldsymbol{\omega}_1, \mathbf{x}), z(\boldsymbol{\omega}_2, \mathbf{x}), \dots, z(\boldsymbol{\omega}_D, \mathbf{x})]^\top$$

where for convenience of notation $z(\boldsymbol{\omega}, \mathbf{x}) = \exp(i\langle \boldsymbol{\omega}, \mathbf{x} \rangle)$. This allows the Gram matrix to be estimated as $\mathbf{K} \simeq \widetilde{\mathbf{K}} = \mathbf{Z}\mathbf{Z}^\top$ where $\mathbf{Z} = [\varphi(\mathbf{x}_1), \varphi(\mathbf{x}_2), \dots, \varphi(\mathbf{x}_D)] \in \mathbb{C}^{n \times D}$ [Rah08, Liu21, HAe16]. To simplify computation in most settings both $p(\cdot)$ and $k(\Delta)$ are real valued functions meaning $\exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle)$ can be replaced with its real component $\cos(\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle)$. The vast majority of literature uses the embeddings Rahimi and Recht provide for $\cos(\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle)$ where $z(\boldsymbol{\omega}, \mathbf{x})$ satisfies equation (36). The first embedding takes the form

$$(38) \quad z(\boldsymbol{\omega}, \mathbf{x}) = [\cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle), \sin(\langle \boldsymbol{\omega}, \mathbf{x} \rangle)]^\top$$

which satisfies (36) since

$$\begin{aligned} z(\boldsymbol{\omega}, \mathbf{x})^\top z(\boldsymbol{\omega}, \mathbf{y}) &= [\cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle), \sin(\langle \boldsymbol{\omega}, \mathbf{x} \rangle)] \begin{bmatrix} \cos(\langle \boldsymbol{\omega}, \mathbf{y} \rangle) \\ \sin(\langle \boldsymbol{\omega}, \mathbf{y} \rangle) \end{bmatrix} \\ &= \cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle) \cos(\langle \boldsymbol{\omega}, \mathbf{y} \rangle) + \sin(\langle \boldsymbol{\omega}, \mathbf{x} \rangle) \sin(\langle \boldsymbol{\omega}, \mathbf{y} \rangle) \\ &= \frac{1}{2} (\cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle + \langle \boldsymbol{\omega}, \mathbf{y} \rangle) + \cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle - \langle \boldsymbol{\omega}, \mathbf{y} \rangle)) + \\ &\quad \frac{1}{2} (\cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle - \langle \boldsymbol{\omega}, \mathbf{y} \rangle) - \cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle + \langle \boldsymbol{\omega}, \mathbf{y} \rangle)) \\ &= \cos(\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle). \end{aligned}$$

The other embedding Rahimi and Recht give is

$$(39) \quad z(\boldsymbol{\omega}, \mathbf{x}) = \sqrt{2} \cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle + b)$$

where $b \sim U[0, 2\pi]$. Using a similar argument we can show that this embedding also satisfies (36). However, Sutherland and Schneider [DJSaJS15] argue that the Gaussian RBF kernel is better suited for

the embedding given in (38). To summarise their argument we denote

$$(40) \quad \varphi_1(\mathbf{x}) = \sqrt{\frac{2}{D}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \cos(\langle \boldsymbol{\omega}_2, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_{D/2}, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \vdots \\ \sin(\langle \boldsymbol{\omega}_{D/2}, \mathbf{x} \rangle) \end{bmatrix}$$

to be the feature map corresponding to embedding in equation (38) and

$$(41) \quad \varphi_2(\mathbf{x}) = \sqrt{\frac{2}{D}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle + b_1) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle + b_D) \end{bmatrix}$$

to be the feature map corresponding to equation (39). They then show that

$$\begin{aligned} \mathbb{V}[\varphi_1(\Delta)] &= \frac{1}{D} (1 + k(2\Delta) - 2k(\Delta)^2) \\ \mathbb{V}[\varphi_2(\Delta)] &= \frac{1}{D} \left(1 + \frac{1}{2}k(2\Delta) - k(\Delta)^2\right) \end{aligned}$$

meaning the variance of φ_1 is smaller whenever

$$\mathbb{V}[\cos(\langle \boldsymbol{\omega}, \Delta \rangle)] = \frac{1}{2} + \frac{1}{2}k(2\Delta) - k(\Delta)^2 \leq \frac{1}{2}.$$

When using the Gaussian kernel,

$$\mathbb{V}[\cos(\langle \boldsymbol{\omega}, \Delta \rangle)] = \frac{1}{2} \left(1 - \exp\left(-\frac{2\|\Delta\|_2^2}{\sigma^2}\right)\right)^2 \leq \frac{1}{2}$$

so that $\varphi_1(\Delta) \leq \varphi_2(\Delta)$ for any $\Delta \in \mathbb{R}^d$. There finding were indeed consistent with our preliminary results. With this in mind, an embedding of φ_1 was always used for our experiments.

Another important result Rahimi and Recht show provides a bound on the sup-norm of the difference between a Gram matrix and its RFF approximation stated in Proposition 25.

Proposition 25. *Let $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y}) = k(\Delta)$ be a continuous shift-invariant, positive definite function defined on compact subset $\mathcal{M} \subset \mathbb{R}^d$ having radius ℓ where $k(0) = 1$ such that $\nabla^2 k(0)$ exists. Then for the feature mapping defined in equation (40) let $\sigma_p^2 = \mathbb{E}_{\boldsymbol{\omega} \sim p(\cdot)} \|\boldsymbol{\omega}\|_2^2 = \text{tr } \nabla^2 k(0)$ then for any $\varepsilon \in \mathbb{R}_{++}$, $\varepsilon \leq \sigma_p \ell$ we have*

$$\mathbb{P} \left[\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{M}} |\langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle_{\mathbb{R}^D} - k(\mathbf{x}, \mathbf{y})| \geq \varepsilon \right] \leq \alpha \left(\frac{\sigma_p \ell}{\varepsilon} \right)^2 \exp \left(-\frac{D\varepsilon^2}{8(d+2)} \right)$$

where $\alpha \in \mathbb{R}_{++}$ does not depend on anything [Rah08, page 3].

Rahimi and Recht prove Proposition 25 for $\alpha = 2^8$ although Sutherland and Schneider improve this to $\alpha = 66$ [DJSaJS15, page 3]. Observe that this bound is somewhat determined by the ratio D/d which is why D is often chosen as a multiple of d .

These results justify the RFF procedure seen in Algorithm 6, which was used to approximate a Gram matrix for the data set \mathbf{X} using the feature map from (40).

Algorithm 6: RFF Algorithm

input : $\mathbf{X} \in \mathbb{R}^{n \times d}$, the dimension of the feature space D .
output: $\widetilde{\mathbf{K}} \simeq \mathbf{K}$ where \mathbf{K} is the Gram matrix corresponding to \mathbf{X} .

- 1 Construct $\mathbf{W} \triangleq [\omega_1, \dots, \omega_D]^\top \in \mathbb{R}^{D \times d}$ where $\omega_i \stackrel{\text{iid}}{\sim} p(\cdot)$
- 2 $\mathbf{Z} = \frac{1}{\sqrt{D}} [\cos(\mathbf{W}\mathbf{X}^\top), \sin(\mathbf{W}\mathbf{X}^\top)]^\top$
- 3 $\widetilde{\mathbf{K}} = \mathbf{Z}\mathbf{Z}^\top$
- 4 **return** $\widetilde{\mathbf{K}}$

Algorithm 6 of course assumes an appropriate construction of \mathbf{W} , commonly called the transformation matrix, and thus has access to a routine which allows one to sample from $p(\cdot)$. When using the RBF Gaussian kernel, the spectral distribution given in equation (35) corresponds to a multivariate Gaussian distribution with mean $\mathbf{0}$ and covariance matrix $\left(\frac{\sigma}{\sqrt{2}}\right)^{-2} \mathbb{1}_{D \times D}$. This means \mathbf{W} can simply be constructed as $\mathbf{W} = \left(\frac{\sigma}{\sqrt{2}}\right)^{-1} [\omega_1, \dots, \omega_D]^\top$ where $\omega_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbb{1}_{D \times D})$. Transformations matrices constructed in this manner are given the notation \mathbf{W}_{RFF} . It can be shown that if \mathbf{W}_{RFF} is used as the transformation matrix in Algorithm 6 then it produces an unbiased estimate, $\widetilde{\mathbf{K}}_{\text{RFF}}$, for the Gram matrix. This is stated more precisely in Lemma 26.

Lemma 26. $\widetilde{\mathbf{K}}_{\text{RFF}}$ is an unbiased estimate of \mathbf{K} , that is

$$\mathbb{E} \left[(\widetilde{\mathbf{K}}_{\text{RFF}})_{ij} \right] = \exp \left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2} \right)$$

[Yu16, page 3].

Proof. We shall prove this for the Gaussian kernel, although proofs for other kernels are analogous. Let $\mathbf{z} = \mathbf{x} - \mathbf{y}$. Recall the kernel is approximated as

$$\sum_{j=1}^D \frac{1}{D} \cos(\langle \omega_j, \mathbf{z} \rangle)$$

where $\omega_i \stackrel{\text{iid}}{\sim} p(\cdot)$. By Bochner's theorem

$$\mathbb{E} [\cos(\langle \omega_j, \mathbf{z} \rangle)] = \exp(-\|\mathbf{z}\|_2^2 / 2\sigma^2)$$

meaning $\widetilde{\mathbf{K}}_{\text{RFF}}$ provides an unbiased estimate of \mathbf{K} . \square

Unfortunately, constructing the transformation matrix using \mathbf{W}_{RFF} does not scale well as the dimension of the feature space increases. Thus the focus of the upcoming sections will be to highlight a few of the more popular alternative methods used in the literature for the construction of the transformation matrix.

3.2. Orthogonal Random Features. In the previous chapter Algorithm 6 assumed some sort of mechanism for producing the transformation matrix \mathbf{W} . The construction presented in 3.1 involved sampling $\omega_i \stackrel{\text{iid}}{\sim} p(\cdot)$. For the Gaussian RBF kernel this meant sampling from the multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbb{1}_{D \times D})$. The transformation matrix constructed in this manner was denoted \mathbf{W}_{RFF} . Recently, there has been a buzz in the literature exploring alternative constructs for the transformation matrix described in Section 3.1 [Liu21]. We shall consider the two methods proposed by Yu *et al.* [Yu16]; the first method here and the second in the following section (3.3). The first method from Yu *et al.* is the Orthogonal Random Features (ORF) method with imposes orthogonality on the transformation matrix. To do this a Gaussian matrix $\mathbf{G} \in \mathbb{R}^{D \times d}$ is first produced, much like in \mathbf{W}_{RFF} . An orthogonal matrix \mathbf{Q} is then created by taking the QR factorisation (see Section 4.2) of \mathbf{G} . However, the random orthogonal matrix, \mathbf{Q} , will not give an unbiased estimate of the kernel matrix. To fix this, the following common probabilistic identity is employed

$$\|z\|_2^2 \sim \chi_k^2, \text{ where } z \sim \mathcal{N}(\mathbf{0}, \mathbb{1}_{k \times k})$$

where χ_k^2 is the chi-squared distribution with k degrees of freedom [Bro91, page 41]. This identity is easily demonstrated by equating a shared moment generating function of $(1 - 2t)^{-\frac{k}{2}}$ for $t < \frac{1}{2}$. Taking the square root of both sides gives $\|z\|_2 \sim \chi_k$ where χ_k is the chi distribution with k degrees of freedom. In the RFF method, each $\omega_i \in \mathbb{R}^D$ was independently taken from the multivariate normal Gaussian distribution meaning that using the identity provided above $\|\omega_i\|_2 \sim \chi_D$. The ORF method augments \mathbf{Q} by scaling its rows by iid χ_D values which can be accomplished through right multiplication with $\mathbf{S} = \text{diag}(\psi_1, \psi_2, \dots, \psi_D)$ where $\psi_i \stackrel{\text{iid}}{\sim} \chi_D$. This means

$$\|(\mathbf{SQ})_{(i)}\|_2 = \|\psi_i \mathbf{Q}_{(i)}\|_2 = \psi_i \sim \chi_D$$

so that the row norms of \mathbf{G} and \mathbf{SQ} have the same distribution. Thus the transformation matrix for the ORF method is

$$(42) \quad \mathbf{W}_{\text{ORF}} = \left(\frac{\sigma}{\sqrt{2}} \right)^{-1} \mathbf{SQ}.$$

The main downside the the ORF method is that the QR factorisation brings a computational cost of $\mathcal{O}(Dd)$. Fortunately when using \mathbf{W}_{ORF} as our transformation matrix in Algorithm 6 the approximate Gram matrix $\widetilde{\mathbf{K}}_{\text{RFF}}$ is an unbiased estimate of \mathbf{K} , stated more formally in Theorem 27.

Theorem 27. $\widetilde{\mathbf{K}}_{\text{ORF}}$ is an unbiased estimate of \mathbf{K} , that is

$$\mathbb{E} \left[(\widetilde{\mathbf{K}}_{\text{ORF}})_{ij} \right] = \exp \left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2} \right)$$

[Yu16, page 3].

Proof. Again, we shall prove this for the Gaussian kernel, although proofs for other kernels are analogous. Let $z = \mathbf{x} - \mathbf{y}$.

$$\mathbb{E} \left[(\widetilde{\mathbf{K}}_{\text{ORF}})_{ij} \right] = \mathbb{E} \left[\frac{1}{D} \sum_{j=1}^D \cos(\langle \omega_j, z \rangle) \right] = \frac{1}{D} \sum_{j=1}^D \mathbb{E} [\cos(\langle \omega_j, z \rangle)].$$

Based on the definition of the ORF method, $\omega_1, \omega_2, \dots, \omega_D$ are D random vectord given by $\omega_i = s_i \mathbf{u}_i$, with $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ a uniformly chosen random orthonormal basis for \mathbb{R}^d , and s_i 's are independently

χ -distributed random variables with d degrees of freedom. It is easy to show that for each i , \mathbf{w}_i is distributed according to $\mathcal{N}(\mathbf{0}, \mathbb{1}_{d \times d})$ and hence by Bochner's theorem

$$\mathbb{E} [\cos (\langle \boldsymbol{\omega}_j, \mathbf{z} \rangle)] = \exp (-\|\mathbf{z}\|_2^2 / 2\sigma)$$

meaning $\widetilde{\mathbf{K}}_{\text{ORF}}$ provides an unbiased estimate of \mathbf{K} . \square

Furthermore, the variance of $(\widetilde{\mathbf{K}}_{\text{ORF}})_{ij}$ is bounded by

$$\mathbb{V} [(\widetilde{\mathbf{K}}_{\text{ORF}})_{ij}] - \mathbb{V} [(\widetilde{\mathbf{K}}_{\text{RFF}})_{ij}] = \frac{1}{D} \left(\frac{g(\tau)}{d} - \frac{(d-1)e^{-\tau^2}\tau^4}{2d} \right)$$

where $\tau = \|\mathbf{x}_i - \mathbf{x}_j\|_2 / \frac{\sigma}{\sqrt{2}}$ and

$$g(\tau) = \frac{e^{\tau^2} (\tau^8 + 6\tau^6 + 7\tau^4 + \tau)}{4} + \frac{e^{\tau^2} \tau^4 (\tau^6 + 2\tau^4)}{2d}$$

[Liu21, page 8]. This shows that there are scenarios for which $\mathbb{V} [(\widetilde{\mathbf{K}}_{\text{ORF}})_{ij}] < \mathbb{V} [(\widetilde{\mathbf{K}}_{\text{RFF}})_{ij}]$, namely when d is large and τ is small. Also, the ratio in variance between $\widetilde{\mathbf{K}}_{\text{ORF}}$ and $\widetilde{\mathbf{K}}_{\text{RFF}}$ for large d can be approximated as

$$\frac{\mathbb{V} [(\widetilde{\mathbf{K}}_{\text{ORF}})_{ij}]}{\mathbb{V} [(\widetilde{\mathbf{K}}_{\text{RFF}})_{ij}]} \simeq 1 - \frac{(s-1)e^{-\tau^2}\tau^4}{d(1-e^{-\tau^2})^2}$$

[Liu21, page 8].

3.3. Random Ortho-Matrices and Structured Orthogonal Random Matrices. The second method we shall consider for producing a transformation matrix also originates from Yu's *et al.* paper, which Choromanski *et al.* [Cho17] generalized as Random Ortho-Matrices (ROM). This second class of methods is underpinned by transformation matrices with the same variance reductions as ORF with the added benefit of time and memory savings. The transformation matrices generated using ROM take the form

$$(43) \quad \mathbf{W}_{\text{ROM}} = \sqrt{d} \prod_{i=1}^k \mathbf{S} \mathbf{D}_i$$

where $\mathbf{S} \in \mathbb{R}^{D \times D}$ has orthogonal rows and $\mathbf{D} = \text{diag}(\delta_1, \dots, \delta_D) \in \mathbb{R}^{D \times D}$ where $\delta_i \stackrel{\text{iid}}{\sim} U(\{-1, 1\})$. This matrix can be forced into a $\mathbb{R}^{D \times d}$ sized matrix by simply extracting the first d columns of \mathbf{D}_1 . The matrix to take the role of \mathbf{S} in virtually every application of ROM is the Hadamard matrix, defined in 28, which facilitates a fast $m \log(n)$ matrix multiplication with a size $m \times n$ and is known as Fast Walsh-Hadamard transform (FWHT) [FaA76].

Definition 28 (Hadamard Matrix). *The Hadamard matrix $\mathbf{H}_i \in \mathbb{R}^{(2^{i-1} \times 2^{i-1})}$ is defined recursively as*

$$\mathbf{H}_i = \begin{cases} [1], & i = 1 \\ \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{H}_{i-1} & \mathbf{H}_{i-1} \\ \mathbf{H}_{i-1} & -\mathbf{H}_{i-1} \end{bmatrix}, & i > 1 \end{cases}.$$

Note that while Hadamard matrices are only defined for dimensions of exact powers of 2, although other sizes can be constructed by removing portions of the matrix given in Definition 28 or by padding with 0. This provides a concrete means for which one can generate a transformation matrix

$$(44) \quad \sqrt{d} \prod_{i=1}^k \mathbf{H} \mathbf{D}_i$$

where \mathbf{H} is an appropriately sized Hadamard matrix. It is easy to check that the matrix generated by equation (44) shares the same expected rows norm lengths as \mathbf{W}_{ORF} and thus enjoys the same variance reduction benefits. Moreover, since matrix multiplication with \mathbf{H} can be performed in $\mathcal{O}(D \log(d))$ time (using FWHT) and multiplication with \mathbf{H} can be performed in $\mathcal{O}(D)$ time, the ROM method has the added benefit of improved run time complexity $\mathcal{O}(D \log(d))$ using only $\mathcal{O}(D)$ extra memory. Table 2 gives a comparison of the time and space complexities for the methods mentioned so far.

TABLE 2. A comparison of various methods for computing a suitable transformation matrix with the Random Fourier Features paradigm. Typically the dimension of the feature space, D , is chosen as some multiple of the dimension of data, d .

<i>Method</i>	<i>Time</i>	<i>Extra Space</i>
RFF [Rah08]	$\mathcal{O}(Dd)$	$\mathcal{O}(Dd)$
ORF [Yu16]	$\mathcal{O}(Dd)$	$\mathcal{O}(Dd)$
ROM (SORF) [Cho17, Yu16]	$\mathcal{O}(D \log(d))$	$\mathcal{O}(D)$

Despite the wide use of the ROM method in various machine learning tasks [Cho17, And15, Cho20, Liu21] a number of high-interest theoretical properties remain unsolved, leaving many aspects of this method shrouded in mystery. Instead, much of what we understand about ROM's estimate capabilities comes from empirical analysis. Nonetheless, we shall still cover a smaller number of important results that have been established.

Choromanski *et al.* [Cho17] show that there are diminishing returns (estimate wise) for choosing larger values of k in equation (44). They also show that choosing odd values of k in (44) provides better estimates than its even-parity $k - 1$ and $k + 1$ counterparts. For this reason a k value of 3 is usually chosen which gives rise to the transformation matrix estimate given in equation (45). The method for constructing transformation matrices in this manner is referred to as Structured Orthogonal Random Features (SORF).

$$(45) \quad \mathbf{W}_{\text{SORF}} = \sqrt{d} \mathbf{H} \mathbf{D}_3 \mathbf{H} \mathbf{D}_2 \mathbf{H} \mathbf{D}_1$$

This is the same transformation matrix estimate that Yu *et al* [Yu16]. provides. Unfortunately using the SORF method in Algorithm 6 does not produce an unbiased estimate of the Gram matrix; however, it does satisfy an asymptotic unbiased property

$$\left| \mathbb{E} \left[(\widetilde{\mathbf{K}}_{\text{SORF}})_{ij} \right] - \mathbb{E} \left[(\widetilde{\mathbf{K}}_{\text{RFF}})_{ij} \right] \right| \leq \frac{6\tau}{\sqrt{d}}$$

where τ is again $\|x_i - x_j\|_2 / \sqrt{2}$ [Liu21, page 8].

Bojarski *et al.* [Boj16, page 4] give an intuitive explanation for the roles of each of the different blocks $\mathbf{H}\mathbf{D}_1$, $\mathbf{H}\mathbf{D}_2$ and $\mathbf{H}\mathbf{D}_3$. The first block can be shown to satisfy

$$\mathbb{P} \left[\|\mathbf{H}\mathbf{D}_1\mathbf{x}\|_\infty > \frac{\log D}{\sqrt{D}} \right] \leq 2d \exp \left(-\frac{\log^2 D}{8} \right), \quad \mathbf{x} \in \mathbb{R}^D$$

[Liu21, page 8] so that it can be thought as a "balancer" leaving no single dimension bearing too much of the l^2 norm. For the second block, the cost of using a structured matrix is the loss of independence. The purpose of the second block is to mitigate this effect by making similar input vectors near-orthogonal. Finally the third block controls the capacity of the entire structure by providing a vector of parameters. Near-independence is now implied by the near-orthogonality (achieved by $\mathbf{H}\mathbf{D}_2$) and the fact that the projections of the Gaussian vector or Radamacher vector onto "almost orthogonal directions" are "close to independent". These roles are portrayed visually in Figure 8.

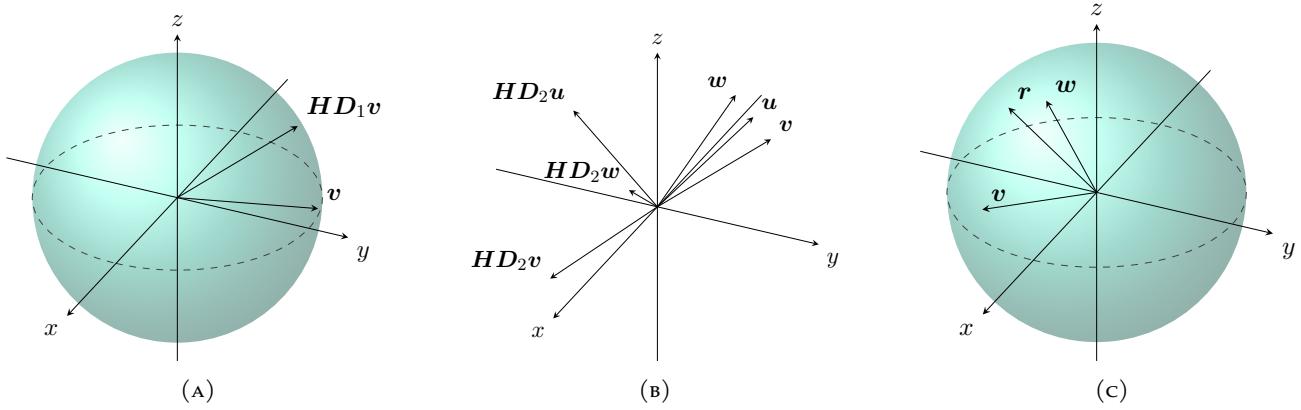


FIGURE 8. A visual representation for the roles of each matrix block in the SORF method. The first block $\mathbf{H}\mathbf{D}_1$ rotates v so that single dimension bears too much of the l^2 norm seen in panel (A). In panel (B) the second block $\mathbf{H}\mathbf{D}_2$ transforms vectors so that their image is near-orthogonal. Panel (C) shows that the projection of a random vector r onto two near-orthogonal vector v, w yields a near-independent vector.

4. KRYLOV SUBSPACE METHODS

In this section we will focus on how iterative methods and, in particular a class of iterative methods known as Krylov Subspace methods, may be used to solve a linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$. While non-iterative methods exist to solve such systems, virtually all of them carry an unwieldy runtime of $\mathcal{O}(n^3)$ for a system of n parameters. Even for current computer systems, this renders many common matrix problems computation intractable. Consequently, the focus of solving linear systems has shifted towards iterative methods. While iterative methods typically demand certain structural properties of the matrices, such as symmetry and positive definiteness, this generally is not a problem since the majority of large matrix problems by nature endow these systems with these desired properties. For example, the Gram matrices used to solve linear systems in Gaussian Processes possess both symmetry and positive definiteness. There are also a number of other properties of iterative methods which make them rather attractive to users. To start, iterative Krylov subspace methods are guaranteed to converge to an exact solution within a finite number of iterations. Even if the method is prematurely stopped before reaching an exact solution, the solution obtained on the final iteration will, in some sense, be a good enough estimate. Additionally, unlike most non-iterative methods, Krylov subspace methods do not require an explicit form of the matrix \mathbf{A} and instead only requires some routine for computing $\mathbf{A}\mathbf{x}$.

4.1. Krylov Subspaces. The main usefulness of Krylov subspaces come from their ability in assisting us to solving linear systems. To this end, consider the problem of solving the linear system

$$(46) \quad \mathbf{A}\mathbf{x}^\star = \mathbf{b}$$

where no explicit form of \mathbf{A} is available and instead one must draw information from \mathbf{A} solely through a routine that can evaluate $\mathbf{A}\mathbf{v}$ for any \mathbf{v} . How could this routine be utilized in such a manner to provide with a solution to (46)? Before answering this, consider the following theorem.

Theorem 29. *For $\mathbf{A} \in \mathbb{C}^{n \times n}$ if $\|\mathbf{A}\| = q < 1$ then $\mathbf{1} - \mathbf{A}$ is invertible and its inverse admits the following representation*

$$(\mathbf{1} - \mathbf{A})^{-1} = \sum_{k=0}^{\infty} \mathbf{A}^k$$

[Ber96, page 287].

Proof. Let us show that the sequence $\mathbf{S}_n = \sum_{k=0}^n \mathbf{A}^k$ (where n is some natural number) of partial sums of the series on the right-hand side of Theorem 29 is fundamental in $\mathbb{C}^{n \times n}$. Indeed, by using the fact that $\|\mathbf{AB}\| \leq \|\mathbf{A}\|\|\mathbf{B}\|$ for any $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$, we conclude that

$$\|\mathbf{S}_{n+p} - \mathbf{S}_n\| \leq \|\mathbf{A}^{n+1}\| + \cdots + \|\mathbf{A}^{n+p}\| \leq q^{n+1} + \cdots + q^{n+p}$$

for all $n, p \in \mathbb{N}$. Since $q < 1$, this yields the desired result. Since $\mathbb{C}^{n \times n}$ is a Banach space, the sequence $\{\mathbf{S}_n\}_n^\infty$ uniformly converges to an operator $\mathbf{S} \in \mathbb{C}^{n \times n}$. Let us show that $\mathbf{S}(\mathbf{1}_{n \times n} - \mathbf{A}) = (\mathbf{1}_{n \times n} - \mathbf{A})\mathbf{S} = \mathbf{1}_{n \times n}$. To this end, we note that

$$\|(\mathbf{1}_{n \times n} - \mathbf{A})\mathbf{S}_n - (\mathbf{1}_{n \times n} - \mathbf{A})\mathbf{S}\| \leq \|\mathbf{1}_{n \times n} - \mathbf{A}\| \cdot \|\mathbf{S} - \mathbf{S}_n\| \rightarrow 0$$

as $n \rightarrow \infty$. Therefore, it suffices to show that the sequence $\{(\mathbb{1}_{n \times n} - \mathbf{A}) \mathbf{S}_n\}_{n=1}^{\infty}$ uniformly converges to the identity operator. We have

$$\mathbf{S}_n (\mathbb{1}_{n \times n} - \mathbf{A}) = (\mathbb{1}_{n \times n} - \mathbf{A}) \mathbf{S}_n = \sum_{k=0}^n \mathbf{A}^k - \sum_{k=1}^{n+1} \mathbf{A}^k = \mathbb{1}_{n \times n} - \mathbf{A}^{n+1},$$

that is, $\|(\mathbb{1}_{n \times n} - \mathbf{A}) \mathbf{S}_n - \mathbb{1}_{n \times n}\| = \|\mathbf{A}^{n+1}\| \leq q^{n+1}$, which vanishes as $n \rightarrow \infty$. Thus $\mathbf{S} = (\mathbb{1}_{n \times n} - \mathbf{A})^{-1}$. \square

Consider a matrix for which $\|\mathbf{A}\| < 2$, it follows that $\|\mathbb{1}_{n \times n} - \mathbf{A}\| < 1$ meaning $\mathbb{1}_{n \times n} - (\mathbb{1}_{n \times n} - \mathbf{A})$ is invertible and $\mathbf{A}^{-1} = (\mathbb{1}_{n \times n} - (\mathbb{1}_{n \times n} - \mathbf{A}))^{-1} = \sum_{k=0}^{\infty} (\mathbb{1}_{n \times n} - \mathbf{A})^k$. Returning to (46) for any $\mathbf{x}_0 \in \mathbb{R}^n$ we have

$$\begin{aligned} \mathbf{x}^* &= \mathbf{A}^{-1} \mathbf{b} = \mathbf{A}^{-1} (\mathbf{A} \mathbf{x}^* - \mathbf{A} \mathbf{x}_0 + \mathbf{A} \mathbf{x}_0) \\ &= \mathbf{x}_0 + \mathbf{A}^{-1} \mathbf{r}_0 \\ &= \mathbf{x}_0 + \sum_{k=0}^{\infty} (\mathbb{1}_{n \times n} - \mathbf{A})^k \mathbf{r}_0 \end{aligned}$$

where $\mathbf{r}_0 = \mathbf{A} \mathbf{x}^* - \mathbf{A} \mathbf{x}_0$. An obvious question that arises is whether a closed form solution of the above equation can be found. To answer this, we need to enlist the help of the Cayley-Hamilton theorem.

Theorem 30 (Cayley-Hamilton). *Let $p_n(\lambda) = \sum_{i=0}^n c_i \lambda^i$ be the characteristic polynomial of the matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, then $p_n(\mathbf{A}) = \mathbf{0}$ [Hor13, page 109].*

The Cayley-Hamilton Theorem implies that

$$\begin{aligned} \mathbf{0} &= c_0 + c_1 \mathbf{A} + \dots + c_{n-1} \mathbf{A}^{n-1} + c_n \mathbf{A}^n \\ \mathbf{0} &= \mathbf{A}^{-1} c_0 + c_1 + \dots + c_{n-1} \mathbf{A}^{n-2} + c_n \mathbf{A}^{n-1} \\ \mathbf{A}^{-1} &= \alpha_1 \mathbb{1}_{n \times n} + \dots + \alpha_{n-1} \mathbf{A}^{n-2} + \alpha_n \mathbf{A}^{n-1} \end{aligned}$$

where $\alpha_i = -c_i/c_0$. The above demonstrates that \mathbf{A}^{-1} can be represented as a matrix polynomial of degree $n-1$. This means that $\sum_{k=0}^{\infty} (\mathbb{1} - \mathbf{A})^k$ indeed possess a closed form solution namely

$$(47) \quad \mathbf{x}^* = \mathbf{x}_0 + \mathbf{A}^{-1} \mathbf{r}_0 = \mathbf{x}_0 + \alpha_1 \mathbf{r}_0 + \dots + \alpha_{n-1} \mathbf{A}^{n-2} \mathbf{r}_0 + \alpha_n \mathbf{A}^{n-1} \mathbf{r}_0.$$

This also shows that $\mathbf{x}^* \in \mathbf{x}_0 + \text{l.s} \{ \mathbf{r}_0, \mathbf{A} \mathbf{r}_0, \mathbf{A}^2 \mathbf{r}_0, \dots, \mathbf{A}^{n-1} \mathbf{r}_0 \}$. One idea for finding a solution to (46) is to use our routine for evaluting \mathbf{Av} to iteratively compute new basis elements for the space generated by $\{ \mathbf{r}_0, \mathbf{A} \mathbf{r}_0, \mathbf{A}^2 \mathbf{r}_0, \dots, \mathbf{A}^{n-1} \mathbf{r}_0 \}$ and at each step carefully choosing a \mathbf{x}_k such that \mathbf{x}_k approaches \mathbf{x}^* , in some form. The subspace constructed using this technique is so important that is has its own name.

Definition 31 (Krylov Subspace). *The Krylov Subspace of order k generated by the matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and the vector $\mathbf{v} \in \mathbb{C}^n$ is defined as*

$$\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \text{l.s} \{ \mathbf{r}_0, \mathbf{A} \mathbf{r}_0, \mathbf{A}^2 \mathbf{r}_0, \dots, \mathbf{A}^{k-1} \mathbf{r}_0 \}$$

for $k \geq 1$ and $\mathcal{K}_0(\mathbf{A}, \mathbf{v}) = \{ \mathbf{0} \}$ [Tre97, page 245].

For the purposes of solving (46) it is important to understand how $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$ grows for larger and larger k since a solution for (46) will be present in a Krylov Subspace that cannot be grown any larger, according

to (47). In other words, an exact solution can be constructed once we have extracted all the information from \mathbf{A} through multiplication of \mathbf{r}_0 . The following theorem provides information on how exactly the Krylov Subspace grows as k increases.

Theorem 32. *There is a positive integer called the grade of \mathbf{v} with respect to \mathbf{A} , denoted $t_{\mathbf{v}, \mathbf{A}}$, where*

$$\dim(\mathcal{K}_k(\mathbf{A}, \mathbf{v})) = \begin{cases} k, & k \leq t_{\mathbf{v}, \mathbf{A}} \\ t_{\mathbf{v}, \mathbf{A}}, & k \geq t_{\mathbf{v}, \mathbf{A}} \end{cases}$$

Theorem 32 essentially tells us for $k \leq t_{\mathbf{v}, \mathbf{A}}$ that $\mathbf{A}^k \mathbf{v}$ is linearly independent to $\mathbf{A}^i \mathbf{v}$ for $0 \leq i \leq k-1$ meaning $\{\mathbf{v}, \mathbf{A}\mathbf{v}, \mathbf{A}^2\mathbf{v}, \dots, \mathbf{A}^{k-1}\mathbf{v}\}$ serves as a basis for $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$ and $\mathcal{K}_{k-1}(\mathbf{A}, \mathbf{v}) \subsetneq \mathcal{K}_k(\mathbf{A}, \mathbf{v})$. Conversely, any new vectors formed beyond $t_{\mathbf{v}, \mathbf{A}}$ will be linearly independent meaning $\mathcal{K}_k(\mathbf{A}, \mathbf{v}) \subsetneq \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{v})$ for $k \geq t_{\mathbf{v}, \mathbf{A}}$. While $t_{\mathbf{v}, \mathbf{A}}$ obviously plays a central role in determining a suitable basis whose span contains $\mathbf{A}^{-1}\mathbf{b}$, its importance is made abundantly clear in the following corollary.

Corollary 33.

$$t_{\mathbf{v}, \mathbf{A}} = \min \{k \mid \mathbf{A}^{-1}\mathbf{v} \in \mathcal{K}_k(\mathbf{A}, \mathbf{v})\}$$

Proof. Recall from Cayley-Hamilton (Theorem 30) that

$$\mathbf{A}^{-1}\mathbf{v} = \sum_{i=0}^{n-1} \alpha_i \mathbf{A}^i \mathbf{v}$$

But since $\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{v})$ for $k \geq t_{\mathbf{v}, \mathbf{A}}$

$$\mathbf{A}^{-1}\mathbf{v} = \sum_{i=0}^{t-1} \beta_i \mathbf{A}^i \mathbf{v}$$

meaning $\mathbf{A}^{-1}\mathbf{v} \in \mathcal{K}_k(\mathbf{A}, \mathbf{v})$ for $k \geq t_{\mathbf{v}, \mathbf{A}}$. Suppose for the sake of contradiction that this also holds for $k = t_{\mathbf{v}, \mathbf{A}} - 1$, that is, $\mathbf{A}^{-1}\mathbf{v} = \sum_{i=0}^{t-2} \gamma_i \mathbf{A}^i \mathbf{v}$. However, this gives

$$\mathbf{v} = \sum_{i=0}^{t-2} \gamma_i \mathbf{A}^{i+1} \mathbf{v} = \sum_{i=0}^{t-1} \gamma_{i-1} \mathbf{A}^i \mathbf{v}$$

implying $\{\mathbf{v}, \mathbf{A}\mathbf{v}, \mathbf{A}^2\mathbf{v}, \dots, \mathbf{A}^{t-1}\mathbf{v}\}$ are linearly dependent which means that $\dim(\mathcal{K}_k(\mathbf{A}, \mathbf{v})) < t$, which provides us with our contradiction. \square

This allows us to make a much stronger statement on the whereabouts of \mathbf{x}^* in relation to the Krylov Subspaces.

Corollary 34. *For any \mathbf{x}_0 we have*

$$\mathbf{x}^* \in \mathbf{x}_0 + \mathcal{K}_{t_{\mathbf{r}_0, \mathbf{A}}}(\mathbf{A}, \mathbf{r}_0)$$

where $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$.

4.2. Gram-Schmidt Process and QR factorisations. Many areas of linear algebra involve studying the column space of matrices. The QR factorisation both provides a powerful tool to better understand the column space of a matrix and serves as an important factorisation mechanism for many numerical methods. Suppose that a matrix $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] \in \mathbb{R}^{n \times n}$ has full rank. The idea of a QR factorisation is to find an alternative orthonormal basis for $\{\mathbf{a}_i\}_{i=1}^n$, say $\{\mathbf{q}_i\}_{i=1}^n$, and to somehow relate the original matrix \mathbf{A} to a new matrix whose columns are $\{\mathbf{q}_i\}_{i=1}^n$. Consider the following procedure that allows us to find an orthonormal basis $\{\mathbf{q}_i\}_{i=1}^n$ for which $\text{l.s } \{\mathbf{a}_i\}_{i=1}^n = \text{l.s } \{\mathbf{q}_i\}_{i=1}^n$. First set $\mathbf{q}_1 = \frac{\mathbf{a}_1}{\|\mathbf{a}_1\|}$, clearly $\text{l.s } \{\mathbf{a}_1\} = \text{l.s } \{\mathbf{q}_1\}$. Next, construct a vector $\mathbf{q}'_2 = \mathbf{a}_2 - r_{1,2} \cdot \mathbf{q}_1$ so that $\mathbf{q}'_2 \perp \mathbf{q}_1$. This means

$$\begin{aligned} 0 &= \langle \mathbf{q}_1, \mathbf{q}'_2 \rangle \\ 0 &= \langle \mathbf{q}_1, \mathbf{a}_2 - r_{1,2} \cdot \mathbf{q}_1 \rangle \\ 0 &= \langle \mathbf{q}_1, \mathbf{a}_2 \rangle - r_{1,2} \cdot \langle \mathbf{q}_1, \mathbf{q}_1 \rangle \\ r_{1,2} &= \langle \mathbf{q}_1, \mathbf{a}_2 \rangle \end{aligned}$$

Since \mathbf{q}'_2 may not be a unit vector we set $\mathbf{q}_2 = \frac{\mathbf{q}'_2}{\|\mathbf{q}'_2\|}$ where $\text{l.s } \{\mathbf{a}_1, \mathbf{a}_2\} = \text{l.s } \{\mathbf{q}_1, \mathbf{q}_2\}$. Next, the vector \mathbf{q}'_3 is constructed so that

$$\mathbf{q}'_3 = \mathbf{a}_3 - r_{1,3} \mathbf{q}_1 - r_{2,3} \mathbf{q}_2$$

is orthogonal to both \mathbf{q}_2 and \mathbf{q}_1 . This amounts to setting $r_{1,3} = \langle \mathbf{q}_1, \mathbf{a}_3 \rangle$ and $r_{2,3} = \langle \mathbf{q}_2, \mathbf{a}_3 \rangle$. Similarly, \mathbf{q}'_3 is normalized so that $\mathbf{q}_3 = \frac{\mathbf{q}'_3}{\|\mathbf{q}'_3\|}$ and $\text{l.s } \{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3\} = \text{l.s } \{\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3\}$. Continuing in this fashion the k^{th} vector in our orthonormal basis is computed as

$$(48) \quad \mathbf{q}_k = \frac{\mathbf{a}_k - \sum_{i=1}^{k-1} r_{i,k} \cdot \mathbf{q}_i}{r_{k,k}}$$

where $r_{i,k} = \langle \mathbf{q}_i, \mathbf{a}_k \rangle$, $r_{k,k} = \|\mathbf{a}_k - \sum_{i=1}^{k-1} r_{i,k} \cdot \mathbf{q}_i\|$ and $\text{l.s } \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\} = \text{l.s } \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k\}$. This procedure is known as the Gram-Schmidt process [Ber96, Tre97, Dem97] and is summarized in the following algorithm.

Algorithm 7: Classical Gram-Schmidt

input : A basis $\{\mathbf{a}_i\}_{i=1}^n$.

output: An orthonormal basis $\{\mathbf{q}_i\}_{i=1}^n$ such that $\text{l.s } \{\mathbf{a}_i\}_{i=1}^n = \text{l.s } \{\mathbf{q}_i\}_{i=1}^n$

```

1 for  $k = 1$  to  $n$  do
2    $\mathbf{q}'_k = \mathbf{a}_k$ 
3   for  $i = 1$  to  $k - 1$  do
4      $r_{i,k} = \langle \mathbf{q}_i, \mathbf{a}_k \rangle$ 
5      $\mathbf{q}'_k = \mathbf{q}'_k - r_{i,k} \mathbf{q}_i$ 
6   end
7    $r_{k,k} = \|\mathbf{q}'_k\|$ 
8    $\mathbf{q}_k = \mathbf{q}'_k / r_{k,k}$ 
9 end
10 return  $\{\mathbf{q}_i\}_{i=1}^n$ 

```

Relating the column space of \mathbf{A} to the orthonormal basis $\{\mathbf{q}_i\}_{i=1}^n$ can be done in a matrix form as

$$[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n] \begin{bmatrix} r_{1,1} & r_{1,2} & \cdots & r_{1,n} \\ & r_{2,2} & & \vdots \\ & & \ddots & \vdots \\ & & & r_{n,n} \end{bmatrix}$$

or more written succinctly as

$$(49) \quad \mathbf{A} = \mathbf{Q}\mathbf{R}$$

where $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n]$ and $R_{i,j} = r_{i,j}$ for $i \leq j$ and $R_{i,j} = 0$ for $i > j$. This is exactly the QR factorisation for a full rank matrix. Note that $\text{Range}(\mathbf{A}) = \text{Range}(\mathbf{Q})$. In general, any square matrix $\mathbf{A} \in \mathbb{K}^{m \times n}$ may be decomposed as $\mathbf{A} = \mathbf{Q}\mathbf{R}$ where $\mathbf{Q} \in \mathbb{K}^{m \times m}$ is an orthogonal matrix and $\mathbf{R} \in \mathbb{K}^{m \times n}$ is an upper triangular matrix. This is known as a full QR factorisation. Since the bottom $(m - n)$ rows of this \mathbf{R} consists entirely of zeros, it is often useful to partition the full QR factorisation in the following manner to shed vacuous entries

$$\mathbf{A} = \mathbf{Q}\mathbf{R} = \mathbf{Q} \begin{bmatrix} \hat{\mathbf{R}} \\ \mathbf{0}_{(m-n) \times n} \end{bmatrix} = [\hat{\mathbf{Q}} \quad \mathbf{Q}'] \begin{bmatrix} \hat{\mathbf{R}} \\ \mathbf{0}_{(m-n) \times n} \end{bmatrix} = \hat{\mathbf{Q}}\hat{\mathbf{R}}.$$

This alternate decomposition is called the reduced (or sometimes the thin) QR factorisation. The following two theorems on the QR factorization are stated without proof.

Theorem 35. Every $\mathbf{A} \in \mathbb{K}^{m \times n}$, ($m \geq n$) has a full QR factorisation, hence also a reduced QR factorisation. [Tre97]

Theorem 36. Each $\mathbf{A} \in \mathbb{K}^{m \times n}$, ($m \geq n$) of full rank has a unique reduced QR factorisation $\mathbf{A} = \hat{\mathbf{Q}}\hat{\mathbf{R}}$ with $r_{k,k} > 0$. [Tre97]

In practice the classical Gram-Schmidt process described in Algorithm 7 is rarely used as the procedure becomes numerically unstable if $\{\mathbf{a}_i\}_{i=1}^n$ are almost linearly dependent. Before looking at ways to resolve these numerical instabilities, a quick recap of projectors is useful. A square matrix \mathbf{P}_G acting on a Hilbert space H that sends $\mathbf{x} \in H$ to its projection onto a subspace G is called the projector onto G . If $\{\mathbf{q}_k\}_{k=1}^m$ is an orthonormal basis in G then

$$\mathbf{P}_G = \mathbf{Q}\mathbf{Q}^*$$

where $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m, \mathbf{0}, \dots, \mathbf{0}] \in \mathbb{R}^{n \times n}$. A special class of projectors which isolates the components of a given vector onto a one dimensional subspace spanned by a single unit vector \mathbf{q} called a rank one orthogonal projector, denoted as \mathbf{P}_q . Each k in the classical Gram-Schmidt process creates \mathbf{q}'_k using the following orthogonal projection

$$(50) \quad \mathbf{q}'_k = \mathbf{P}_{A_k^\perp} \mathbf{a}_k$$

where $A_k = \text{l.s.} \{\mathbf{a}_i\}_{i=1}^k$ and $\mathbf{P}_{A_1^\perp} = \mathbb{1}_{n \times n}$ for convenience. A modified version of the Gram-Schmidt process performs the same orthogonal projection broken up as $k - 1$ orthogonal projections of rank $n - 1$ as so

$$\mathbf{q}'_k = \mathbf{P}_{A_k^\perp} \mathbf{a}_k$$

$$\begin{aligned}
&= (\mathbb{1}_{n \times n} - \mathbf{Q}_k \mathbf{Q}_k^*) \mathbf{a}_k \\
&= \left(\prod_{i=1}^{k-1} (\mathbb{1}_{n \times n} - \mathbf{q}_i \mathbf{q}_i^*) \right) \mathbf{a}_k \\
&= (\mathbb{1}_{n \times n} - \mathbf{q}_1 \mathbf{q}_1^*) (\mathbb{1}_{n \times n} - \mathbf{q}_2 \mathbf{q}_2^*) \cdots (\mathbb{1}_{n \times n} - \mathbf{q}_{k-1} \mathbf{q}_{k-1}^*) \mathbf{a}_k \\
&= \mathbf{P}_{\mathbf{q}_k^\perp} \cdots \mathbf{P}_{\mathbf{q}_1^\perp} \mathbf{a}_k
\end{aligned}$$

where $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k, 0, \dots, 0]$. While its clear that $\mathbf{P}_{A_k^\perp} \mathbf{a}$ and $\mathbf{P}_{\mathbf{q}_k^\perp} \cdots \mathbf{P}_{\mathbf{q}_1^\perp} \mathbf{a}_k$ used for computing \mathbf{q}'_k are algebraically equivalent, they differ arithmetically as the latter expression evaluates \mathbf{q}'_k using the follow procedure

$$\begin{aligned}
\mathbf{q}_k^{(1)} &= \mathbf{a}_k \\
\mathbf{q}_k^{(2)} &= \mathbf{P}_{\mathbf{q}_1^\perp} \mathbf{q}_k^{(1)} \\
\mathbf{q}_k^{(3)} &= \mathbf{P}_{\mathbf{q}_2^\perp} \mathbf{q}_k^{(2)} \\
&\vdots \\
\mathbf{q}'_k &= \mathbf{q}_k^{(k)} = \mathbf{P}_{\mathbf{q}_{k-1}^\perp} \mathbf{q}_k^{(k-1)}.
\end{aligned}$$

Applying projections sequentially in this manner produces smaller numerical errors. The modified Gram-Schmidt process [Tre97, Dem97] is summarized in the following algorithm.

Algorithm 8: Modified Gram-Schmidt

```

input : A basis  $\{\mathbf{a}_i\}_{i=1}^n$ .
output: An orthonormal basis  $\{\mathbf{q}_i\}_{i=1}^n$  such that  $\text{l.s } \{\mathbf{a}_i\}_{i=1}^n = \text{l.s } \{\mathbf{q}_i\}_{i=1}^n$ 

1 for  $k = 1$  to  $n$  do
2    $\mathbf{q}'_k = \mathbf{a}_k$ 
3 end
4 for  $k = 1$  to  $n$  do
5    $r_{k,k} = \|\mathbf{q}'_k\|$ 
6    $\mathbf{q}_k = \mathbf{q}'_k / r_{k,k}$ 
7   for  $i = k + 1$  to  $n$  do
8      $r_{i,k} = \langle \mathbf{q}_k, \mathbf{q}'_i \rangle$ 
9      $\mathbf{q}_i = \mathbf{q}_i - r_{i,k} \mathbf{q}_k$ 
10  end
11 end
12 return  $\{\mathbf{q}_i\}_{i=1}^n$ 

```

4.3. Arnoldi and Lanczos Algorithm. As a quick reminder, we are in search of an iterative process to solve the linear system $\mathbf{Ax}^* = \mathbf{b}$ where no explicit form of \mathbf{A} is available and we may only rely on a routine that computes \mathbf{Av} for any \mathbf{v} to extract information on \mathbf{A} . In Section 4.1 it was shown that

$\mathbf{x}^* \in \mathcal{K}_{t_{r_0}, A}(\mathbf{A}, \mathbf{r}_0)$. With many iterative methods, computing an exact value for \mathbf{x}^* is out the question with the view that $t_{r_0, A}$ is impractically large. We must instead resort to approximating \mathbf{x}^* by \mathbf{x}_k for which $\mathbf{x}^k \in \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$ where $k \ll t_{r_0}$. To find an appropriate value for \mathbf{x}_k , a good start would be to find a basis $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$. Definition 31 indicated that $\{\mathbf{A}^{i-1}\mathbf{r}_0\}_{i=1}^k$ serves as a basis for $\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$. However, for numerical reasons this is a poor choice of basis since each consecutive term becomes closer and closer to being linearly dependent. From now on, for more convenient notation we shall set $n = t_{r_0, A}$ so that $\mathbf{x}^* \in \mathcal{K}_n(\mathbf{A}, \mathbf{r}_0)$. To search for a more appropriate basis let $\mathbf{K} \in \mathbb{K}^{n \times n}$ be the invertible matrix

$$\mathbf{K} = [\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{n-1}\mathbf{r}_0].$$

Since \mathbf{K} is invertible we can compute $\mathbf{c} = -\mathbf{K}^{-1}\mathbf{A}^n\mathbf{r}_0$ so that

$$\begin{aligned}\mathbf{AK} &= [\mathbf{Ar}_0, \mathbf{A}^2\mathbf{r}_0, \dots, \mathbf{A}^n\mathbf{r}_0] \\ \mathbf{AK} &= \mathbf{K} \cdot [\mathbf{e}_2, \mathbf{e}_3, \dots, \mathbf{e}_n, -\mathbf{c}] \triangleq \mathbf{KC}\end{aligned}$$

or, in other terms

$$\mathbf{K}^{-1}\mathbf{AK} = \mathbf{C} = \begin{bmatrix} 0 & 0 & \cdots & 0 & -c_1 \\ 1 & 0 & \cdots & 0 & -c_2 \\ 0 & 1 & \cdots & 0 & \vdots \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -c_n \end{bmatrix}.$$

Note here that \mathbf{C} is upper Hessenberg. While this form is simple, it is of little practical use since the matrix \mathbf{K} is very likely to be ill-conditioned. To remedy this we can replace \mathbf{K} with an orthogonal matrix which spans the same space. These are exactly the properties that the \mathbf{V} matrix offers in the QR factorisation of \mathbf{K} . With this in mind let $\mathbf{K} = \mathbf{VR}$ be the full QR factorisation of \mathbf{K} . Then

$$\begin{aligned}\mathbf{AVR} &= \mathbf{AK} \\ \mathbf{AV} &= \mathbf{AKR}^{-1} \\ \mathbf{AV} &= \mathbf{KCR}^{-1} \\ \mathbf{AV} &= \mathbf{VRCR}^{-1} \\ \mathbf{AV} &\triangleq \mathbf{VH}.\end{aligned}$$

Since \mathbf{R} and \mathbf{R}^{-1} and both upper triangular and \mathbf{C} is upper Hessenberg, \mathbf{H} is also upper Hessenberg. This form provides us with a \mathbf{V} such that the range of \mathbf{V} is $\mathcal{K}_n(\mathbf{A}, \mathbf{r}_0)$ and

$$(51) \quad \mathbf{V}^\top \mathbf{AV} = \mathbf{H}.$$

Again, in practice, it may be very difficult to compute this entire expression forcing us to search for approximative alternatives. Consider (51) for which the only first k columns of \mathbf{V} have been computed. Let $\mathbf{V}_k = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$ and $\mathbf{V}_u = [\mathbf{v}_{k+1}, \mathbf{v}_{k+2}, \dots, \mathbf{v}_n]$. Then

$$\begin{aligned}\mathbf{V}^\top \mathbf{AV} &= \mathbf{H} \\ [\mathbf{V}_k, \mathbf{V}_u]^\top \mathbf{A} [\mathbf{V}_k, \mathbf{V}_u] &= \begin{bmatrix} \mathbf{H}_k & \mathbf{H}_{u,k} \\ \mathbf{H}_{k,u} & \mathbf{H}_u \end{bmatrix}\end{aligned}$$

$$\begin{bmatrix} \mathbf{V}_k^\top \mathbf{A} \mathbf{V}_k & \mathbf{V}_k^\top \mathbf{A} \mathbf{V}_u \\ \mathbf{V}_u^\top \mathbf{A} \mathbf{V}_k & \mathbf{V}_u^\top \mathbf{A} \mathbf{V}_u \end{bmatrix} = \begin{bmatrix} \mathbf{H}_k & \mathbf{H}_{u,k} \\ \mathbf{H}_{k,u} & \mathbf{H}_u \end{bmatrix}$$

where \mathbf{H}_k , $\mathbf{H}_{u,k}$, $\mathbf{H}_{k,u}$ and \mathbf{H}_u are the relevant sub matrices. This provides us with the equality

$$(52) \quad \mathbf{V}_k^\top \mathbf{A} \mathbf{V}_k = \mathbf{H}_k$$

noting that \mathbf{H}_k is upper Hessenberg for the same reason that \mathbf{H} is. We know that when $n = t_{\mathbf{r}_0, \mathbf{A}}$ we can find a $\mathbf{V} \in \mathbb{K}^{n \times n}$ and $\mathbf{H} \in \mathbb{K}^{n \times n}$ that satisfies $\mathbf{AV} = \mathbf{VH}$. However, in general, we may not be so fortunate in finding a $\mathbf{V}_k \in \mathbb{K}^{n \times k}$ and $\mathbf{H}_k \in \mathbb{K}^{n \times k}$ to satisfy $\mathbf{AV}_k = \mathbf{V}_k \mathbf{H}_k$ for any $k < n$. Instead we can adjust this equality by adding an error $\mathbf{E}_k \in \mathbb{K}^{n \times k}$ to force an equality. Our expression now becomes

$$(53) \quad \mathbf{V}_k^\top \mathbf{A} \mathbf{V}_k = \mathbf{H}_k + \mathbf{E}_k.$$

A judicious choice of \mathbf{E}_k must be made to also retain equality in (52), meaning $\mathbf{V}_k^\top \mathbf{E}_k = \mathbf{0}$. Since $\{\mathbf{v}_i\}_{i=1}^k$ forms an orthonormal basis for $\mathcal{K}_n(\mathbf{A}, \mathbf{r}_0)$, consider the following choice of \mathbf{E}_k ,

$$\mathbf{E}_k = \mathbf{v}_{k+1} \mathbf{h}_k^\top$$

where \mathbf{h}_k is any vector in \mathbb{K}^k . Notice that

$$\mathbf{V}_k^\top \mathbf{E} = \mathbf{V}^\top (\mathbf{v}_{k+1} \mathbf{h}_k) = (\mathbf{V}^\top \mathbf{v}_{k+1}) \mathbf{h}_k^\top = \mathbf{0}.$$

Since this holds for any $\mathbf{h}_k \in \mathbb{K}^k$, to preserve sparsity and to keep this form as simple as possible we can set $\mathbf{h}_k = [0, 0, \dots, h_{k+1,k}]^\top$. This means \mathbf{AV}_k can be written as

$$(54) \quad \mathbf{AV}_k = \mathbf{V}_k \mathbf{H}_k + \mathbf{v}_{k+1} \mathbf{h}_k^\top$$

where

$$\mathbf{V}_k \mathbf{H}_k = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k] \begin{bmatrix} h_{1,1} & \cdots & \cdots & \cdots & h_{1,k} \\ h_{2,1} & \cdots & \cdots & \cdots & \vdots \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & h_{k,k-1} & h_{k,k} \\ 0 & \cdots & 0 & 0 & 0 \end{bmatrix}.$$

Equating the j^{th} columns of (54) yields

$$\mathbf{Av}_j = \sum_{i=1}^{j+1} h_{i,j} \mathbf{v}_i.$$

Again since $\{\mathbf{v}_i\}_{i=1}^n$ form an orthonormal basis, multiplying both sides by \mathbf{v}_m for $1 \leq m \leq j$ gives

$$\mathbf{v}_m^\top \mathbf{Av}_j = \sum_{i=1}^{j+1} h_{i,j} \mathbf{v}_m^\top \mathbf{v}_i = h_{m,j}$$

and so

$$(55) \quad h_{j+1,j} \mathbf{v}_{j+1} = \mathbf{Av}_j - \sum_{i=1}^j h_{i,j} \mathbf{v}_i.$$

From (55) we find that \mathbf{v}_{j+1} can be computed using a recurrence involving its previous Krylov factors. This bears a striking resemblance to (48) having a virtually identical setup to computing an orthonormal

basis using the modified Gram-Schmidt process (Algorithm 8). As such, values for \mathbf{v}_{j+1} and $h_{j+1,j}$ can be evaluated using a procedure very similar to the modified Gram-Schmidt process better known as the Arnoldi Algorithm [Tre97, Dem97], presented in Algorithm 9.

Algorithm 9: Arnoldi Algorithm

```

input :  $\mathbf{A}$ ,  $r_0$  and  $k$ , the number of columns of  $\mathbf{V}$  to compute.
output:  $\mathbf{V}_k$ ,  $\mathbf{H}_k$ .

1  $\mathbf{v}_1 = \mathbf{r}_0 / \|\mathbf{r}_0\|$ 
2 for  $j = 1$  to  $k$  do
3    $\mathbf{z} = \mathbf{A}\mathbf{v}_j$ 
4   for  $i = 1$  to  $j$  do
5      $h_{i,j} = \langle \mathbf{v}_i, \mathbf{z} \rangle$ 
6      $\mathbf{z} = \mathbf{z} - h_{i,j}\mathbf{v}_i$ 
7   end
8    $h_{j+1,j} = \|\mathbf{z}\|$ 
9   if  $h_{j+1,j} = 0$  then
10    | return  $\mathbf{V}_k$ ,  $\mathbf{H}_k$ 
11   end
12    $\mathbf{v}_{j+1} = \mathbf{z} / h_{j+1,j}$ 
13 end
14 return  $\mathbf{V}_k$ ,  $\mathbf{H}_k$ 

```

When \mathbf{A} is symmetric then $\mathbf{H} = \mathbf{T}$ becomes a tridiagonal matrix, simplifying a large amount of the Arnoldi algorithm since the matrix elements from \mathbf{T} can be written as

$$\mathbf{T} = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \ddots & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & \beta_{n-1} \\ & & & \beta_{n-1} & \alpha_n \end{bmatrix}.$$

As before, equating the j^{th} columns of $\mathbf{AV} = \mathbf{VT}$ yields

$$(56) \quad \mathbf{Av}_j = \beta_{j-1}\mathbf{v}_{j-1} + \alpha_j\mathbf{v}_j + \beta_j\mathbf{v}_{j+1}.$$

Again since $\{\mathbf{v}_i\}_{i=1}^n$ form an orthonormal basis, multiplying both sides of (56) by \mathbf{v}_j gives $\mathbf{v}_j \mathbf{Av}_j = \alpha_j$. This simplified version of the Arnoldi algorithm used for computing $\{\mathbf{v}_i\}_{i=1}^n$ and \mathbf{T} for symmetric matrices is better known as the Lanczos algorithm [Dem97], presented more explicitly in Algorithm 10.

Algorithm 10: Lanczos Algorithm

```

input :  $\mathbf{A}, \mathbf{r}_0$  and  $k$ , the number of columns of  $\mathbf{V}$  to compute.
output:  $\mathbf{V}_k, \mathbf{T}_k$ .

1  $\mathbf{v}_1 = \mathbf{r}_0 / \|\mathbf{r}_0\|, \beta_0 = 0, \mathbf{v}_0 = 0$ 
2 for  $j = 1$  to  $k$  do
3    $\mathbf{z} = \mathbf{A}\mathbf{v}_j$ 
4    $\alpha_j = \langle \mathbf{v}_j, \mathbf{z} \rangle$ 
5    $\mathbf{z} = \mathbf{z} - \alpha_j \mathbf{v}_j - \beta_{j-1} \mathbf{v}_{j-1}$ 
6    $\beta_j = \|\mathbf{z}\|$ 
7   if  $\beta_j = 0$  then
8     return  $\mathbf{V}_k, \mathbf{T}_k$ 
9   end
10   $\mathbf{v}_{j+1} = \mathbf{z} / \beta_j$ 
11 end
12 return  $\mathbf{V}_k, \mathbf{T}_k$ 

```

For the Lanczos algorithm, (54) can be re-written in the a more compact form as

$$(57) \quad \mathbf{AV}_k \triangleq \mathbf{V}_k \mathbf{T}_{k+1,k}$$

where $\mathbf{T}_{k+1,k} = \mathbf{T}_k + \mathbf{v}_{k+1} \mathbf{t}_k^\top$.

4.4. Optimality Conditions. So far we have shown that $\mathbf{x}^* \in \mathcal{K}_{t_{\mathbf{r}_0}, \mathbf{A}}(\mathbf{A}, \mathbf{r}_0)$ where $n = t_{\mathbf{r}_0}$ is the grade of \mathbf{r}_0 with respect to \mathbf{A} . Moreover from Section 4.3 we found ways to construct a basis for $\mathcal{K}_{t_{\mathbf{r}_0}, \mathbf{A}}(\mathbf{A}, \mathbf{r}_0)$ allowing us to generate vectors with these affine spaces, namely the Arnoldi algorithm (Algorithm 9) and Lanczos algorithm (Algorithm 10) for non-symmertic and symmertic systems, respectively. From now on $\mathcal{K}_{t_{\mathbf{r}_0}, \mathbf{A}}(\mathbf{A}, \mathbf{r}_0)$ will be abbreviated to $\mathcal{K}_{t_{\mathbf{r}_0}, \mathbf{A}}$ when the context is clear. The question still remains however, how should one choose an \mathbf{x}_k that best approximates \mathbf{x}^* satisfying (46)? The following are some of the most well known methods for selecting a suitable \mathbf{x}_k .

- (a) Select an $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k$ which minimizes $\|\mathbf{x}_k - \mathbf{x}^*\|_2$. While this method seems like the most intuitive and natural way to select \mathbf{x}_k , it is unfortunately of no practical use since there is not enough information in the Krylov Subspace to find an \mathbf{x}_k which matches this profile.
- (b) Select an $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k$ which minimizes $\|\mathbf{r}_k\|_2$ (recall this is the residual of \mathbf{x}_k , that is, $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$). This method is possible to implement. Two well known algorithms stem from this class of methods, namely MINRES (minimum residual) and GMRES (general minimum residual) which solve linear systems for symmetric and non-symmertic \mathbf{A} respectively.
- (c) When \mathbf{A} is a positive definite matrix it defines a norm $\|\mathbf{r}\|_{\mathbf{A}} = (\mathbf{r}^\top \mathbf{A} \mathbf{r})^{1/2}$, called the energy norm. Select an $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k$ which minimizes $\|\mathbf{r}_k\|_{\mathbf{A}^{-1}}$ which is equivalent to minimizing $\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{A}}$. This technique is known as the CG (conjugate gradient) algorithm.

- (d) Select an $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k$ for which $\mathbf{r}_k \perp \mathcal{W}_k$ where \mathcal{W}_k is some k -dimensional subspace. Two well known algorithms that belong to this family of methods are SYMMLQ (Symmetric LQ Method) and a variant of GMRES used for solving symmetric and non-symmetric methods respectively.

Interestingly, when \mathbf{A} is symmetric positive definite and $\mathcal{W}_k = \mathcal{K}_k$ the last two selection methods are equivalent. This is stated more precisely in Theorem 37.

Theorem 37. *In the context of the above selection method, if $\mathbf{A} \succ \mathbf{0}$ and $\mathcal{W}_k = \mathcal{K}_k$ in method (4) then it produces the same \mathbf{x}_k in method (3) [Dem97].*

Proof. Since $\mathcal{W}_k = \mathcal{K}_k$, any basis for \mathcal{W}_k is a basis for \mathcal{K}_k . In fact, we can write $\mathbf{W}_k = \mathbf{K}_k \mathbf{B}$ for some non-singular $\mathbf{B} \in \mathbb{R}^{k \times k}$. This means

$$\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k = \mathbf{B}^\top \mathbf{K}_k^\top \mathbf{A} \mathbf{K}_k$$

and since \mathbf{A} is symmetric positive definite then $\mathbf{K}_k^\top \mathbf{A} \mathbf{K}_k$ is also symmetric positive definite, and hence the entire product is non-singular. \square

In fact the last method can be used to bring together a number of different analytical aspects and unify them in a general framework known as projection methods. Selecting an \mathbf{x}_k from our Krylov subspace allows k degrees of freedom meaning k constraints must be used to determine a unique \mathbf{x}_k for selection. As seen in method (4) already, typically orthogonality constraints are imposed on the residual \mathbf{r}_k . Specifically we would like to find a $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{K}_k$ where $\mathbf{r}_k \perp \mathcal{W}_k$. This is sometimes referred to as the Petrov-Galerkin (or just Galerkin) conditions. Projection methods for which $\mathcal{W}_k = \mathcal{K}_k$ are known as orthogonal projections while methods for which $\mathcal{W}_k = \mathbf{A}\mathcal{K}_k$ are known as oblique projections. If we set $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{z}_k$ for some $\mathbf{z}_k \in \mathcal{K}_k$ then the Petrov-Galerkin conditions imply $\mathbf{r}_0 - \mathbf{A}\mathbf{z}_k \perp \mathcal{W}_k$, or alternatively $\langle \mathbf{r}_0 - \mathbf{A}\mathbf{z}_k, \mathbf{w} \rangle = 0$ for every $\mathbf{w} \in \mathcal{W}_k$. To impose these conditions it will help to have an appropriate basis for \mathcal{K} and \mathcal{W} . Suppose we have access to such a basis where $\{\mathbf{v}_i\}_{i=1}^k$ and $\{\mathbf{w}_i\}_{i=1}^k$ are basis elements for \mathcal{K} and \mathcal{W} respectively. Let

$$\begin{aligned} \mathbf{K}_k &\triangleq [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k] \in \mathbb{R}^{n \times k} \\ \mathbf{W}_k &\triangleq [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k] \in \mathbb{R}^{n \times k} \end{aligned}$$

then the Petrov-Galerkin conditions can be imposed as follows

$$\begin{aligned} \mathbf{K}_k \mathbf{y}_k &= \mathbf{z}_k, \quad \text{for some } \mathbf{y}_k \in \mathbb{R}^k \\ \mathbf{W}_k^\top (\mathbf{r}_0 - \mathbf{A}\mathbf{K}_k \mathbf{y}_k) &= \mathbf{0}. \end{aligned}$$

Moreover if $\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k$ is invertible then \mathbf{x}_k can be expressed as

$$(58) \quad \mathbf{x}_k = \mathbf{x}_0 + \mathbf{K}_k (\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k)^{-1} \mathbf{W}_k \mathbf{r}_0.$$

This justifies a general form of the projection method algorithm presented in Algorithm 11.

Algorithm 11: General Projection Method

output: An approximation of \mathbf{x}^* , \mathbf{x}_k .

```

1 for  $k = 1, \dots$  until convergence do
2   Select  $\mathcal{K}_k$  and  $\mathcal{W}_k$ 
3   Form  $\mathbf{K}_k$  and  $\mathbf{W}_k$ 
4   Solve  $(\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k) \mathbf{y}_k = \mathbf{W}_k^\top \mathbf{r}_0$ 
5    $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{K}_k \mathbf{y}_k$ 
6 end
7 return  $\mathbf{x}_k$ 

```

4.5. Conjugate Gradient Algorithm. From Section 4.4 that the Petrov-Galerkin conditions for the CG algorithm used an orthogonal projection and the matrix \mathbf{A} was assumed to be positive definite. To derive the CG algorithm we can start by using some machinery that the Lanczos algorithm provides us with. Recall, the Lanczos algorithm produces the form $\mathbf{A}\mathbf{V}_k = \mathbf{V}_k \mathbf{T}_k + \mathbf{v}_{k+1} \mathbf{t}_k^\top$ where $\mathbf{t}_k \triangleq [0, 0, \dots, 0, \beta_k]^\top \in \mathbb{K}^k$ and the columns of \mathbf{V}_k span \mathcal{K}_k . Recall that \mathbf{x}_k can be expressed as $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{K}_k (\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k)^{-1} \mathbf{W}_k \mathbf{r}_0$ ((58)) when $\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k$ is invertible. For the CG algorithm $\mathcal{K} = \mathcal{W}$ and $\mathbf{A} \succ 0$. Under these conditions we can easily show that $\mathbf{W}_k^\top \mathbf{A} \mathbf{K}_k$ is indeed invertible. Thus the approximate vector can be expressed as $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{z}_k$ where $\mathbf{z}_k \in \mathcal{K}_k$. In terms of the Petrov-Galerkin conditions this means that \mathbf{z}_k must satisfy $\mathbf{r}_0 - \mathbf{A}\mathbf{z}_k \perp \mathcal{W}_k$. Furthermore since $\mathcal{K}_k = \text{Range}(\mathbf{V}_k)$ where \mathbf{V}_k has full column rank then \mathbf{z}_k can be represented as $\mathbf{z}_k = \mathbf{V}_k \mathbf{y}$ for a unique $\mathbf{y} \in \mathbb{R}^k$ so that

$$(59) \quad \mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{y}.$$

Coupling this with the Petrov-Galerkin conditions ensures

$$\begin{aligned}
(60) \quad & \mathbf{V}_k^\top (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_k \mathbf{y}) = \mathbf{0} \\
& \mathbf{V}_k^\top \mathbf{A} \mathbf{V}_k \mathbf{y} = \mathbf{V}_k^\top \mathbf{r}_0 \\
& \mathbf{T}_k \mathbf{y} = \|\mathbf{r}_0\| e_1.
\end{aligned}$$

In the CG algorithm \mathbf{x}_{k+1} is computed as the recurrence of the following three sets of vectors

- (a) The approximate solutions \mathbf{x}_k
- (b) The residual vectors \mathbf{r}_k
- (c) The conjugate gradient vectors \mathbf{p}_k

The conjugate gradient vectors include the word 'gradient' since the attempt to find the direction of steepest descent that minimizes $\|\mathbf{r}_k\|_{\mathbf{A}^{-1}}$. It also includes the word 'conjugate' since $\langle \mathbf{p}_k, \mathbf{A} \mathbf{p}_j \rangle = 0$ for $i \neq j$, that is, vectors \mathbf{p}_i and \mathbf{p}_j are mutually A -conjugate (shown in Lemma 38).

Since \mathbf{A} is symmetric positive definite then so is $\mathbf{T}_k = \mathbf{V}_k \mathbf{A} \mathbf{V}_k$. We can take the Cholesky decomposition of \mathbf{T}_k to get

$$(61) \quad \mathbf{T}_k = \mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^\top$$

where \mathbf{L}_k is a unit lower bidiagonal matrix and \mathbf{D}_k is diagonal written as

$$\mathbf{L}_k = \begin{bmatrix} 1 & & & \\ l_1 & \ddots & & \\ & \ddots & \ddots & \\ & & l_{k-1} & 1 \end{bmatrix}, \quad \mathbf{D}_k = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_k \end{bmatrix}.$$

Combining equations (59), (60) and (61)

$$\begin{aligned} \mathbf{x}_k &= \mathbf{x}_0 + \mathbf{V}_k \mathbf{y} \\ \mathbf{x}_k &= \mathbf{x}_0 + \|r_0\| \mathbf{V}_k \mathbf{T}_k^{-1} \mathbf{e}_1 \\ \mathbf{x}_k &= \mathbf{x}_0 + \|r_0\| \mathbf{V}_k (\mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^\top)^{-1} \mathbf{e}_1 \\ \mathbf{x}_k &= \mathbf{x}_0 + (\mathbf{V}_k \mathbf{L}_k^{-\top}) (\|r_0\| \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} \mathbf{e}_1) \\ \mathbf{x}_k &\triangleq \mathbf{x}_0 + \tilde{\mathbf{P}}_k \tilde{\mathbf{y}}_k \end{aligned}$$

where $\tilde{\mathbf{P}}_k = \mathbf{V}_k \mathbf{L}_k^{-\top}$ and $\tilde{\mathbf{y}}_k = \|r_0\| \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} \mathbf{e}_1$. The matrix $\tilde{\mathbf{P}}_k$ can be written as $\tilde{\mathbf{P}}_k = [\tilde{\mathbf{p}}_1, \tilde{\mathbf{p}}_2, \dots, \tilde{\mathbf{p}}_k]$. Lemma 38 shows that the columns of $\tilde{\mathbf{P}}_k$ are A -conjugate.

Lemma 38. *The columns of $\tilde{\mathbf{P}}_k$ are A -conjugate, in otherwords $\tilde{\mathbf{P}}_k^\top A \tilde{\mathbf{P}}_k$ is diagonal.*

Proof. We compute

$$\begin{aligned} \tilde{\mathbf{P}}_k^\top A \tilde{\mathbf{P}}_k &= (\mathbf{V}_k \mathbf{L}_k^{-\top})^\top A (\mathbf{V}_k \mathbf{L}_k^{-\top}) \\ &= \mathbf{L}_k^{-1} (\mathbf{V}_k^\top A \mathbf{V}_k) \mathbf{L}_k^{-\top} \\ &= \mathbf{L}_k^{-1} (\mathbf{T}_k) \mathbf{L}_k^{-\top} \\ &= \mathbf{L}_k^{-1} (\mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^\top) \mathbf{L}_k^{-\top} \\ &= \mathbf{D}_k \end{aligned} \tag{61}$$

as wanted. \square

Since \mathbf{L}_k is a lower bidiagonal, setting $\mathbf{a} \triangleq l_{k-1} \mathbf{e}_{k-1}$, it can be written in the form

$$\mathbf{L}_k = \begin{bmatrix} \mathbf{L}_{k-1} & \mathbf{0} \\ \mathbf{a}^\top & 1 \end{bmatrix}$$

meaning

$$\mathbf{L}_k^{-1} = \begin{bmatrix} \mathbf{L}_{k-1}^{-1} & \mathbf{0} \\ \star & 1 \end{bmatrix}$$

where the elements of \star are of no importance. With this a recurrence for the columns of $\tilde{\mathbf{P}}_k$ can now be derived in terms of \mathbf{y}_k . To start we can show that the first $k - 1$ entries of $\tilde{\mathbf{y}}_k$ shares the first $k - 1$ entires with $\tilde{\mathbf{y}}_{k-1}$ and that $\tilde{\mathbf{P}}_k$ and $\tilde{\mathbf{P}}_{k-1}$ share the same first $k - 1$ columns. We can begin by computing a recurrence for $\tilde{\mathbf{y}}_k$ as follows

$$\tilde{\mathbf{y}}_k = \|r_0\| \mathbf{D}_k^{-1} \mathbf{L}_k^{-1} \mathbf{e}_1^k$$

$$\begin{aligned}
&= \|\mathbf{r}_0\| \begin{bmatrix} \mathbf{D}_{k-1}^{-1} & \mathbf{0} \\ \mathbf{0} & d_k^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{k-1}^{-1} & \mathbf{0} \\ * & 1 \end{bmatrix} \mathbf{e}_1^k \\
&= \|\mathbf{r}_0\| \begin{bmatrix} \mathbf{D}_{k-1}^{-1} \mathbf{L}_{k-1}^{-1} & \mathbf{0} \\ * & d_k^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{e}_1^k \\ 0 \end{bmatrix} \\
&= \begin{bmatrix} \tilde{\mathbf{y}}_{k-1} \\ \eta_k \end{bmatrix}
\end{aligned}$$

where \mathbf{e}_i^k is the i^{th} unit vector with k dimensions. To get a recurrence for the columns of $\tilde{\mathbf{P}}_{k-1} = [\tilde{\mathbf{p}}_1, \tilde{\mathbf{p}}_2, \dots, \tilde{\mathbf{p}}_k]$ since \mathbf{L}_{k-1}^\top is upper triangular then so is $\mathbf{L}_{k-1}^{-\top}$, thus forming the leading $(k-1) \times (k-1)$ submatrix of $\mathbf{L}_k^{-\top}$. In effect, $\tilde{\mathbf{P}}_{k-1}$ is identical to the leading $k-1$ columns of

$$\tilde{\mathbf{P}}_k = \mathbf{V}_k \mathbf{L}_k^{-\top} = [\mathbf{V}_{k-1}, \mathbf{v}_k] \begin{bmatrix} \mathbf{L}_{k-1}^{-1} & \mathbf{0} \\ * & 1 \end{bmatrix} = [\mathbf{V}_{k-1} \mathbf{L}_{k-1}^{-1}, \tilde{\mathbf{p}}_k] = [\tilde{\mathbf{P}}_{k-1}, \tilde{\mathbf{p}}_k].$$

Moreover rearranging $\tilde{\mathbf{P}}_k = \mathbf{V}_k \mathbf{L}_k^{-\top}$ we get $\tilde{\mathbf{P}}_k \mathbf{L}_k^\top = \mathbf{V}_k$. Equating the k^{th} column yields

$$(62) \quad \tilde{\mathbf{p}}_k = \mathbf{v}_k - l_{k-1} \tilde{\mathbf{p}}_{k-1}.$$

Finally we can use

$$(63) \quad \mathbf{x}_k = \mathbf{x}_0 + \tilde{\mathbf{P}}_k \tilde{\mathbf{y}}_k = \mathbf{x}_0 + [\tilde{\mathbf{P}}_{k-1}, \tilde{\mathbf{p}}_k] \begin{bmatrix} \tilde{\mathbf{y}}_{k-1} \\ \eta_k \end{bmatrix} = \mathbf{x}_0 + \tilde{\mathbf{P}}_{k-1} \tilde{\mathbf{y}}_{k-1} + \eta_k \tilde{\mathbf{p}}_k = \mathbf{x}_{k-1} + \eta_k \tilde{\mathbf{p}}_k$$

as a recurrence for \mathbf{x}_k . A recurrence for \mathbf{r}_k is easily computed as

$$(64) \quad \mathbf{r}_k = b - \mathbf{A}\mathbf{x}_k = b - \mathbf{A}(\mathbf{x}_{k-1} + \eta_k \tilde{\mathbf{p}}_k) = (b - \mathbf{A}\mathbf{x}_{k-1}) - \eta_k \mathbf{A}\tilde{\mathbf{p}}_k = \mathbf{r}_{k-1} - \eta_k \mathbf{A}\tilde{\mathbf{p}}_k$$

Altogether we are left with recurrences for \mathbf{v}_k from Lanczos, $\tilde{\mathbf{p}}_k$ (62), the residual \mathbf{r}_k (64), and for the approximate solution \mathbf{x}_k (63). However, further simplifications can be made to bring about a more efficient algorithm. Recall from Section 4.3 that $\mathbf{A}\mathbf{V}_k = \mathbf{V}_k \mathbf{T}_k + \mathbf{v}_{k+1} \mathbf{t}_k^\top$ where $\mathbf{t}_k = [0, 0, \dots, 0, \beta_k]^\top \in \mathbb{R}^k$ meaning

$$\mathbf{r}_k = \mathbf{r}_0 - \mathbf{A}\mathbf{V}_k \mathbf{y}_k = \mathbf{r}_0 - \mathbf{V}_k \mathbf{T}_k \mathbf{y}_k - \langle \mathbf{t}_k, \mathbf{y} \rangle \mathbf{v}_{k+1} = -\beta_k y_k \mathbf{v}_{k+1}.$$

This tells us that \mathbf{r}_k is parallel to \mathbf{v}_{k+1} and orthogonal to all \mathbf{v}_i , $1 \leq i \leq k$. This further implies that \mathbf{r}_k is orthogonal to all \mathbf{r}_i , $1 \leq i \leq k-1$ since they are just \mathbf{v}_i scaled by some constant factor. So replacing \mathbf{r}_{k-1} with \mathbf{v}_k/η_k and defining $\mathbf{p}_k \triangleq \tilde{\mathbf{p}}_k/\gamma_k$ gives us a new set of recurrences

$$\begin{aligned}
\mathbf{x}_k &= \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k \\
\mathbf{r}_k &= \mathbf{r}_{k-1} - \alpha_k \mathbf{A}\mathbf{p}_k \\
\mathbf{p}_k &= \mathbf{r}_{k-1} + \beta_k \mathbf{p}_{k-1}
\end{aligned}$$

where $\alpha_k = \eta_k/\gamma_k$. From Lemma 38 we have shown that the columns of $\tilde{\mathbf{P}}_k$ are A -conjugate (that is $\langle \tilde{\mathbf{p}}_i, \mathbf{A}\tilde{\mathbf{p}}_j \rangle = 0$, $i \neq j$) and that $\tilde{\mathbf{P}}_k^\top \mathbf{A} \tilde{\mathbf{P}}_k = \mathbf{D}_k$. This also means that $\langle \mathbf{r}_i, \mathbf{r}_j \rangle = 0$, $i \neq j$. Note that from our recurrence for $\mathbf{p}_k = \mathbf{r}_{k-1} + \beta_k \mathbf{p}_{k-1}$ we have

$$\langle \mathbf{A}\mathbf{p}_k, \mathbf{p}_k \rangle = \langle \mathbf{A}\mathbf{p}_k, \mathbf{r}_{k-1} + \beta_k \mathbf{p}_{k-1} \rangle = \langle \mathbf{A}\mathbf{p}_k, \mathbf{r}_{k-1} \rangle.$$

We can now find an expression for α_k as

$$\langle \mathbf{r}_{k-1}, \mathbf{r}_k \rangle = \langle \mathbf{r}_{k-1}, \mathbf{r}_{k-1} - \alpha_k \mathbf{A}\mathbf{p}_k \rangle$$

$$\begin{aligned}\langle \mathbf{r}_{k-1} - \mathbf{r}_k, \mathbf{r}_{k-1} \rangle &= \langle \mathbf{r}_{k-1}, \mathbf{r}_{k-1} \rangle - \alpha_k \langle \mathbf{p}_k, \mathbf{A}\mathbf{p}_k \rangle \\ \alpha_k &= \frac{\langle \mathbf{r}_{k-1}, \mathbf{r}_{k-1} \rangle}{\langle \mathbf{p}_k, \mathbf{A}\mathbf{p}_k \rangle}.\end{aligned}$$

Similarly, using the recurrence for \mathbf{p}_k , an expression for β_k can be computed as

$$\begin{aligned}\langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{p}_k \rangle &= \langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{r}_{k-1} + \beta_k \mathbf{p}_{k-1} \rangle \\ \langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{p}_k \rangle &= \langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{r}_{k-1} \rangle + \beta_k \langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{p}_{k-1} \rangle \\ \beta_k &= -\frac{\langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{r}_{k-1} \rangle}{\langle \mathbf{A}\mathbf{p}_{k-1}, \mathbf{p}_{k-1} \rangle}.\end{aligned}$$

This formula requires an additional dot product which was not present before. Fortunately, this dot product can be eliminated using our recurrence for \mathbf{r}_k

$$\begin{aligned}\langle \mathbf{r}_k, \mathbf{r}_k \rangle &= \langle \mathbf{r}_k, \mathbf{r}_{k-1} - \alpha_k \mathbf{A}\mathbf{p}_k \rangle \\ \langle \mathbf{r}_k, \mathbf{r}_k \rangle &= \langle \mathbf{r}_k, \mathbf{r}_{k-1} \rangle - \alpha_k \langle \mathbf{r}_k, \mathbf{A}\mathbf{p}_k \rangle \\ \alpha_k &= -\frac{\langle \mathbf{r}_k, \mathbf{r}_k \rangle}{\langle \mathbf{r}_k, \mathbf{A}\mathbf{p}_k \rangle}.\end{aligned}$$

Equating the two expressions for α_k affords

$$\begin{aligned}-\frac{\langle \mathbf{r}_k, \mathbf{r}_k \rangle}{\langle \mathbf{r}_k, \mathbf{A}\mathbf{p}_k \rangle} &= \frac{\langle \mathbf{r}_{k-1}, \mathbf{r}_{k-1} \rangle}{\langle \mathbf{p}_k, \mathbf{A}\mathbf{p}_k \rangle} \\ -\frac{\langle \mathbf{r}_k, \mathbf{r}_k \rangle}{\langle \mathbf{r}_{k-1}, \mathbf{r}_{k-1} \rangle} &= \frac{\langle \mathbf{r}_k, \mathbf{A}\mathbf{p}_k \rangle}{\langle \mathbf{p}_k, \mathbf{A}\mathbf{p}_k \rangle}.\end{aligned}$$

This means that

$$\beta_k = \frac{\langle \mathbf{r}_{k-1}, \mathbf{r}_{k-1} \rangle}{\langle \mathbf{r}_{k-2}, \mathbf{r}_{k-2} \rangle}.$$

These recurrences are computed iteratively to form the basis of the CG algorithm, seen in Algorithm 12.

Algorithm 12: CG Algorithm

input : $A \succ 0$, b and an initial guess x_0 .
output: An approximation of x^* , x_k .

```

1  $r_0 = b - Ax_0$ ,  $p_1 = r_0$ 
2 for  $k = 1, \dots$  until  $\|r_{k-1}\| \leq \tau$  do
3    $\alpha_k = \frac{\langle r_{k-1}, r_{k-1} \rangle}{\langle p_k, A p_k \rangle}$ 
4    $x_k = x_{k-1} + \alpha_k p_k$ 
5    $r_k = r_{k-1} - \alpha_k A p_k$ 
6    $\beta_{k+1} = \frac{\langle r_k, r_k \rangle}{\langle r_{k-1}, r_{k-1} \rangle}$ 
7    $p_{k+1} = r_k + \beta_{k+1} p_k$ 
8 end
9 return  $x_k$ 

```

The next few theorems pertain to the rate of convergence of the CG algorithm.

Theorem 39. Let the CG iteration (Algorithm 12) be applied to a symmetric positive definite matrix problem $\mathbf{A}\mathbf{x} = \mathbf{b}$. As long as the iteration has not yet converged (that is, $\mathbf{r}_{n-1} \neq 0$), the algorithm proceeds without divisions by zero, and we have the following identities of subspaces:

$$(65) \quad \mathcal{K}_n = \text{l.s } \{\mathbf{x}_i\}_{i=1}^n = \text{l.s } \{\mathbf{p}_i\}_{i=0}^{n-1} = \text{l.s } \{\mathbf{r}_i\}_{i=0}^{n-1} = \text{l.s } \{\mathbf{A}^i \mathbf{b}\}_{i=0}^{n-1}.$$

Moreover the residuals are orthogonal,

$$\langle \mathbf{r}_n, \mathbf{r}_j \rangle = 0 \quad (j < n)$$

and the search directions are \mathbf{A} -conjugate [Tre97, page 295].

Proof. The orthogonality of the residuals has already been shown within the derivation of the CG algorithm. Seeing that the search directions are \mathbf{A} -conjugate is a very straight forward application of Lemma 38. To show (65) we can induct on k . From the initial guess \mathbf{x}_0 and the expression $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k$, using an induction argument it follows that \mathbf{x}_k belongs to $\text{l.s } \{\mathbf{p}_i\}_{i=0}^{k-1}$. Similarly, from $\mathbf{p}_k = \mathbf{r}_{k-1} + \beta_k \mathbf{p}_{k-1}$ it follows that this belongs to $\text{l.s } \{\mathbf{r}_i\}_{i=0}^{n-1}$. Finally, from $\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_k \mathbf{A} \mathbf{p}_k$ that this belongs to $\text{l.s } \{\mathbf{A}^i \mathbf{b}\}_{i=0}^{n-1}$. \square

It is fairly straight forward to confirm that CG minimizes error at each step.

Theorem 40. Let the CG iteration be applied to a symmetric positive definite matrix $\mathbf{A}\mathbf{x} = \mathbf{b}$. If the iteration has not already converged (that is, $\|\mathbf{r}_{k-1}\| > 0$), then \mathbf{x}_k is the unique point in \mathcal{K}_k that minimizes $\|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{A}}$. The convergence is the monotonic,

$$\|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{A}} \leq \|\mathbf{x}^* - \mathbf{x}_{k-1}\|_{\mathbf{A}}$$

and $\|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{A}} = 0$ is achieved for some $k \leq n$ [Tre97, page 296].

Proof. From Theorem 39 we know that \mathbf{x}_k belongs to \mathcal{K}_k . To show that it is the unique point in \mathcal{K}_k that minimizes the error, consider an arbitrary point $\mathbf{z} = \mathbf{x}_k - \Delta \mathbf{z} \in \mathcal{K}_k$, with error $\mathbf{e} = \mathbf{x}^* - \mathbf{z} = \mathbf{e}_k + \Delta \mathbf{x}$. Note in this context $\mathbf{e}_k = \mathbf{x}^* - \mathbf{x}_k$, and does not represent standard basis vectors. We calculate

$$\begin{aligned} \|\mathbf{e}\|_{\mathbf{A}}^2 &= (\mathbf{e}_k + \Delta \mathbf{x})^\top \mathbf{A} (\mathbf{e}_k + \Delta \mathbf{x}) \\ &= \mathbf{e}_k^\top \mathbf{A} \mathbf{e}_k + (\Delta \mathbf{x})^\top \mathbf{A} (\Delta \mathbf{x}) + 2\mathbf{e}_k^\top \mathbf{A} (\Delta \mathbf{x}). \end{aligned}$$

The final term of this expression is $2\mathbf{e}_k^\top (\Delta \mathbf{x})$, an inner product of \mathbf{e}_k with a vector in \mathcal{K}_k , and by Theorem 40, any such inner product is zero. Effectively, this means

$$\|\mathbf{e}\|_{\mathbf{A}}^2 = \mathbf{e}_k^\top \mathbf{A} \mathbf{e}_k + (\Delta \mathbf{x})^\top \mathbf{A} (\Delta \mathbf{x}).$$

Only the second of these terms depends on $\Delta \mathbf{x}$, since \mathbf{A} is positive definite, that term is greater than or equal to 0, attaining the value of 0 if and only if $\Delta \mathbf{x} = \mathbf{0}$, that is, $\mathbf{x}_k = \mathbf{x}$. Thus $\|\mathbf{e}\|_{\mathbf{A}}$ is minimal if and only if $\mathbf{x}_k = \mathbf{x}$, as claimed. \square

As seen in (47) Krylov subspace methods, in some sense, aim to find a polynomial $p_n \in P_n$ which minimizes $\|p_n(\mathbf{A})\| = \|\sum_{k=0}^{n-1} \alpha_{k+1} \mathbf{A}^k \mathbf{r}_0\|$, where P_n is the set of polynomials with degree less than or equal to n with $p(\mathbf{0}) = \mathbf{1}$. In the context of the CG algorithm, this is equivalent to finding a polynomial $p_n \in P_n$ which instead minimizes

$$(66) \quad \|p_n(\mathbf{A}) \mathbf{e}_0\|_{\mathbf{A}}.$$

From Theorem 39, we can derive the following convergence theorem.

Theorem 41. *If the CG iteration has not already converged before step k (that is, $\|\mathbf{r}_{k-1}\| > 0$), then (66) has a unique solution $p_n \in P_n$, and the iterate \mathbf{x}_k has error $\mathbf{e}_k = p_n(\mathbf{A})\mathbf{e}_0$ for this same polynomial p_n . Consequently we have*

$$(67) \quad \frac{\|\mathbf{e}_k\|_{\mathbf{A}}}{\|\mathbf{e}_0\|_{\mathbf{A}}} = \inf_{p \in P_n} \frac{\|p(\mathbf{A})\mathbf{e}_0\|_{\mathbf{A}}}{\|\mathbf{e}_0\|_{\mathbf{A}}} \leq \inf_{p \in P_n} \max_{\lambda \in \Lambda(\mathbf{A})} |p(\lambda)|$$

[Tre97, page 298].

Proof. From Theorem 39 it follows that $\mathbf{e}_k = p(\mathbf{A})\mathbf{e}_0$ for some $p \in P_n$. The equality in (67) is a consequence of this and Theorem 40. As for the inequality in (67), if $\mathbf{e}_0 = \sum_{j=1}^m a_j \mathbf{v}_j$ is an expansion of \mathbf{e}_0 in orthonormal eigenvectors of \mathbf{A} , then we have $p(\mathbf{A})\mathbf{e}_0 = \sum_{j=1}^m a_j p(\lambda_j) \mathbf{v}_j$ and thus

$$\|\mathbf{e}_0\|_{\mathbf{A}}^2 = \sum_{j=1}^m a_j^2 \lambda_j, \quad \|p(\mathbf{A})\mathbf{e}_0\|_{\mathbf{A}}^2 = \sum_{j=1}^m a_j^2 \lambda_j (p(\lambda_j))^2.$$

These identities indicate $\|p(\mathbf{A})\mathbf{e}_0\|_{\mathbf{A}}^2 / \|\mathbf{e}_0\|_{\mathbf{A}}^2 \leq \max_{\lambda \in \Lambda(\mathbf{A})} |p(\lambda)|^2$, which implies the inequality in question. \square

4.6. Minimum Residual. In contrast to CG, MINRES is applicable to a wider range of linear systems, namely symmetric indefinite systems. From Section 4.4 we saw that MINRES minimizes the residual \mathbf{r}_k with respect to the Euclidean norm at each iteration (hence the name), that is \mathbf{x}_k is chosen so that

$$(68) \quad \mathbf{x}_k = \arg \min_{\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k} \|\mathbf{Ax} - \mathbf{b}\|_2.$$

It can be shown that this is equivalent setting $\mathcal{W} = \mathbf{A}\mathcal{K}_k$ in the Petrov-Galerkin conditions where $\det(\mathbf{A}) \neq 0$, in other words MINRES is an oblique projection method. We can also show that when $\det(\mathbf{A}) \neq 0$ and $\mathcal{W} = \mathbf{A}\mathcal{K}_k$ the matrix $\mathbf{W}^\top \mathbf{A} \mathbf{K}$ is non-singular, similar to Theorem 37. So under the conditions of MINRES, using a similar argument for CG, \mathbf{x}_k can be expressed as $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{y}_k$. Employing (57) produced by the Lanczos algorithm in 4.3 we can manipulate our optimality condition (68) as follows

$$\begin{aligned} \mathbf{x}_k &= \arg \min_{\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k} \|\mathbf{Ax} - \mathbf{b}\|_2 \\ \iff \mathbf{y}_k &= \arg \min_{\mathbf{y} \in \mathbb{R}^k} \|\mathbf{A}(\mathbf{x}_0 + \mathbf{V}_k \mathbf{y}) - \mathbf{b}\|_2 \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^k} \|-(\mathbf{b} - \mathbf{Ax}_0) + \mathbf{AV}_k \mathbf{y}\|_2 \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^k} \|\mathbf{V}_k \mathbf{T}_{k+1,k} \mathbf{y} - \mathbf{r}_0\|_2 \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^k} \|\mathbf{T}_{k+1,k} \mathbf{y} - \mathbf{V}_k^\top \mathbf{r}_0\|_2 \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^k} \|\mathbf{T}_{k+1,k} \mathbf{y} - \beta_0 \mathbf{e}_1\|_2 \end{aligned}$$

[Gre97, page 43]. Under the general project method Algorithm 11, this affords the following high-level description of the MINRES algorithm [Tre97, page 268].

Algorithm 13: High-Level MINRES

output: An approximation of x^*, x_k .

```

1 for  $k = 1, \dots$  until convergence do
2   Lanczos-Step ( $\mathbf{A}, \mathbf{v}_k, \mathbf{v}_{k-1}, \beta_k$ )  $\rightarrow \alpha_k, \beta_{k+1}, \mathbf{v}_{k+1}$ 
3   Find  $\mathbf{y}$  to minimize  $\|\mathbf{T}_{k+1,k}\mathbf{y} - \beta_0 \mathbf{e}_1\|_2$ 
4    $\mathbf{x}_k = \mathbf{V}_k \mathbf{y}$ 
5 end
6 return  $\mathbf{x}_k$ 

```

At each step solving \mathbf{y} is a matter of solving a $(n+1) \times n$ least squares problem with Hessenberg structure. This can be done by performing a QR factorisation on the successive $\mathbf{T}_{k+1,k}$ matrices. Whats more, using a clever bit of thinking, we can actually compute the QR factorisation of $\mathbf{T}_{k+1,k}$ from the QR factorisation of $\mathbf{T}_{k,k-1}$ using an inexpensive $\mathcal{O}(n)$ Householder reflection [Tre97, page 268]. Computing the QR factorisation of $\mathbf{T}_{k,k-1}$ yields

$$(69) \quad \mathbf{Q}_k \mathbf{T}_{k,k+1} = \begin{bmatrix} \mathbf{R}_k \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \gamma_1^{(1)} & \delta_2^{(1)} & \varepsilon_3^{(1)} & & \\ & \gamma_1^{(2)} & \delta_3^{(2)} & \varepsilon_4^{(1)} & \\ & & \ddots & \ddots & \ddots \\ & & & \ddots & \ddots & \varepsilon_k^{(1)} \\ & & & & \ddots & \delta_k^{(2)} \\ & & & & & \gamma_k^{(2)} \\ & & & & & 0 \end{bmatrix}$$

where $\mathbf{Q}_k = \prod_{i=1}^k \mathbf{Q}_{i,i+1}$ is the product of Householder reflections designed to annihilate the β_i s in the subdiagonal of $\mathbf{T}_{k,k+1}$ [Cho07, page 25] where each $\mathbf{Q}_{i,i+1}$ is defined as

$$\mathbf{Q}_{i,i+1} \triangleq \begin{bmatrix} \mathbb{1}_{(i-1) \times (i-1)} & & & \\ & c_i & s_i & & \\ & s_i & -c_i & & \\ & & & \mathbb{1}_{(k-i) \times (k-i)} & \end{bmatrix}$$

[Cho07, page 22]. As mentioned for each i , $\mathbf{Q}_{i,i+1}$ is orthogonal, symmetric and constructed to annihilate β_{k+1} that is the bottom right element of $\mathbf{T}_{k,k+1}$. To see this better an alternative way of writing $\mathbf{Q}_k \mathbf{T}_{k,k+1}$ is $\mathbf{Q}_{k,k+1} \cdot \mathbf{Q}_{2,3} \mathbf{Q}_{1,2} \mathbf{T}_{k,k+1}$. Notice however that $\mathbf{Q}_{k,k+1}$ is the only rotation matrix that will involve β_{k+1} in its matrix multiplication with $\mathbf{T}_{k,k+1}$. To study the influence of $\mathbf{Q}_{k,k+1}$ in a more compact way we only need to consider the matrix product

$$(70) \quad \begin{bmatrix} c_k & s_k \\ s_k & -c_k \end{bmatrix} \begin{bmatrix} \gamma_k^{(1)} & \delta_{k+1}^{(1)} & 0 \\ \beta_{k+1} & \alpha_{k+1} & \beta_{k+1} \end{bmatrix}.$$

To annihilate β_{k+1} , we find the appropriate choice of c_k and s_k are

$$\rho_k = \sqrt{\gamma_k^{(1)2} + \beta_{k+1}^2}, \quad c_k \triangleq \frac{\gamma_k^{(1)}}{\rho_k}, \quad s_k \triangleq \frac{\beta_{k+1}}{\rho_k}$$

so that (70) becomes

$$\begin{bmatrix} c_k & s_k \\ s_k & -c_k \end{bmatrix} \begin{bmatrix} \gamma_k^{(1)} & \delta_{k+1}^{(1)} & 0 \\ \beta_{k+1} & \alpha_{k+1} & \beta_{k+1} \end{bmatrix} = \begin{bmatrix} \gamma_k^{(2)} & \delta_{k+1}^{(2)} & \varepsilon_{k+2}^{(1)} \\ 0 & \gamma_{k+1}^{(1)} & \delta_{k+2}^{(1)} \end{bmatrix}.$$

Furthermore, if we set

$$\begin{aligned} \mathbf{Q}_k(\beta_0 \mathbf{e}_1) &= \prod_{i=1}^k \mathbf{Q}_{i,i+1}(\beta_0 \mathbf{e}_1) \\ &= \beta_0 \mathbf{Q}_{k,k+1} \cdots \mathbf{Q}_{2,3} \begin{bmatrix} c_1 \\ s_1 \\ \mathbf{0}_{k-1} \end{bmatrix} \\ &= \beta_0 \mathbf{Q}_{k,k+1} \cdots \mathbf{Q}_{3,4} \begin{bmatrix} c_1 \\ s_1 c_1 \\ s_1 s_2 \\ \mathbf{0}_{k-2} \end{bmatrix} \\ &= \beta_0 \mathbf{Q}_{k,k+1} \begin{bmatrix} c_1 \\ s_1 c_1 \\ \vdots \\ s_1 \cdots s_{k-1} \\ 0 \end{bmatrix} \\ &= \mathbf{Q}_{k,k+1} \begin{bmatrix} \mathbf{t}_{k-1} \\ \phi_{k-1} \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{t}_k \\ \phi_{k-1} \end{bmatrix} \end{aligned}$$

where $\mathbf{t}_k = [\tau_1, \tau_2, \dots, \tau_k]^\top$ then our optimality condition becomes

$$\mathbf{y}_k = \arg \min_{\mathbf{y} \in \mathbb{R}^k} \left\| \begin{bmatrix} \mathbf{R}_k \\ \mathbf{0} \end{bmatrix} \mathbf{y} - \begin{bmatrix} \mathbf{t}_k \\ \phi_{k-1} \end{bmatrix} \right\|_2$$

which can be used to formulate subproblems [Cho07, page 25]. From the new optimality condition it becomes obvious that the optimal solution will satisfy $\mathbf{R}_k \mathbf{y}_k = \mathbf{t}_k$. However, instead of solving for \mathbf{y}_k directly, MINRES solves

$$(71) \quad \mathbf{R}_k^\top \mathbf{D}_k^\top = \mathbf{V}_k^\top$$

where $\mathbf{D}_k = [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_k] \triangleq \mathbf{V}_k \mathbf{R}_k^{-1}$. This is done by forward substitution obtaining the last column \mathbf{d}_k of \mathbf{D}_k at iteration k . At the same time \mathbf{x}_k is updated as

$$\begin{aligned} \mathbf{x}_k &= \mathbf{V}_k \mathbf{y}_k = \mathbf{D}_k \mathbf{R}_k \mathbf{y}_k = \mathbf{D}_k \mathbf{t}_k \\ &= [\mathbf{D}_{k-1}, \mathbf{d}_k] \begin{bmatrix} \mathbf{t}_{k-1} \\ \tau_t \end{bmatrix} \\ &= \mathbf{x}_{k-1} + \tau_k \mathbf{d}_k. \end{aligned} \tag{72}$$

The d_j in (72) can be found as

$$\begin{cases} \mathbf{d}_1 = \mathbf{v}_1 / \gamma_1 \\ \mathbf{d}_2 = (\mathbf{v}_2 - \delta_2 \mathbf{d}_1) / \gamma_2^{(2)} \\ \mathbf{d}_j = (\mathbf{v}_j - \delta_j^{(2)} \mathbf{d}_{j-1} - \varepsilon_j \mathbf{d}_{j-2}) / \gamma_j^{(2)}, \quad j = 3, \dots, k \end{cases}$$

from (71) [CHO11, page 4]. The final form of the MINRES is presented in Algorithm 14.

Algorithm 14: MINRES Algorithm

```

input :  $\mathbf{A}$  where  $\mathbf{A} = \mathbf{A}^\top$ ,  $\mathbf{b}$  and an initial guess  $\mathbf{x}_0$ .
output: An approximation of  $\mathbf{x}^*$ ,  $\mathbf{x}_k$ .
1  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ ,  $\tau_0 = \beta_0 = \|\mathbf{r}_0\|$ ,  $\mathbf{v}_1 = \mathbf{r}_0 / \beta_0$ ,  $\delta_1^{(1)} = 0$ ,  $\mathbf{v}_0 = \mathbf{d}_0 = \mathbf{d}_{-1} = \mathbf{0}$ 
2 for  $k = 1, \dots$  until convergence do
3   Lanczos-Step ( $\mathbf{A}, \mathbf{v}_k, \mathbf{v}_{k-1}, \beta_k$ )  $\rightarrow \alpha_k, \beta_{k+1}, \mathbf{v}_{k+1}$ 
4   
$$\begin{bmatrix} \delta_k^{(2)} & \varepsilon_{k+1}^{(1)} \\ \gamma_k^{(1)} & \delta_{k+1}^{(1)} \end{bmatrix} = \begin{bmatrix} c_{k-1} & s_{k-1} \\ s_{k-1} & -c_{k-1} \end{bmatrix} \begin{bmatrix} \delta_k^{(1)} & 0 \\ \alpha_k & \beta_{k+1} \end{bmatrix}$$

5    $\gamma_k^{(2)} = \sqrt{\gamma_k^{(1)2} + \beta_{k+1}^2}$ ,  $c_k = \frac{\gamma_k^{(1)}}{\rho_k}$ ,  $s_k = \frac{\beta_{k+1}}{\rho_k}$ 
6    $\tau_k = c_k \phi_{k-1}$ ,  $\phi_k = s_k \phi_{k-1}$ 
7    $\mathbf{d}_k = (\mathbf{v}_k - \delta_k^{(2)} \mathbf{d}_{k-1} - \varepsilon_k \mathbf{d}_{k-2}) / \gamma_k^{(2)}$ 
8    $\mathbf{x}_k + \tau_k \mathbf{d}_k$ 
9 end
10 return  $\mathbf{x}_k$ 

```

5. APPLICATIONS AND RESULTS

So far we have discussed a number of different approximation techniques that can be used to speed up some of the bottle necks within our Gaussian Process regression and classification algorithms. In this chapter, we shall see how they can be used within the naive implementations from Section 1 to provide faster processing, at the expense of a small amount of accuracy. We shall also go through an extensive treatment of how the various methods reviewed in previous chapters were implemented and how experiments were set up so that comparisons between them can be made and we can determine the most successful model.

5.1. Gaussian Processes Prediction Reviewed. In Chapter 1, a naive implementation for Gaussian Process prediction was presented in Algorithm 1. While this does provide a simple and convenient way to produce predictions for a regression task, it is not all smooth sailing after this. Unfortunately, there are a number of problems with this algorithm in terms of scalability where various bottlenecks are present on lines 1, 2 and 4.

Algorithm 1: Unoptimized GPR

```

input : Observations  $\mathbf{X}, \mathbf{y}$  and a test input  $\mathbf{x}_*$ .
output: A prediction  $\bar{f}_*$  with its corresponding variance  $\mathbb{V}[f_*]$ .
1  $\mathbf{L} = \text{cholesky}(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbf{1}_{n \times n})$ 
2  $\boldsymbol{\alpha} = \text{lin-solve}(\mathbf{L}^\top, \text{lin-solve}(\mathbf{L}, \mathbf{y}))$ 
3  $\bar{f}_* = \mathbf{K}_{\mathbf{x}_*\mathbf{X}} \boldsymbol{\alpha}$ 
4  $\mathbf{v} = \text{lin-solve}(\mathbf{L}, \mathbf{K}_{\mathbf{x}_*\mathbf{X}})$ 
5  $\mathbb{V}[f_*] = \mathbf{K}_{\mathbf{x}_*\mathbf{x}_*} - \mathbf{v}^\top \mathbf{v}$ 
6 return  $\bar{f}_*, \mathbb{V}[f_*]$ 

```

To start computing the kernel matrix, $\mathbf{K}_{\mathbf{X}\mathbf{X}}$, on line 1 carries an $\mathcal{O}(n^2)$ runtime since $\mathcal{O}(n^2)$ pairwise kernel evaluations must be made. Moreover, solving linear systems using a Cholesky decomposition (seen on lines 1, 2 and 4) will incur a runtime of $\mathcal{O}(n^3)$. While this quadratic scaling is not a significant issue for smaller datasets, this will scale very poorly as most modern desktops are unable to feasibly perform training with datasets having more than 10^5 samples [WJMaSKaAGW21, page 2]. The Gaussian Process classifier from Algorithm 2 suffers for the exact same reasons as GPR. To mitigate the computational burden of these bottle necks, we can replace these procedures with their inexact counterparts discussed in Chapters 2, 3 and 4.

To start, computing the kernel matrix can be done using either the Nystrom method or the RFF technique, both of which provide better asymptotic runtimes. Similarly, the Cholesky decomposition and linear solves can be replaced with either CG or MINRES to improve runtime performance. This makes sense when approximating the kernel matrix as the Cholesky decomposition provides an almost-exact solution to solving these linear system, which is rather counterintuitive if approximations are used earlier on within the algorithm.

5.2. Experimental Setup. Each method was implemented in python3 and run using the python3.8.5 interpreter distributed by the official python website. Various scientific computing python libraries, such

as numpy, pandas and scipy, were used to implement methods that enhanced runtime performance. A just-in-time compiler was also employed to improve the performance of the FWHT and Krylov-Subspace methods as fast runtime is critically important for the success of both these methods. Experiments were carried out on a variety of different datasets listed in Table 3. Code pertinent to experimentation can be found at the following GitHub link: <https://github.com/Michae1CC/Optimising-GPs>.

TABLE 3. Descriptions and sources for each of the datasets used in experiments.

Name	Description	d	n	Source
3DSN	The 3D Spatial Network (3DSN) dataset was constructed by adding elevation information to a 2D road network in North Jutland, Denmark (covering a region of $185 \times 135 \text{ km}^2$).	2	2000	UCI
Abalone	Physical measurements of abalones.	7	4177	UCI
magic04	Simulated registration of high energy gamma particles in a ground-based atmospheric Cherenkov gamma telescope.	10	19020	UCI
Wine	Chemical measurements of wine.	11	4898	UCI
Temp	Weather station data from rural Queensland.	1	11324	Qld Gov
Stocks	Daily stock prices spanning 2000 to 2019.	4	4904	tiskw
quadratic	A dataset artificially constructed where samples were drawn from a Gaussian process with a mean function of $(x, y) \mapsto x^2 + y^2$ using a RBF kernel.	2	5000	NA

5.2.1. *Kernel Matrix Approximation Testing*. To test the various kernel matrix approximation techniques, each method was used to construct an estimate of the actual kernel matrix with varying sample sizes. The definition of samples changes depending of which family of approximation technique is considered. For the Nystrom technique, the number of samples refers to s , that is the number of columns sampled. In the case of the RFF technique the samples refer to D , that being the dimension of the constructed feature space. The Nystrom technique was implemented according to [PDAWM05] and RFF according to [Rah08] and [Liu21]. Methods were then made to compute approximations for kernel matrices for each of the above data sets for a fixed σ and sample size. This was repeated 30 times for each dataset. For every method the relative Frobenius error $\|K_{XX} - \widehat{K}_{XX}\|_F / \|K_{XX}\|_F$ and relative infinity error $\|K_{XX} - \widehat{K}_{XX}\|_\infty / \|K_{XX}\|_\infty = \|\widehat{K}_{XX}\|_\infty$ (in the case of the RBF kernel) was recorded, using the notation that \widehat{K}_{XX} represents an approximated kernel matrix. The values of σ employed here were 0.1, 1.0 and 10.0. A σ value of 2.1 was also used to compare with results from [PDAWM05]. Additionally, computation time and memory usage for each kernel approximation was also recorded. It should be noted, however, that the time and memory spent on constructing probabilites for the Nystrom methods are *not* included in the kernel construction time and have instead been recorded separately in Figures 45-50. This is because the probabilites are typically known prior to kernel matrix construction and are simply used to provide an approximation to the kernel matrix instead on needing to save the

entire data kernel matrix for a new prediction. Thus, they act as an efficient means to quickly save and rebuild the data kernel matrix. The errors, time and memory were averaged across the repetitions for each experiment.

To understand how well each kernel matrix approximation method performs in terms of providing predictions, each method was required to train a Gaussian process with $4/5^{ths}$ of each data set and then predict the remaining $1/5^{th}$. Predictions for each data set were repeated 15 times across various sample sizes. The mean square error (MSE) and classification error was captured for regression and classification tasks respectively. To make comparisons between different approximation methods as fair as possible, the matrix $(\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbb{1}_{n \times n})^{-1}$ was directly inverted to eliminate any error that an in-exact linear system solver might introduce. To do this efficiently, since both the Nystrom and RFF methods produce decompositions to their approximations of the form $\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} = \mathbf{U}\mathbf{W}^\dagger\mathbf{U}^\top$ we can make use of the matrix inversion lemma (see Lemma 12) to compute

$$\begin{aligned} (\widehat{\mathbf{K}}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbb{1}_{n \times n})^{-1} &= (\mathbf{U}\mathbf{W}\mathbf{U}^\top + \sigma_n^2 \mathbb{1}_{n \times n})^{-1} \\ &= \sigma_n^{-2} \mathbb{1}_{n \times n} + \sigma_n^{-2} \mathbb{1}_{n \times n} \mathbf{U} (\mathbf{W}^\dagger + \mathbf{U}^\top \sigma_n^{-2} \mathbb{1}_{n \times n} \mathbf{U})^{-1} \mathbf{U}^\top \sigma_n^{-2} \mathbb{1}_{n \times n}. \end{aligned}$$

Only a single value of σ was used to create predictions for each dataset. The value of σ chosen corresponds to the best value of σ out of 0.1, 2.1, 1.0 and 10.0 when an exact GP algorithm was used to provide predictions (a value of 2.1 was used to help compare with the results from [PDAWMW05]). The time and memory of performing the matrix inverse using the above procedure was recorded. Again, errors, time and memory were averaged across the repetitions for each different sample.

5.2.2. Krylov Subspace Methods Approximation Testing. Similar to the kernel matrix approximation setup, to see how well CG and MINRES compare in providing predictions, each method was used to solve the linear system

$$(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbb{1}_{n \times n}) \boldsymbol{\alpha} = \mathbf{y}$$

within the GPR algorithm and

$$(\mathbf{K}_{\mathbf{X}\mathbf{X}} + \mathbf{W}^{-1}) \boldsymbol{\alpha} = \mathbf{K}_{\mathbf{x}_* \mathbf{X}}^\top$$

within the GPC algorithm in terms of $\boldsymbol{\alpha}$. For consistency, predictions for each data set was repeated 15 times across various samples. The mean square error (MSE) and classification error was captured for regression and classification tasks respectively. In an effort to provide a fair comparison between the different linear solvers, the kernel matrix was constructed in a manner that ensured no error was introduced through any other part of the prediction process. As with the kernel matrix approximation prediction set up, only a single value of σ was used to create predictions for each dataset, that being the value of σ that provided the best predictions results using exacts methods. Each method was used to train a Gaussian process with $4/5^{ths}$ of each data set and required to predict the remaining $1/5^{th}$ over all the different datasets and was repeated 15 times for each dataset with a varying number of maximum iterations. The usual metrics such as prediction error, time and memory ususage were recorded for every experiment and averaged across repetitions for each dataset and method.

5.3. Discussion.

5.3.1. Kernel Matrix Approximation. Starting with the performance of the Nystrom methods, assessing Figures 33 to 44, overall the rls method is superior for sampling distributions as it is virtually always among the top three methods for any combination of dataset, k or σ . Interestingly, non-uniform sampling techniques performed better with datasets for which the spectrum of the corresponding kernel matrix was also non-uniform. For example, the spectrum for the kernel matrices of the magic04 and Stocks dataset (see Figures 17 and 18) are relatively uniform and from Figures 37 and 39 we find that the non-uniform sampling techniques give rather poor approximations and generally behave just as poorly, or even worse, than the naive uniform sampling distribution. In contrast, non-uniform methods very much shine in the case of the 3DSN and Wine datasets where some of the better methods could provide almost exact kernel matrix approximations (refer to Figures 33, 34, 43 and 44) after a few thousand samples. A good comparison of these two scenarios is seen in Figure 9.

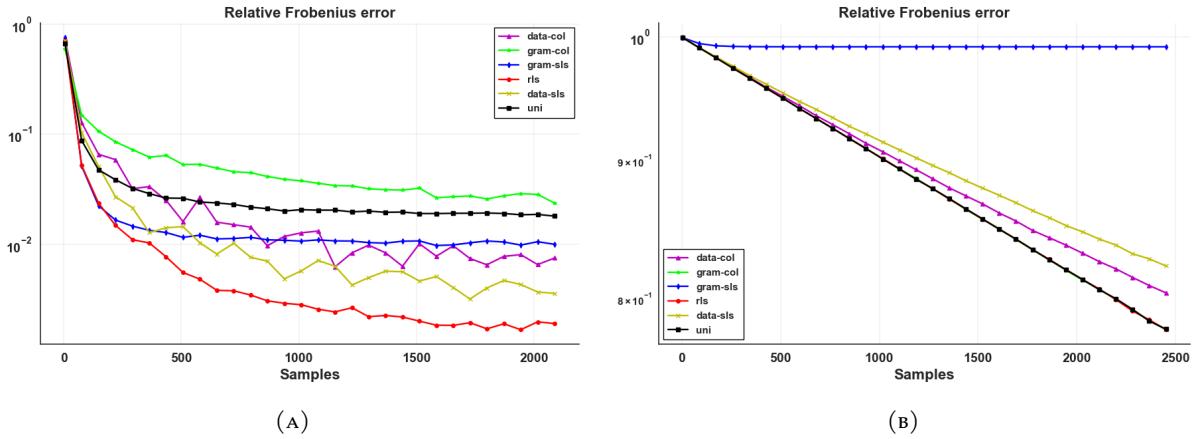


FIGURE 9. Non-uniform sampling distributions perform much better when the corresponding kernel matrix also has a non-uniform spectrum, as seen in the left panel were various sampling distributions are used to construct the Abalone kernel matrix. In contrast, there is not much difference between sampling distributions for kernel matrices with highly uniform spectrums, as depicted in the right panel.

This phenomena is likely due to the fact that the spectrum of a matrix decays faster when it is less uniform, meaning that the matrix can be better expressed as a low-rank approximation. Evidently, the majority of non-uniform sampling distributions are intelligent enough to select columns that are best representative in this low rank approximation, most of the time. What is especially interesting about the Abalone dataset is, despite the rls sampling distribution being almost uniform (compared to other sampling distributions, see Figures 22 and 28), this method seems to perform significantly better than uniform selection, shown in Figure 35. The small amount of variation in the rls sampling distribution presumably gives it an edge over the uniform method in selecting columns that better provides information to the lower rank approximation.

Turning our attention to the RFF methods, we find that the ORF and SORF methods do not really live up to their acclaimed theoretical error bounds. While the ORF and SORF methods do occasionally produce better errors over the standard transformation matrix sampling (see Figures 63–74), it is hard to justify

using either of these methods as neither offer any obvious advantage in terms of time and memory saving (see Figures 75-80 and 81-86). This may come as a surprise given the SORF method, even with a JIT compiler, was used to speed up various parts of the algorithm. This dissonance between the theoretical bounds presented in Table 2 and these "real world" experiments is that generating Gaussian variables is still, in practice, quicker than the asymptotically more efficient construction of SORF's transform matrix. When fed large enough datasets, the SORF method is likely to overtake the "basic" sampling method proposed by Rami and Recht. Unfortunately, due to the memory constraints of the computers used for this project's experimental, using such large data sets could not be conducted to test this hypothesis.

To make a fair comparison in kernel matrix approximation between the RFF and Nystrom methods, the time that each method used to build an approximation was graphed alongside the corresponding relative error. In this setup, we can think about each method as having some sort of time budget and that its task is to produce the best approximation it can within this limit, regardless of how many samples it requires or how big of a feature space it uses. Comparing the methods in this manner the Nystrom family is superior in lowering the Frobenius error in the approximations it produces, while the RFF methods provide approximations with smaller infinity errors (see Figures 87-98). In other words, Nystrom methods produce approximations where entries are, *on average*, closer to the true matrix compared to RFF methods. RFF methods, on the otherhand, produce approximations where entries are *in the worst case* closer to the true matrix compared to the Nystrom methods. This makes sense since much of the theory behind the Nystrom methods is aimed toward lowering Frobenius errors, while RFF methods focuses on lowering infinity errors. An example of this is seen in Figure 10.

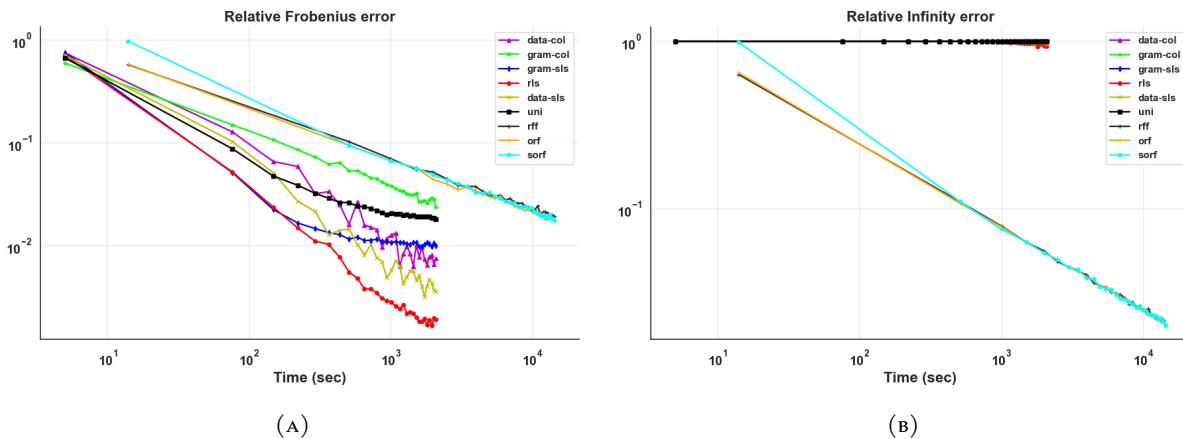


FIGURE 10. This gives a comparison of errors for constructing the kernel matrix associated with the Abalone dataset. Nystrom methods (graphed with primary colors) generally provided better Frobenius errors (left panel) while RFF methods (graphed with composite colors) gave better infinity errors (right panel).

One final point worth mentioning is that RFF methods seemed to require much larger amounts of memory to produce their approximations, sometimes double or even triple the amount of memory the of Nystrom methods used shown in Figure 11 (also see Figures 57-62 and Figures 81-86).

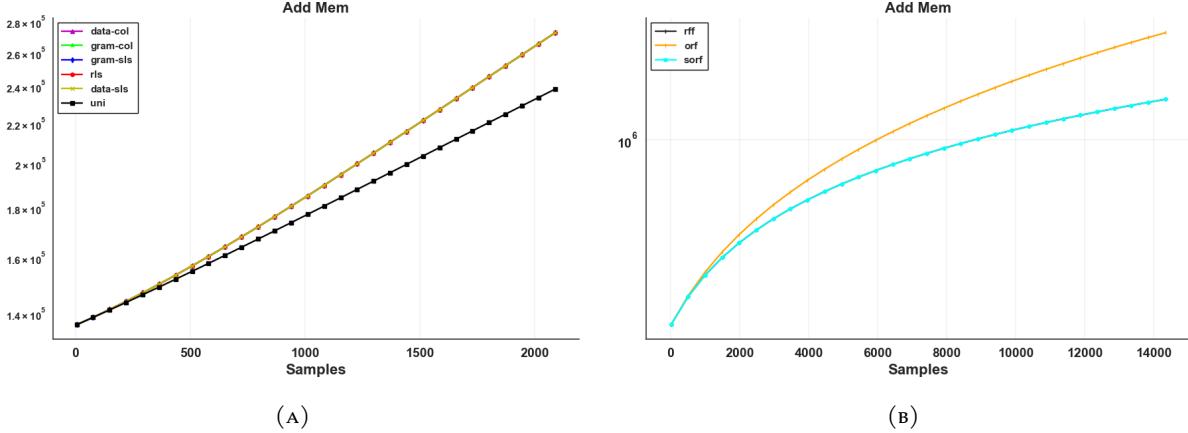


FIGURE 11. A comparison of memory usage for constructing the kernel matrix associated with the Abalone dataset. RFF methods (right panel) typically had much larger memory footprints compared to Nystrom methods (left panel). Keep in mind that the interpretation of "samples" (on the horizontal axis) changes depending on the method.

Next, two methods from the Nystrom class and one method from the RFF class, that we deemed to be the best among their respective families, were used within the Gaussian Process algorithm (replacing exact kernel matrix creation) to form predictions. The rls and data columns methods were chosen to represent the best of the Nystrom family while the Rami and Recht's "basic" RFF method was selected to represent the best of the RFF techniques. Similar to relative error comparisons, the time each method required to complete an approximation was graphed along side the corresponding prediction error. The results, suggest that there was no clear winner. Generally speaking, Rami and Recht's RFF method delivered the best predictions with smaller time budgets were the rls generally caught up and overtook it when a larger time budget was allowed. The data columns method, tends to yield the poorest prediction errors among the selected methods. Thus, the RFF method would probably be the best method to use in practice for two main reasons. First, if one was to select an inexact method for kernel matrix production, it most likely means that there is not a large time budget. In this case RFF's ability to produce the best approximation within shorter timeframes makes it the best option. Second, the time taken to construct probability distributions used in the Nystrom methods are not included in Figures 87-98 (and can instead be found in Figures 45-50) making it harder to justify using the rls method if the probability distribution is not already available since a much larger time budget would be required to deliver the same level of prediction accuracy.

5.3.2. Krylov Subspace Methods. The CG and MINRES methods where also used within the Gaussian Process algorithm to see which one of the two would provide better predictions. Historically CG is preferred over MINRES for Gaussian process prediction since it is typically a superior linear system solver when dealing with PSD matrices. There are, however, a number of circumstances where MINRES can actually outperform CG. Take for example the work of Pleiss *et al.* [Ple20] which leverages a variant of MINRES to assist in the computation kernel matrix square roots. We were also wondering if MINRES is better suited as a linear solver within the context of Gaussian processes. In our experimentation, the number of

iterations that each method used to form an approximation was graphed along side the corresponding prediction error produced. Differences in execution time per iteration between the two methods were nominal. Looking at the results in Figures 99-104, it appears that the MINRES method performs just as well and, occasionally, even better than CG. This was perhaps most apparent with the prediction errors produced from the abalone and 3DSN datasets, see Figure 12. These differences became more noticeable

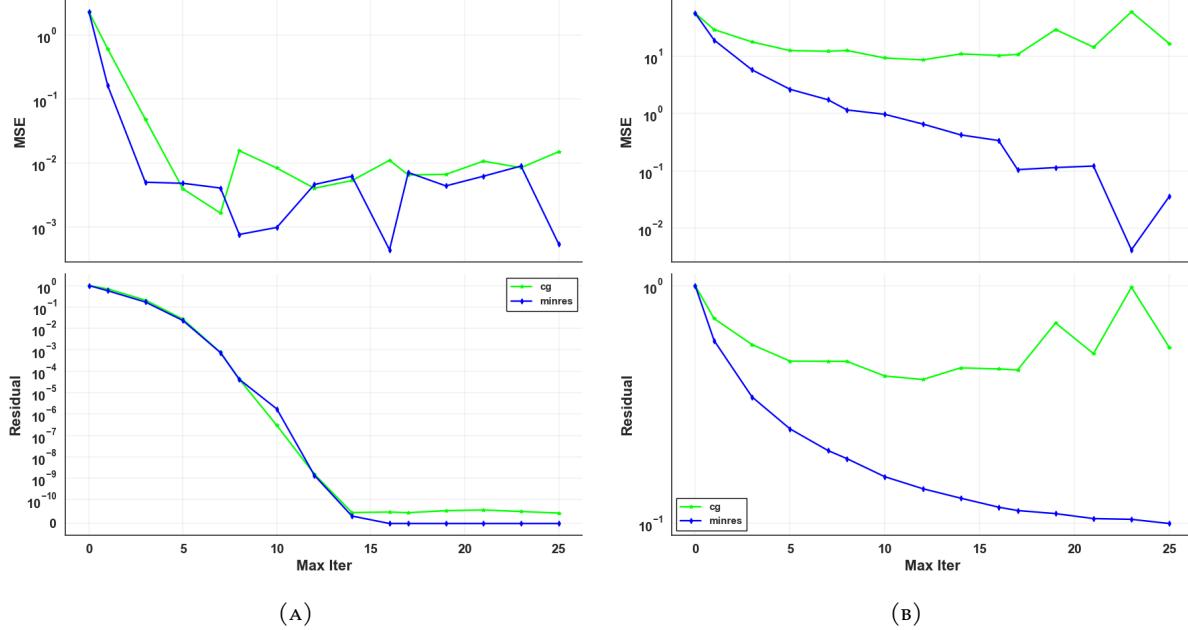


FIGURE 12. Prediction errors for the 3DSN (left panel) and Abalone (right panel) datasets make it most apparent that MINRES is the superior solver for GPR.

when each of these linear solvers where paired with the RFF method (instead of using an exact kernel matrix) to provide predictions shown in Figure 13 (also see Figures 105-110). In particular, the MINRES method affords better results even on the quadratic dataset. The quadratic dataset is an artificial dataset constructed using a Gaussian process with a quadratic mean function $x^2 + y^2$ and some amount of variance. Its main purpose was to ascertain if the applied method was prone to overfitting. The fact that, overall, MINRES yielded lower errors indicates that it is also more robust against overfitting compared to CG. To understand why MINRES may produce smaller errors for regression tasks, recall that the mean square error of ρ is computed as

$$\|K_{X^*X}\rho - y^*\|_2^2$$

where ρ is our best estimate for a solution to the linear system $[K_{XX} + \sigma_n^2 \mathbb{1}_{n \times n}] \rho = y$. The MINRES method seeks to solve this linear system by solving an empirical version of the MSE up to a noise value of sigma, that is, ρ is chosen so that

$$\|[K_{XX} + \sigma_n^2 \mathbb{1}_{n \times n}] \rho - y\|_2^2$$

is minimized. In contrast, CG seeks vectors that successively minimizes the power norm of ρ ; something which we are not so interested in when comparing and contrasting test accuracies across different

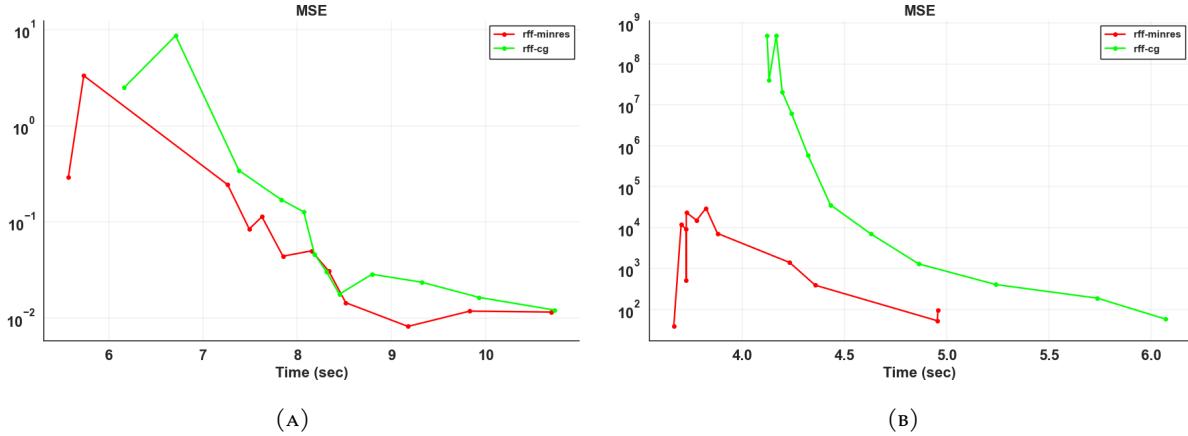


FIGURE 13. The difference in predictive performance between MINRES and CG is even more noticeable when the two methods are paired with a RFF kernel estimate. These differences are best seen when comparing prediction errors provided for the 3DSN (left panel) and Abalone (right panel) datasets.

methods. While both MINRES and CG will eventually yield the same solution, the success of MINRES is likely due to its ability find early estimates of ρ that more directly lowers the empirical MSE. For small values of σ_n^2 , will also correspond in a lower true MSE. If the above analysis is correct, then MINRES is really only advantageous when the value of sigma is small so that less weight is placed on the model's prior. Indeed, when the variance was artificially increased for both the abalone and quadratic dataset prediction, the errors between the CG and MINRES where much closer.

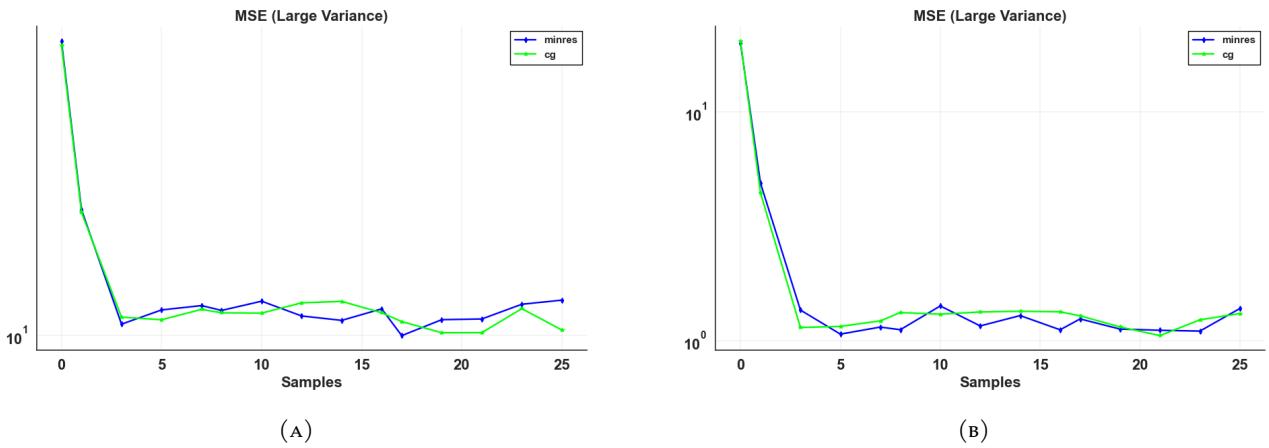


FIGURE 14. Large variances place greater emphasis on the regression model's prior. This makes the true and empirical MSE more dissimilar causing MINRES to lose its edge over CG. The effect is best seen on the abalone (right panel) and the quadratic dataset (left panel).

5.4. Conclusion. The aim of this thesis was to explore how various approximation techniques could be used within the Gaussian Process algorithm to speed up prediction time. Chapter 1 focuses on what kernel reproducing Hilbert spaces and how the kernel trick can be employed as a clever means to afford non-linear estimations while avoiding non-linear solvers. Next, Gaussian processes were defined and shown how they could be used as an instrument for machine learning to facilitate classification and regression tasks. However, when implemented naively, Gaussian process scale very poorly with the number of samples on account of two predominate bottlenecks within the prediction process. The first bottleneck arises from having to produce the kernel matrix and the second from solving linears involving said kernel matrix.

The remainder of the thesis focuses on popular estimation techniques to accelerate the prediction process. Chapter 2 looked at the Nystrom method to speed up kernel matrix approximation. Specifically, a common sampling technique variant was studied in which various sampling distributions where used to construct a sketching matrix. We also devised novel sampling distributions that depended only on the data matrix and did not require any information about the kernel matrix. In Chapter 3, another (and more recent) approach to kernel matrix approximation was reviewed, that being the Random Fourier Feature (RFF) method. In this chapter various transformation matrices were considered for the construction of a Gaussian random matrix. Finally, in Chapter 4, a number of Krylov subspace methods were surveyed to seek alternative ways of solving linear systems within Gaussian process predictions. The particular Krylov subspace methods studied were the conjugate gradient (CG) and minimum residual (MINRES) algorithms.

Chapter 5, summarized our extensive experimental analysis to determine which methods delivered the best results. For the kernel estimation techniques, each method was required to create a its own approximation for the kernel matrices of various datasets where the relative infinity and Frobenius error was captured to measure the quality of approximations. Moreover, selected kernel methods were used within the Gaussian Process algorithm to provided predictions for the same datasets. Similarly, Krylov subspace methods were used in place of the Cholesky decomposition in Algorithm 1 to afford a comparative prediction.

When looking at a kernel matrix construct in isloation, 5 different Nystrom column sampling distributions where compared with 3 RFF transformation matrix construction methods. For Nystrom methods, it was found that some of the more sophisticated sampling distributions provided superior results for kernel matrices with non-uniform spectrums. Surprisingly, for RFF methods, the use of sophisticated techniques used to construct the transformation did not provide much of an advantage over basic methods. In terms of Gaussian process predictions, RFF methods delivered better errors with small time budgets over Nystrom methods.

When comparing CG and MINRES, we found that MINRES was consistently more competitive, if not better, than CG, even when paired with RFF for regression tasks. This may seem surprising given the ubiquitous preference of CG over any other linear solver in many Gaussian process software libraries. We believe MINRES's superior performance is due to its ability to more directly lower the Euclidean distance between the predictions and the true outputs.

Looking forward, some of the findings discovered in this thesis have, to our knowledge, not been reported elsewhere. As a result, we intend to publish these findings in a reputable journal given their obvious appeal to the wider scientific computing community. In terms of research, many of the datasets used in this thesis were far too small to observe the asymptotic benefits of using approximation methods over most of the exact methods. It would be interesting to determine how well these techniques scale on very large datasets. In particular, our discoveries are likely to benefit the agricultural sector as more data is collected over the forthcoming years to perform crop analysis on. Another, direction future research could be taken is the application of the approximation techniques applied to multi-output or multi-task Gaussian process models. In many machine learning scenarios we may want to predict multiple outputs using the same set of inputs. As an example, remote sensing researchers attempt to predict the intensity of multiple bands of light reflecting off farm land to give an indication of crop growth. Different versions of the Gaussian processes algorithm exist to predict multiple outputs simultaneously [Bon07]; however, they also suffer from the same bottlenecks as single output Gaussian processes. It would be intriguing to learn whether the approximation techniques studied in this thesis could potentially improve the prediction time of multi-task Gaussian Processes.

REFERENCES

- [Ras06] Carl Edward and Williams Rasmussen Christopher K. I, *Gaussian processes for machine learning / Carl Edward Rasmussen, Christopher K.I. Williams.*, Adaptive computation and machine learning, MIT Press, Cambridge, Mass., 2006 (eng).
- [HHF73] H. Howard Frisinger, *Aristotle's legacy in meteorology*, Bulletin of the American Meteorological Society **54** (1973), no. 3, 198–204.
- [Yul27] G. Udny Yule, *On a Method of Investigating Periodicities in Disturbed Series, with Special Reference to Wolfer's Sunspot Numbers*, Philosophical transactions of the Royal Society of London. Series A, Containing papers of a mathematical or physical character **226** (1927), no. 636-646, 267–298 (eng).
- [Box08] George E. P. and Jenkins Box Gwilym M and Reinsel, *Time series analysis : forecasting and control / George E.P. Box, Gwilym M. Jenkins, Gregory C. Reinsel.*, 4th ed., Wiley series in probability and statistics, John Wiley, Hoboken, N.J., 2008 (eng).
- [VdW19] Mark Van der Wilk, *Sparse Gaussian process approximations and applications*, University of Cambridge, 2019.
- [Cao18] Yanshuai Cao, *Scaling Gaussian Processes*, University of Toronto (Canada), 2018.
- [SD22] Matías and Estévez Salinero-Delgado José and Pipia, *Monitoring Cropland Phenology on Google Earth Engine Using Gaussian Process Regression*, Remote Sensing **14** (2022), no. 1, DOI 10.3390/rs14010146.
- [Pot13] Andries and Lawson Potgieter Kenton and Huete, *Determining crop acreage estimates for specific winter crops using shape attributes from sequential MODIS imagery*, International Journal of Applied Earth Observation and Geoinformation **23** (2013), DOI 10.1016/j.jag.2012.09.009.
- [NdF13] Nando de Freitas, *University of British Columbia CPSC 540, Lecture notes in Machine Learning*, University of British Columbia, 2013.
- [Mur12] Kevin P. Murphy, *Machine learning : a probabilistic perspective / Kevin P. Murphy.*, Adaptive computation and machine learning, MIT Press, Cambridge, MA, 2012 (eng).
- [Ber96] Z.G. Sheftel Berezansky G.F, *Functional analysis. Volume 1 / Y.M. Berezansky, Z.G. Sheftel, G.F. Us ; translated from the Russian by Peter V. Malyshov.*, 1st ed. 1996., Operator

- Theory: Advances and Applications, 85, Basel ; Boston ; Berlin : BirkhaIuser Verlag, Basel ; Boston ; Berlin, 1996 (eng).
- [Tre97] Lloyd N. (Lloyd Nicholas) and Bau Trefethen David, *Numerical linear algebra / Lloyd N. Trefethen, David Bau.*, SIAM Society for Industrial and Applied Mathematics, Philadelphia, 1997 (eng).
- [Dem97] James W Demmel, *Applied numerical linear algebra / James W. Demmel.*, Society for Industrial and Applied Mathematics, Philadelphia, Pa., 1997 (eng).
- [Ste08] Ingo and Christmann Steinwart Andreas, *Support Vector Machines*, 1st ed. 2008., Information Science and Statistics, Springer New York, New York, NY, 2008 (eng).
- [Ber03] Alain and Thomas-Agnan Berlinet Christine, *Reproducing Kernel Hilbert Spaces in Probability and Statistics*, Springer, SpringerLink (Online service), Boston, MA, 2003 (eng).
- [Ste99] Michael L Stein, *Interpolation of Spatial Data Some Theory for Kriging / by Michael L. Stein.*, 1st ed. 1999., Springer Series in Statistics, Springer New York : Imprint: Springer, New York, NY, 1999 (eng).
- [Bos92] Bernhard and Guyon Boser Isabelle and Vapnik, *A training algorithm for optimal margin classifiers*, Proceedings of the fifth annual workshop on computational learning theory, 1992, pp. 144–152 (eng).
- [Cor95] Corinna Cortes, *Support-Vector Networks*, Machine learning **20** (1995), no. 3, 273 (eng).
- [Kro14] Dirk P and C.C. Chan Kroese Joshua, *Statistical Modeling and Computation by Dirk P. Kroese, Joshua C.C. Chan.*, 1st ed. 2014., Springer New York : Imprint: Springer, New York, NY, 2014 (eng).
- [Fle00] R Fletcher, *Practical Methods of Optimization*, John Wiley and Sons, Incorporated, New York, 2000 (eng).
- [Bis06] Christopher M Bishop, *Pattern recognition and machine learning / Christopher M. Bishop.*, Information science and statistics, Springer, New York, 2006 (eng).
- [Spi90] David J and Lauritzen Spiegelhalter Steffen L, *Sequential updating of conditional probabilities on directed graphical structures*, Networks **20** (1990), no. 5, 579–605.
- [WJMaSKaAGW21] Wesley J. Maddox and Sanyam Kapoor and Andrew Gordon Wilson, *When are Iterative Gaussian Processes Reliably Accurate?*, CoRR **abs/2112.15246** (2021), available at [2112.15246](https://arxiv.org/abs/2112.15246).

- [MWM11] Michael W. Mahoney, *Randomized algorithms for matrices and data*, CoRR **abs/1104.5557** (2011).
- [Hal11] Nathan and Martinsson Halko Per-Gunnar and Tropp, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM review **53** (2011), no. 2, 217–288.
- [PGMaJT21] Per-Gunnar Martinsson and Joel Tropp, *Randomized Numerical Linear Algebra: Foundations and Algorithms*, arXiv, 2021.
- [FR20] Fred Roosta, *University of Queensland MATH3204, Lecture notes in Numerical Linear Algebra and Optimisation*, University of Queensland, 2020.
- [Dri06] Petros and Kannan Drineas Ravi and Mahoney, *Fast Monte Carlo Algorithms for Matrices I: Approximating Matrix Multiplication*, SIAM Journal on Computing **36** (2006), no. 1, 132-157, DOI 10.1137/S0097539704442684, available at <https://doi.org/10.1137/S0097539704442684>.
- [PDaMWM17] Petros Drineas and Michael W. Mahoney, *Lectures on Randomized Numerical Linear Algebra*, arXiv, 2017.
- [PDaMWM05] Petros Drineas and Michael W. Mahoney, *On the Nystrom Method for Approximating a Gram Matrix for Improved Kernel-Based Learning*, Journal of Machine Learning Research **6** (2005), no. 72, 2153-2175.
- [AGaMWM13] Alex Gittens and Michael W. Mahoney, *Revisiting the Nystrom Method for Improved Large-Scale Machine Learning*, CoRR **abs/1303.1849** (2013), available at [1303.1849](https://doi.org/10.4236/jmlr.v14i1.3430).
- [CMaCM17] Cameron Musco and Christopher Musco, *Recursive Sampling for the Nystrom Method*, arXiv, 2017.
- [PDe11] Petros Drineas etal., *Fast approximation of matrix coherence and statistical leverage*, CoRR **abs/1109.3843** (2011).
- [MBCaCMaCM15] Michael B. Cohen and Cameron Musco and Christopher Musco, *Ridge Leverage Scores for Low-Rank Approximation*, CoRR **abs/1511.07263** (2015).
- [Kum09] Sanjiv and Mohri Kumar Mehryar and Talwalkar, *Sampling techniques for the nystrom method*, Artificial intelligence and statistics, 2009, pp. 304–311.
- [Hoa78] David C and Welsch Hoaglin Roy E, *The Hat Matrix in Regression and ANOVA*, The American statistician **32** (1978), no. 1, 17–22 (eng).

- [Ala15] Ahmed and Mahoney Alaoui Michael W, *Fast Randomized Kernel Ridge Regression with Statistical Guarantees*, Advances in Neural Information Processing Systems, 2015.
- [Kar10] Noureddine El Karoui, *THE SPECTRUM OF KERNEL RANDOM MATRICES*, The Annals of statistics **38** (2010), no. 1, 1–50 (eng).
- [Pre92] William H. (William Henry) Press, *Numerical recipes in C : the art of scientific computing / William H. Press ... [et al.]*, 2nd ed., Cambridge University Press, Cambridge, 1992 (eng).
- [Wan] Guorong and Wei Wang Yimin and Qiao, *Generalized Inverses: Theory and Computations*, Developments in Mathematics, vol. 53, Springer Singapore, Singapore (eng).
- [Gre97] Anne Greenbaum, *Iterative methods for solving linear systems* Anne Greenbaum., Frontiers in applied mathematics ; 17, Society for Industrial and Applied Mathematics SIAM, 3600 Market Street, Floor 6, Philadelphia, PA 19104, Philadelphia, Pa., 1997 (eng).
- [Cho07] Sou-Cheng (Terrya) Choi, *Iterative methods for singular linear equations and least - squares problems*, ProQuest Dissertations Publishing, 2007 (eng).
- [CHO11] Sou-Cheng T and PAIGE CHOI Christopher C and SAUNDERS, *MINRES-QLP: A KRYLOV SUBSPACE METHOD FOR INDEFINITE OR SINGULAR SYMMETRIC SYSTEMS*, SIAM journal on scientific computing **33** (2011), no. 3-4, 1810–1836 (eng).
- [Rah08] Ali and Recht Rahimi Benjamin, *Random Features for Large-Scale Kernel Machines*, Advances in Neural Information Processing Systems, 2008.
- [Hor13] Roger A and Johnson Horn Charles R, *Matrix analysis / Roger A. Horn, Charles R. Johnson.*, 2nd ed., Cambridge University Press, New York, 2013 (eng).
- [Pot21] Andres and Wu Potapczynski Luhuan and Biderman, *Bias-Free Scalable Gaussian Processes via Randomized Truncations* (2021) (eng).
- [Hah33] Hans Hahn, *S. Bochner, Vorlesungen über Fouriersche Integrale: Mathematik und ihre Anwendungen, Bd. 12.) Akad. Verlagsges., Leipzig 1932, VIII. u. 229S. Preis brosch. RM 14,40, geb. RM16*, Monatshefte für Mathematik **40** (1933), no. 1, A27–A27 (ger).
- [Liu21] Fanghui and Huang Liu Xiaolin and Chen, *Random Features for Kernel Approximation: A Survey on Algorithms, Theory, and Beyond*, IEEE transactions on pattern analysis and machine intelligence **PP** (2021) (eng).

- [HAe16] Haim Avron et al, *Quasi-Monte Carlo Feature Maps for Shift-Invariant Kernels*, Journal of Machine Learning Research **17** (2016), no. 120, 1-38.
- [DJSaJS15] Danica J. Sutherland and Jeff Schneider, *On the Error of Random Fourier Features*, 2015.
- [Yu16] Felix X and Suresh Yu Ananda Theertha and Choromanski, *Orthogonal Random Features* (2016) (eng).
- [Bro91] Peter J and Davis Brockwell Richard A, *Time Series: Theory and Methods*, Second Edition., Springer Series in Statistics, Springer New York, SpringerLink (Online service), New York, NY, 1991 (eng).
- [Cho17] Krzysztof and Rowland Choromanski Mark and Weller, *The Unreasonable Effectiveness of Structured Random Orthogonal Embeddings* (2017) (eng).
- [FaA76] Fino and Algazi, *Unified Matrix Treatment of the Fast Walsh-Hadamard Transform*, IEEE transactions on computers **C-25** (1976), no. 11, 1142–1146 (eng).
- [And15] Alexandr and Indyk Andoni Piotr and Laarhoven, *Practical and Optimal LSH for Angular Distance* (2015) (eng).
- [Cho20] Krzysztof and Likhoshevstov Choromanski Valerii and Dohan, *Rethinking Attention with Performers* (2020) (eng).
- [Boj16] Mariusz and Choromanska Bojarski Anna and Choromanski, *Structured adaptive and random spinners for fast machine learning computations* (2016) (eng).
- [Ple20] Geoff and Jankowiak Pleiss Martin and Eriksson, *Fast Matrix Square Roots with Applications to Gaussian Processes and Bayesian Optimization*, arXiv, 2020.
- [Bon07] Edwin V and Chai Bonilla Kian and Williams, *Multi-task Gaussian Process Prediction* (J. Platt and D. Koller and Y. Singer and S. Roweis, ed.), Vol. 20, Curran Associates, Inc., 2007.

APPENDIX A. SUPPLEMENTARY RESULTS

Additional results that may have not been included in the main text for conciseness.

A.1. Gram Matrix Spectral Values.

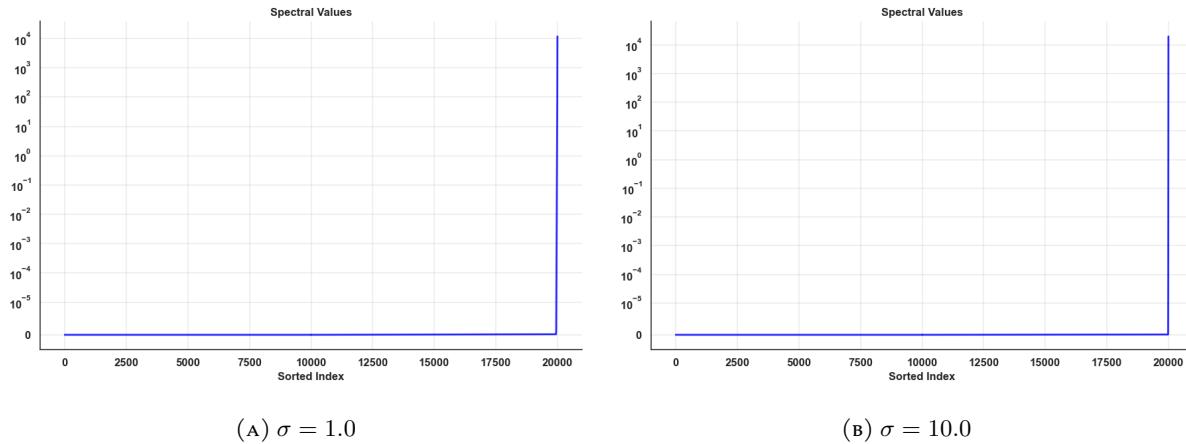


FIGURE 15. Spectral values for 3D-spatial network data.

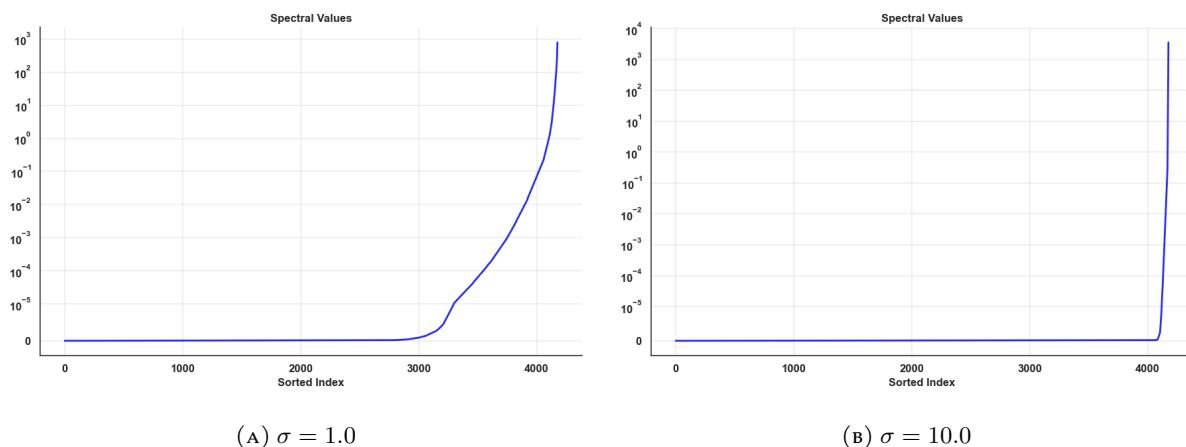


FIGURE 16. Spectral values for abalone data.

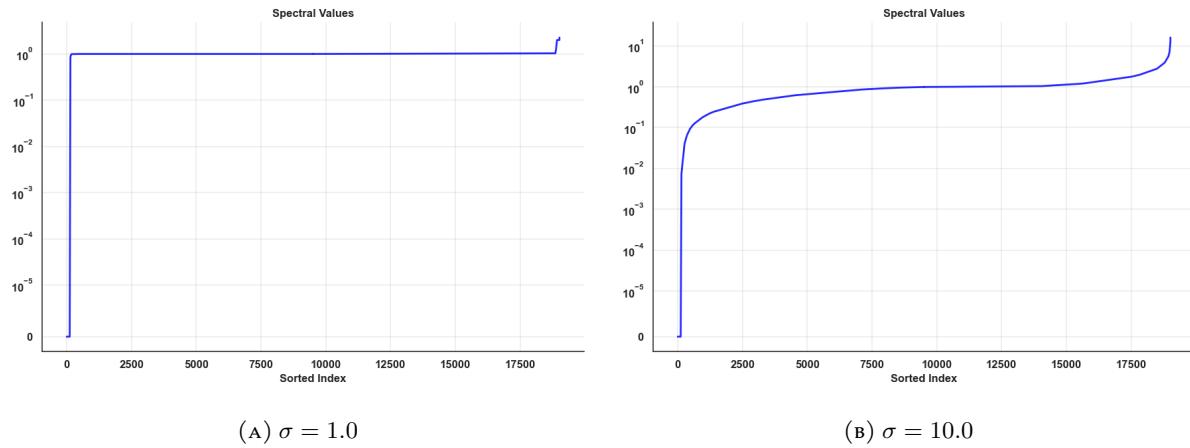


FIGURE 17. Spectral values for magic04 data.

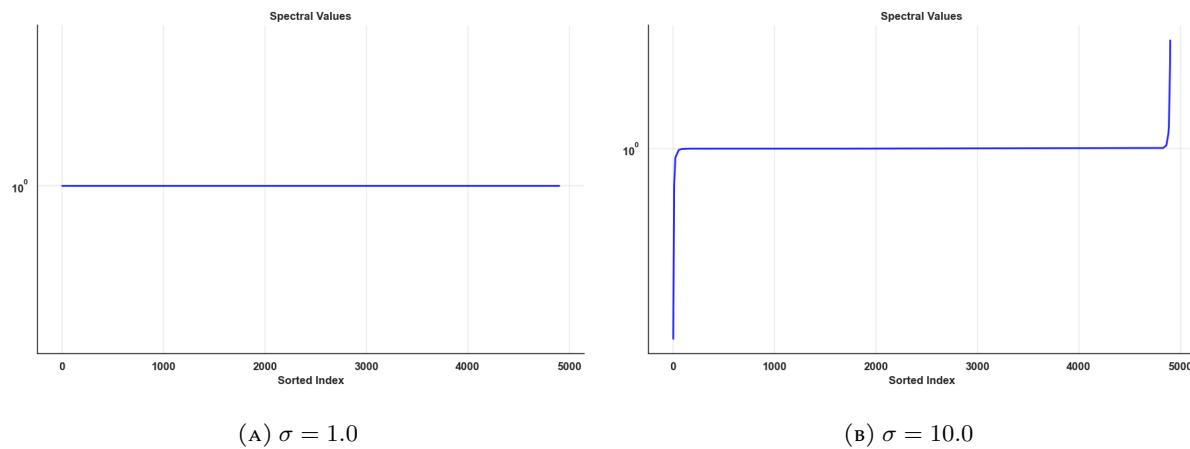


FIGURE 18. Spectral values for stock market data.

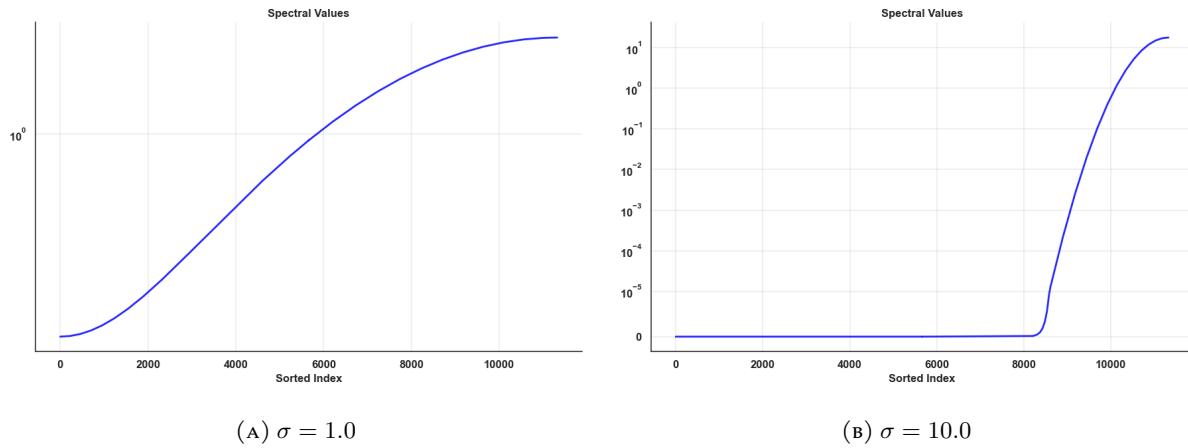


FIGURE 19. Spectral values for temperature data.

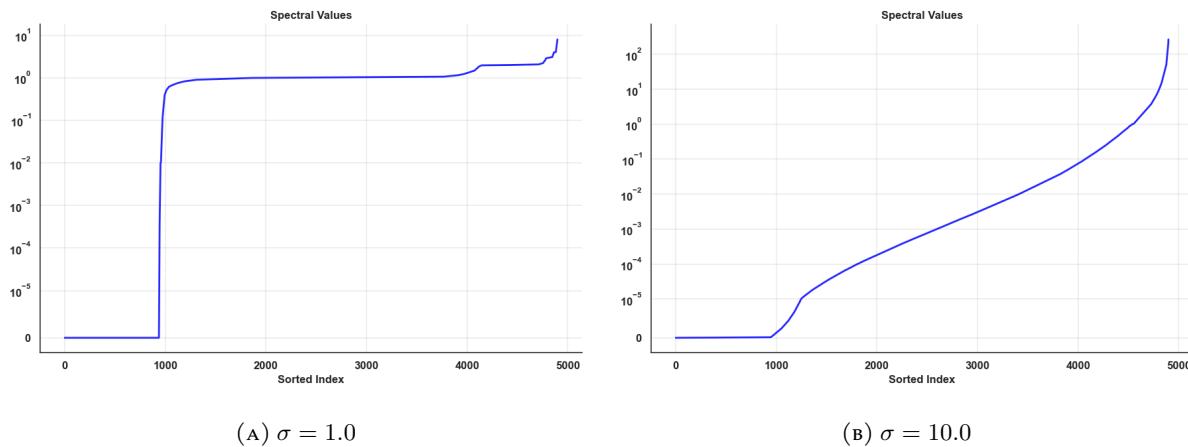


FIGURE 20. Spectral values for wine data.

A.2. Nystrom Scores.

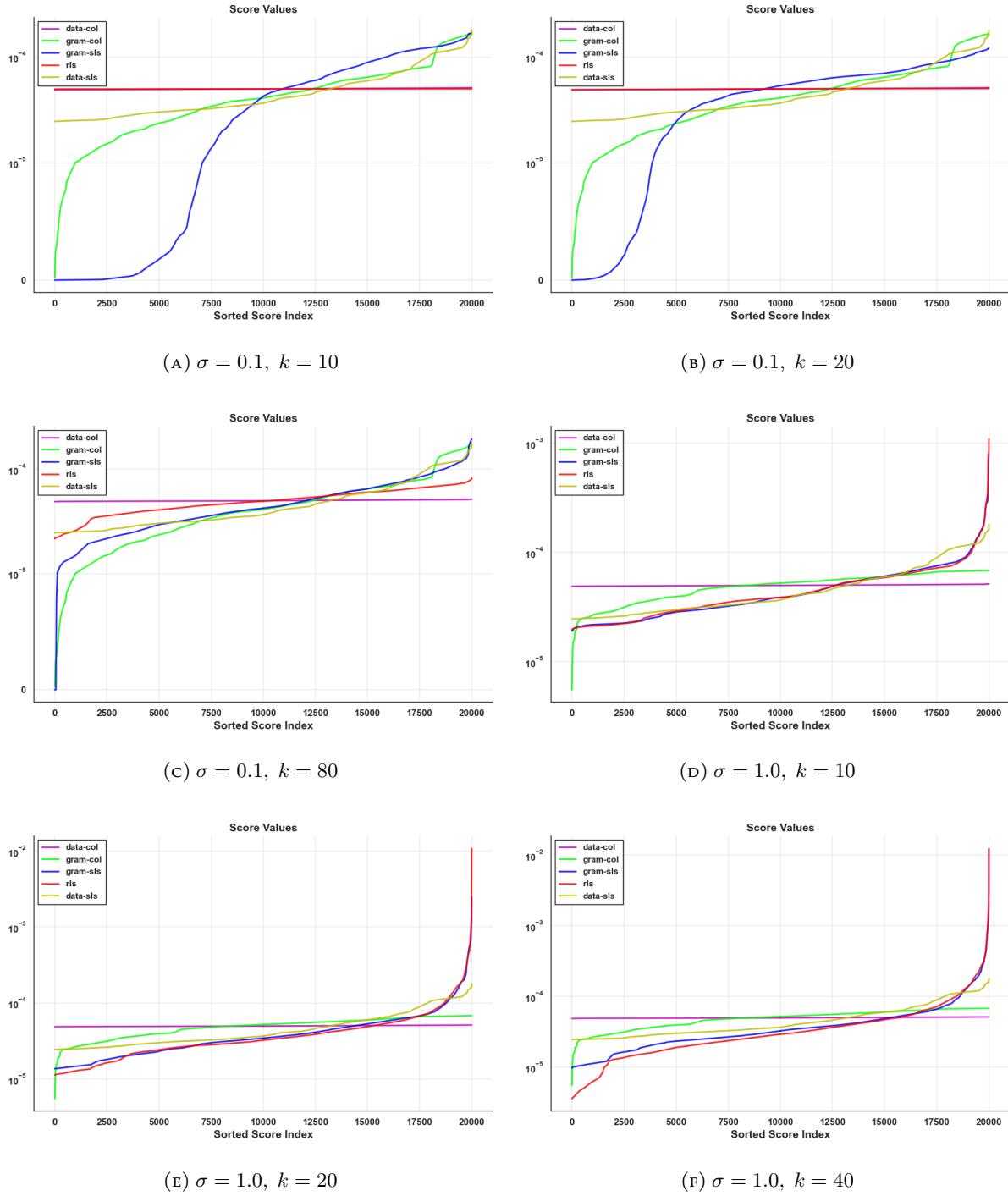


FIGURE 21. Nystrom scores for the 3D-spatial network data set.

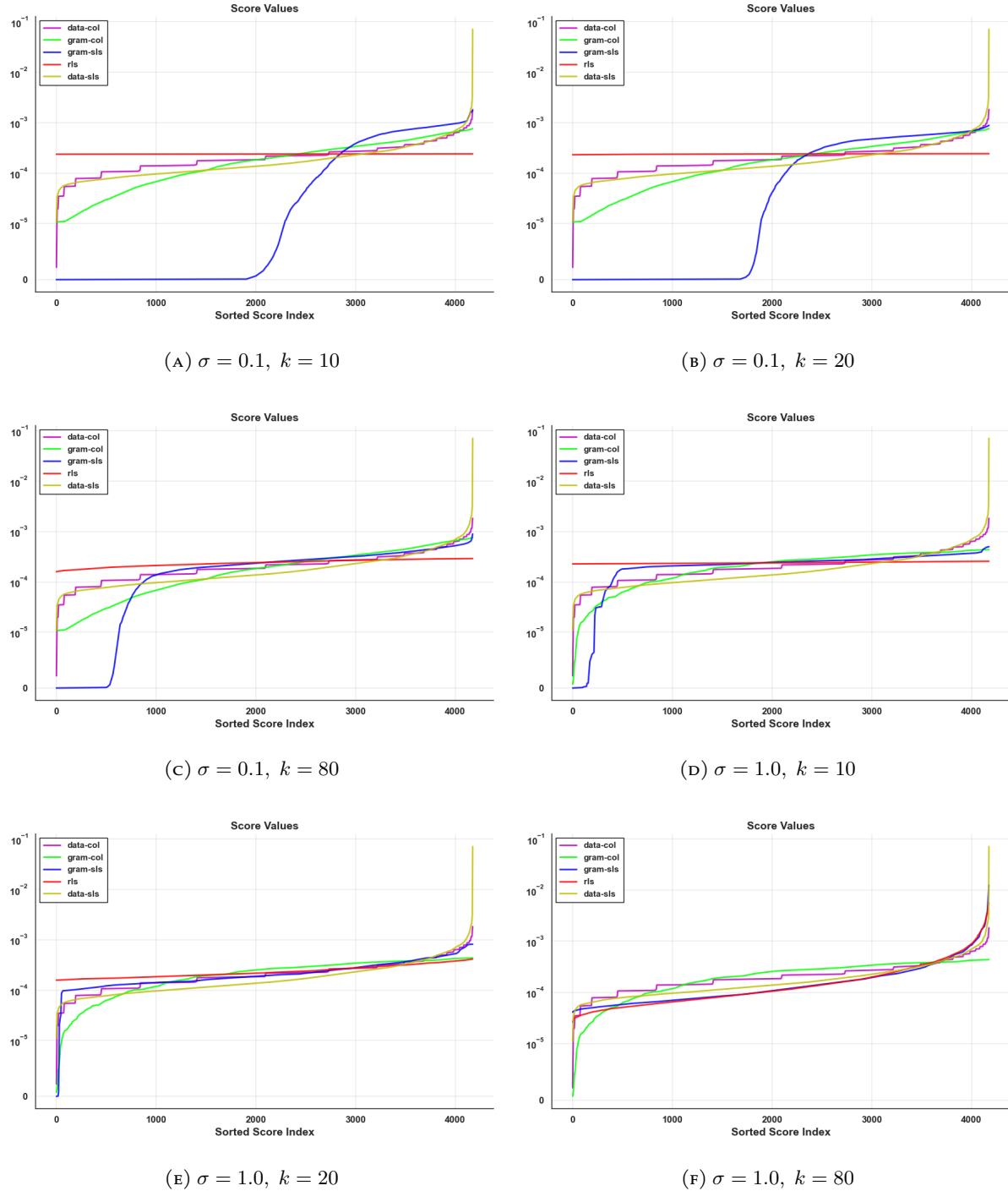


FIGURE 22. Nystrom scores for the Abalone data set.

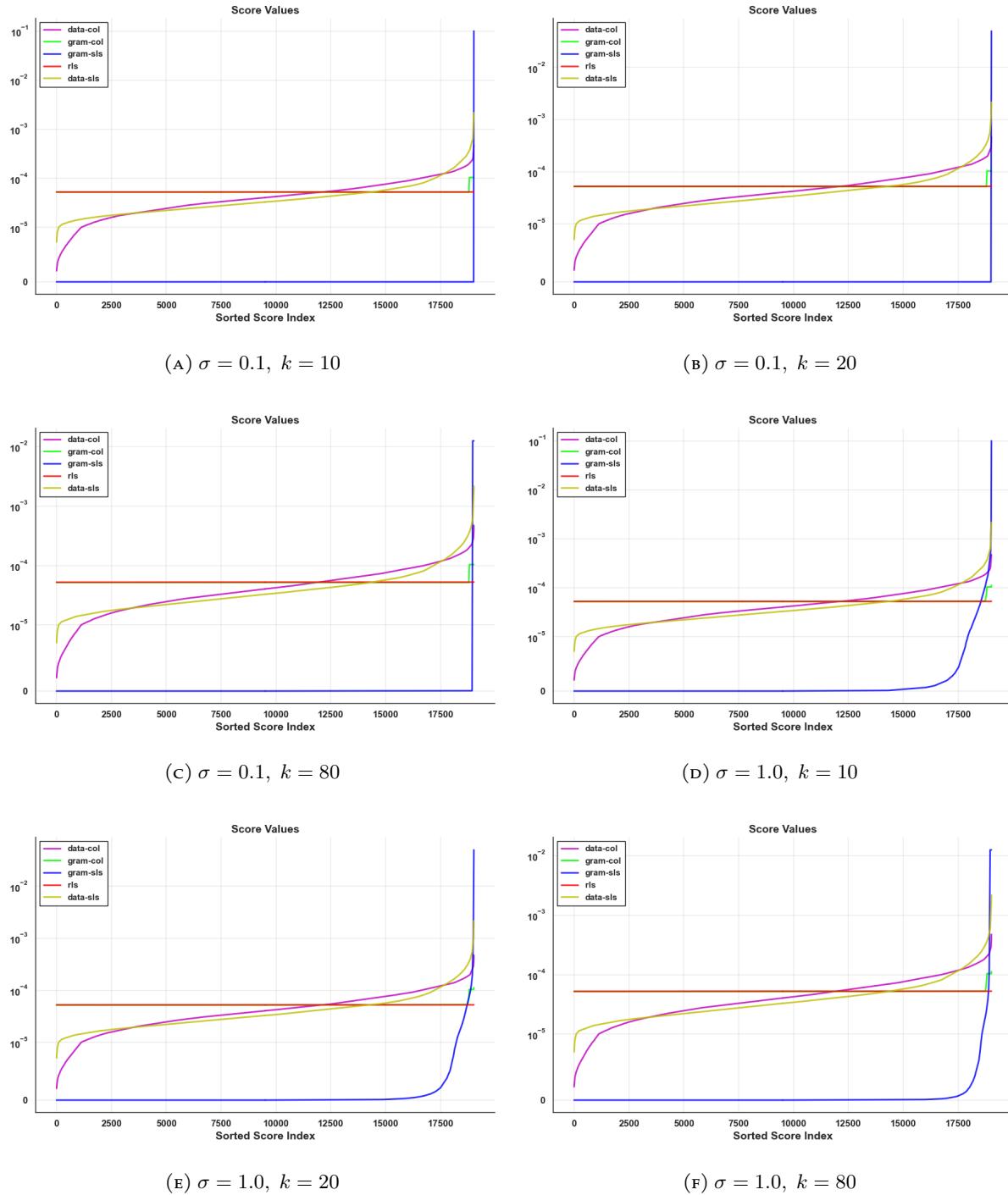


FIGURE 23. Nystrom scores for the Magic data set.

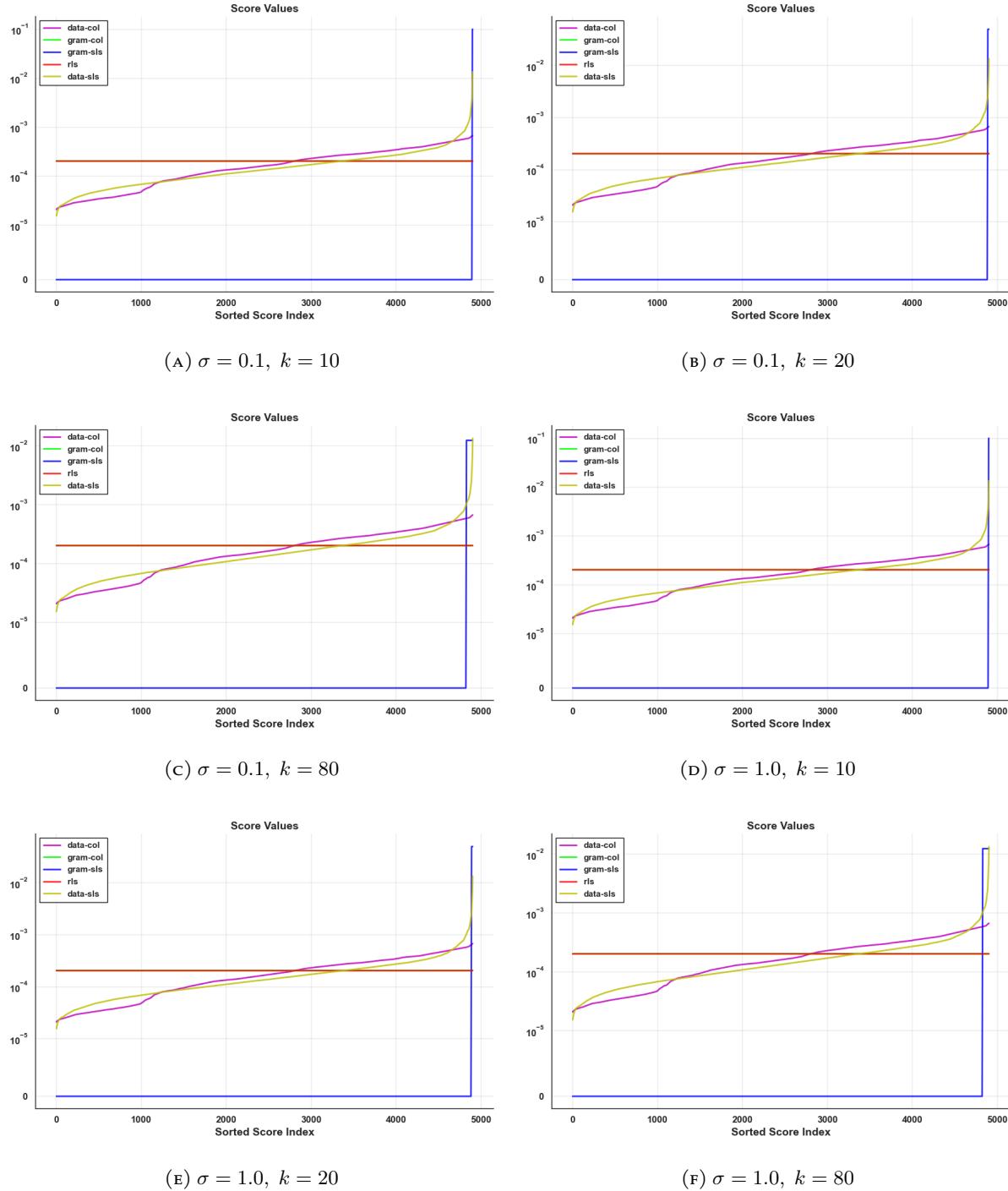


FIGURE 24. Nystrom scores for the stock market data set.

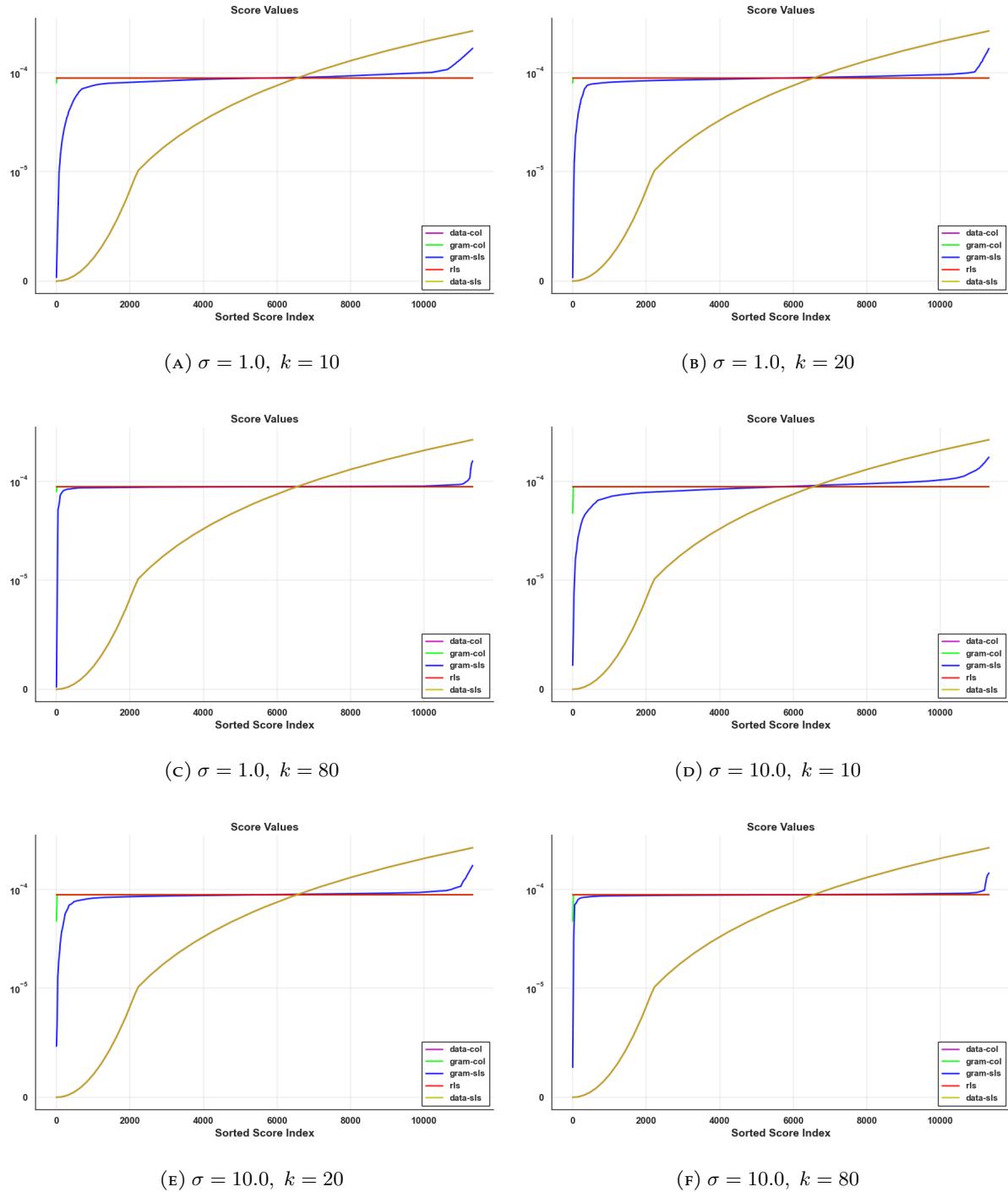


FIGURE 25. Nystrom scores for the temperature data set.

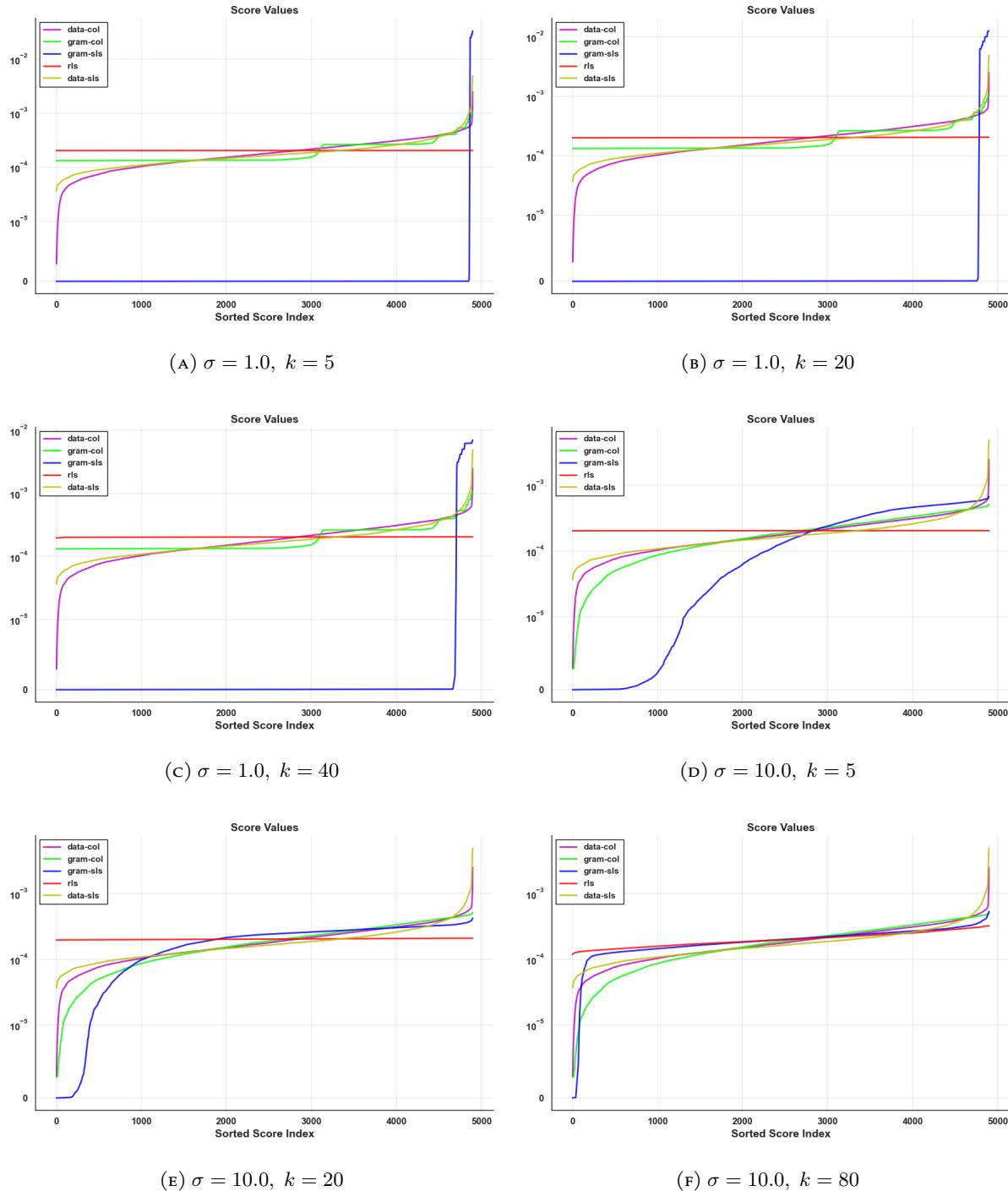


FIGURE 26. Nystrom scores for the wine data set.

A.3. Ridge Leverage Scores.

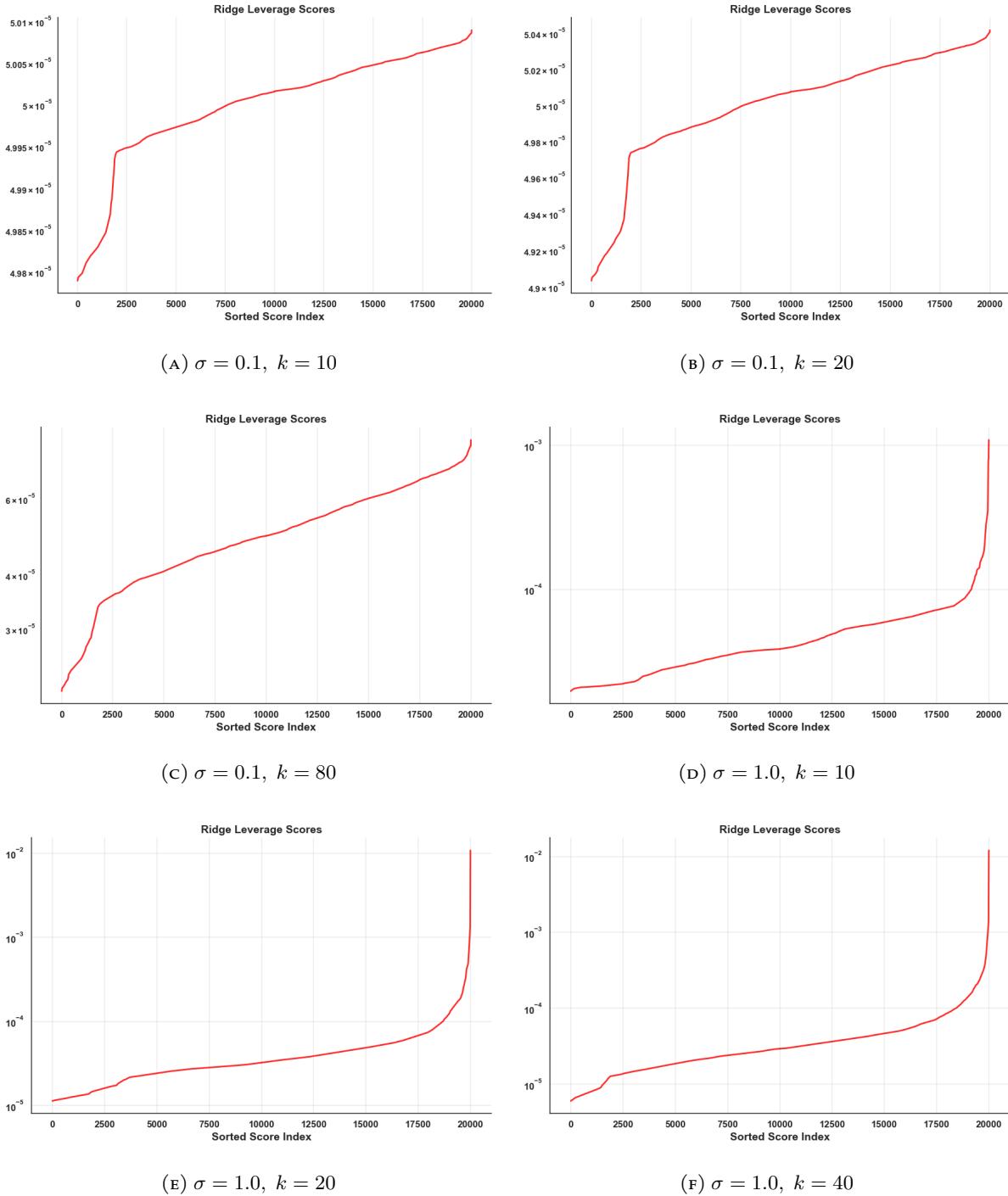


FIGURE 27. Ridge Leverage scores for the 3D-spatial network data set.

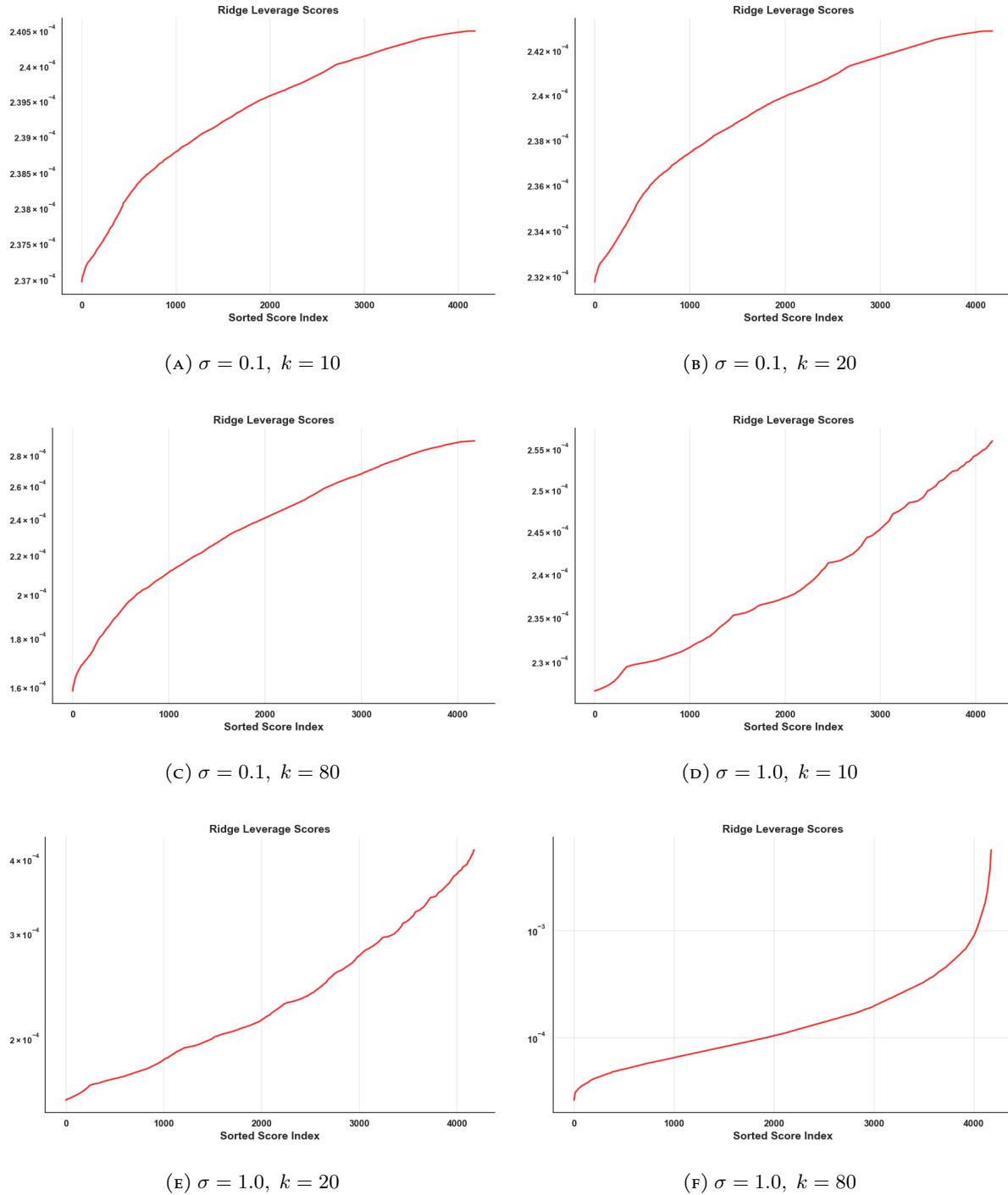
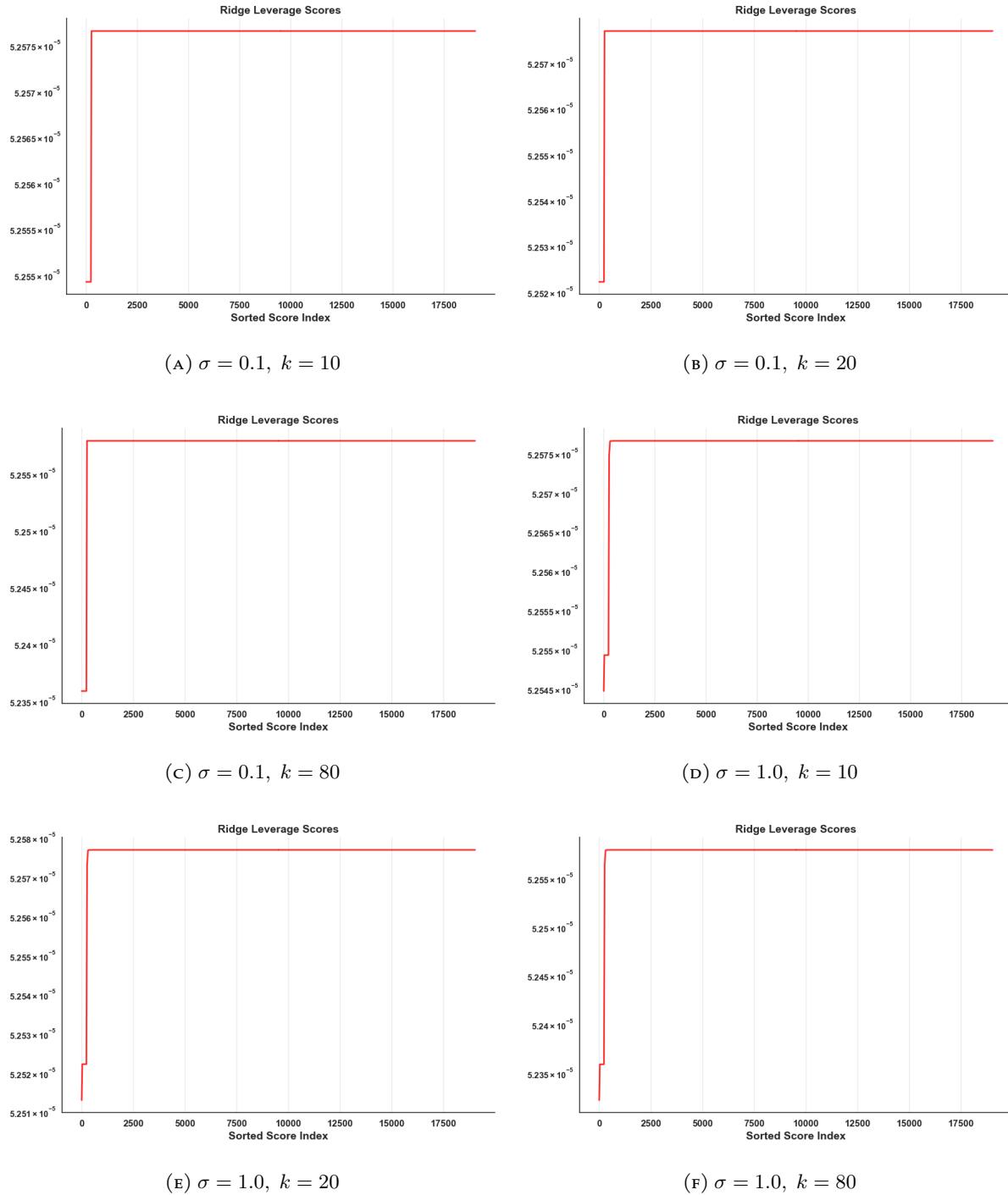


FIGURE 28. Ridge Leverage scores for the Abalone data set.



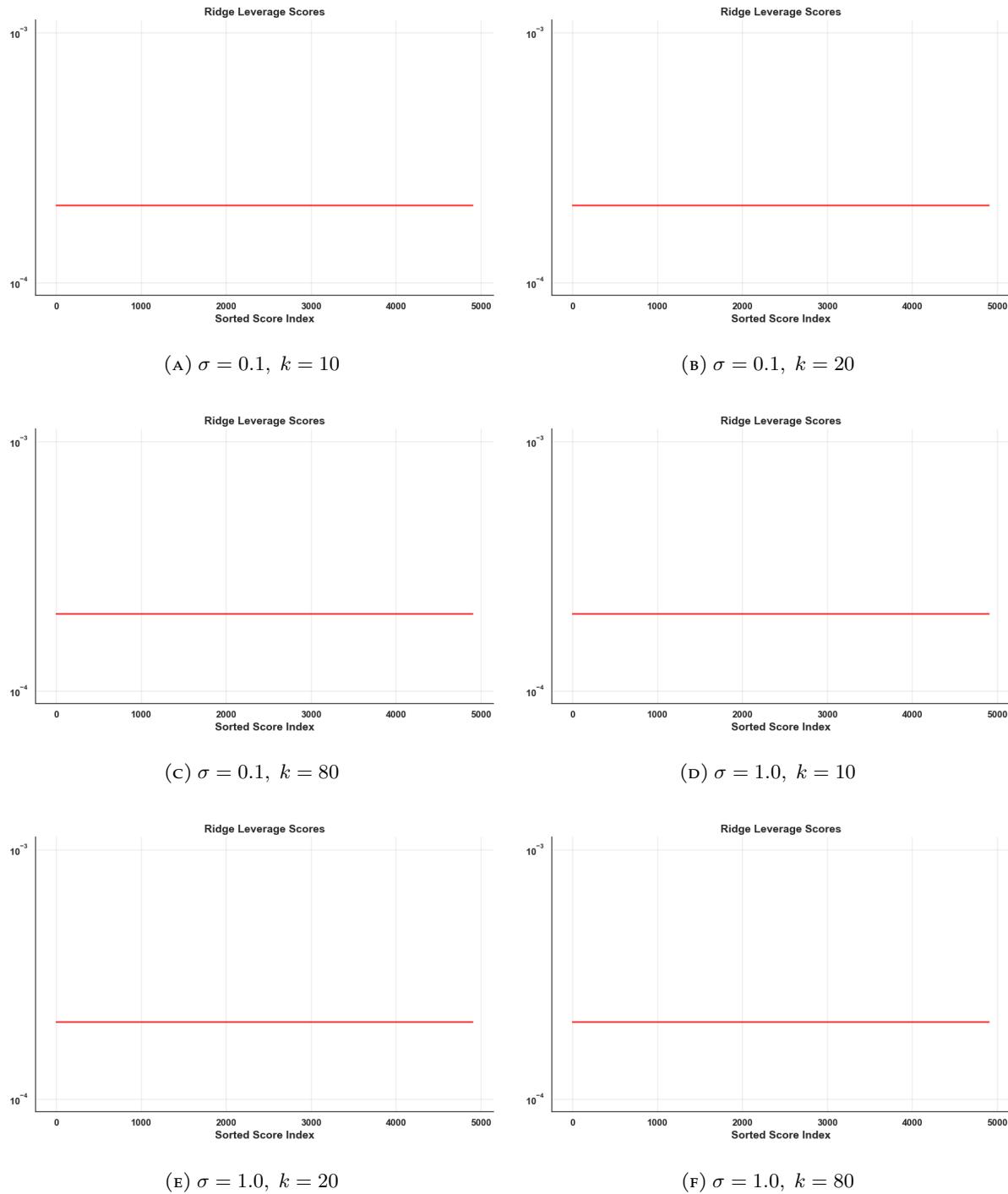


FIGURE 30. Ridge Leverage scores for the stock market data set.

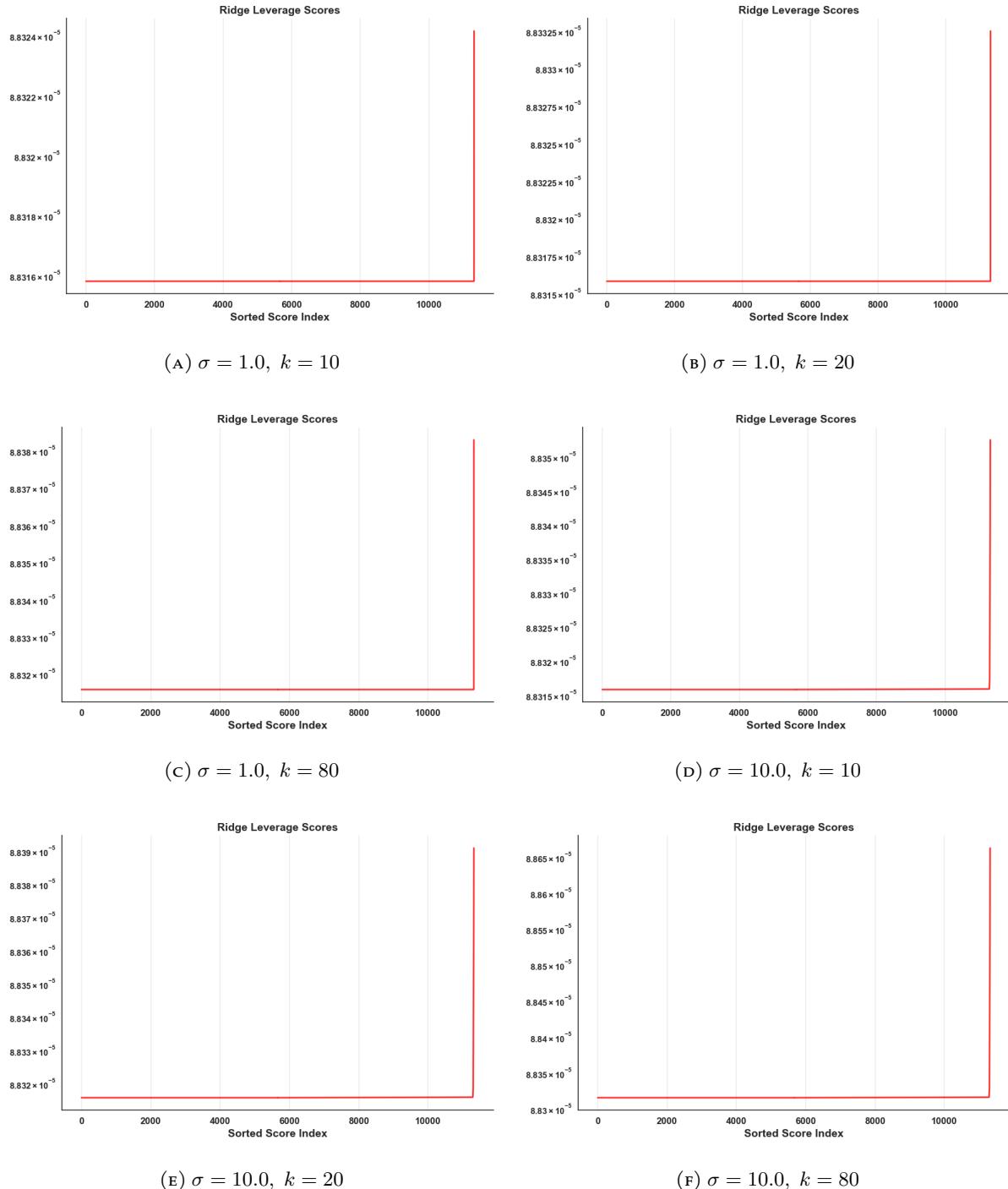


FIGURE 31. Ridge Leverage scores for the temperature data set.

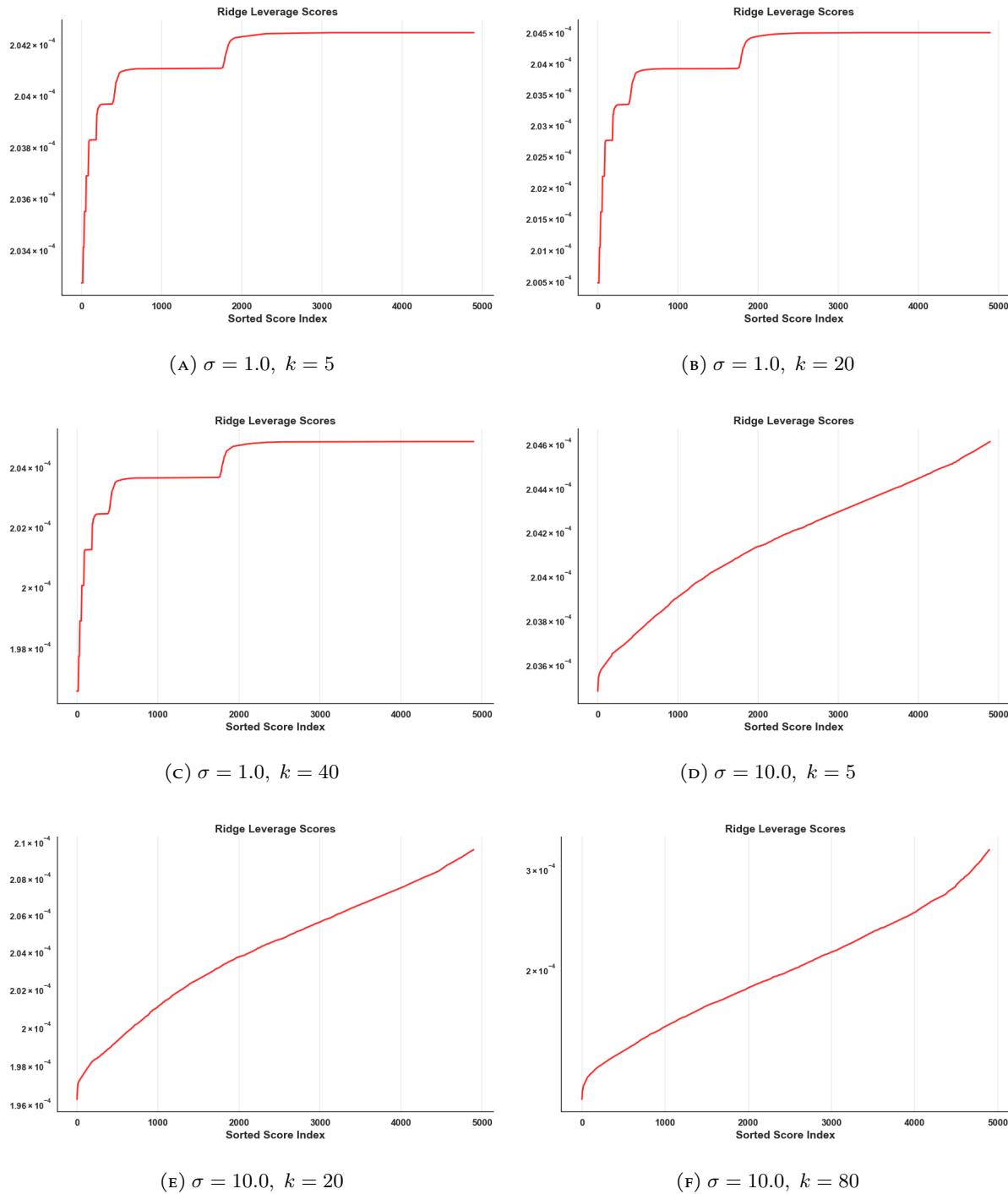


FIGURE 32. Ridge Leverage scores for the wine data set.

A.4. Nystrom Errors.

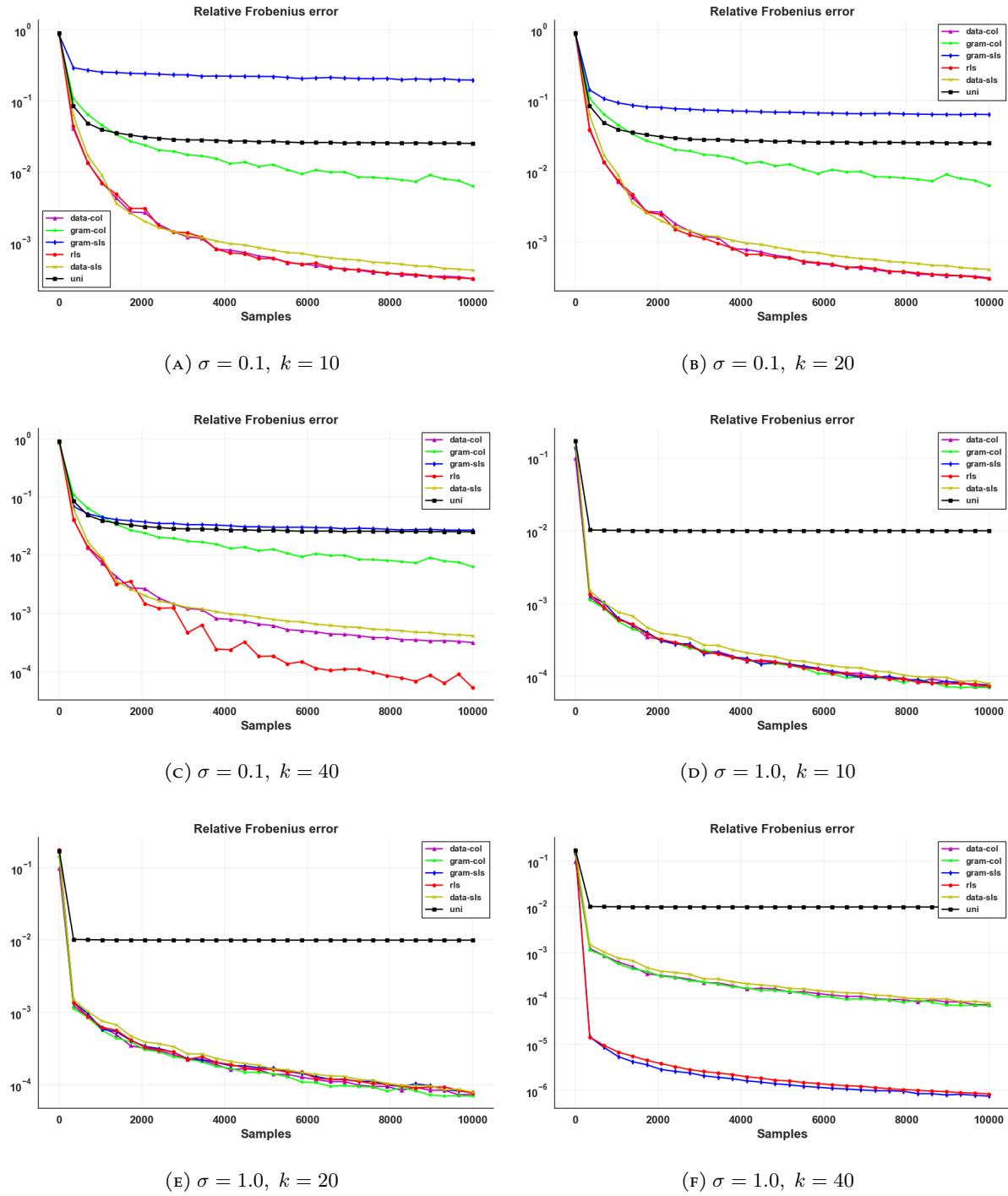


FIGURE 33. Nystrom Frobenius errors for the 3DSN data set.

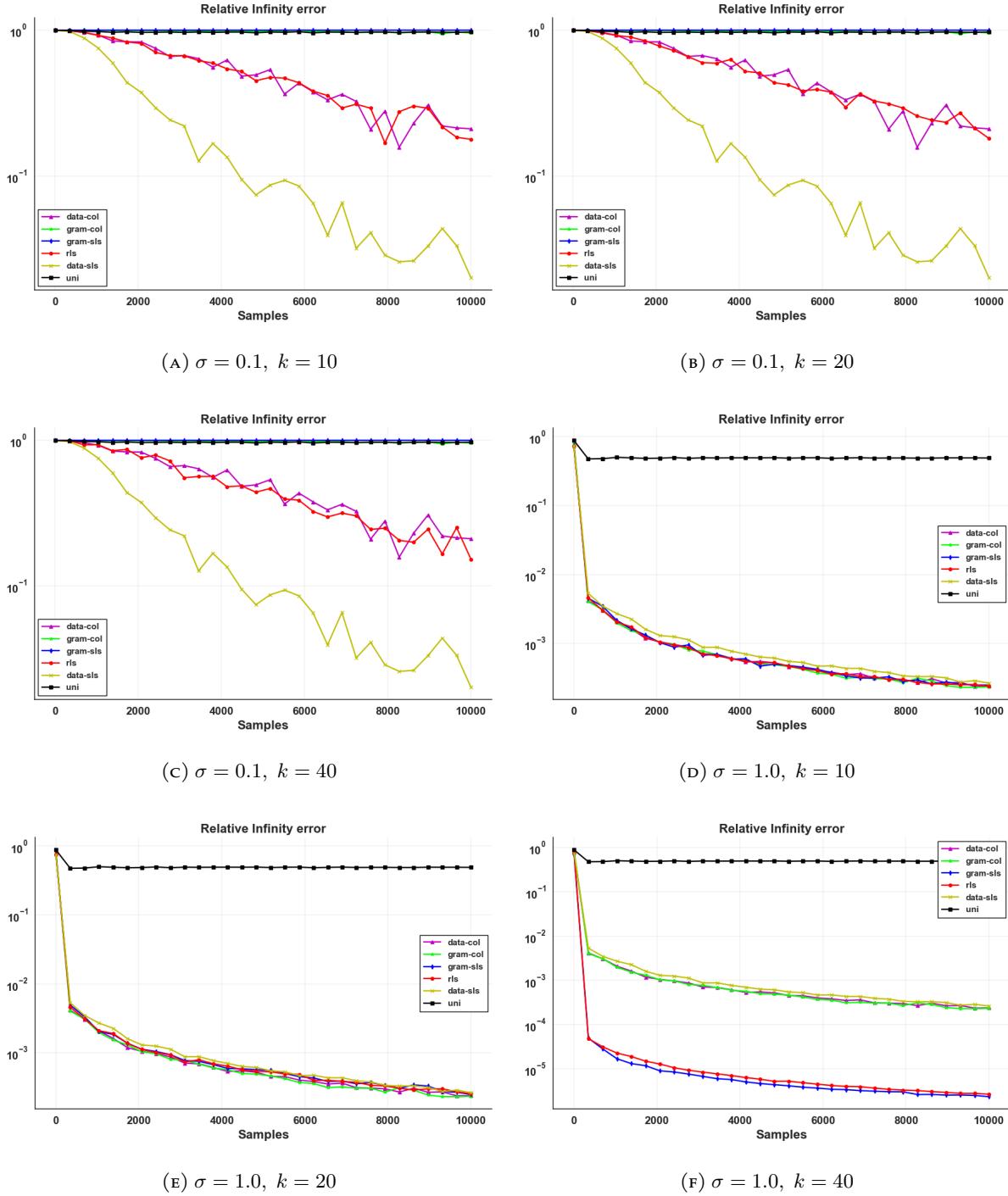


FIGURE 34. Nystrom infinity errors for the 3DSN data set.

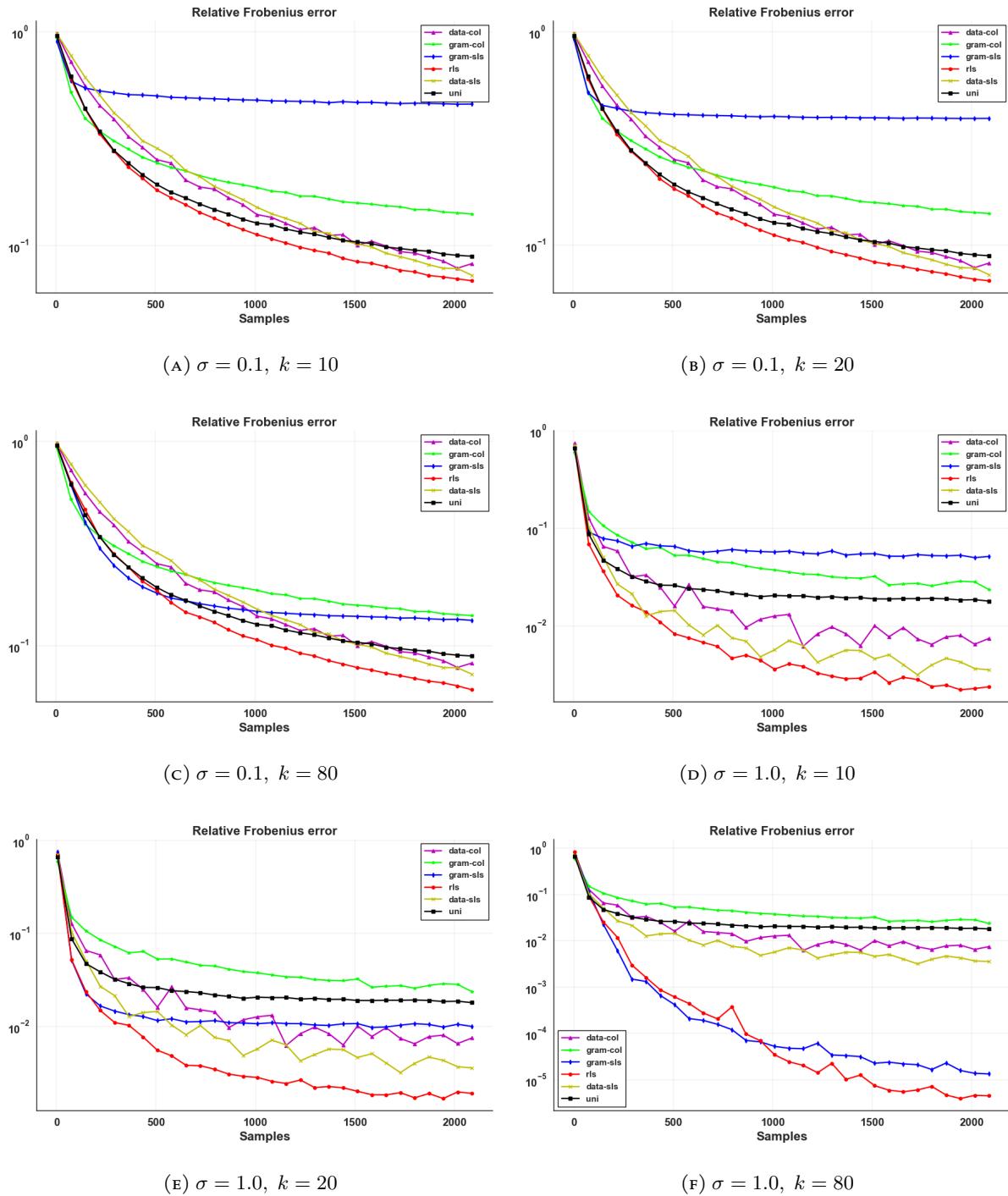


FIGURE 35. Nystrom Frobenius errors for the Abalone data set.

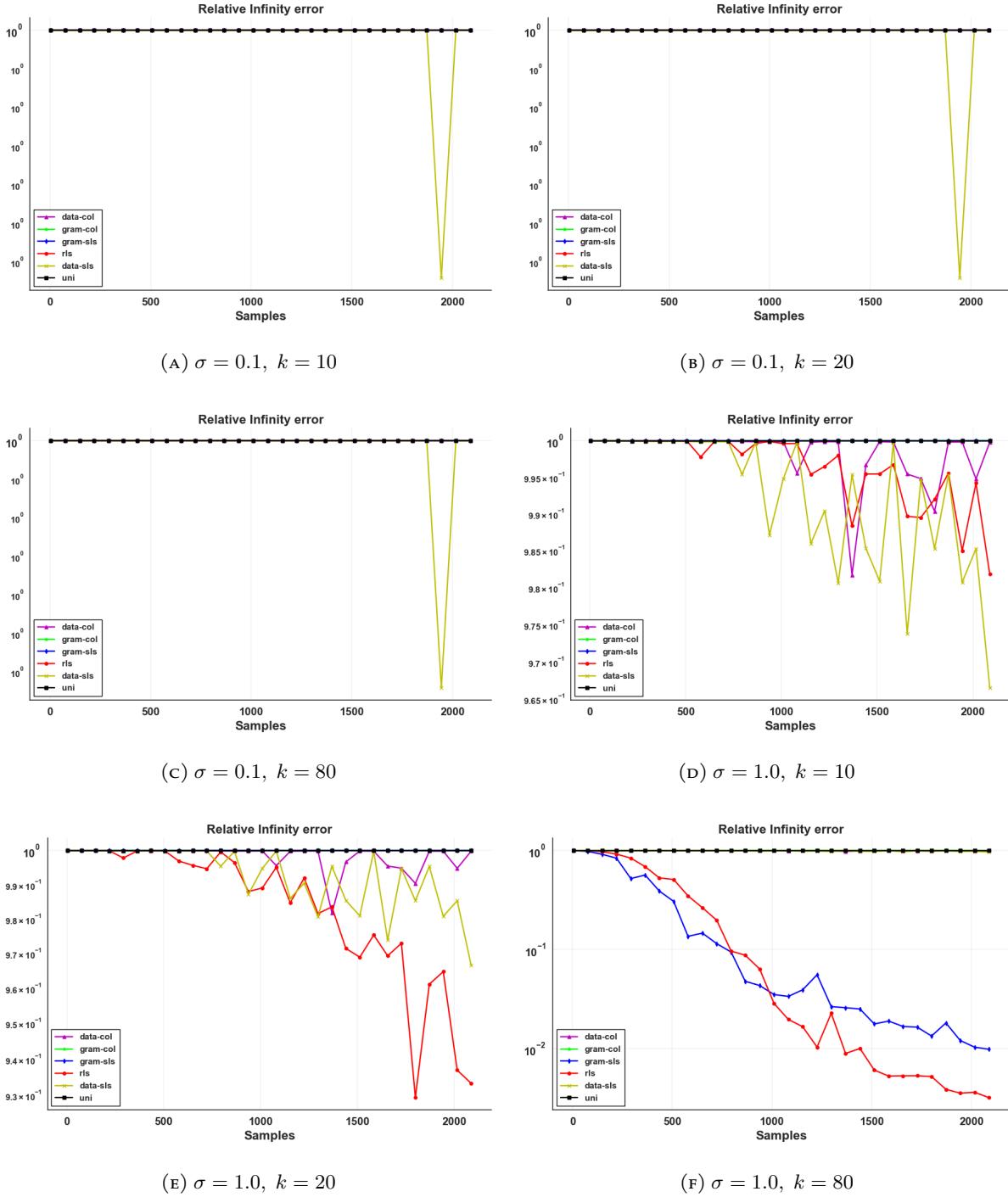


FIGURE 36. Nystrom infinity errors for the Abalone data set.

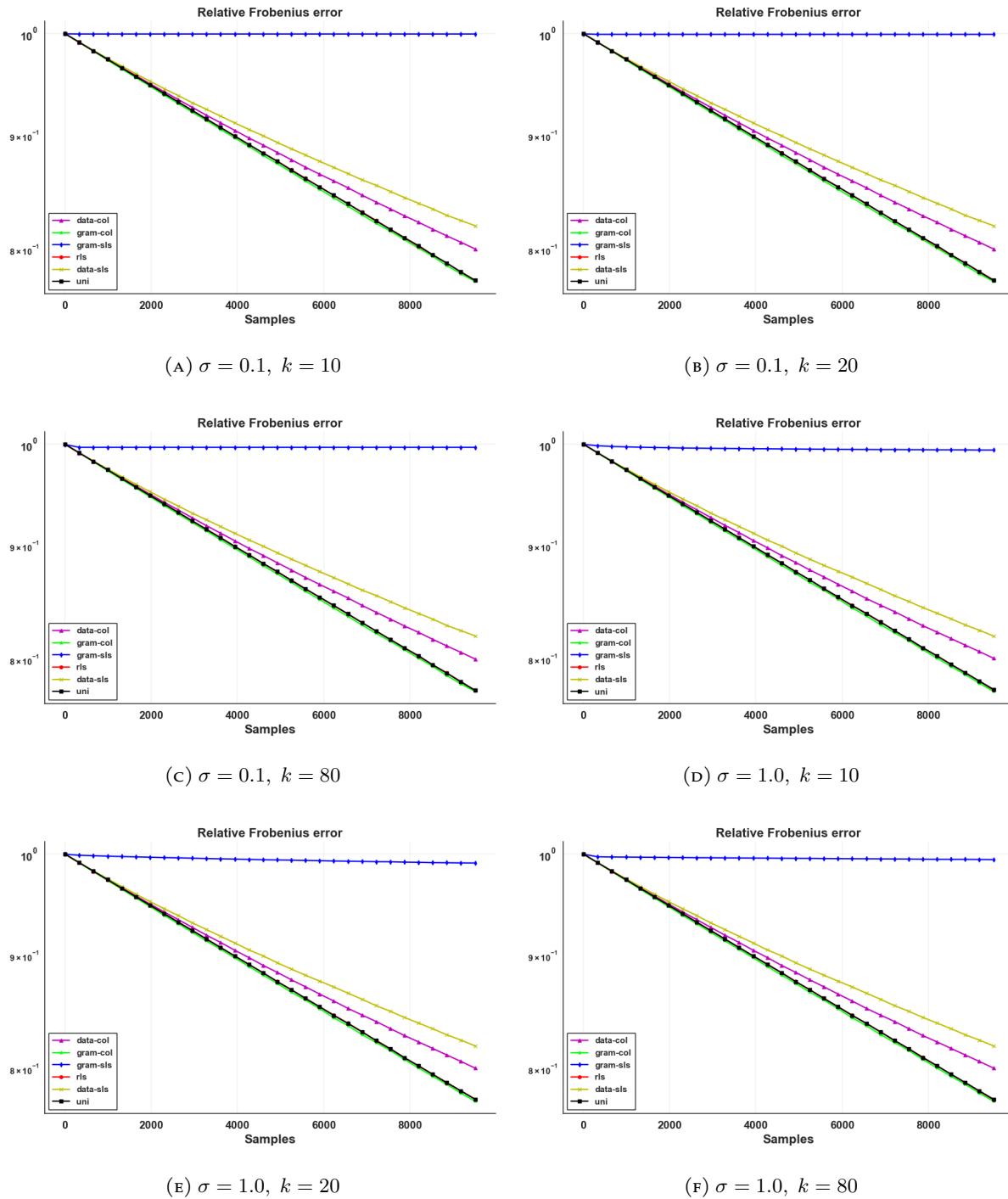


FIGURE 37. Nystrom Frobenius errors for the Magic data set.

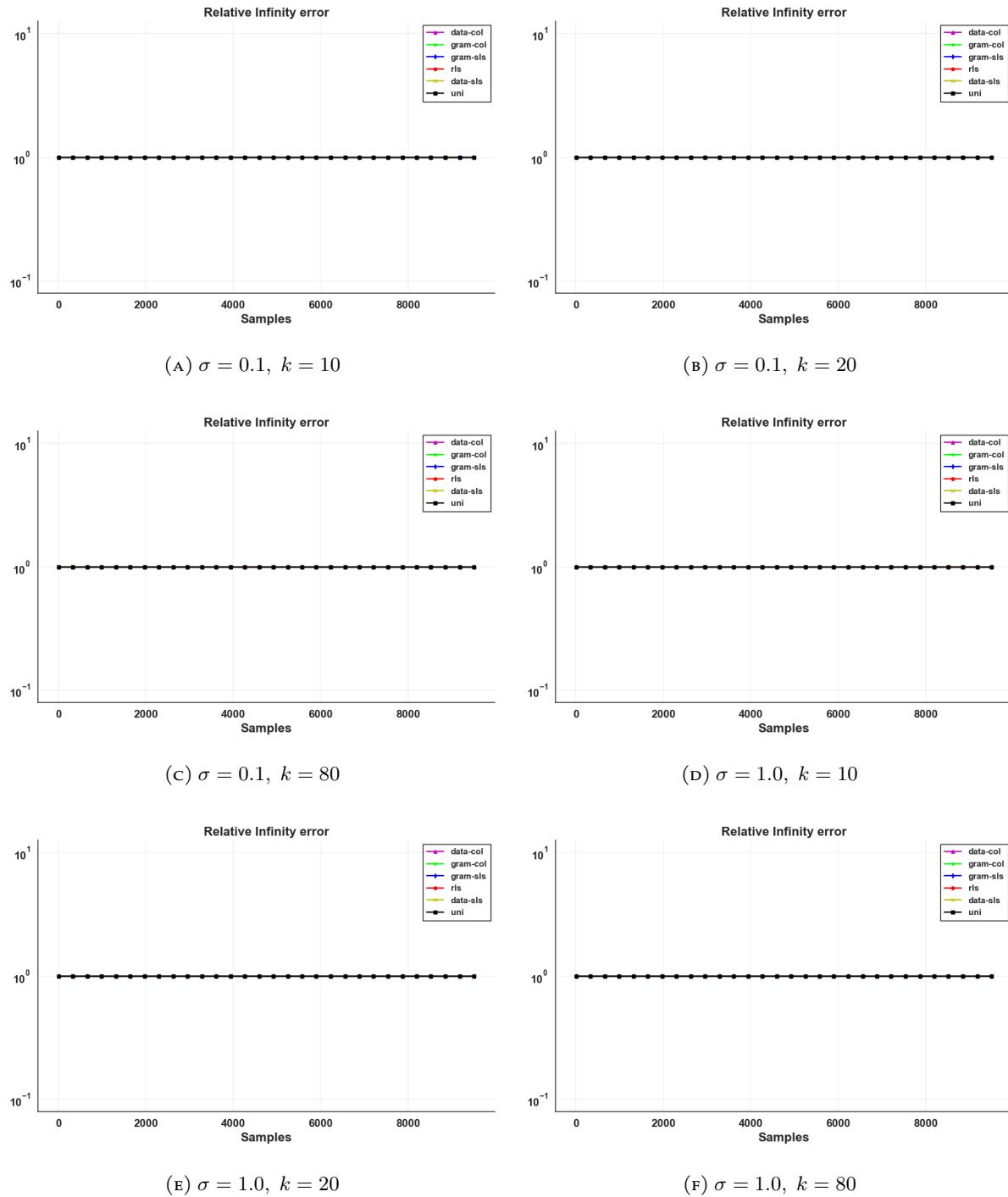


FIGURE 38. Nystrom infinity errors for the Magic data set.

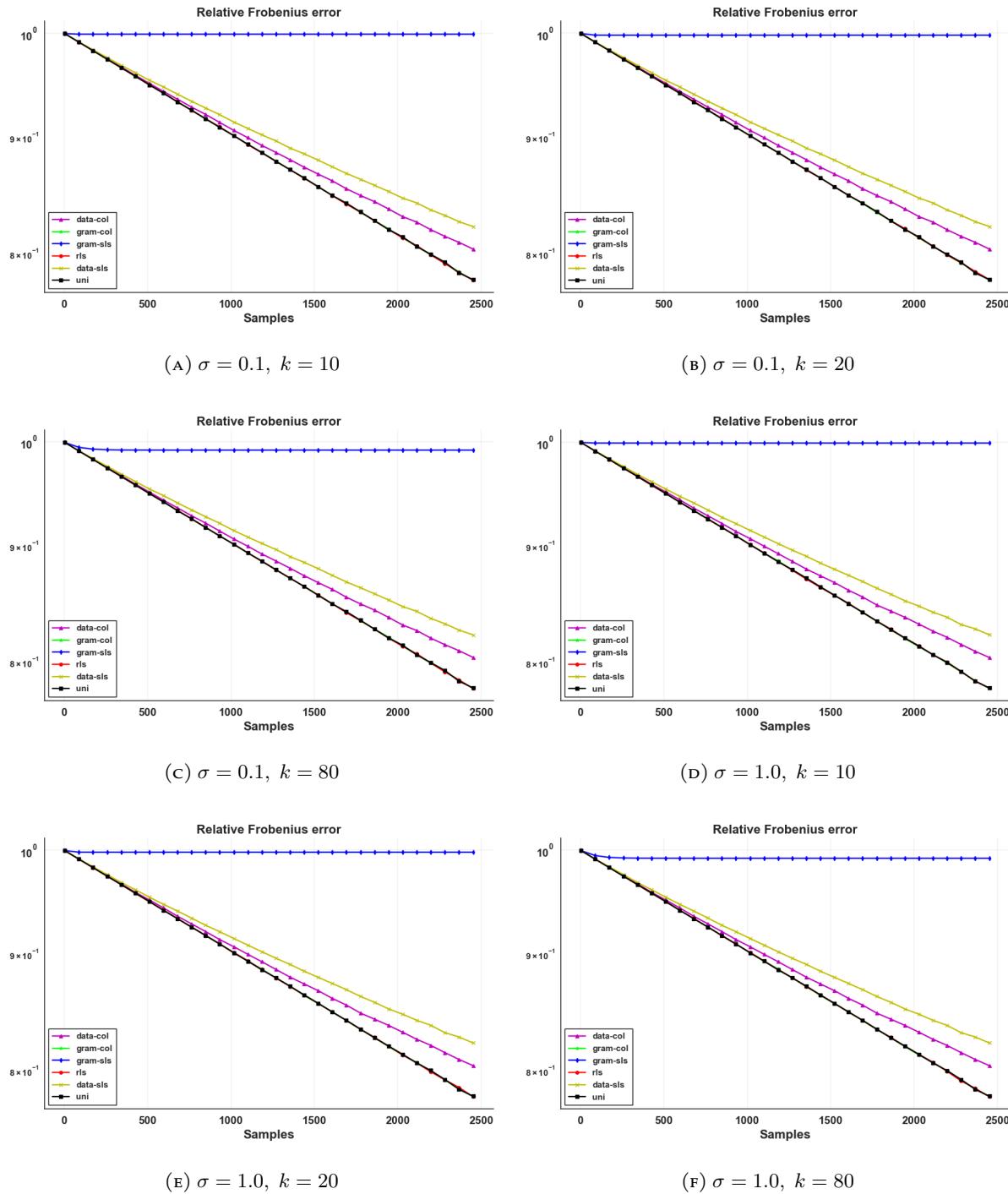


FIGURE 39. Nystrom Frobenius errors for the Stocks data set.

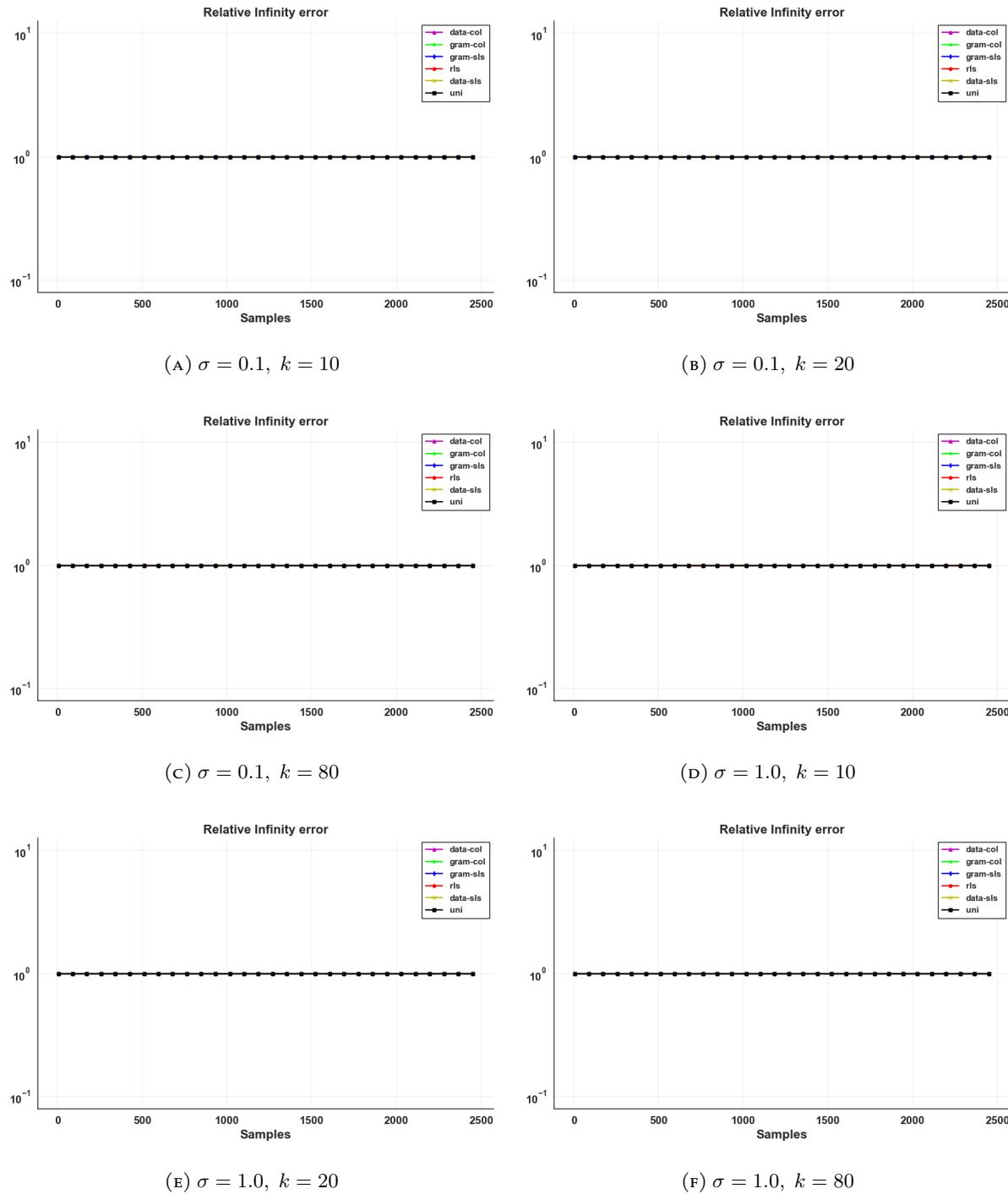


FIGURE 40. Nystrom infinity errors for the Stocks data set.

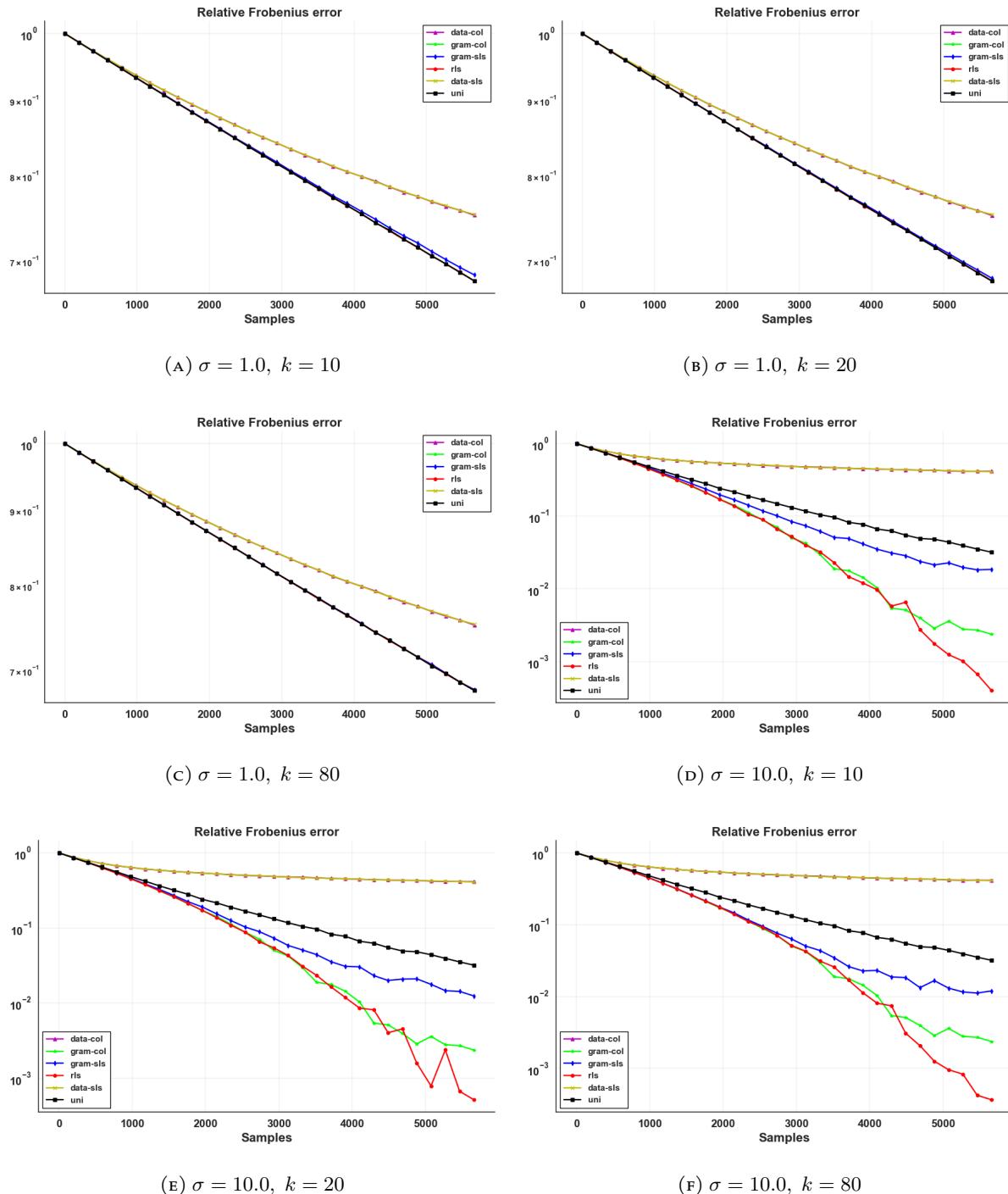


FIGURE 41. Nystrom Frobenius errors for the Temp data set.

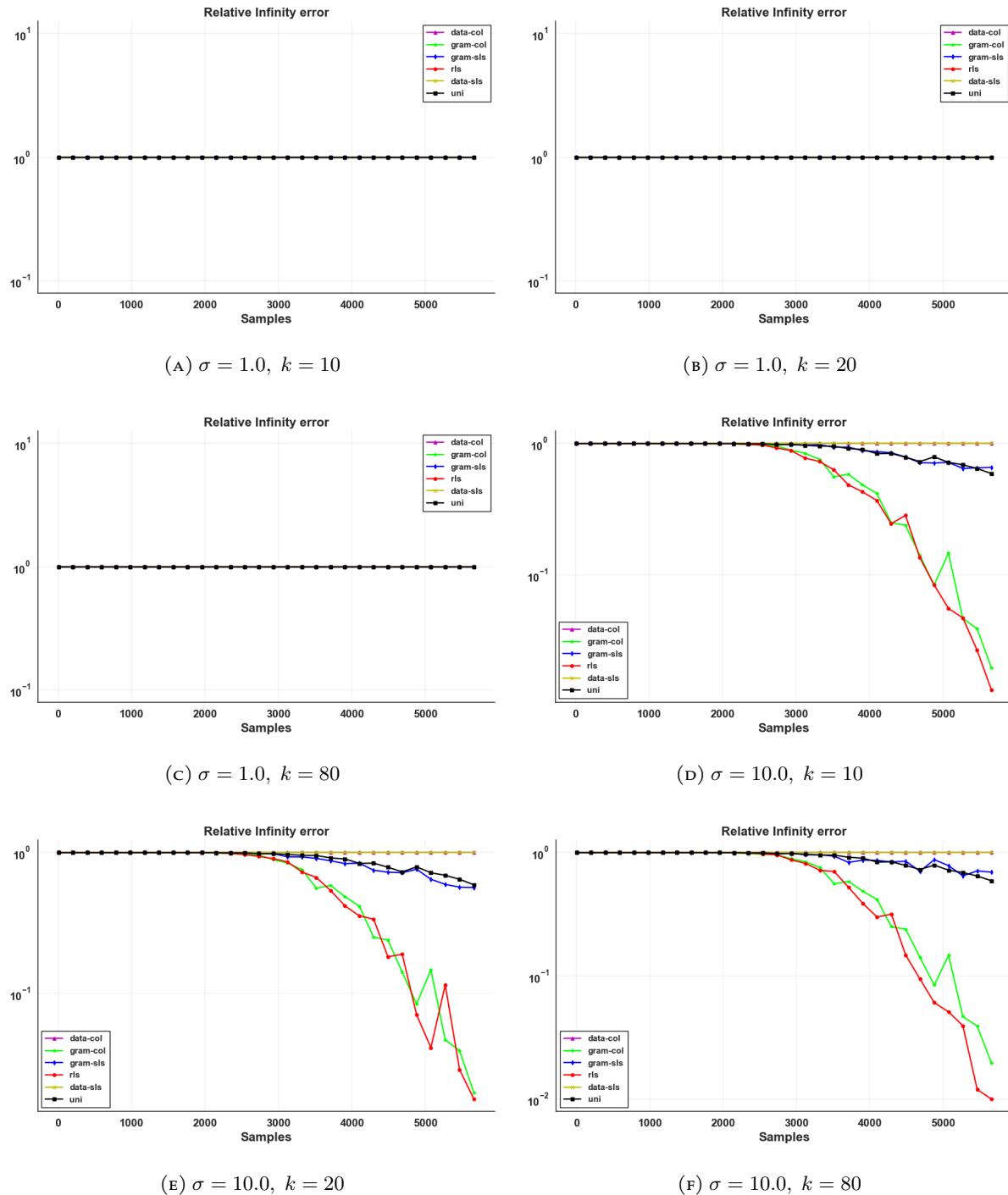


FIGURE 42. Nystrom infinity errors for the Temp data set.

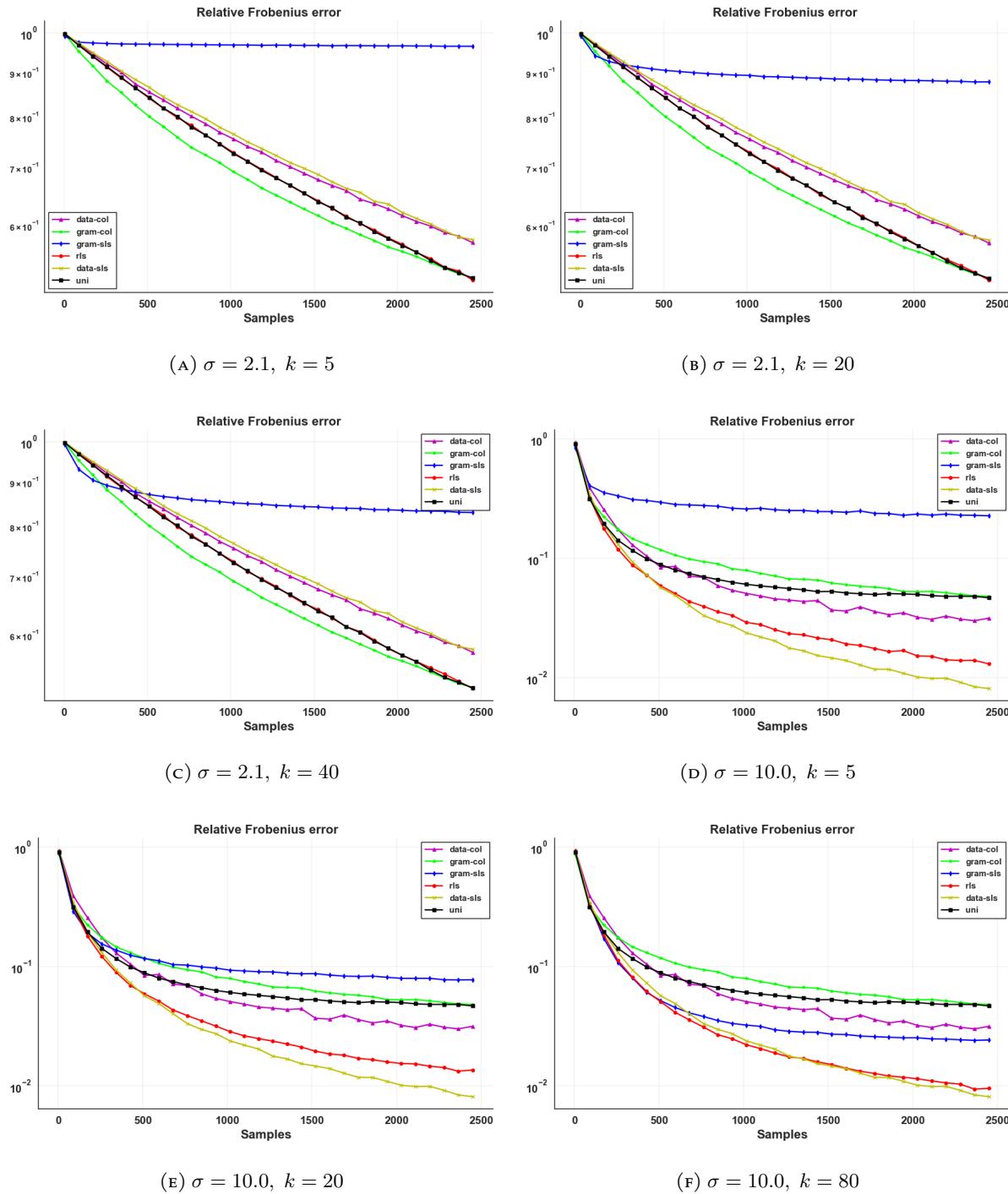


FIGURE 43. Nystrom Frobenius errors for the wine data set.

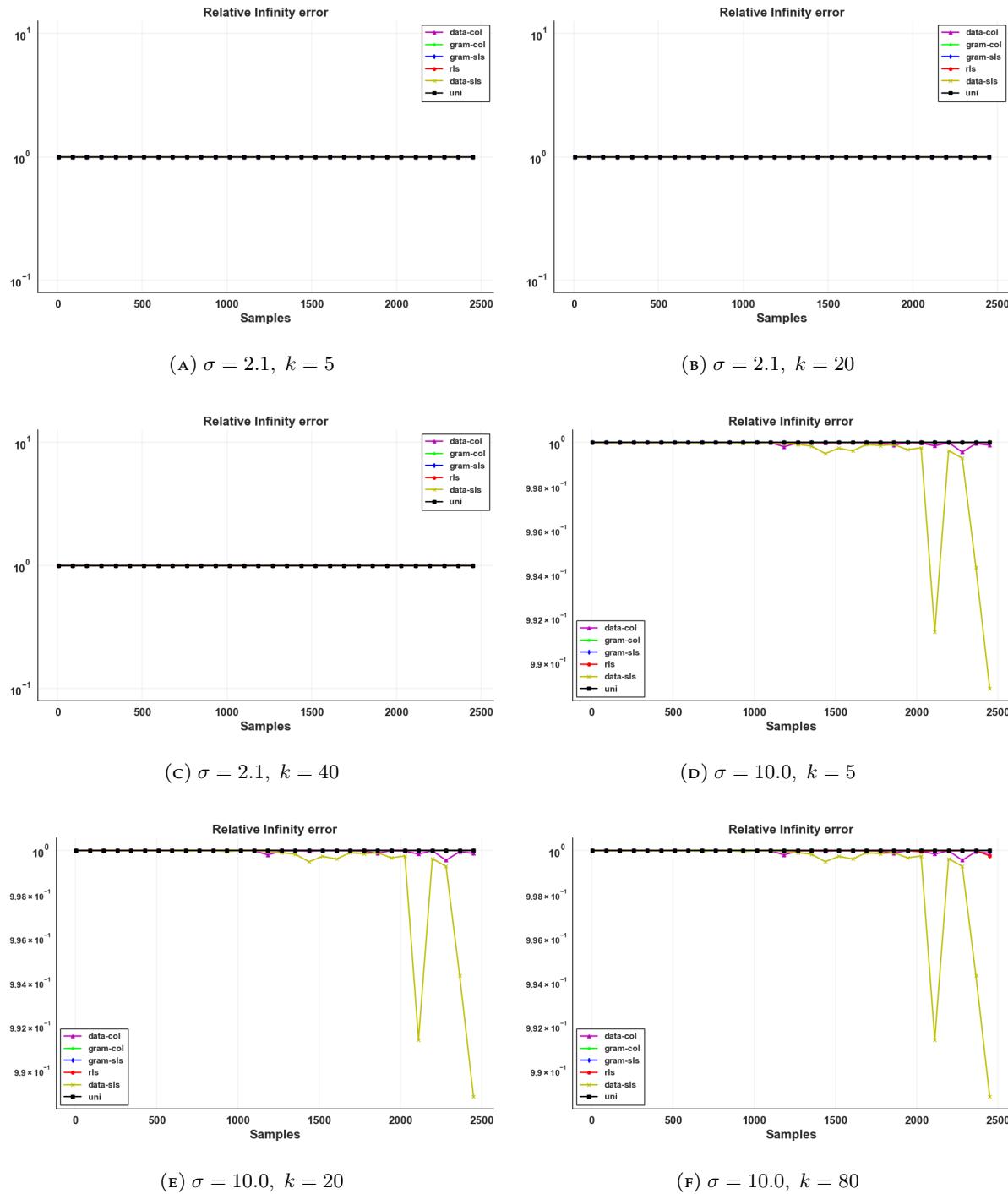


FIGURE 44. Nystrom infinity errors for the wine data set.

A.5. Nystrom Sampling Distribution Construction Times.

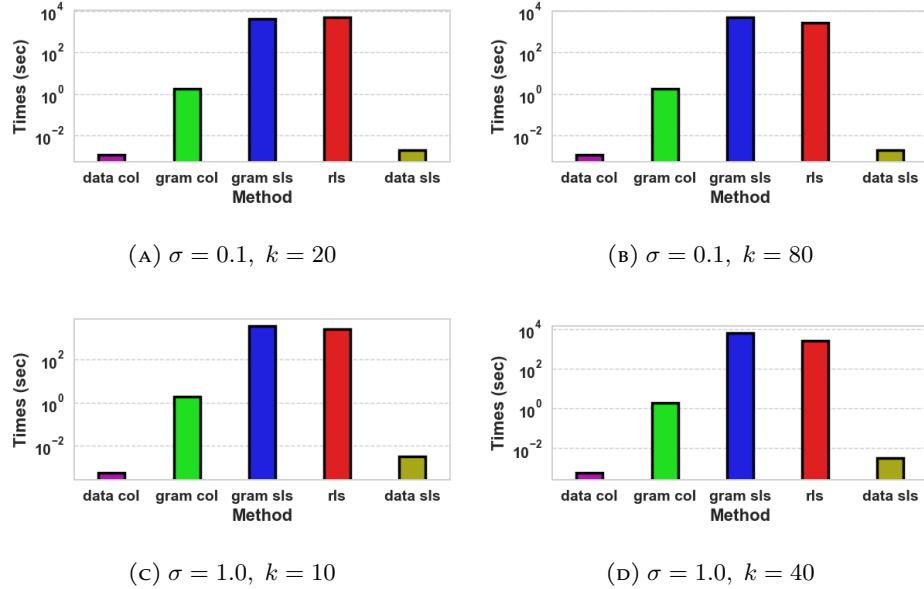


FIGURE 45. Nystrom sampling distribution construction times for the 3D-spatial network data set.

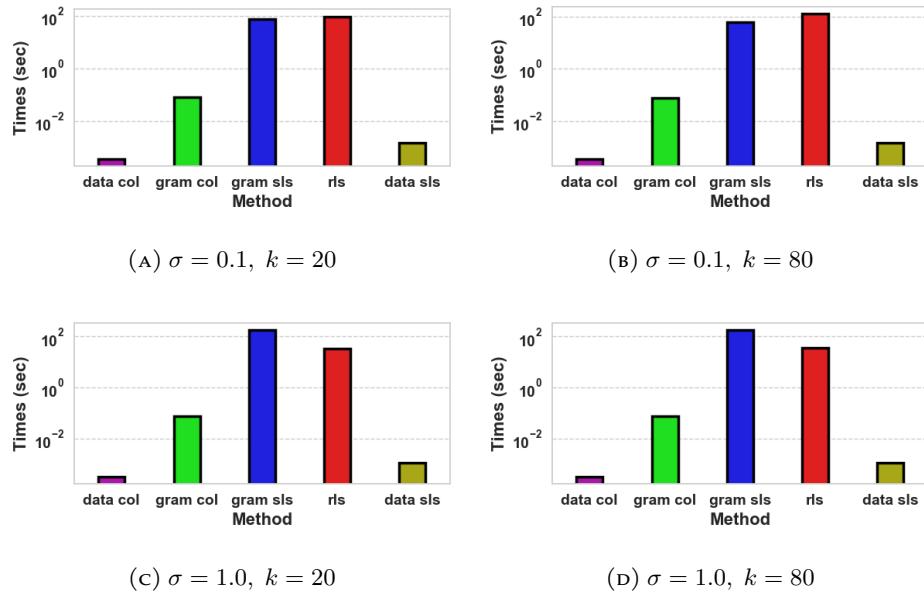


FIGURE 46. Nystrom sampling distribution construction times for the Abalone data set.

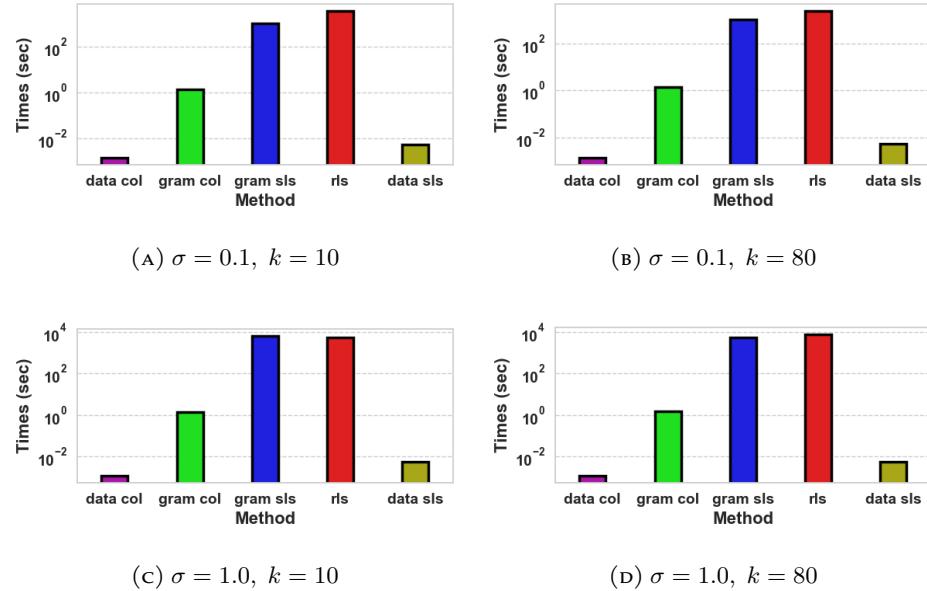


FIGURE 47. Nystrom sampling distribution construction times for the Magic data set.

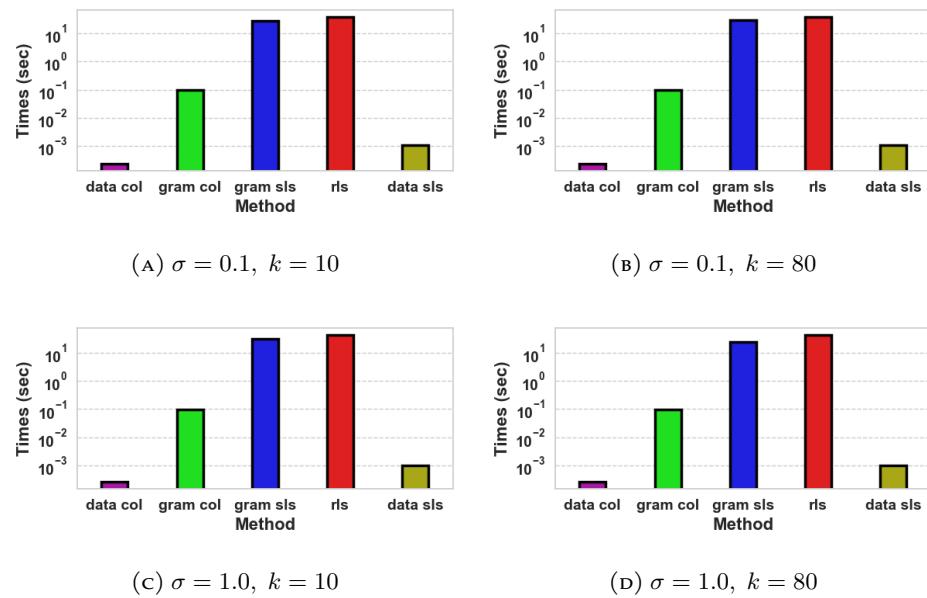


FIGURE 48. Nystrom sampling distribution construction times for the stock market data set.

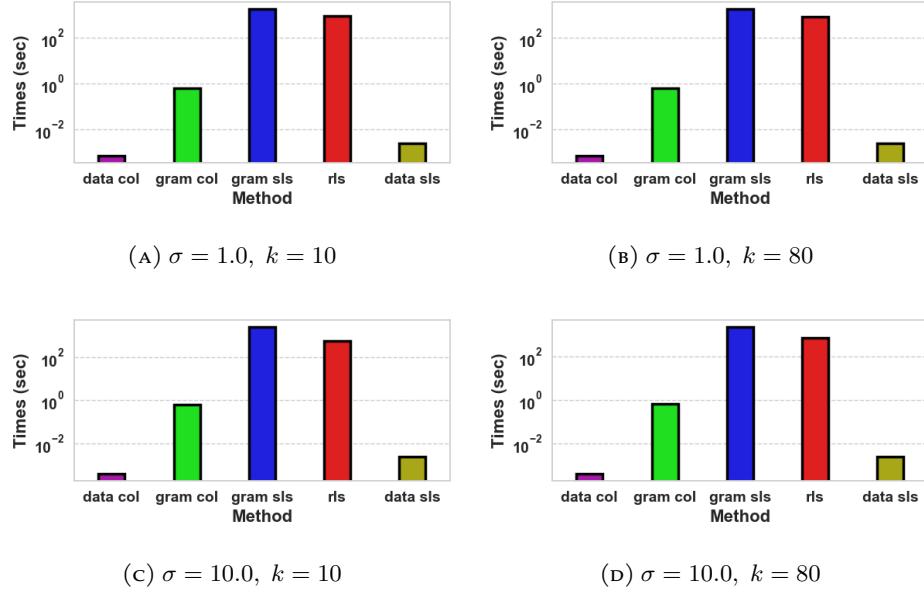


FIGURE 49. Nystrom sampling distribution construction times for the temperature data set.

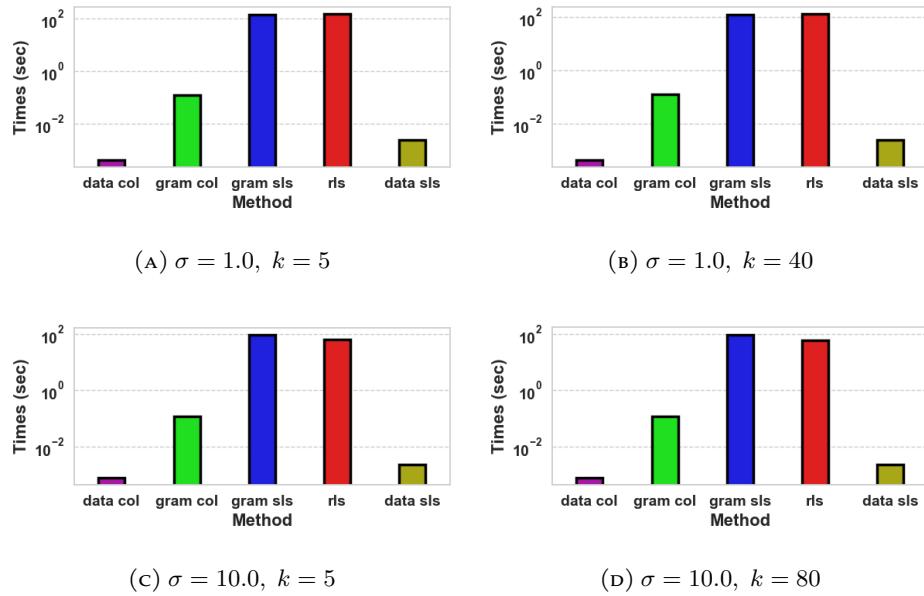


FIGURE 50. Nystrom sampling distribution construction times for the wine data set.

A.6. Nystrom Sketch Times.

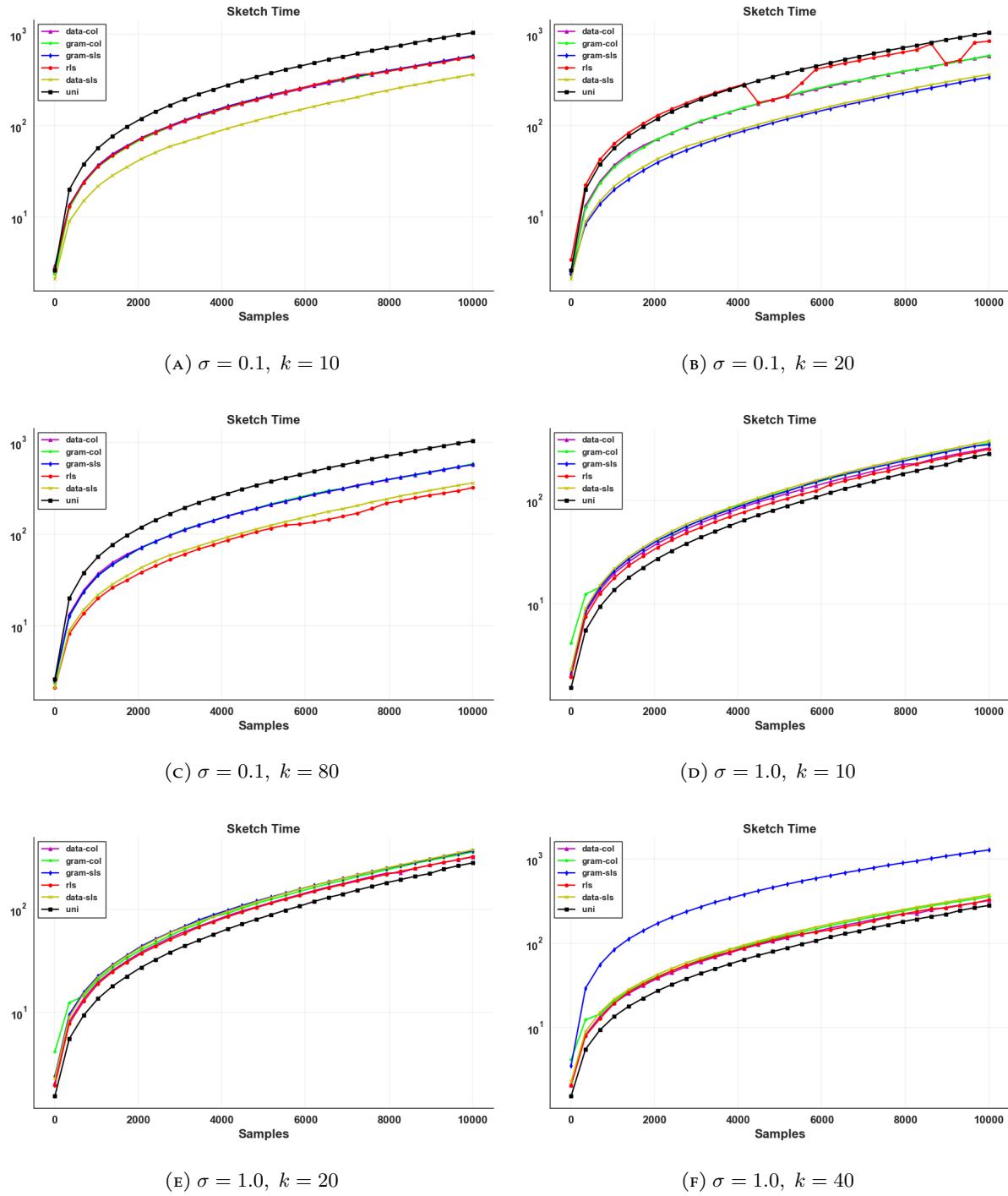


FIGURE 51. Nystrom sketch times for the 3D-spatial network data set.

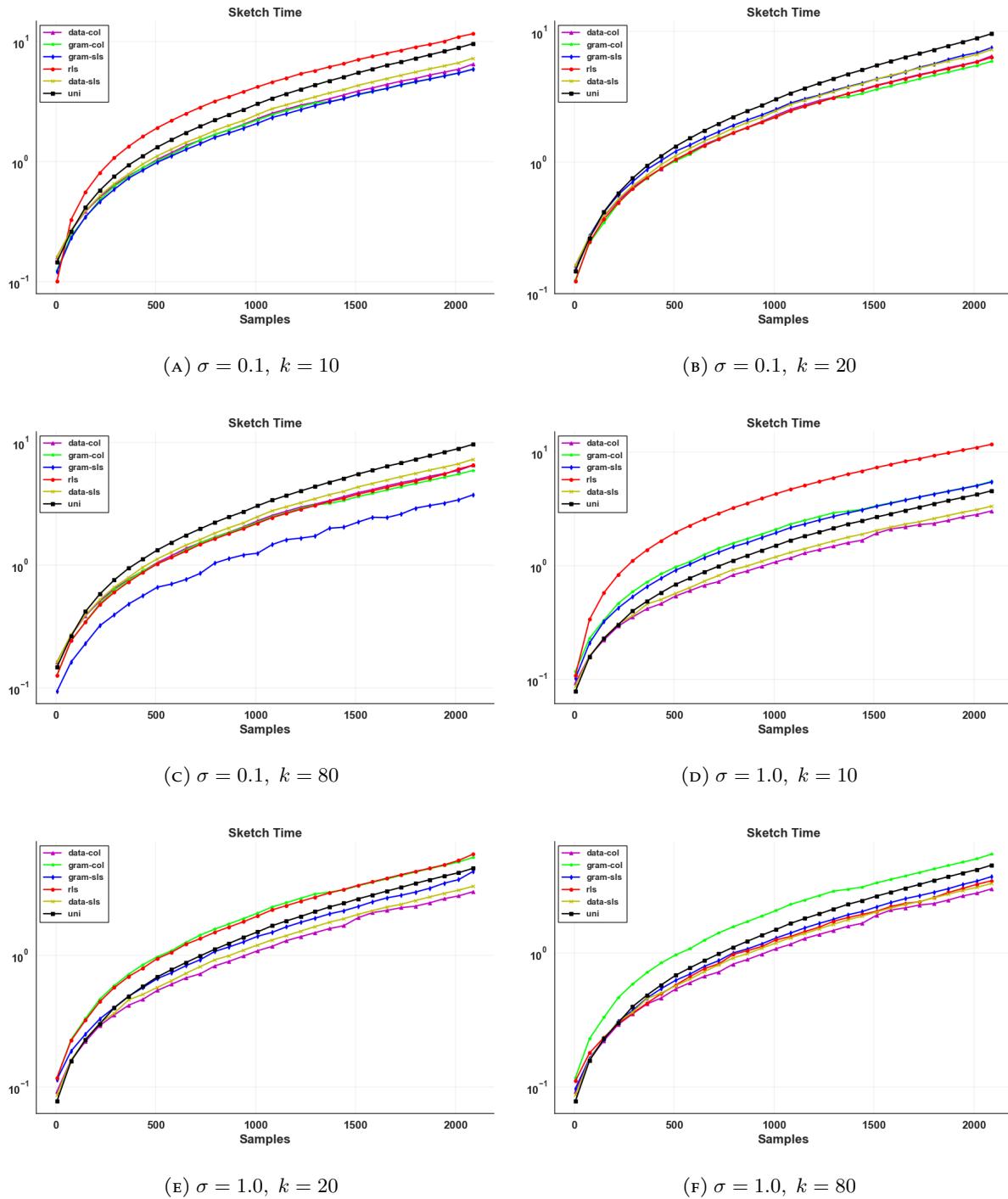


FIGURE 52. Nystrom sketch times for the Abalone data set.

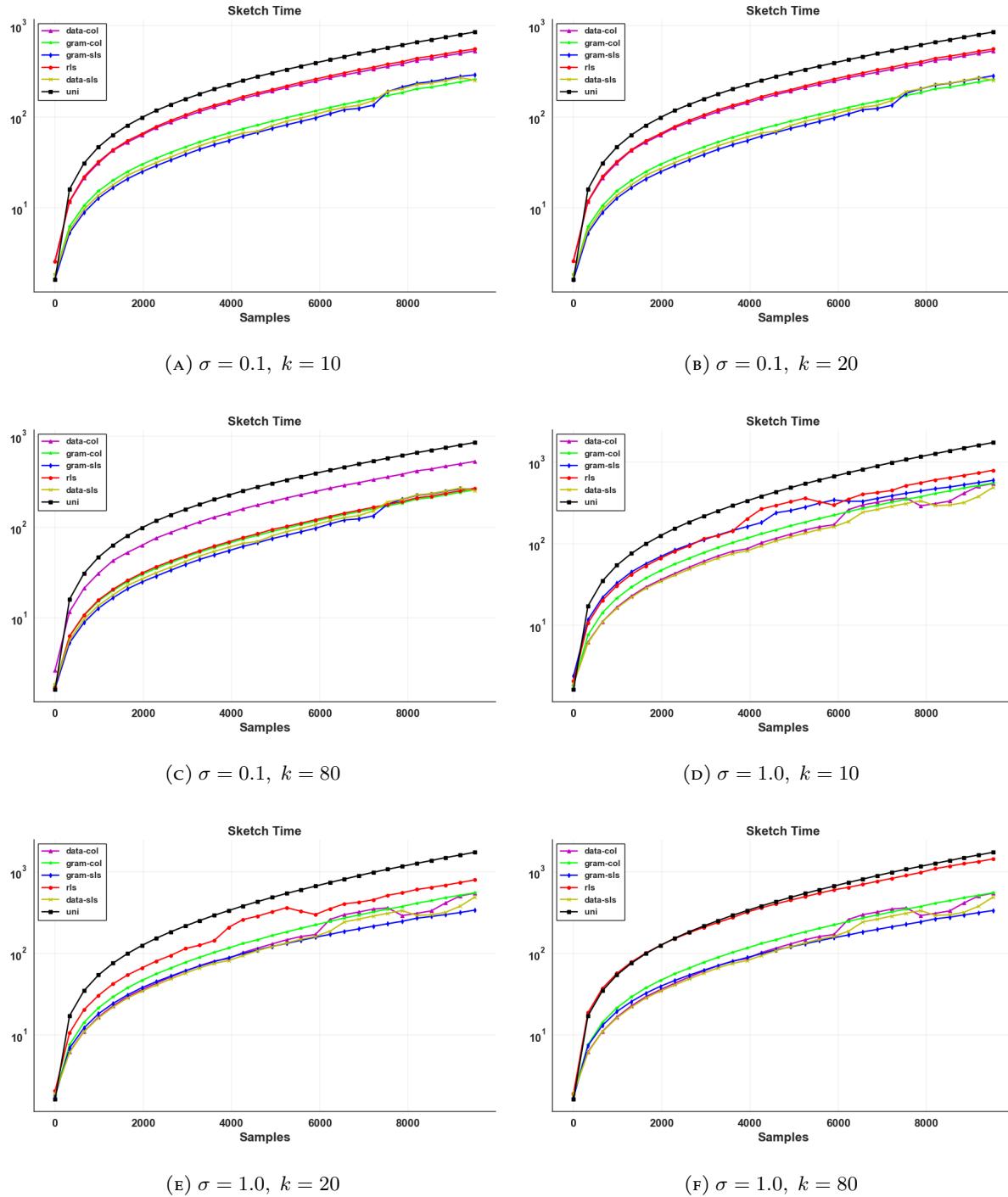


FIGURE 53. Nystrom sketch times for the Magic data set.

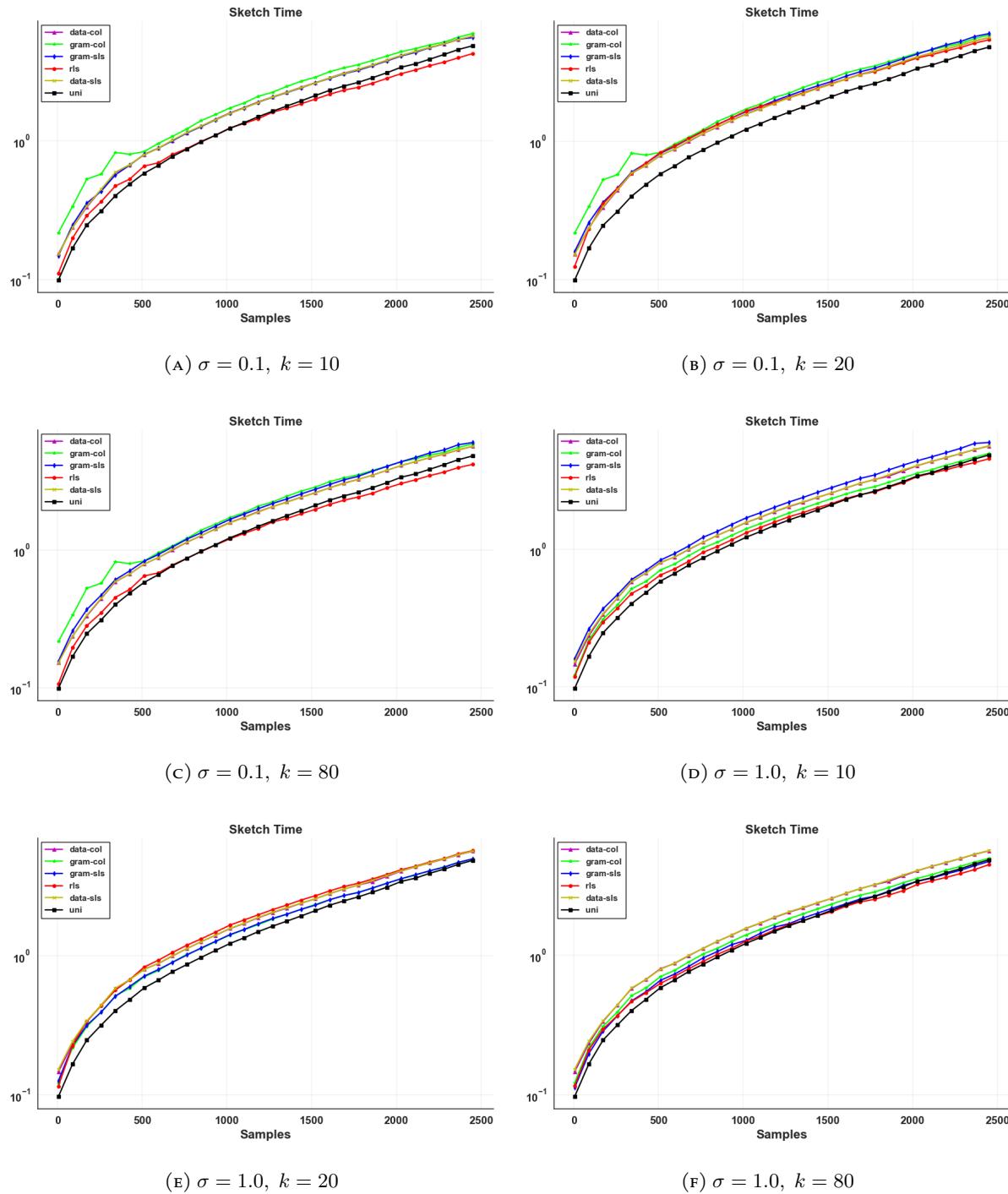


FIGURE 54. Nystrom sketch times for the stock market data set.

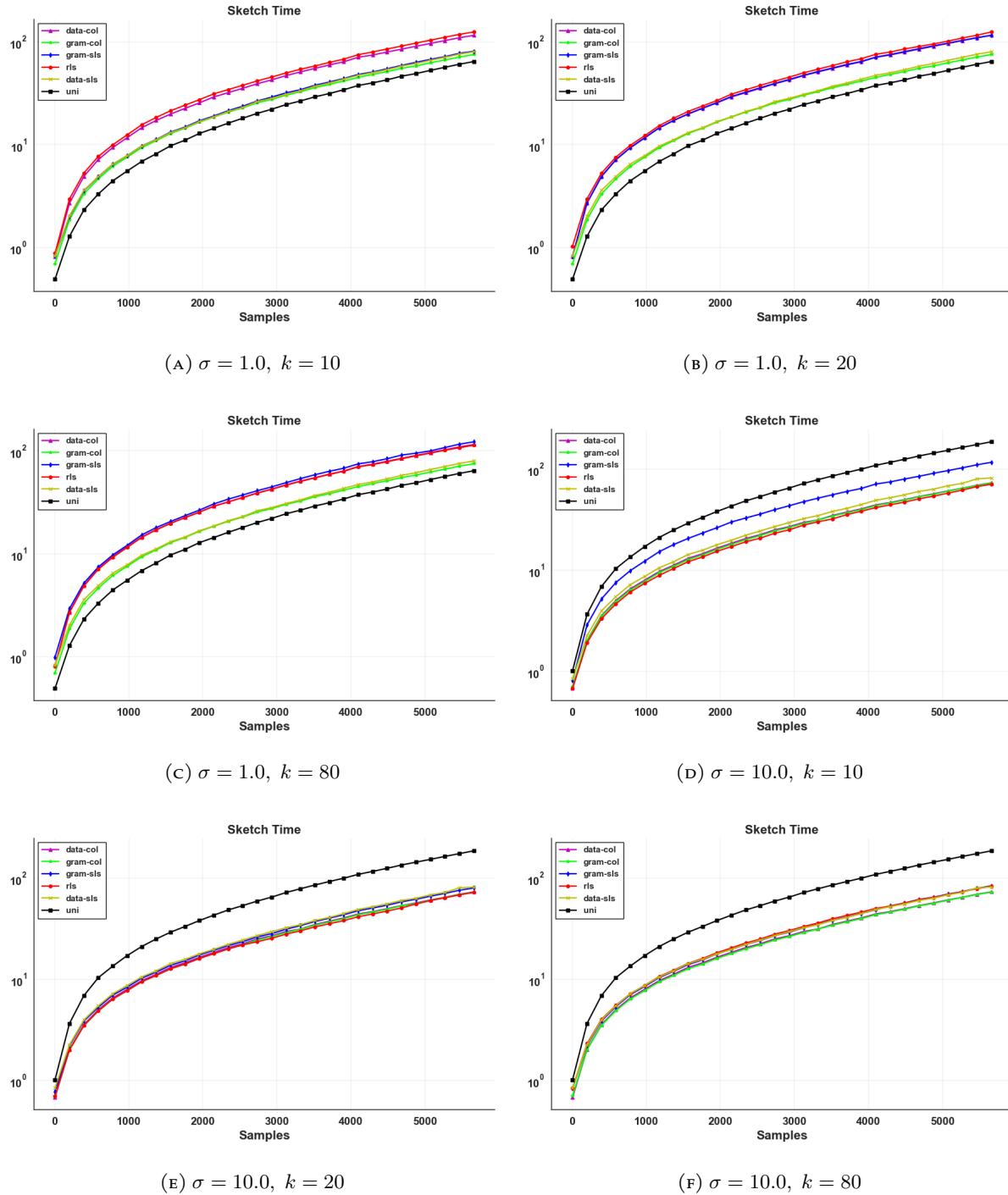


FIGURE 55. Nystrom sketch times for the temperature data set.

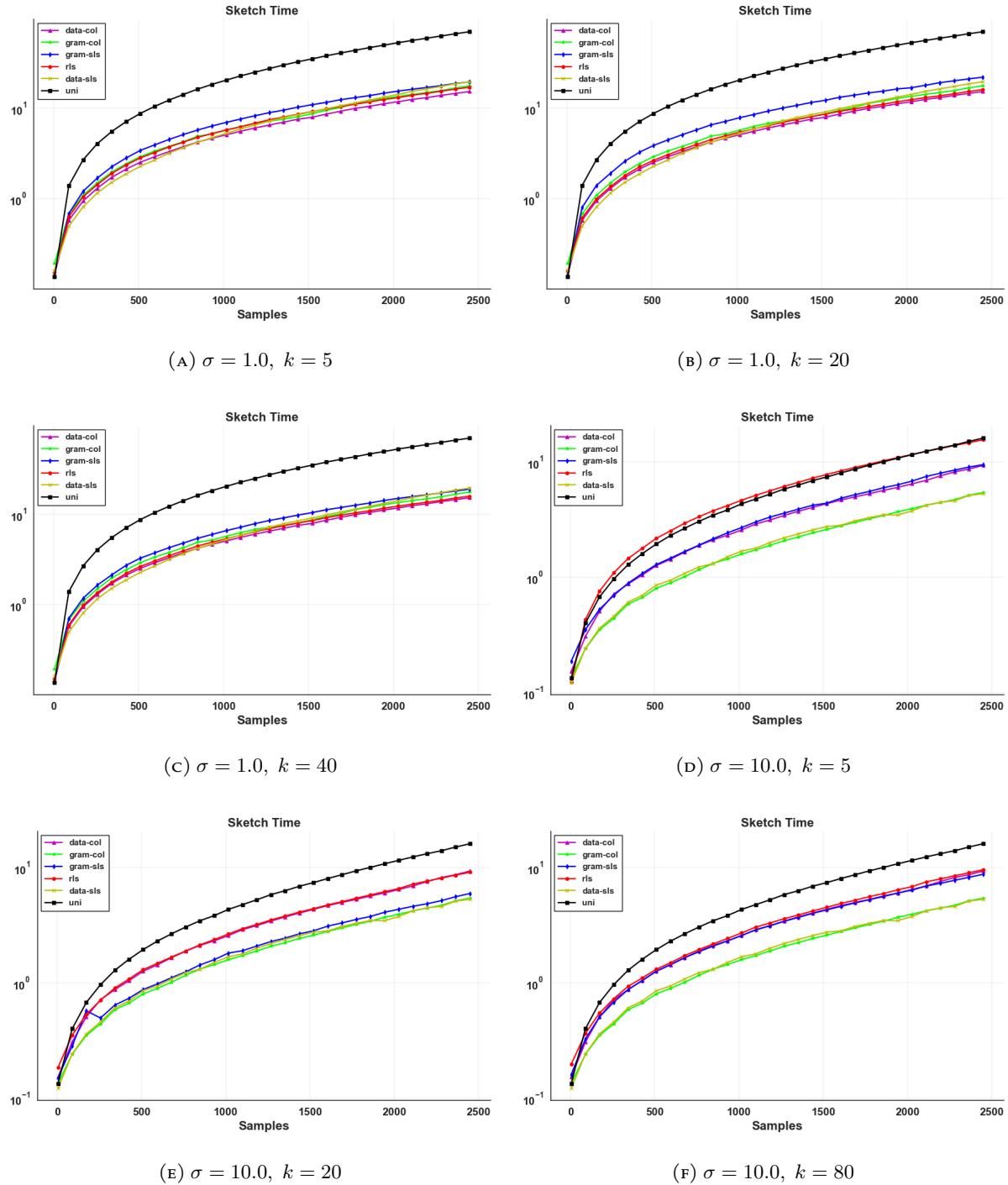


FIGURE 56. Nystrom sketch times for the wine data set.

A.7. Nystrom Memory Usage.

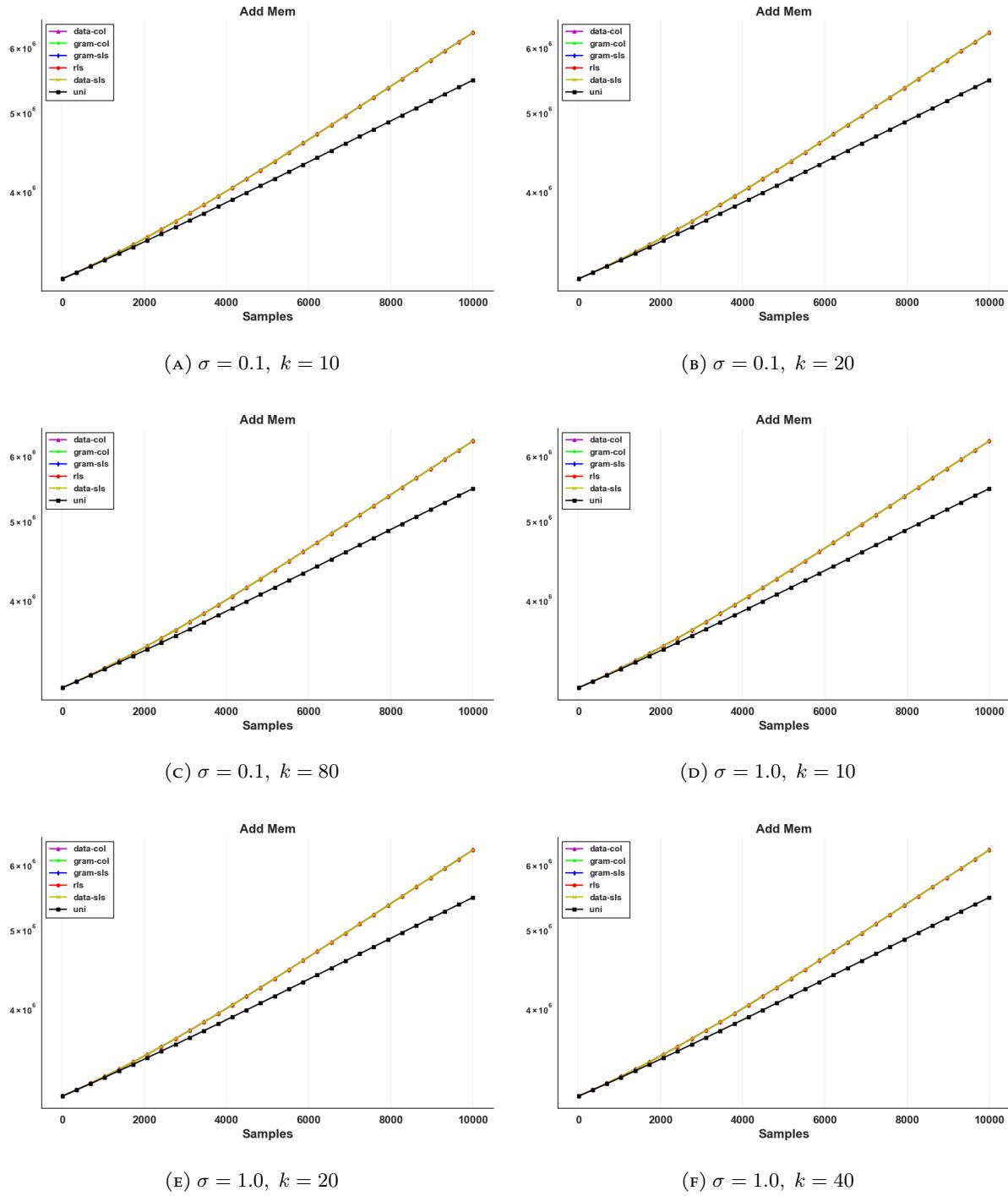


FIGURE 57. Nystrom sketch times for the 3D-spatial network data set.

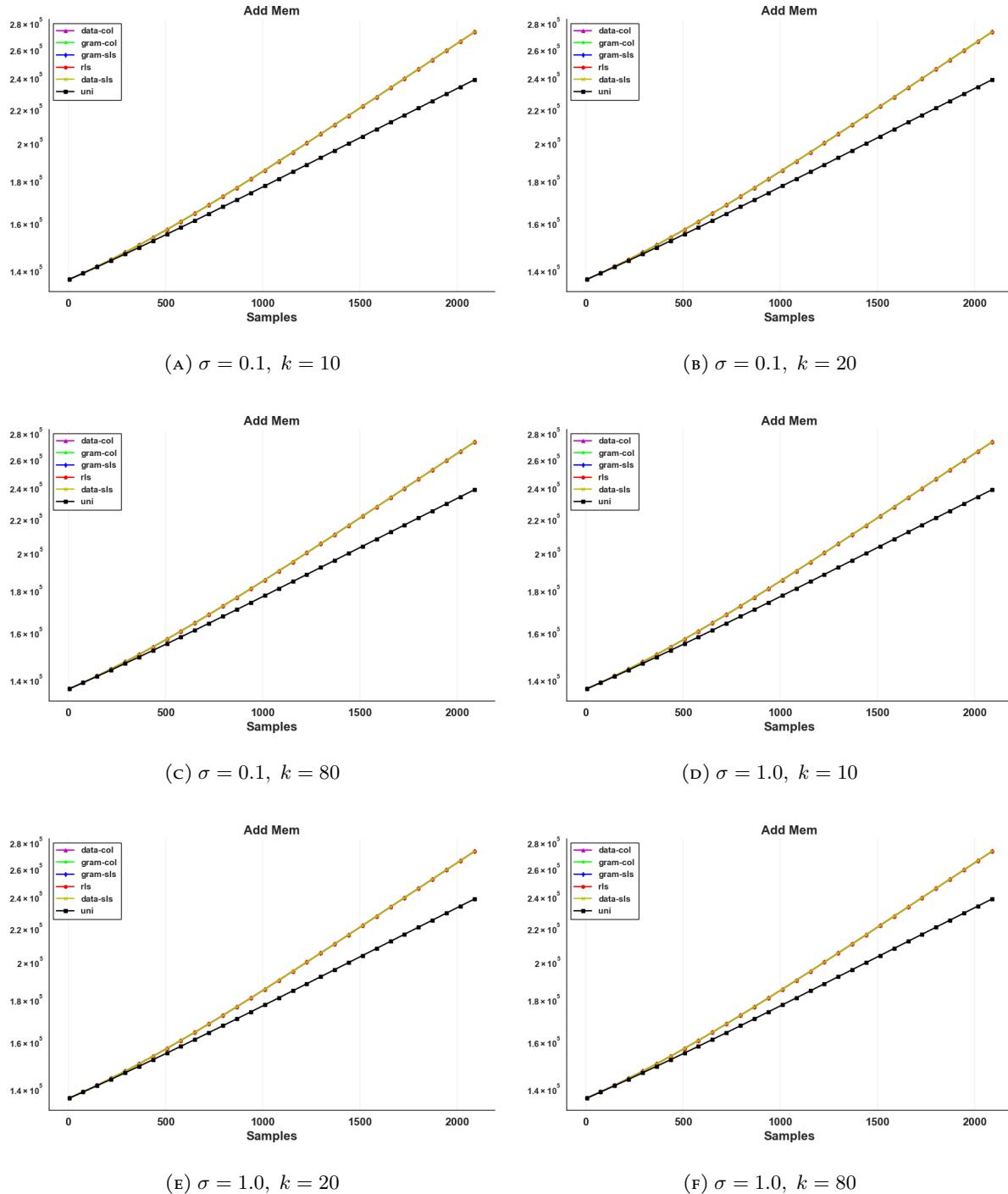


FIGURE 58. Nystrom sketch times for the Abalone data set.

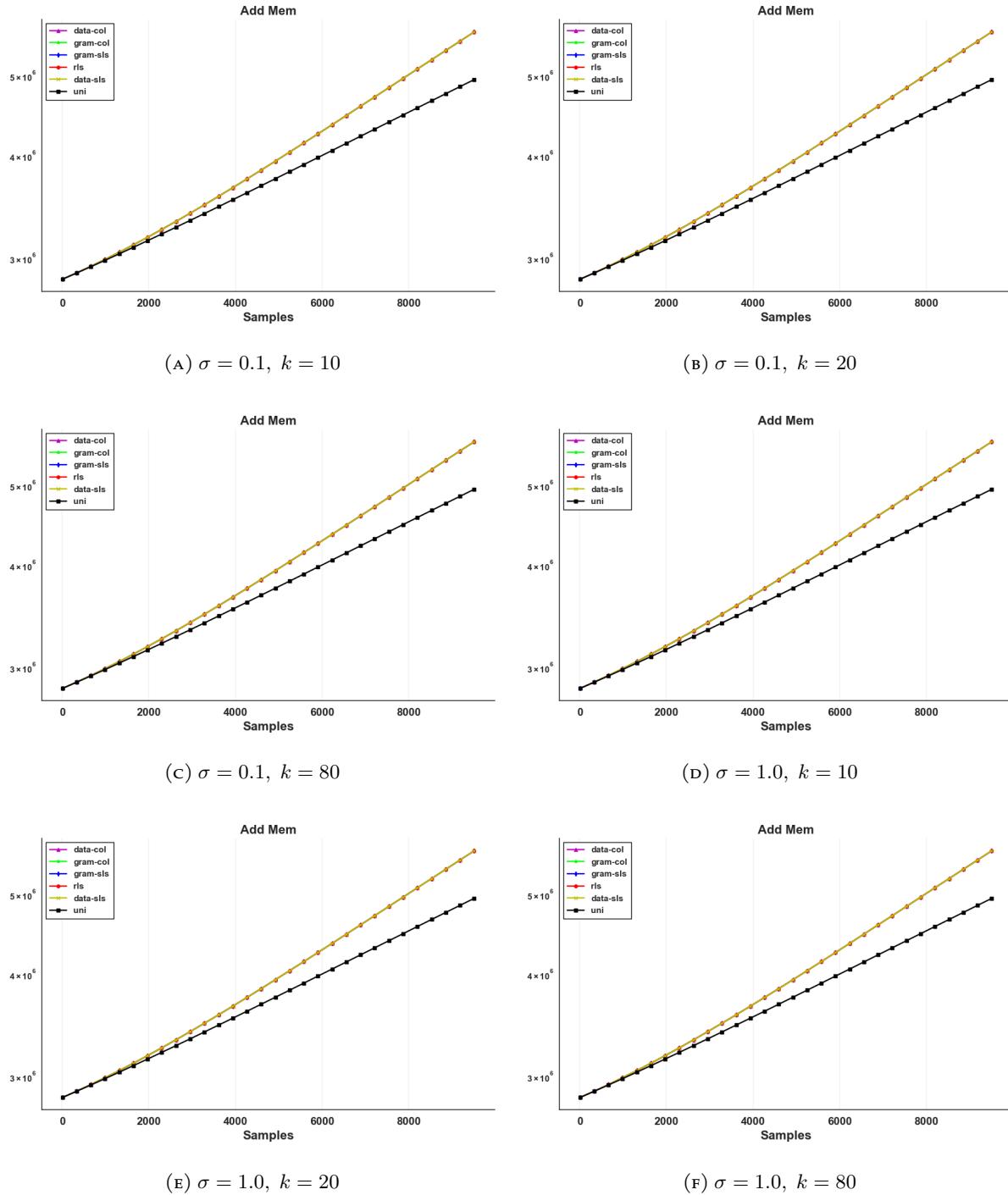


FIGURE 59. Nystrom sketch times for the Magic data set.

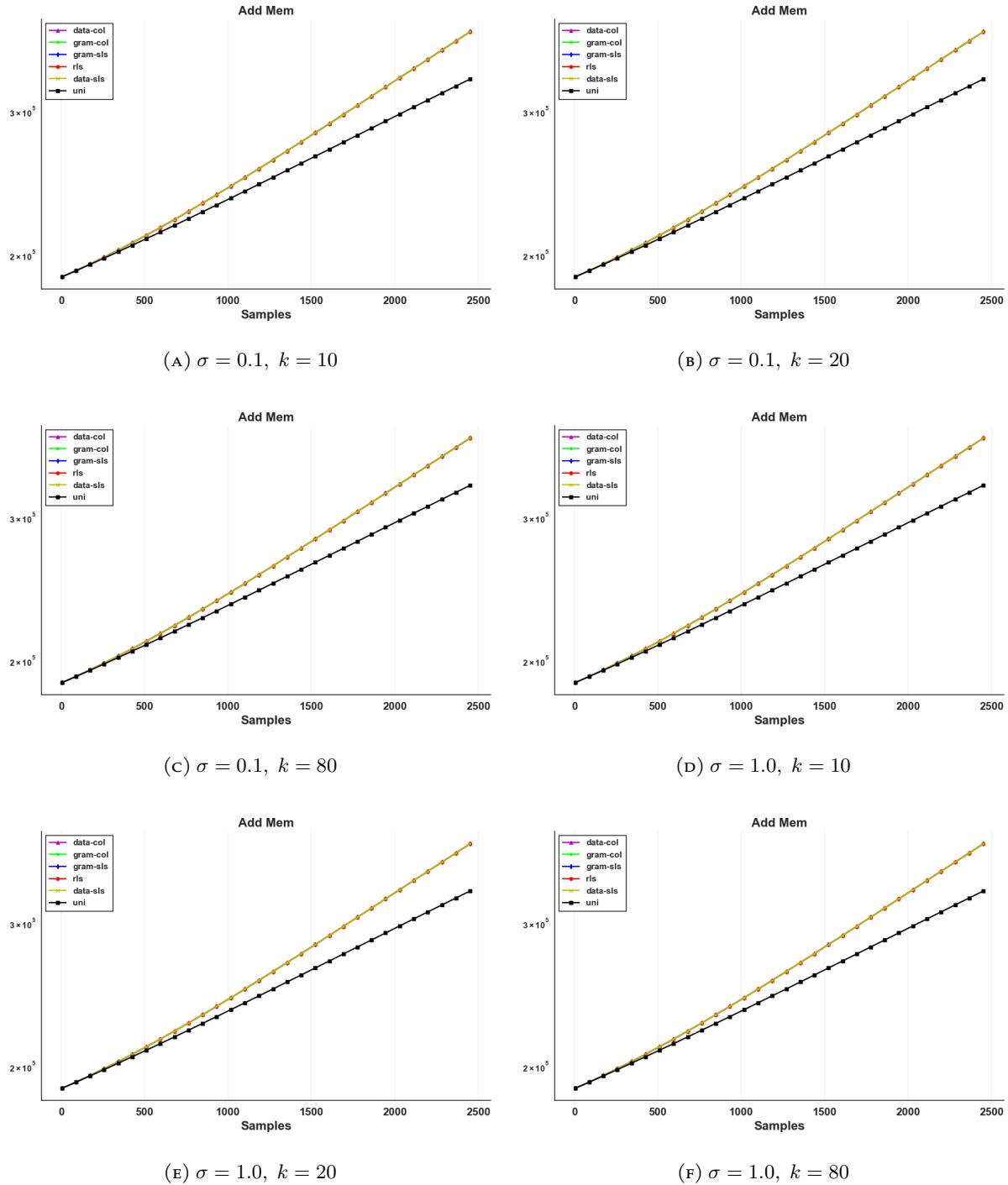


FIGURE 60. Nystrom sketch times for the stock market data set.

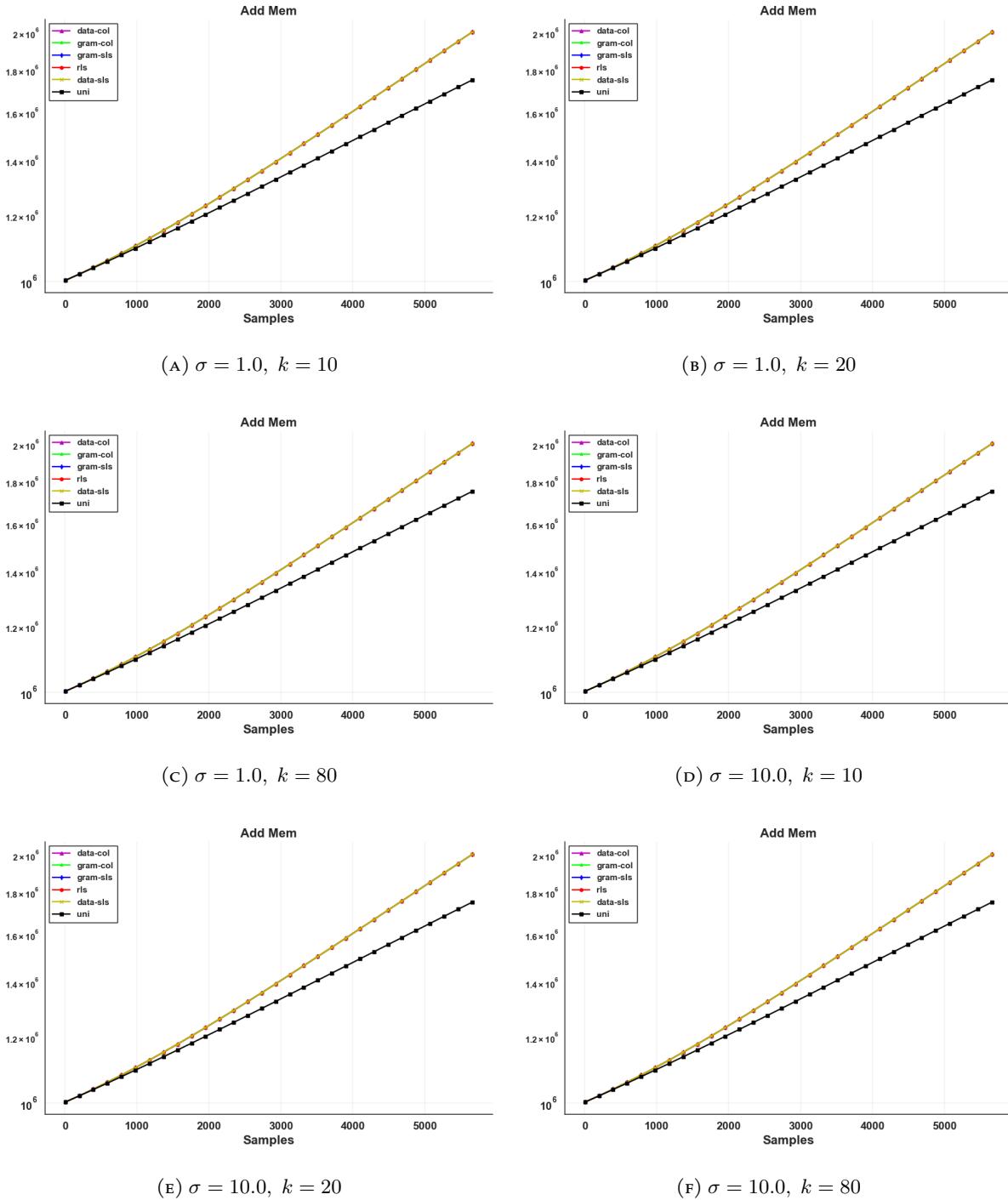


FIGURE 61. Nystrom sketch times for the temperature data set.

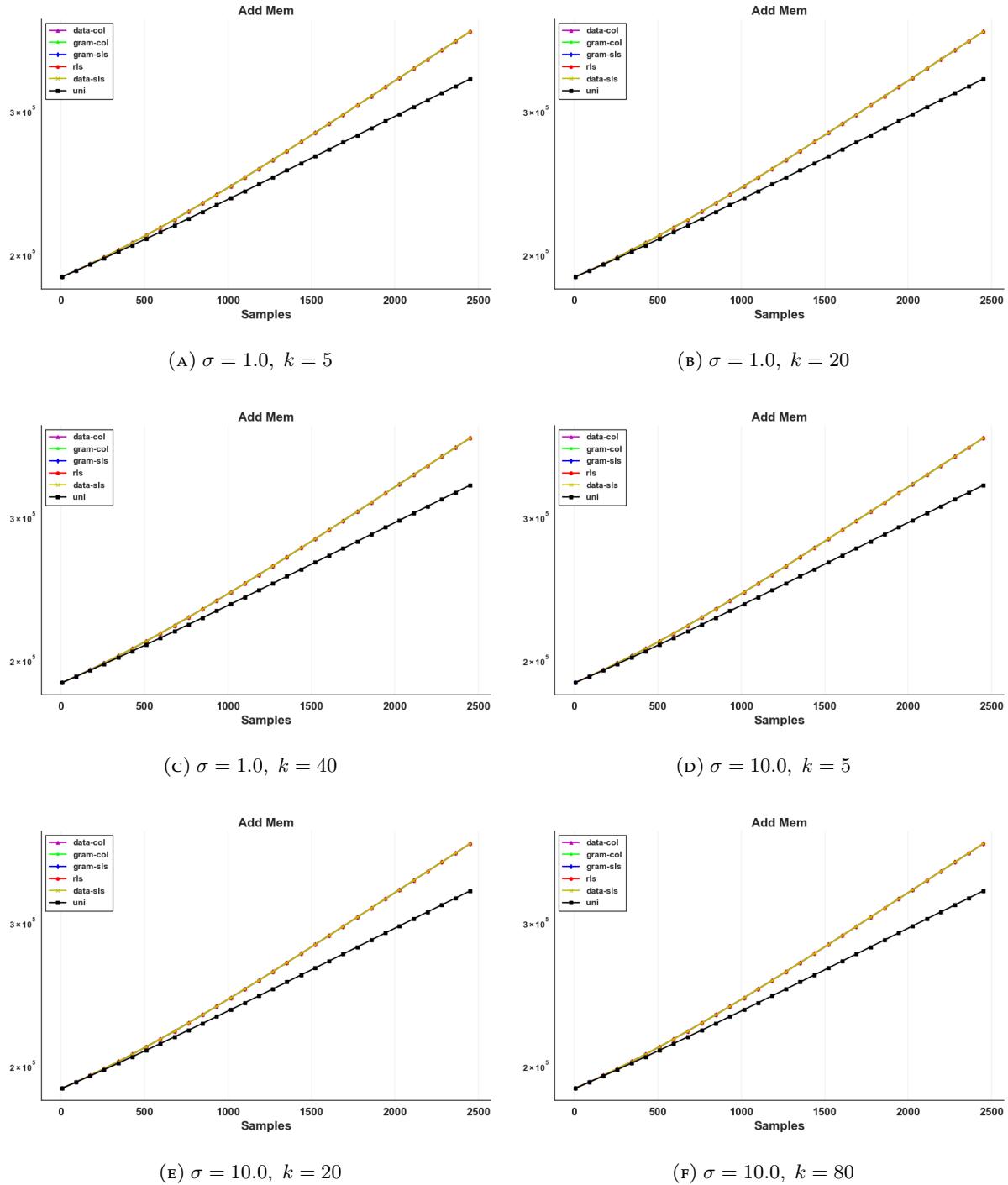


FIGURE 62. Nystrom sketch times for the wine data set.

A.8. RFF Errors.

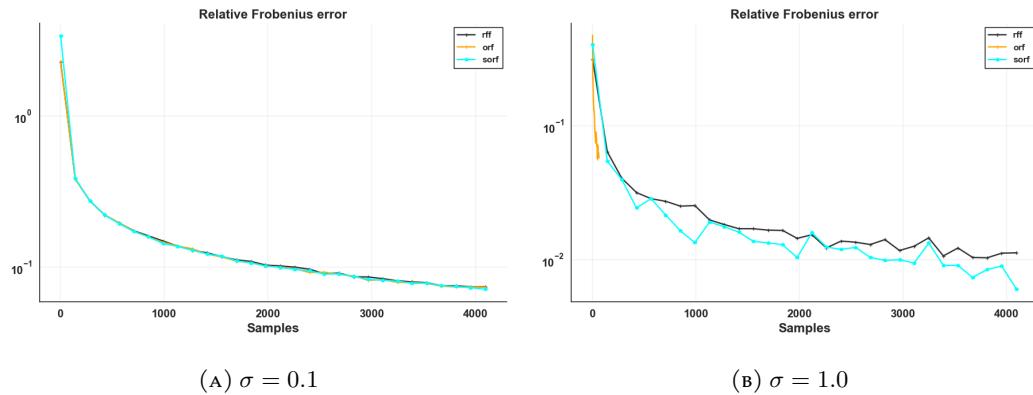


FIGURE 63. RFF Frobenius errors for the 3DSN data set.

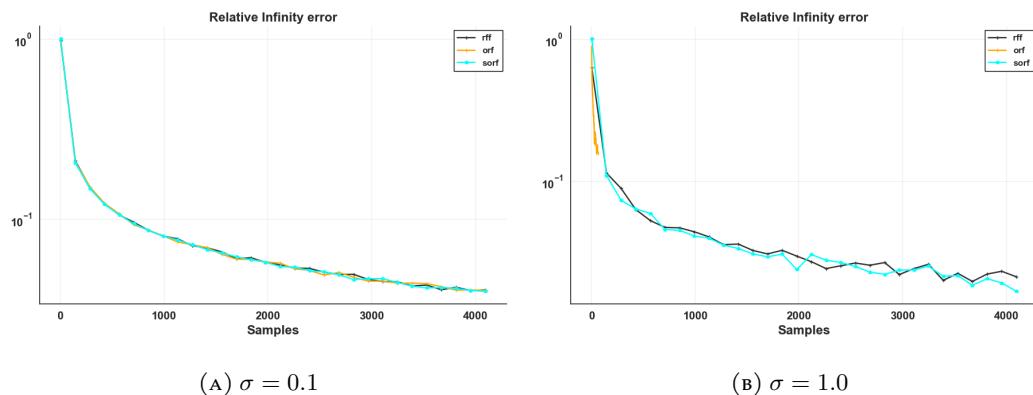


FIGURE 64. RFF absolute errors for the 3DSN data set.

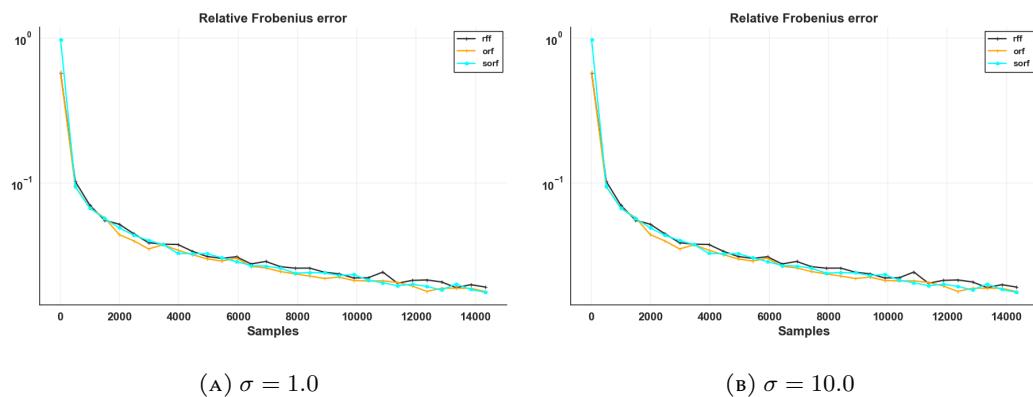


FIGURE 65. RFF Frobenius errors for the Abalone data set.

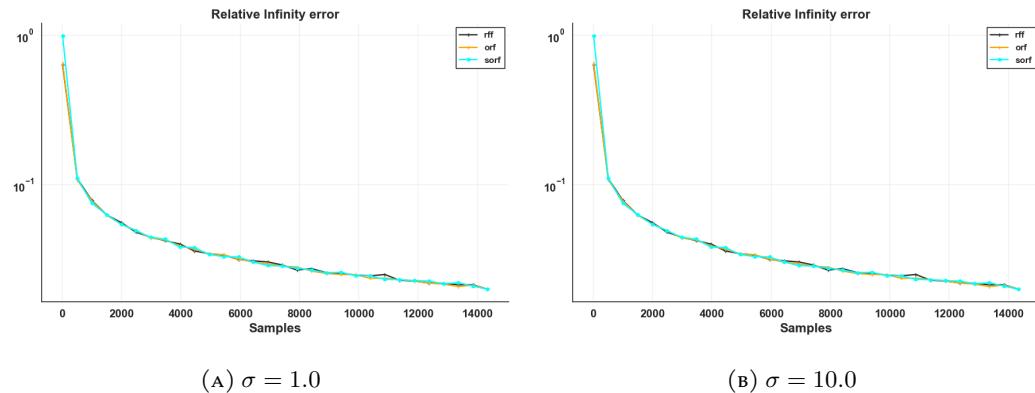


FIGURE 66. RFF infinity errors for the Abalone data set.

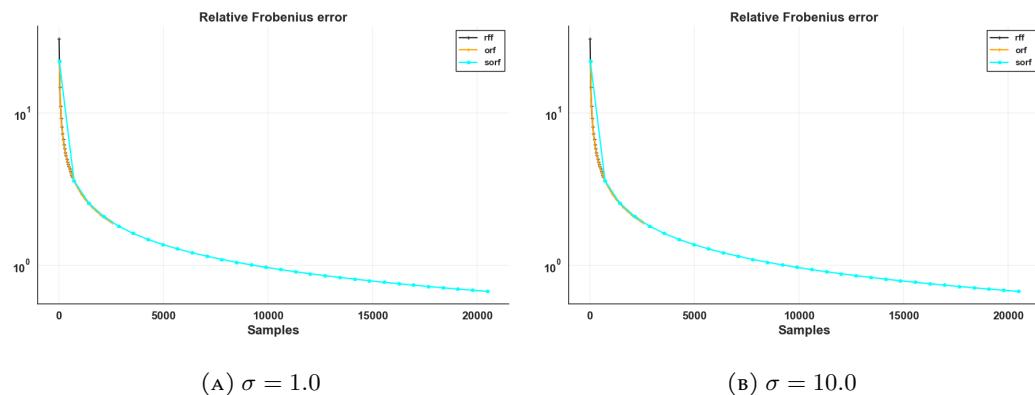


FIGURE 67. RFF Frobenius errors for the magic04 data set.

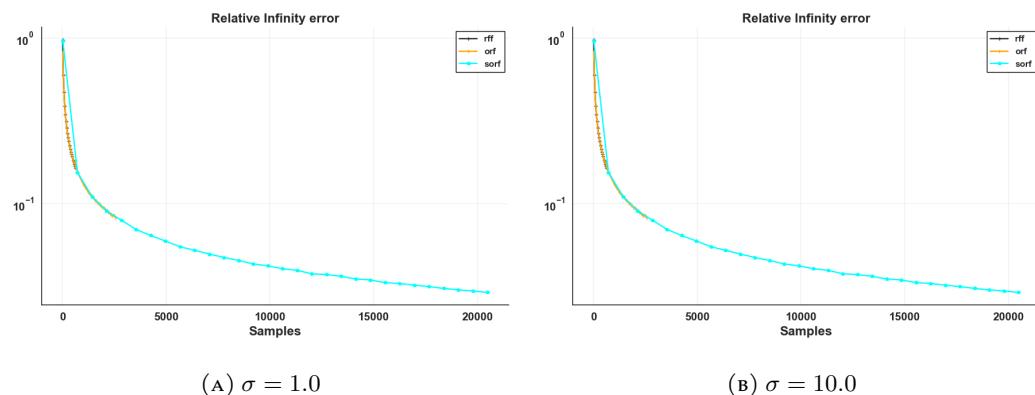


FIGURE 68. RFF infinity errors for the magic04 data set.

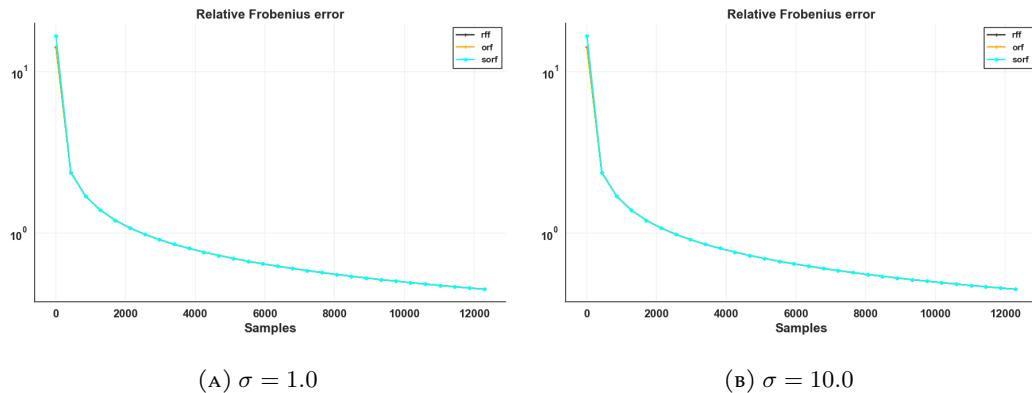


FIGURE 69. RFF Frobenius errors for the Stocks data set.

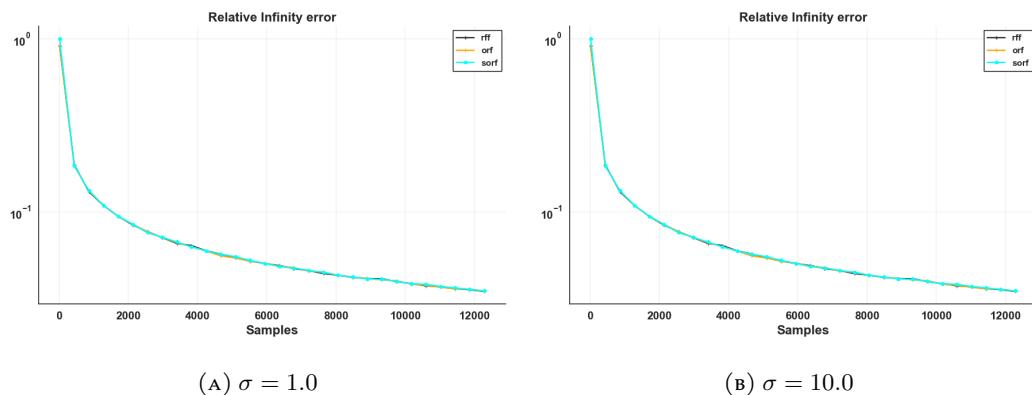


FIGURE 70. RFF infinity errors for the Stocks data set.

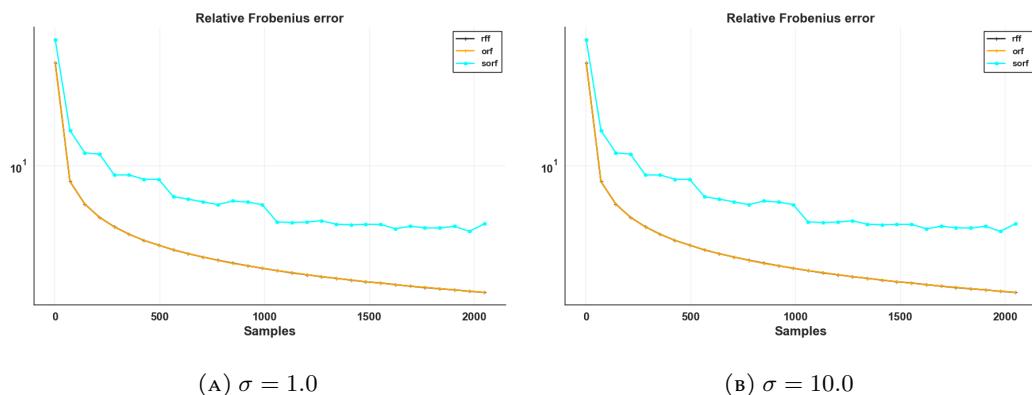


FIGURE 71. RFF Frobenius errors for the Temp data set.

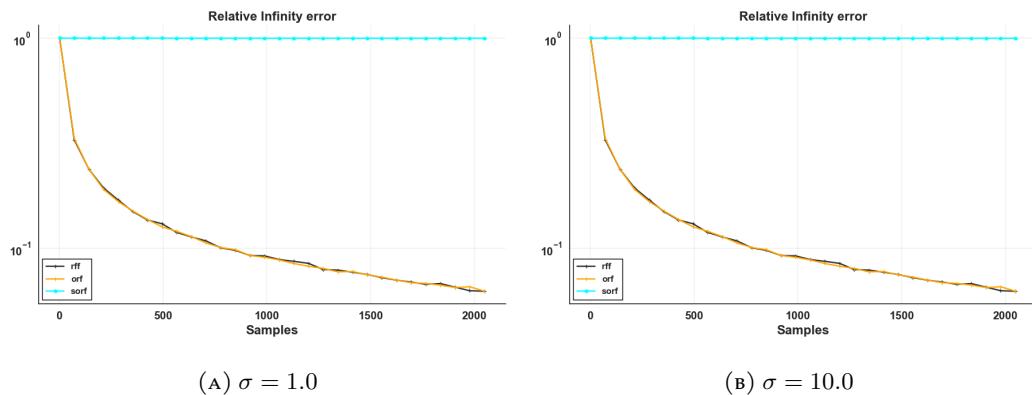


FIGURE 72. RFF infinity errors for the Temp data set.

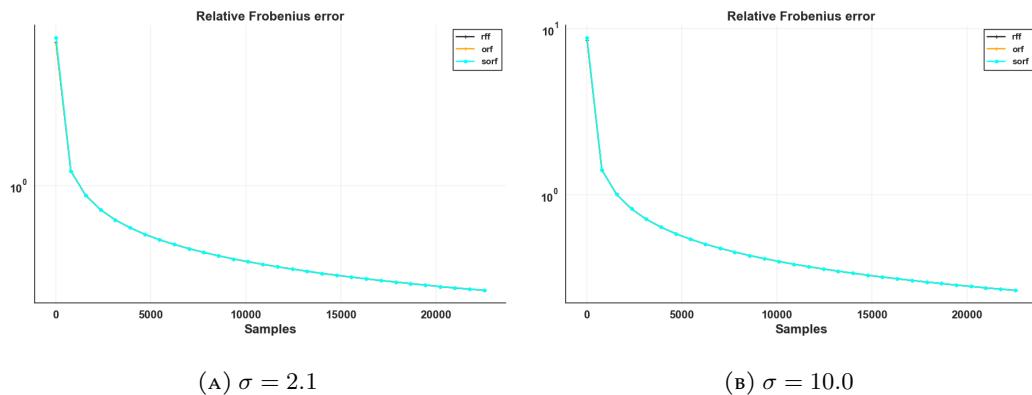


FIGURE 73. RFF Frobenius errors for the Wine data set.

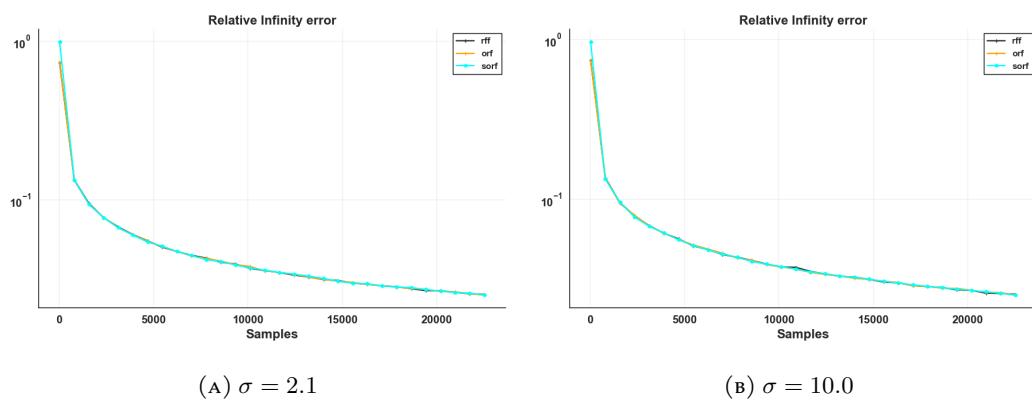


FIGURE 74. RFF infinity errors for the Wine data set.

A.9. RFF Times.

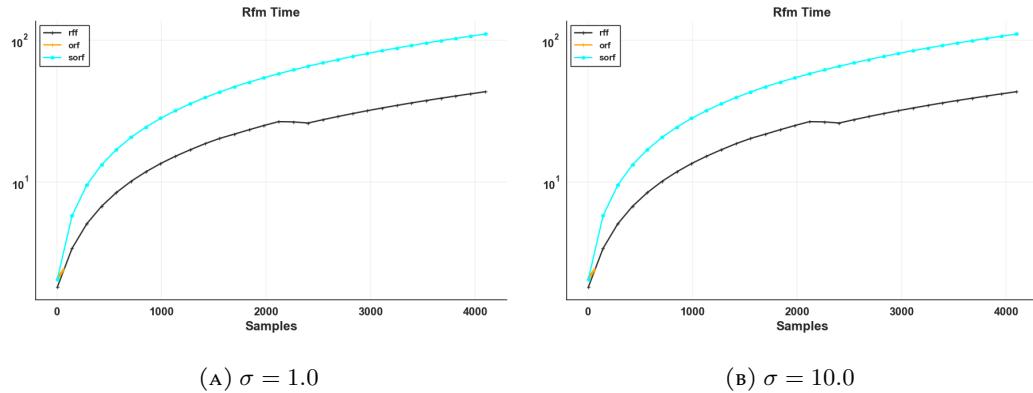


FIGURE 75. RFF approximation times for the 3DSN data set.

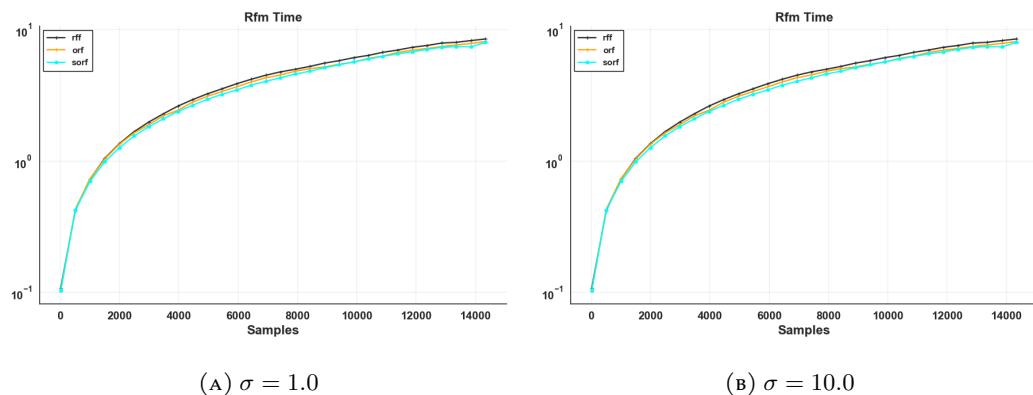


FIGURE 76. RFF approximation times for the Abalone data set.

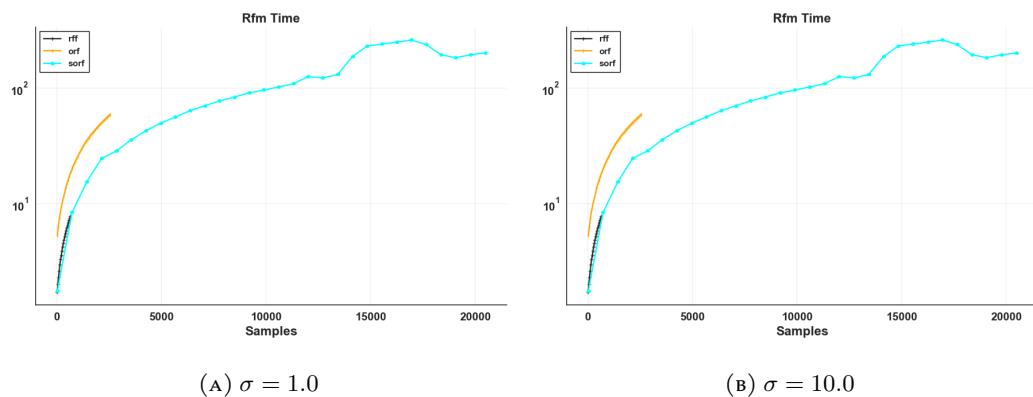


FIGURE 77. RFF approximation times for the magic04 data set.

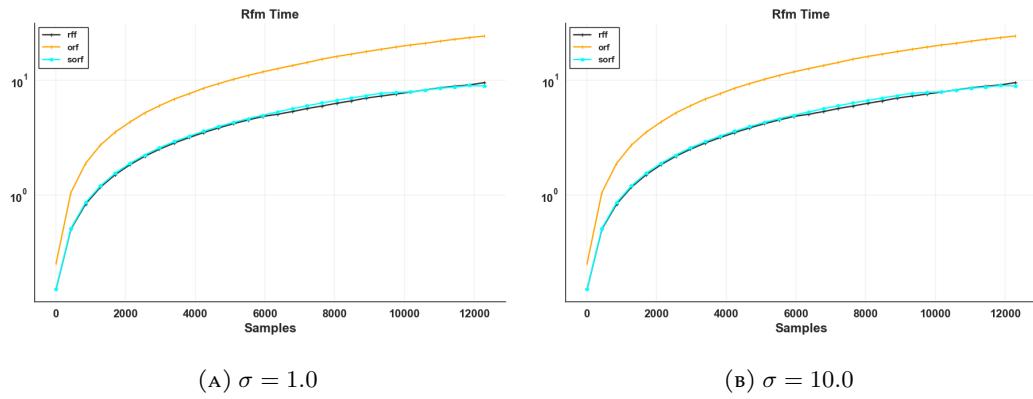


FIGURE 78. RFF approximation times for the Stocks data set.

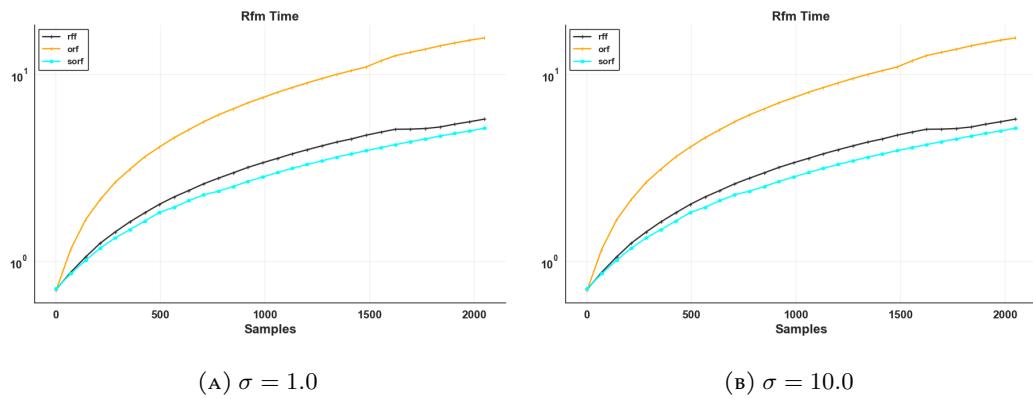


FIGURE 79. RFF approximation times for the Temp data set.

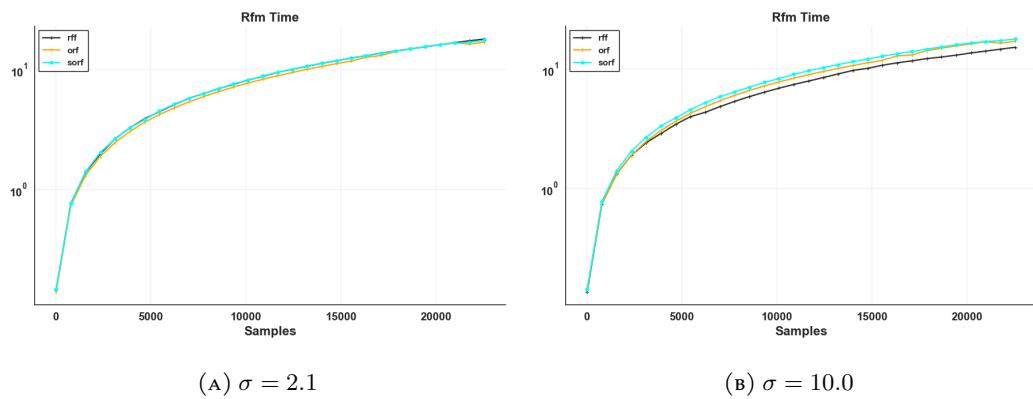


FIGURE 8o. RFF approximation times for the Wine data set.

A.10. RFF Additional Memory Usage.

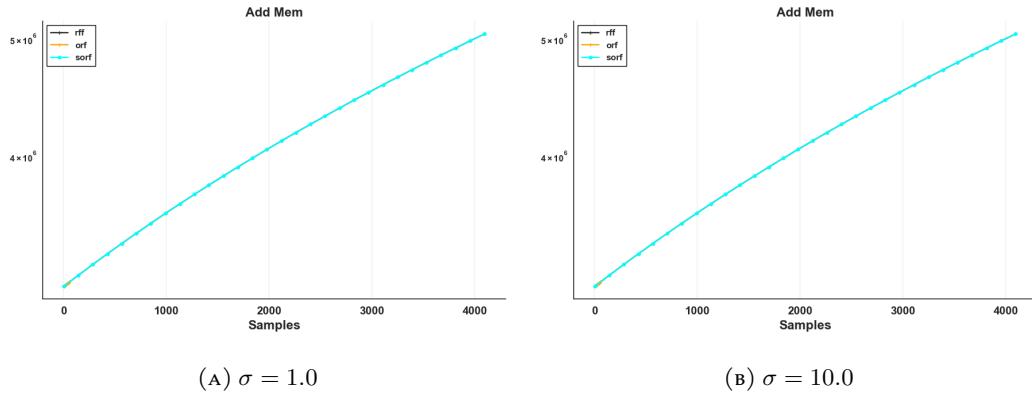


FIGURE 81. RFF additional memory usage for the 3DSN data set.

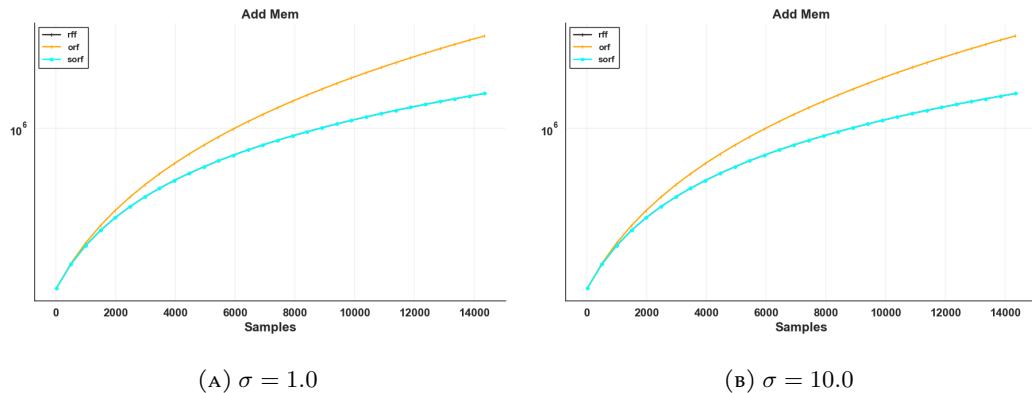


FIGURE 82. RFF additional memory usage for the Abalone data set.

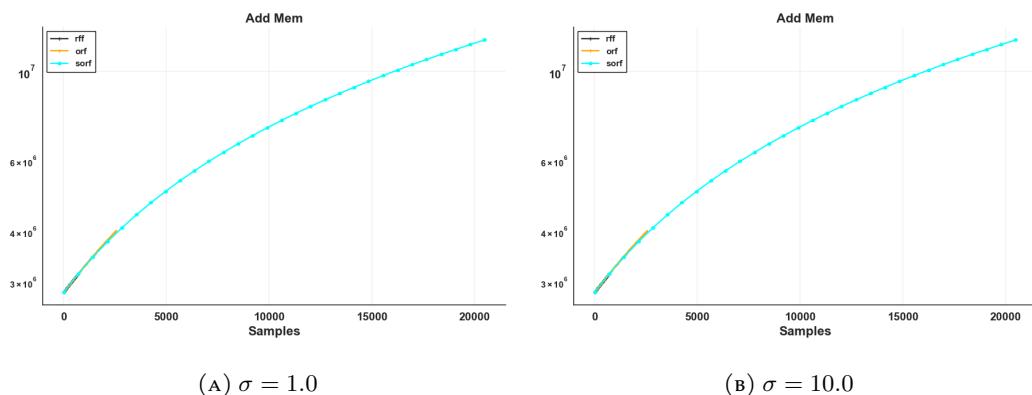


FIGURE 83. RFF additional memory usage for the magic04 data set.

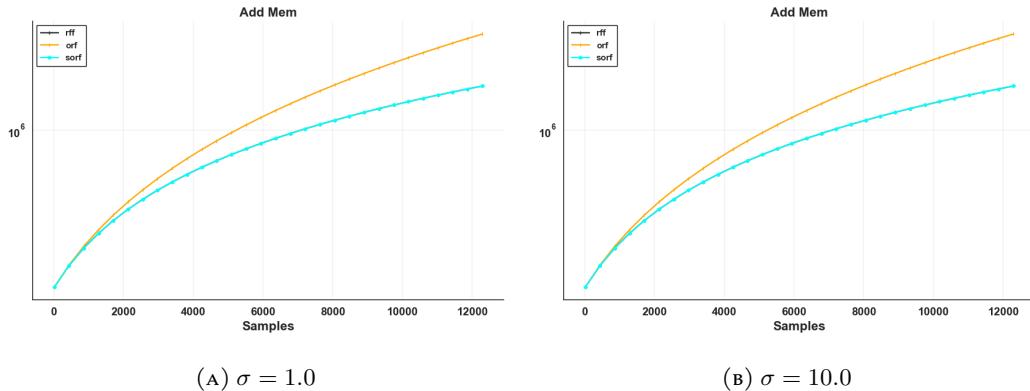


FIGURE 84. RFF additional memory usage for the Stocks data set.

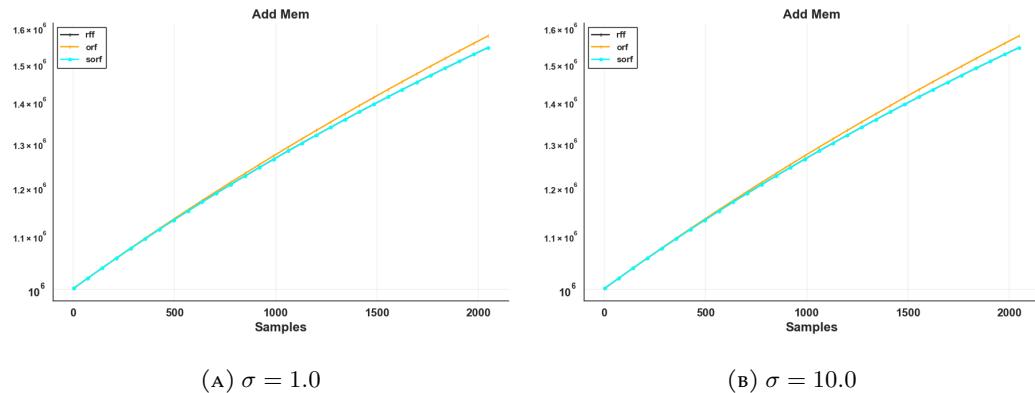


FIGURE 85. RFF additional memory usage for the Temp data set.

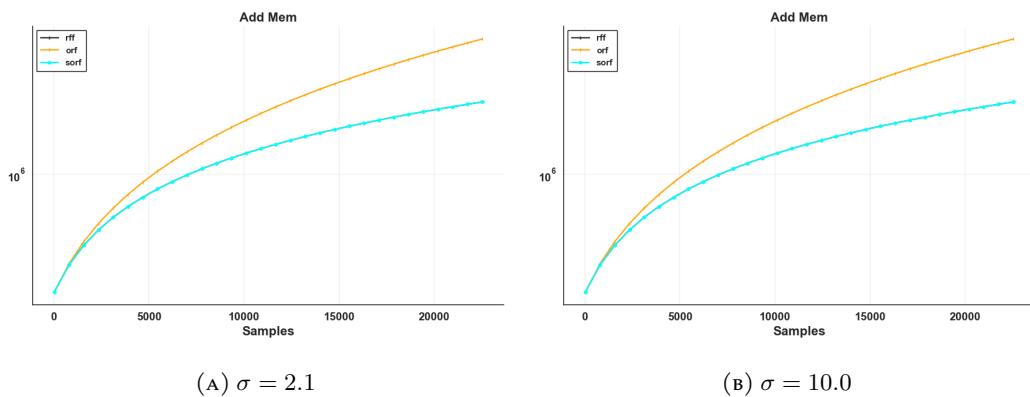


FIGURE 86. RFF additional memory usage for the Wine data set.

A.11. Kernel Comparisons.

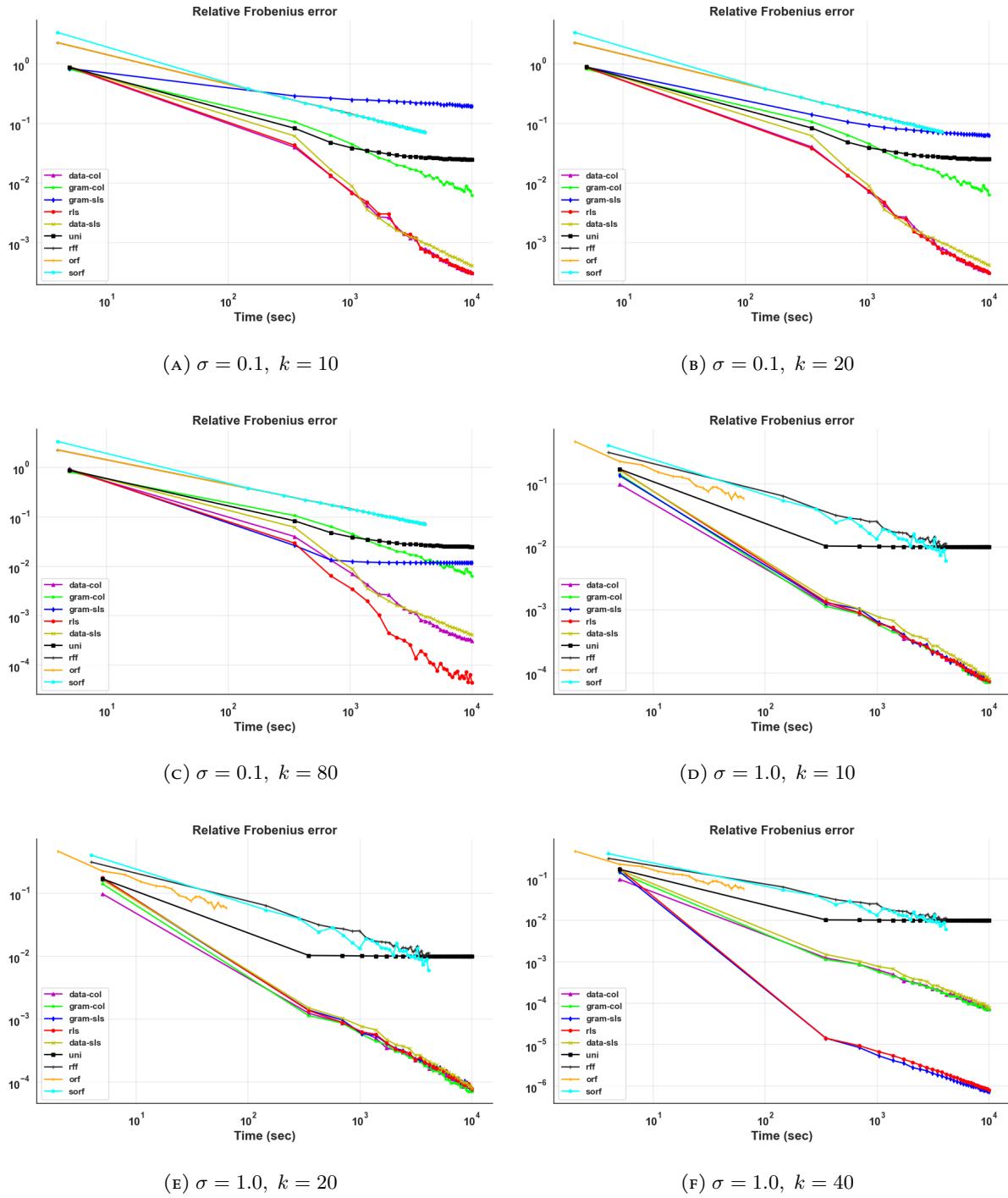


FIGURE 87. Comparing Frobenius errors of kernel methods for the 3DSN data set.

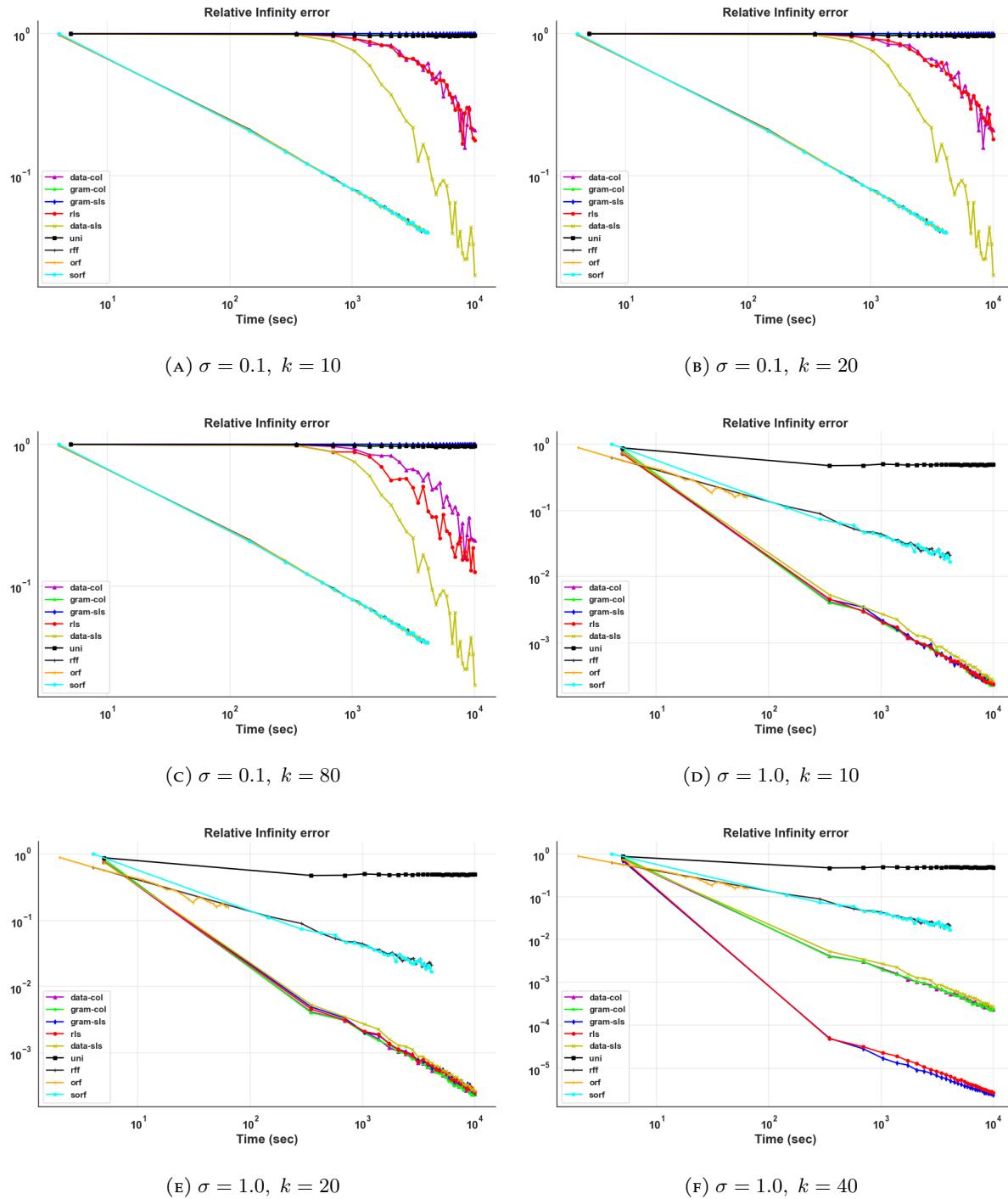


FIGURE 88. Comparing infinity errors of kernel methods for the 3DSN data set.

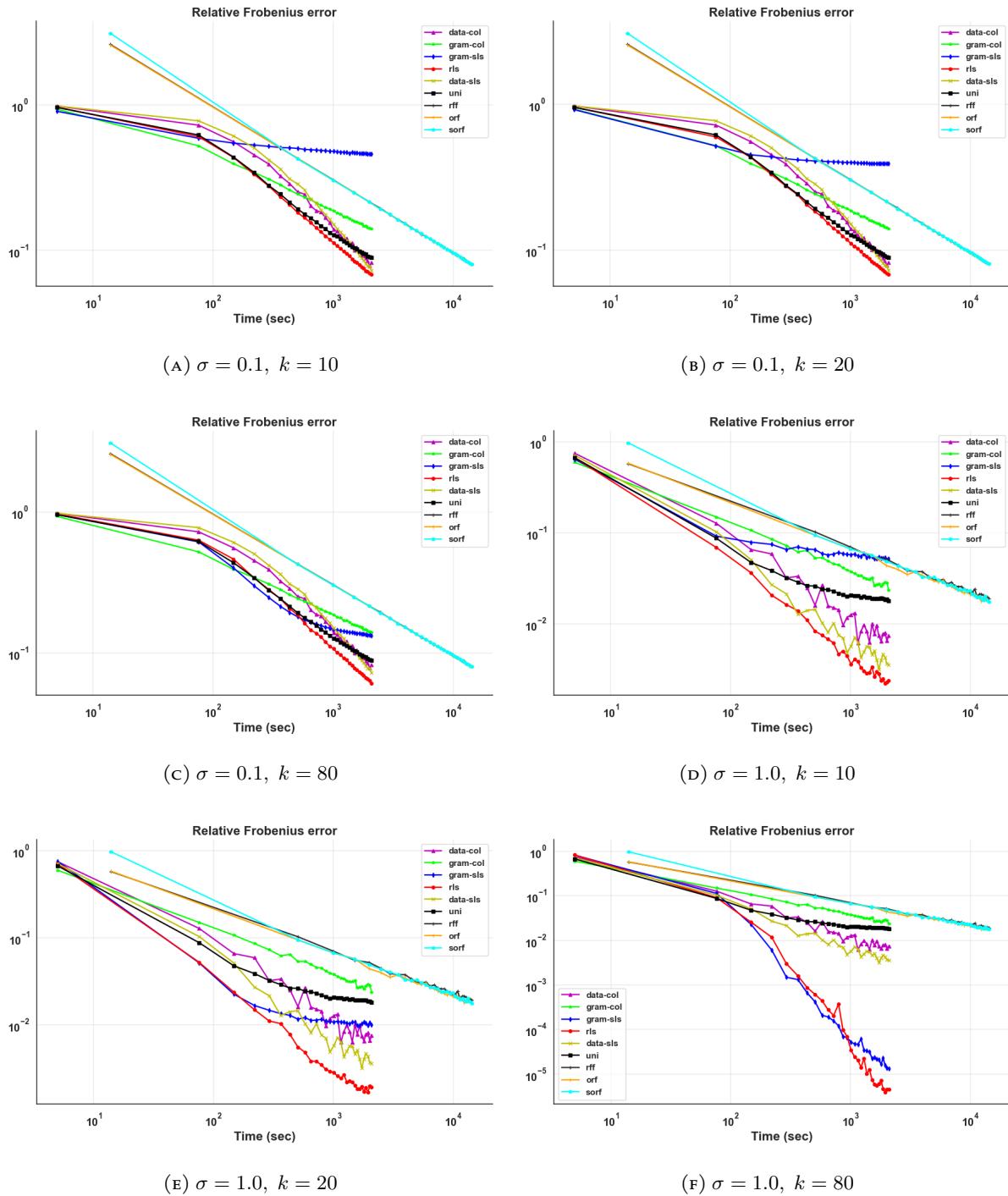


FIGURE 89. Comparing Frobenius errors of kernel methods for the Abalone data set.

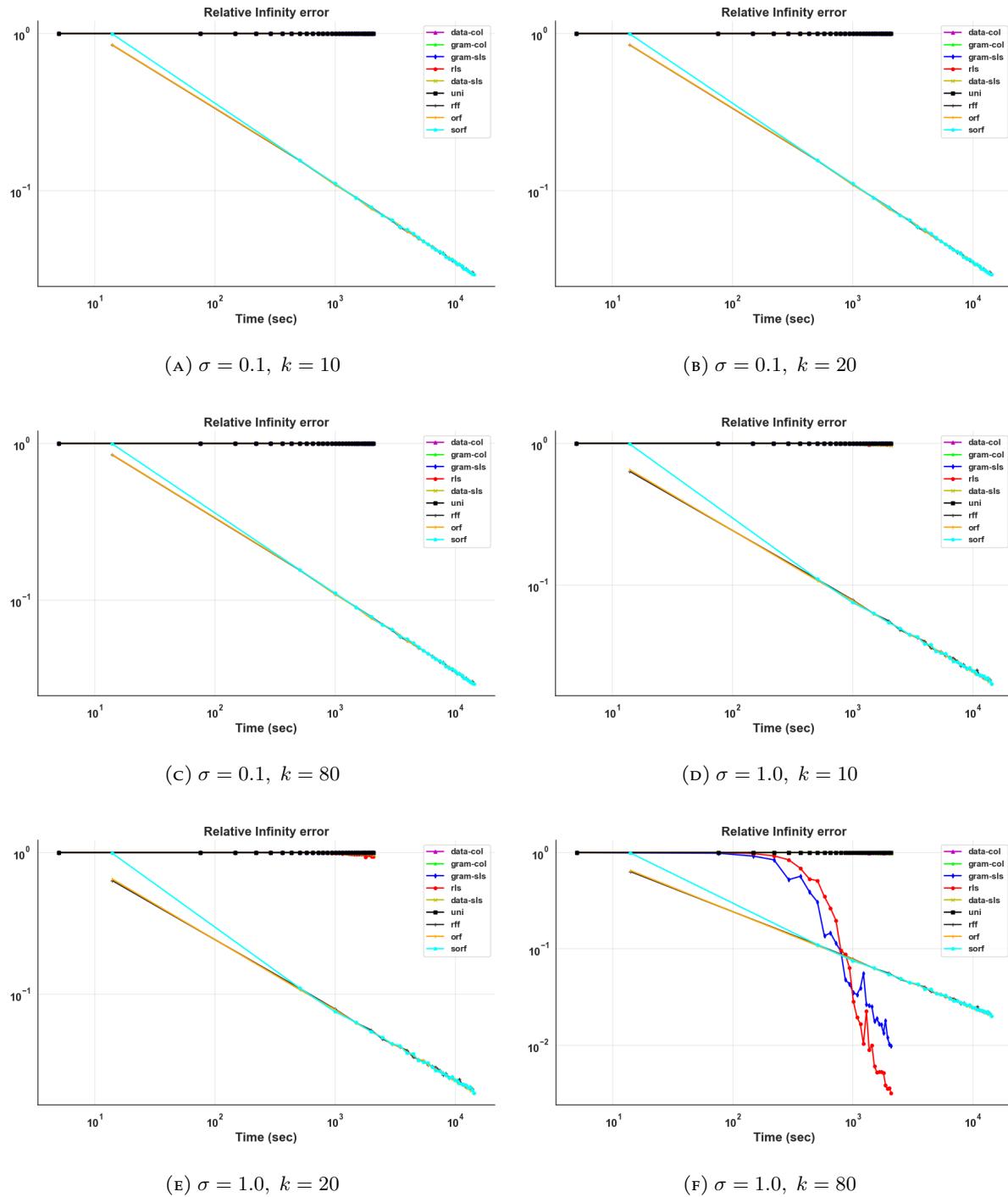


FIGURE 90. Comparing infinity errors of kernel methods for the Magic data set.

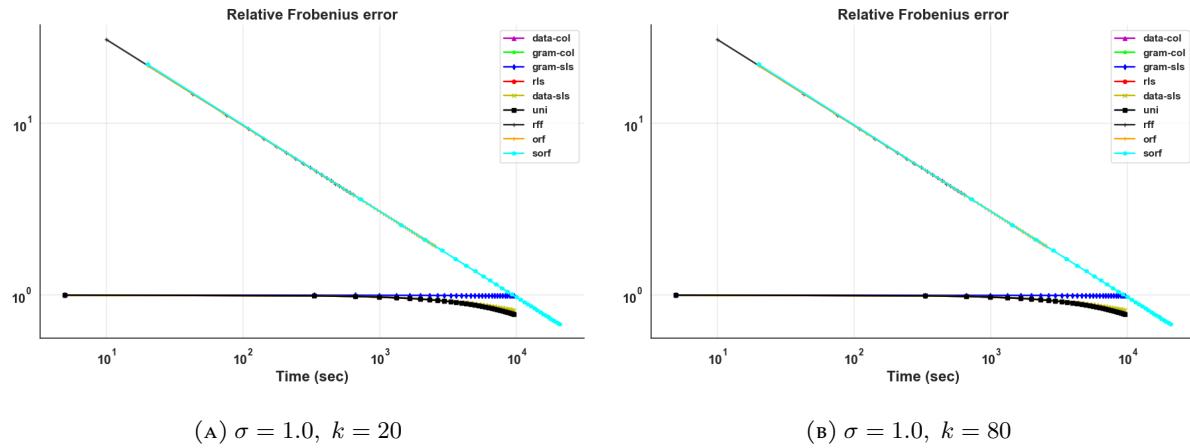


FIGURE 91. Comparing Frobenius errors of kernel methods for the Magic data set.

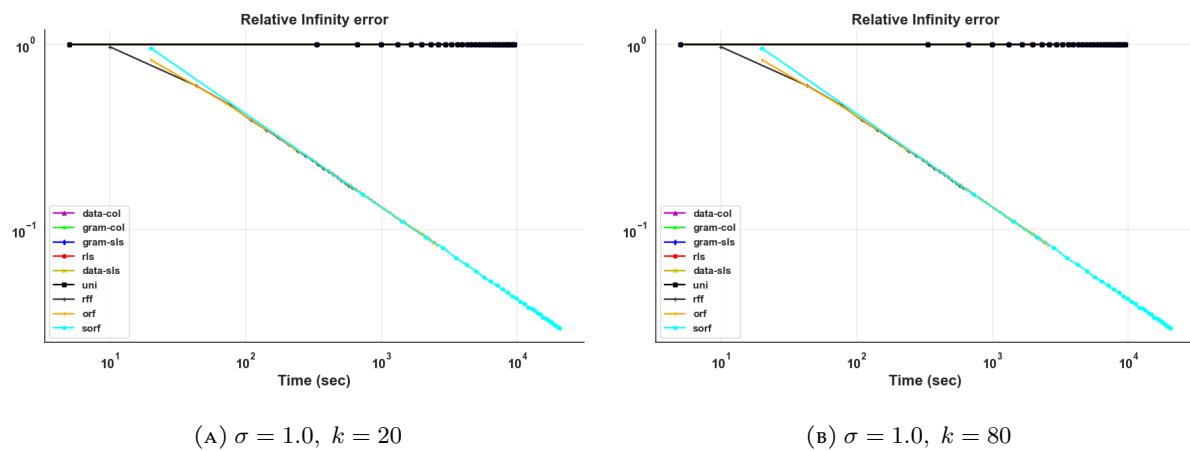


FIGURE 92. Comparing infinity errors of kernel methods for the Magic data set.

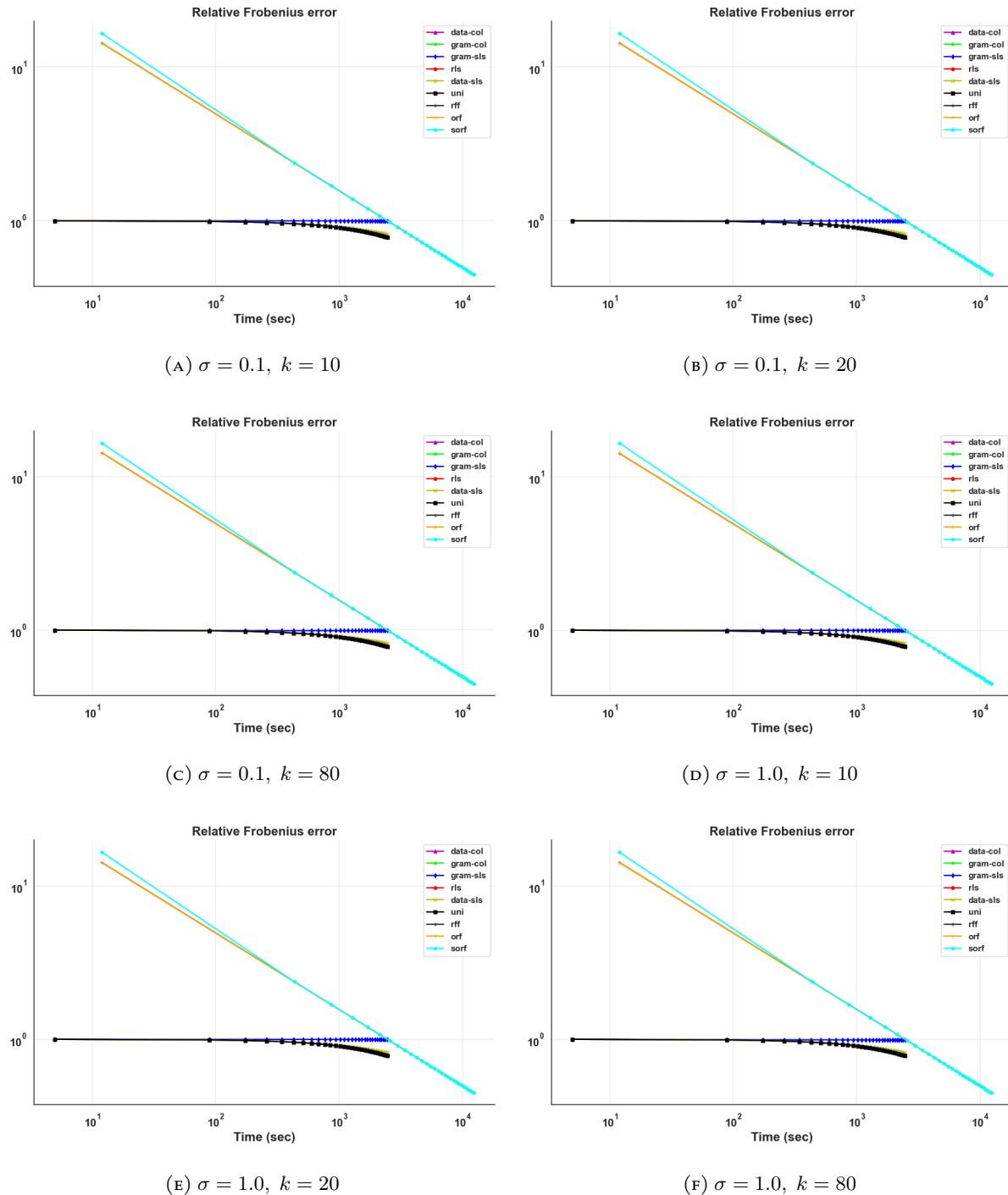


FIGURE 93. Comparing Frobenius errors of kernel methods for the Stocks data set.

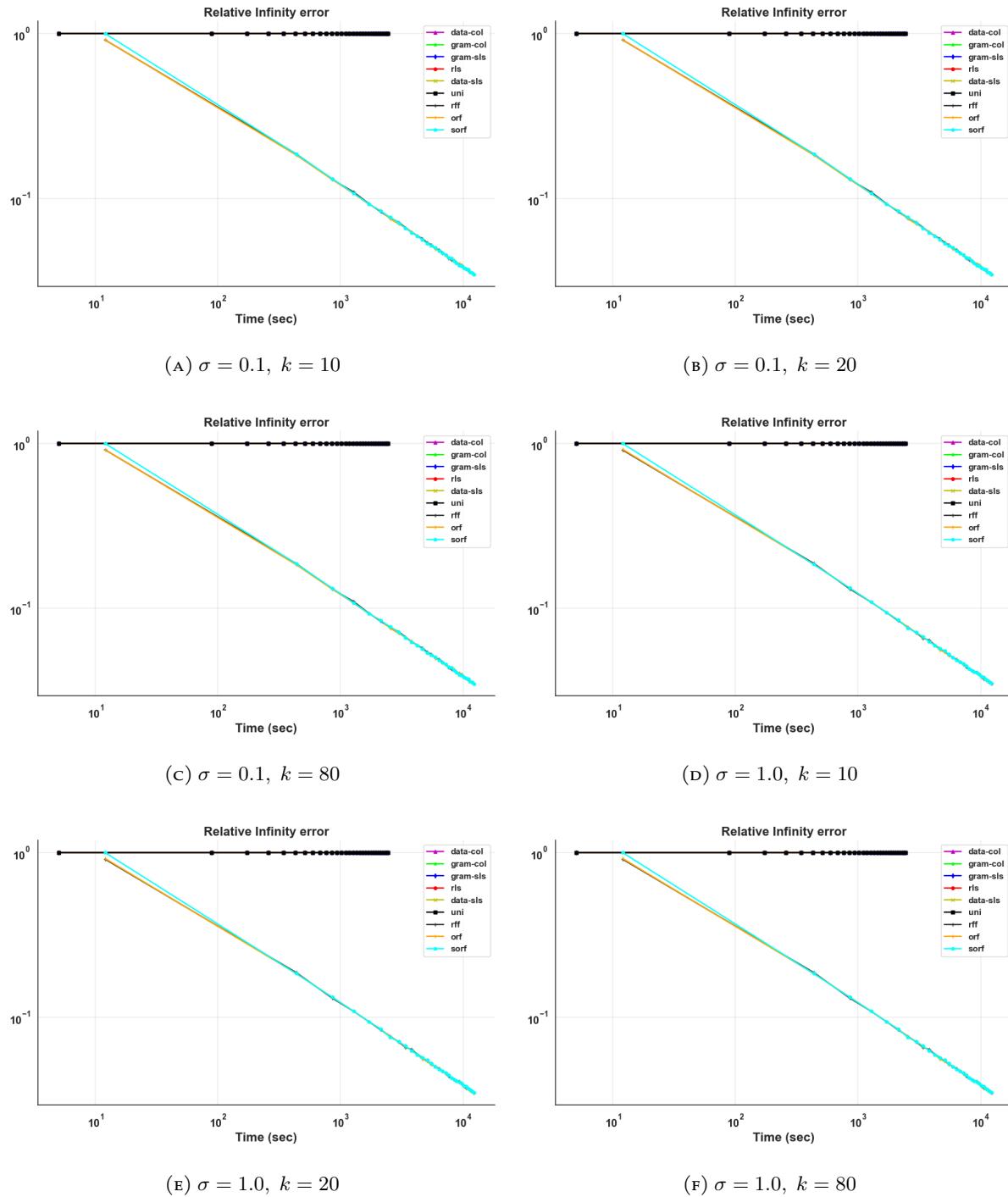


FIGURE 94. Comparing infinity errors of kernel methods for the Stocks data set.

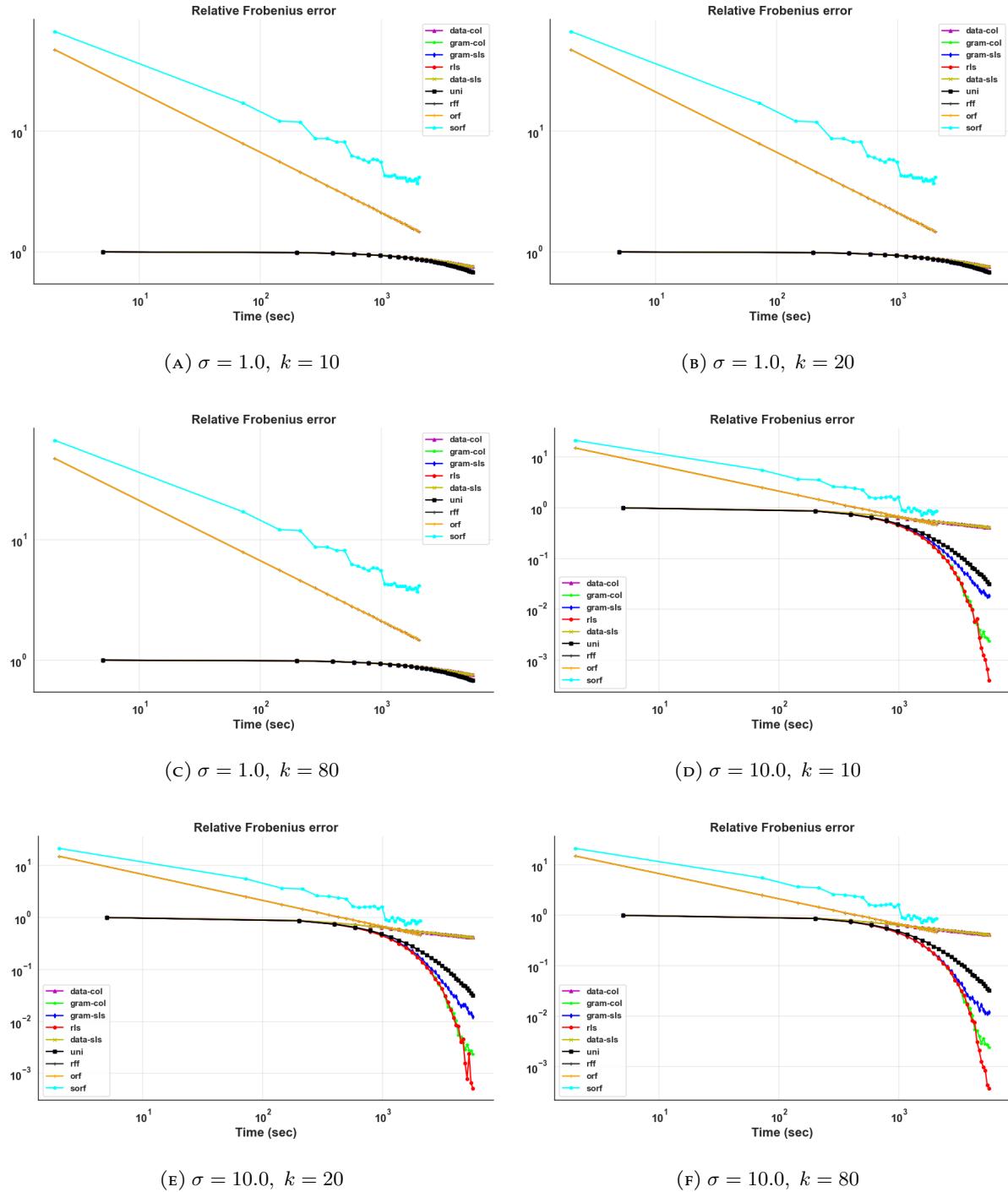


FIGURE 95. Comparing Frobenius errors of kernel methods for the Temp data set.

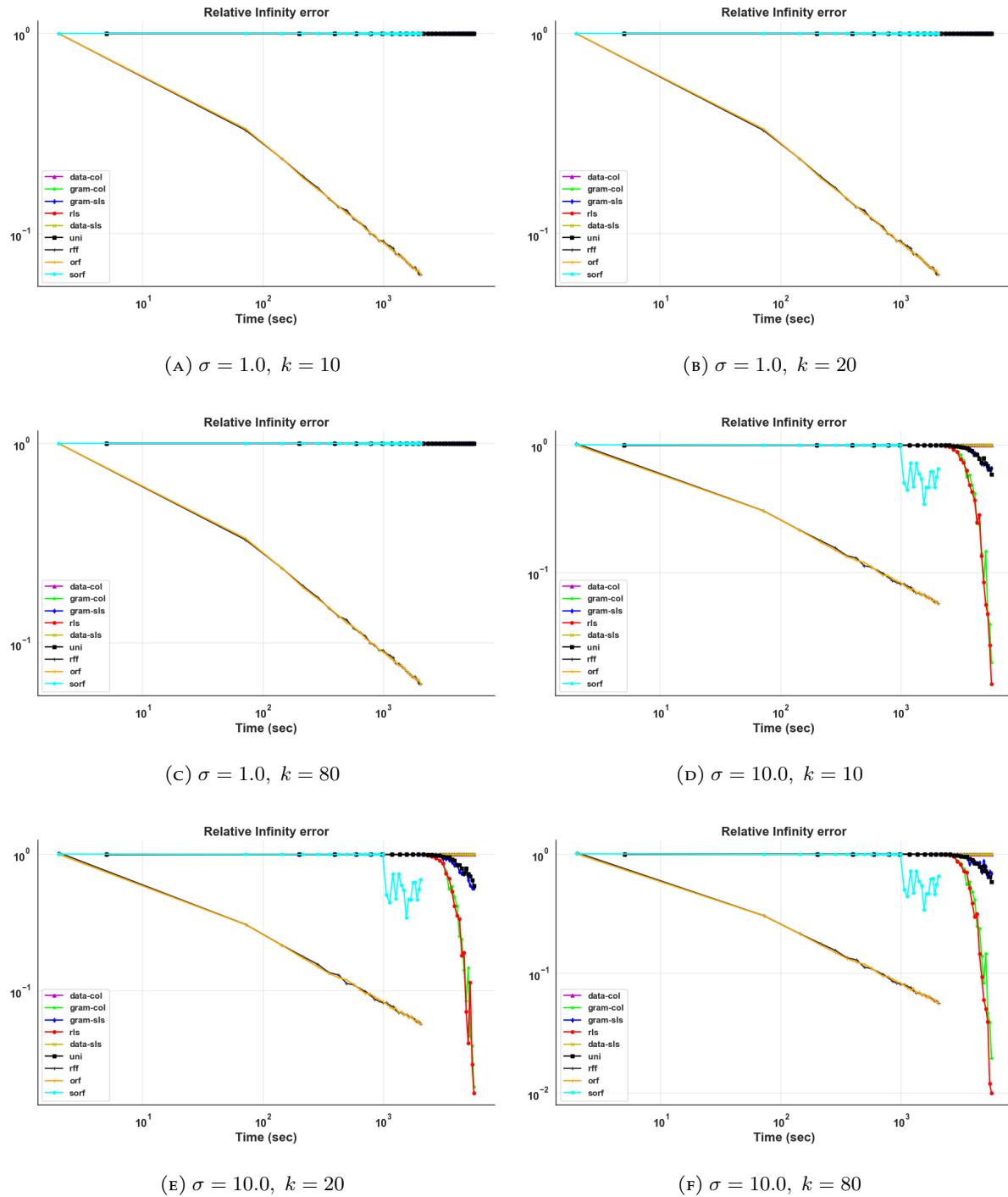


FIGURE 96. Comparing infinity errors of kernel methods for the Temp data set.

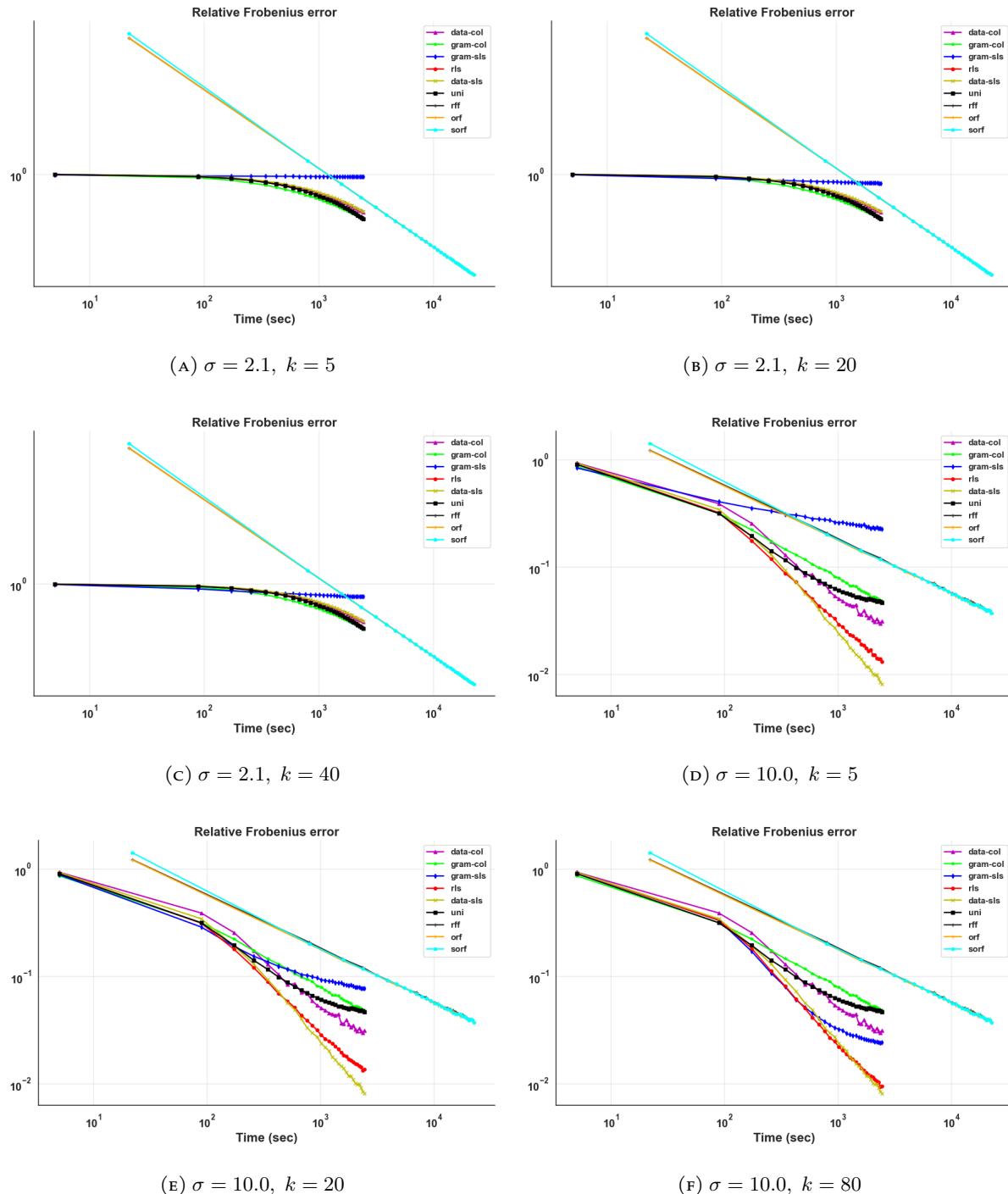


FIGURE 97. Comparing Frobenius errors of kernel methods for the wine data set.

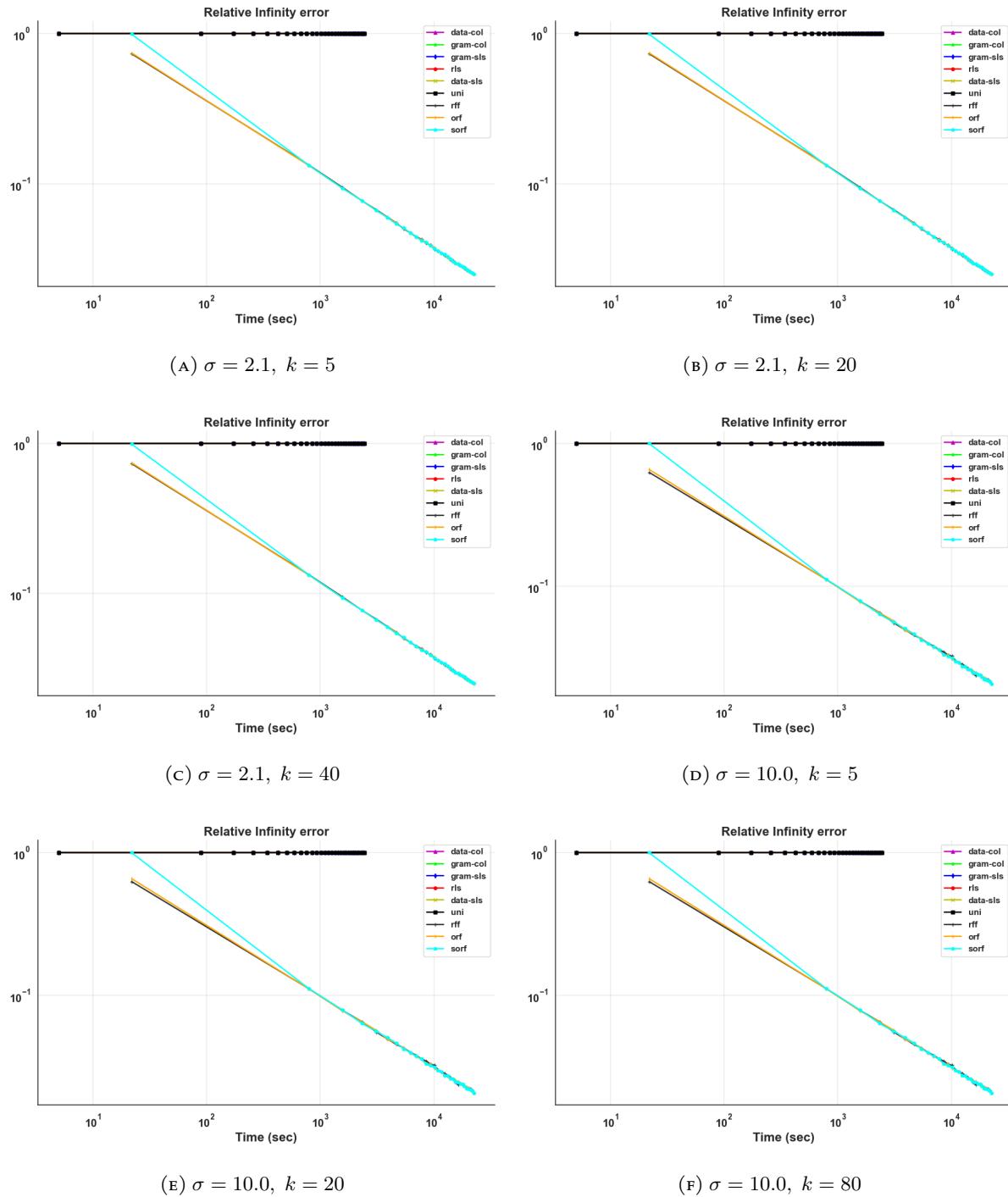


FIGURE 98. Comparing infinity errors of kernel methods for the wine data set.

A.12. prediction errors.

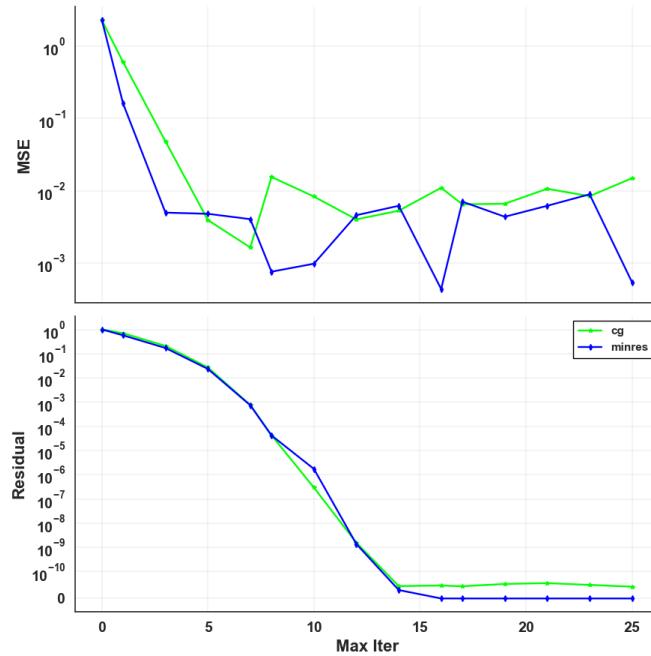


FIGURE 99. Comparison of prediction errors for the 3DSN data set, here $\sigma = 1.0$.

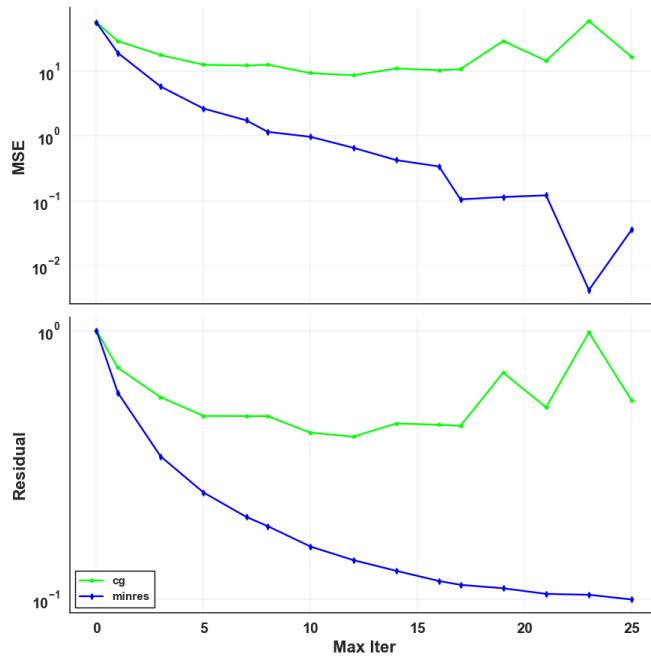
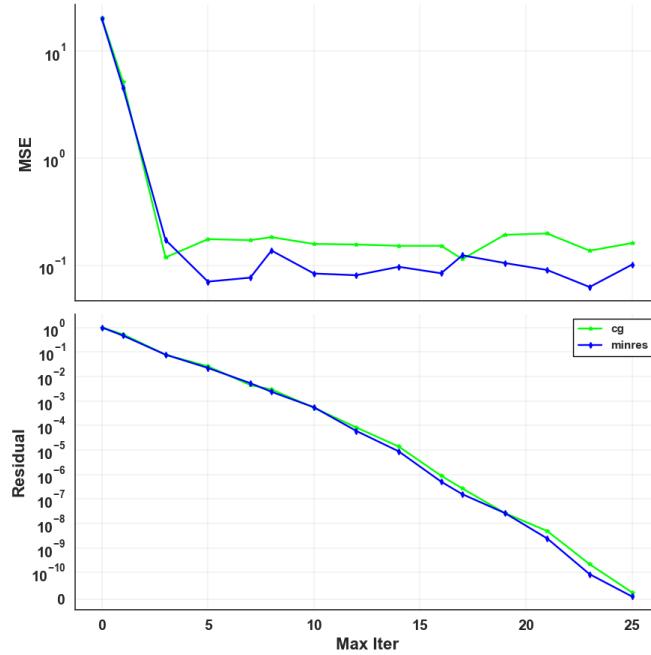
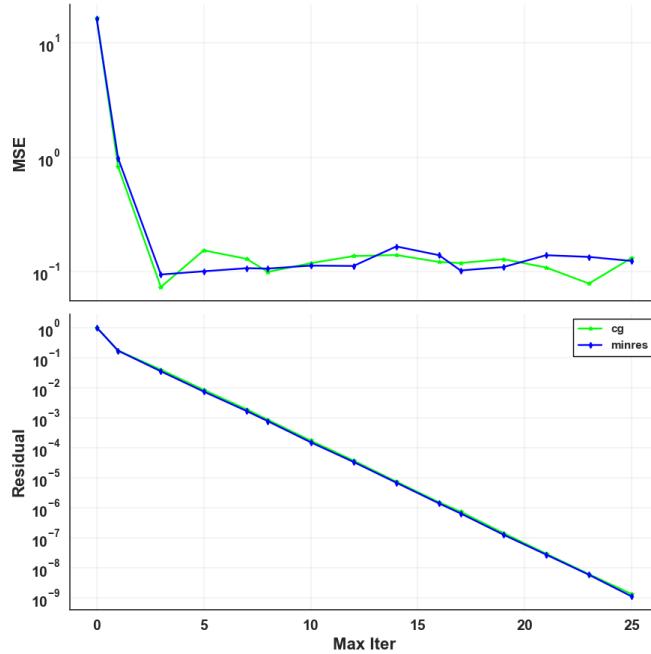


FIGURE 100. Comparison of prediction errors for the Abalone data set, here $\sigma = 1.0$.

FIGURE 101. Comparison of prediction errors for the quadratic data set, here $\sigma = 2.1$.FIGURE 102. Comparison of prediction errors for the Stock Market data set, here $\sigma = 10.0$.

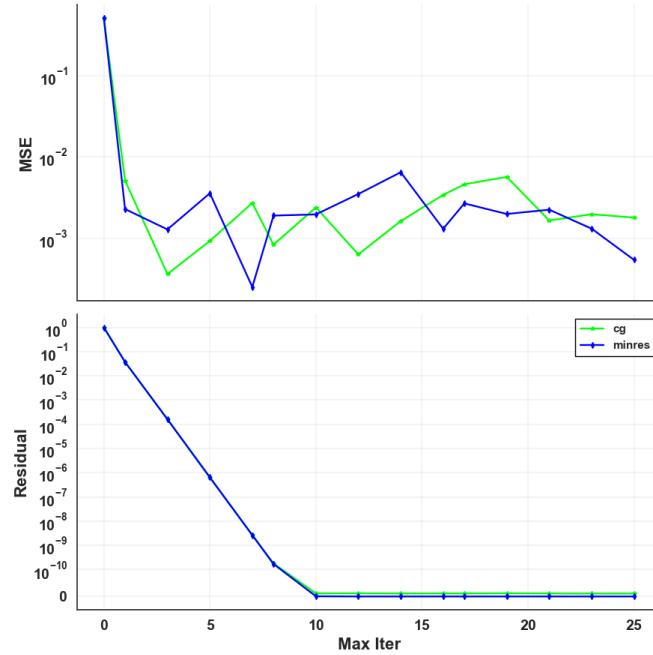


FIGURE 103. Comparison of prediction errors for the temperature data set, here $\sigma = 10.0$.

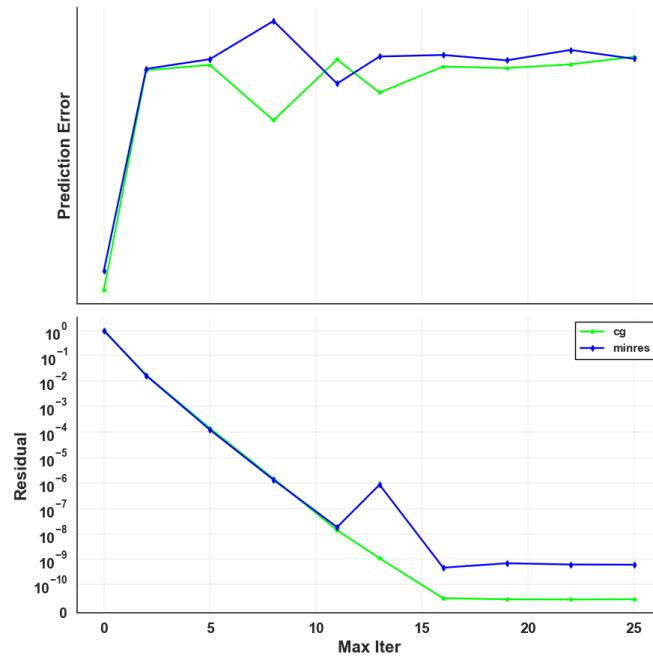


FIGURE 104. Comparison of prediction errors for the wine data set, here $\sigma = 2.1$.

A.13. Comparison of Combining Techniques.

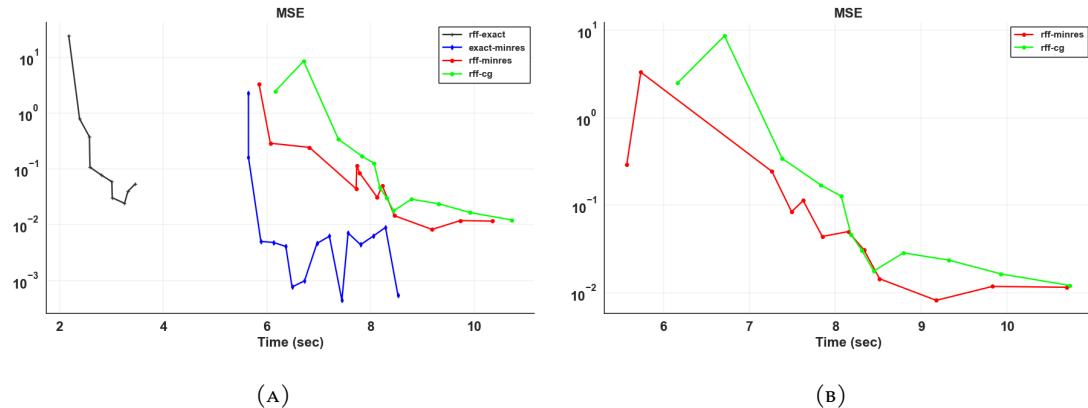


FIGURE 105. Comparing the combination of various methods on the 3DSN data set, here $\sigma = 1.0$.

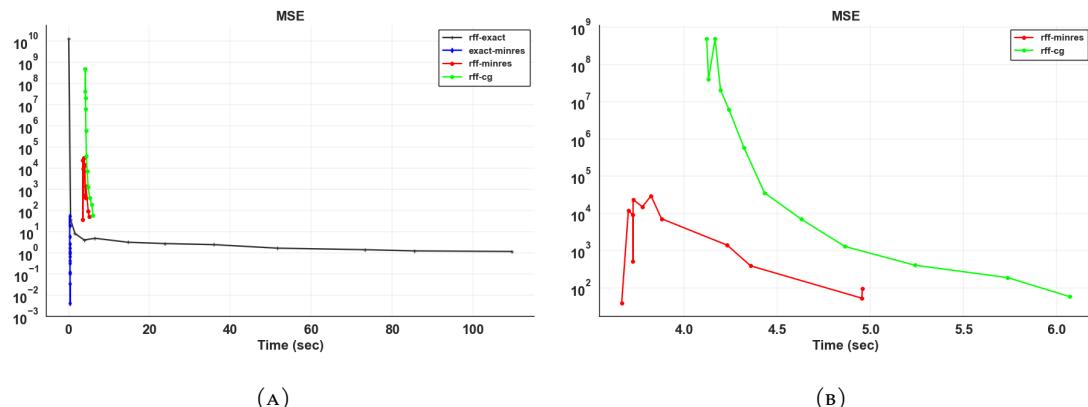


FIGURE 106. Comparing the combination of various methods on the Abalone data set, here $\sigma = 1.0$.

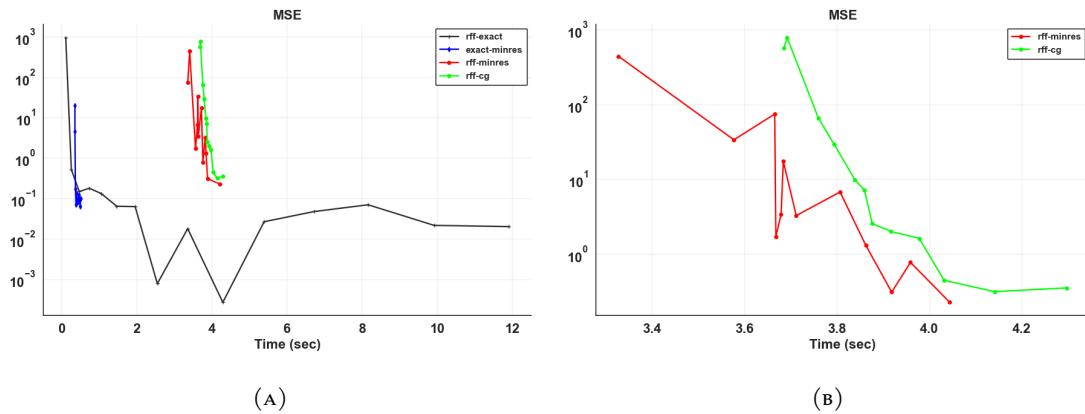


FIGURE 107. Comparing the combination of various methods on the quadratic data set, here $\sigma = 2.1$.

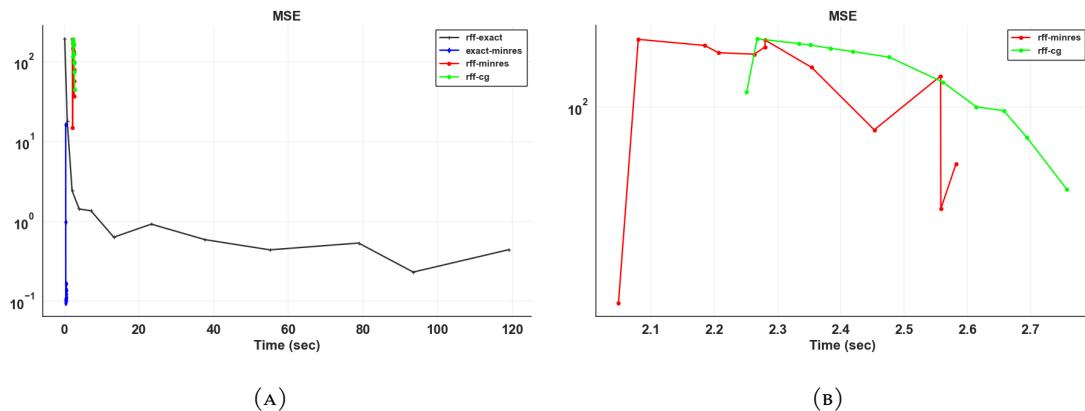


FIGURE 108. Comparing the combination of various methods on the Stock Market data set, here $\sigma = 10.0$.

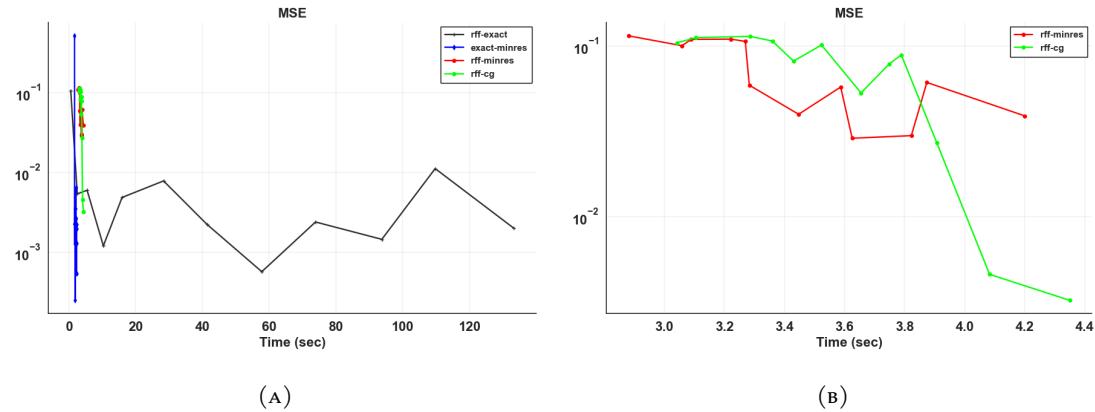


FIGURE 109. Comparing the combination of various methods on the temperature data set, here $\sigma = 10.0$.

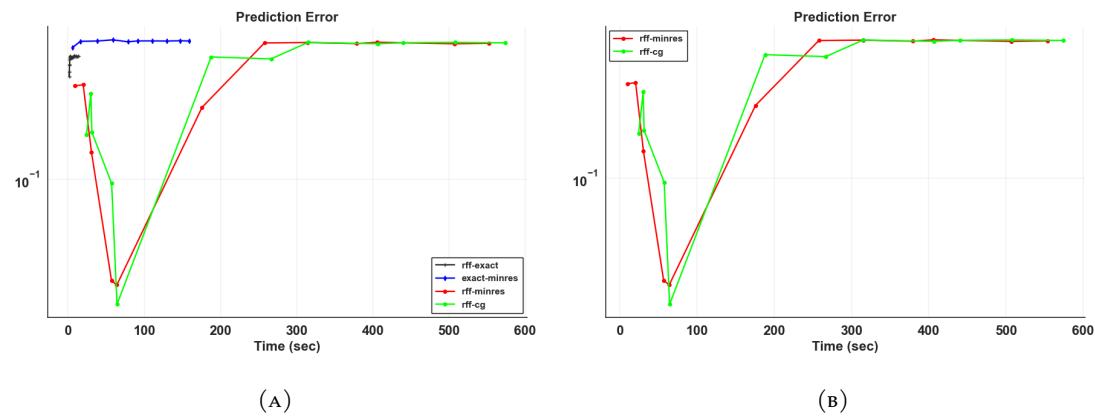


FIGURE 110. Comparing the combination of various methods on the temperature data set, here $\sigma = 2.1$.