



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

BACHELOR OF MATHEMATICS (HONOURS)
JUNE 2022

THE UNIVERSITY OF QUEENSLAND
SCHOOL OF MATHEMATICS AND PHYSICS

CONTENTS

ACKNOWLEDGEMENTS	iii
SYMBOLS AND NOTATION	iv
INTRODUCTION	1
1. RANDOM FOURIER FEATURES.....	4
1.1. THEORY AND COMPUTATION	4
1.2. ORTHOGONAL RANDOM FEATURES	8
1.3. RANDOM ORTHO-MATRICES AND STRUCTURED ORTHOGONAL RANDOM MATRICES	9
2. APPLICATIONS AND RESULTS	12
2.1. GAUSSIAN PROCESSES PREDICTION REVIEWED	12
2.2. EXPERIMENTAL SETUP	12
2.2.1. KERNEL MARTIX APPROXIMATION TESTING	13
2.2.2. KRYLOV SUBSPACE METHODS APPROXIMATION TESTING	14
2.3. DISCUSSION	14
2.3.1. KERNEL MATRIX APPROXIMATION	15
REFERENCES	17
APPENDIX A. SUPPLEMENTARY RESULTS	22
A.1. GRAM MATRIX SPECTRAL VALUES	22
A.2. NYSTROM SCORES	24
A.3. RIDGE LEVERAGE SCORES	27
A.4. NYSTROM ERRORS	30
A.5. NYSTROM ERRORS	36
A.6. KERNEL COMPARISONS	40

ACKNOWLEDGEMENTS

I would like to deeply thank my supervisor Dr. Masoud Kamgarpour for his advice and all of his time spent with me. I consider myself lucky and am glad to have been his student for my honours year. I would also like to thank my co-supervisor Dr. Anna Puskás for the same reasons. A special thanks to Dr. Valentin Buciumas for his time spent teaching me while he was at The University of Queensland.

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).
∇	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.

$\mathbf{A}_{(.,.)}$	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 submatrix $\mathbb{R}^{3 \times (j'-j+1)}$ matrix 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.
$\text{cholesky}(\mathbf{A})$	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 $\mathbf{X} = \mathbf{A}^{-1} \mathbf{B}$ in 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 $f \in \mathcal{O}(g)$ 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 conditional 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 $x_i \in \mathbb{R}_{>0}$ 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 limitations 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 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

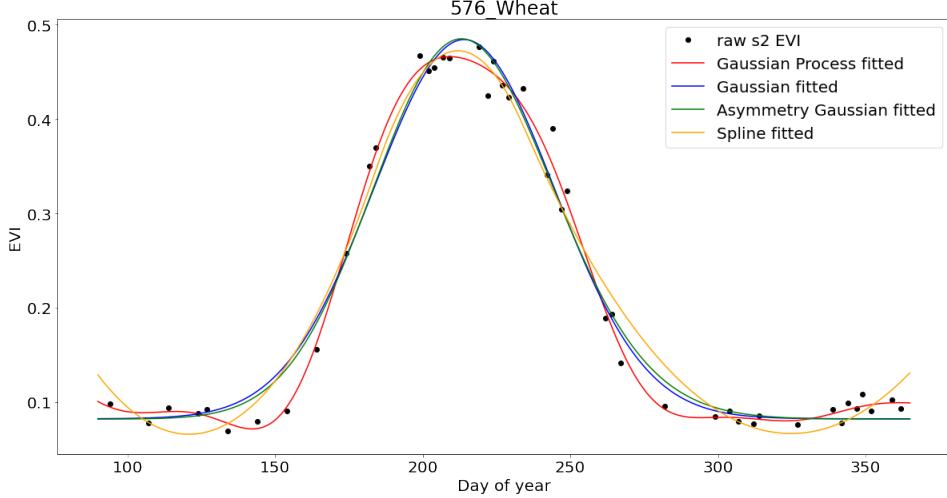


FIGURE 1. Potgieter’s team found that GPs were superior in terms of predicting a phenological timeline for a number of common seasonal crops over other parametric models.

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 ?? 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 ?? and 1 will cover techniques for approximating a large matrix used with GPs that provides information on how similar each observation is to one other. chapter ?? then gives alternative methods for solving linear systems, an essential component required for the GP algorithm to work.

1. RANDOM FOURIER FEATURES

As seen in ?? GPs rely heavily on the Gram matrix (see ??) 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].

1.1. 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 1, which characterises positive definite functions.

Theorem 1 (Bochner's). *A continuous and shift-invariant function $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y}) = k(\Delta)$ is positive definite (see ??) 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 $k(0) = 1$ normalizes μ_k to a probability distribution $p(\cdot)$. For instance, the spectral distribution of the Gaussian RBF kernel is

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

[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 2.

Definition 2 (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 $k(0) = 1$ can be computed

as

$$(2) \quad k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{C}^d} \exp(i\langle \boldsymbol{\omega}, \mathbf{x} - \mathbf{y} \rangle) p(\boldsymbol{\omega}) d\boldsymbol{\omega}.$$

The main idea of RFF is to approximate the integral in 2 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

$$(3) \quad \varphi(\mathbf{x}) = \frac{1}{\sqrt{D}} [z(\boldsymbol{\omega}_1, \mathbf{x}), z(\boldsymbol{\omega}_2, \mathbf{x}), \dots, z(\boldsymbol{\omega}_D, \mathbf{x})]^\top$$

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 2. The first embedding takes the form

$$(4) \quad z(\boldsymbol{\omega}, \mathbf{x}) = [\cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle), \sin(\langle \boldsymbol{\omega}, \mathbf{x} \rangle)]^\top$$

which satisfies 2 since

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

The other embedding Rahimi and Recht give is

$$(5) \quad z(\boldsymbol{\omega}, \mathbf{x}) = \sqrt{2} \cos(\langle \boldsymbol{\omega}, \mathbf{x} \rangle + b)$$

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

embedding given in 4. To summarise their argument we denote

$$(6) \quad \varphi_1(\mathbf{x}) = \sqrt{\frac{2}{D}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \cos(\langle \boldsymbol{\omega}_2, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_{D/2}, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \vdots \\ \sin(\langle \boldsymbol{\omega}_{D/2}, \mathbf{x} \rangle) \end{bmatrix}$$

to be the feature map corresponding to embedding in equation 4 and

$$(7) \quad \varphi_2(\mathbf{x}) = \sqrt{\frac{2}{D}} \begin{bmatrix} \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle + b_1) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle + b_D) \end{bmatrix}$$

to be the feature map corresponding to equation 5. 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 3.

Proposition 3. *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 6 let $\sigma_p^2 = \mathbb{E}_{\boldsymbol{\omega} \sim p(\cdot)} \|\boldsymbol{\omega}\|_2^2 = \text{tr } \nabla^2 k(0)$ then for any $\varepsilon \in \mathbb{R}_{>0}$, $\varepsilon \leq \sigma_p \ell$ 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 $\alpha \in \mathbb{R}_{>0}$, $\alpha < \infty$ does not depend on anything [Rah08, page 3].

Rahimi and Recht prove Proposition 3 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 1, which was used to approximate a Gram matrix for the data set \mathbf{X} using the feature map from 6.

Algorithm 1: 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 1 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 1 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 1 then it produces an unbiased estimate, $\widetilde{\mathbf{K}}_{\text{RFF}}$, for the Gram matrix. This stated more precisely in Lemma 4.

Lemma 4. $\widetilde{\mathbf{K}}_{\text{RFF}}$ 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} . □

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.

1.2. Orthogonal Random Features. In the previous chapter Algorithm 1 assumed some sort of mechanism for producing the transformation matrix \mathbf{W} . The construction presented in 1.1 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 1.1 [Liu21]. We shall consider the two methods proposed by Yu *et al.* [Yu16]; the first method here and the second in the following section (1.3). 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-factorization (see ??) 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

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

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 $\|\mathbf{z}\|_2 \sim \chi_k$ 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

$$(8) \quad \mathbf{W}_{\text{ORF}} = \left(\frac{\sigma}{\sqrt{2}} \right)^{-1} \mathbf{SQ}.$$

The main downside the the ORF method is that the QR-factorization brings a computational cost of $\mathcal{O}(Dd)$. Fortunately when using \mathbf{W}_{ORF} as our transformation matrix in Algorithm 1 the approximate Gram matrix $\widetilde{\mathbf{K}}_{\text{RFF}}$ is an unbiased estimate of \mathbf{K} , stated more formally in Theorem 5.

Theorem 5. $\widetilde{\mathbf{K}}_{\text{ORF}}$ 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].

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

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

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].

1.3. 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

$$(9) \quad \mathbf{W}_{\text{ROM}} = \sqrt{d} \prod_{i=1}^k \mathbf{S} \mathbf{D}_i$$

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 6, 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 6 (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 6 or by padding with 0. This provides a concrete means for which one can generate a transformation matrix

$$(10) \quad \sqrt{d} \prod_{i=1}^k \mathbf{H} \mathbf{D}_i$$

where \mathbf{H} is an appropriately sized Hadamard matrix. It is easy to check that the matrix generated by equation 10 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.

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

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]	$\mathcal{O}(D \log(d))$	$\mathcal{O}(D)$

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 10. They also show that choosing odd values of k in 10 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 11. The method for constructing transformation matrices in this manner is referred to as Structured Orthogonal Random Features (SORF).

$$(11) \quad \mathbf{W}_{\text{SORF}} = \sqrt{d} \mathbf{H} \mathbf{D}_3 \mathbf{H} \mathbf{D}_2 \mathbf{H} \mathbf{D}_1$$

This is the same transformation matrix estimate that Yu *et al.* provides. Unfortunately using the SORF method in Algorithm 1 does not produce an unbiased estimate of the Gram matrix; however, it does satisfy an asymptotic unbiased property

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

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 2.

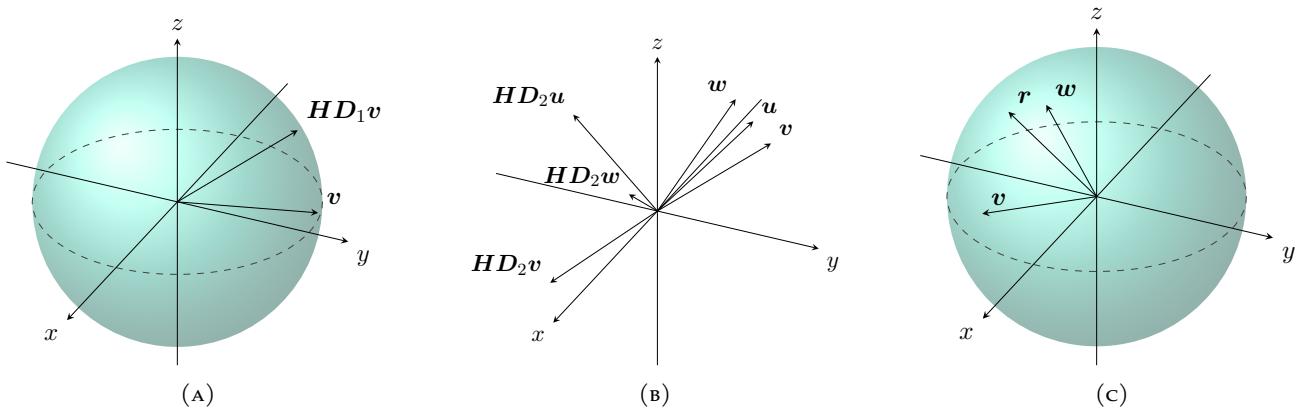


FIGURE 2. A visual representation for the roles of each matrix block in the SORF method. The first block $\mathbf{H}\mathbf{D}_1$ rotates \mathbf{v} so that single dimension bears too much of the l^2 norm seen in panel (A). In panel (B) the second block $\mathbf{H}\mathbf{D}_2$ transforms vectors so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

2. APPLICATIONS AND RESULTS

So far we have discussed a number of different approximation techniques to use to speed up some of the bottle necks within our Gaussian Process regression and classification algorithms. Here, we shall see how they can be used within the naive implementations from ?? to provide faster methods, at the expense of small amounts of accuracy. We shall also go through an extensive treatment of how the various methods reviewed in previous chapters were implemented and the setup of the experiments used to make comparisons to determine the most successful model.

2.1. Gaussian Processes Prediction Reviewed. In chapter ??, a naive implementation for Gaussian Process prediction was presented in ?? . While this does provide a simple and convenient way to produce predictions for a regression task, it is not quite all smooth sailing from here. Unfortunately, there are a number of problems with this algorithm in terms of scalability. For convenience, this algorithm has been restated below but this time with its various bottlenecks highlighted in red.

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 may not be much of an issue for smaller datasets, this will scale very poorly where most modern desktops could not handle training datasets any larger than 10^5 with samples [WJMaSKaAGW21, page 2]. The Gaussian Process classifier from ?? suffers for the exact same reasons as GPR. To mitigate the computational burden of these bottle necks, we can replace these procedures with their approximative counterparts discussed in Chapters ??, 1 and ??.

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, again, improve runtime performance. This especially makes sense approximating the kernel matrix as the Cholesky decomposition will provide an almost-exact solution to solving these linear system, as this is rather counterintuitive if approximations are used earlier on within the algorithm.

2.2. 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 help implement methods to enhance runtime performance. A just-in-time compiler was also used 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.

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

<i>Name</i>	<i>Description</i>	<i>d</i>	<i>n</i>	<i>Source</i>
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 $185 \times 135 km^2$).	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

2.2.1. *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. For 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]. Each method was used 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 $\|K_{XX} - \widehat{K}_{XX}\|_\infty / \|K_{XX}\|_\infty = \left\| K_{XX} - \widehat{K}_{XX} \right\|_\infty$ (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 used were 0.1, 1.0 and 10.0, although a σ value of 2.1 was also used for the Wine data set to compare with [PDAWM05] results. 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 TODO. 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 new prediction. Thus they act as an efficient means to quickly save and rebuild the data kernel matrix. The errors, time and memory were average across the repetitions for each experiment.

To understand how well each kernel matrix approximation method performs in terms of providing predictions, each method was used to train a Gaussian process with $4/5^{ths}$ of each data set and required to predict the remaining $1/5^{th}$. Predictions for each data set was 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 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 ??) 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, 1.0 and 10.0 (as well as 2.1 for the Wine dataset) when an exact GP algorithm was used to provide predictions. The time and memory of performing the matrix inverse using the above was recorded. Again, errors, time and memory were average across the repetitions for each different sample.

2.2.2. 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 $\boldsymbol{\alpha}$. As before, 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. To allow for a fair comparison between the different linear solvers, the kernel matrix was constructed in an exact manner to ensure 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.

2.3. Discussion.

2.3.1. *Kernel Matrix approximation.* Starting with the performance of the Nystrom methods, assessing Figures 21 to 32, overall the rls method is certainly the best methods among the sampling distributions as it virtually always is 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 5 and 6) are relatively uniform and from Figures 25 and 27 we find that the non-uniform sampling techniques give rather poor approximations and generally behave just as bad, or even worse, than the naive uniform sampling distribution. In contrast, non-uniform methods really shone for datasets such as the 3DSN and Wine datasets where some of the better methods could provide almost exact kernel matrix approximations (refer to Figures 21, 22, 31 and 32) after a few thousand samples. This phenomena is likely due to the fact that when the spectrum of a matrix is less uniform, it decays faster, 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 10 and 16), the rls methods seems to perform significantly better than uniform selection, shown in Figure 23. The small amount of variation in the rls sampling distribution presumably provides it with a large enough advantage over the uniform method and selects columns that better provides information to the lower rank approximation.

Looking at 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 33 to 44), it is hard to argue using either of these methods as neither them offer no obvious advantage in terms of time and memory saving (see Figures TODO). This may come as a surprise for most with the SORF method under deliverance, 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 the single line of code to produce the i.i.d. Gaussian matrix for the standard transformation matrix using the numpy python library has been incredibly well engineered and optimized to produce matrix samples brilliantly fast and designed to cope with large scaling. When given large enough datasets, the SORF method likely will overtake the naive sampling method. Unfortunately, due to the memory constraints of the computers used for this project, such large data sets could not be experimented with.

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 along side the corresponding relative error. In this setup, we can think about each method having some sort of time budget and that its task is to produce the best approximation it can within this time 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 far better at lowering the Frobenius error in the approximations it produces, while the RFF methods provide approximations with smaller infinity errors (see Figures 45 to 56). In other words, the entries of the approximations produced by Nystrom methods are *on average* better than those from RFF methods, while *in the worst case* the entries of the approximations produced by RFF are better than those from Nystrom methods. This makes sense since much of the theory in the Nystrom methods was aimed

towards lowering the Frobenius errors, while the theory behind the RFF methods was aimed towards lowering the infinity errors. 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 Nystrom methods used (refer to Figures TODO).

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.
- [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 etal, *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 Likhoshesterov 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).

APPENDIX A. SUPPLEMENTARY RESULTS

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

A.1. Gram Matrix Spectral Values.

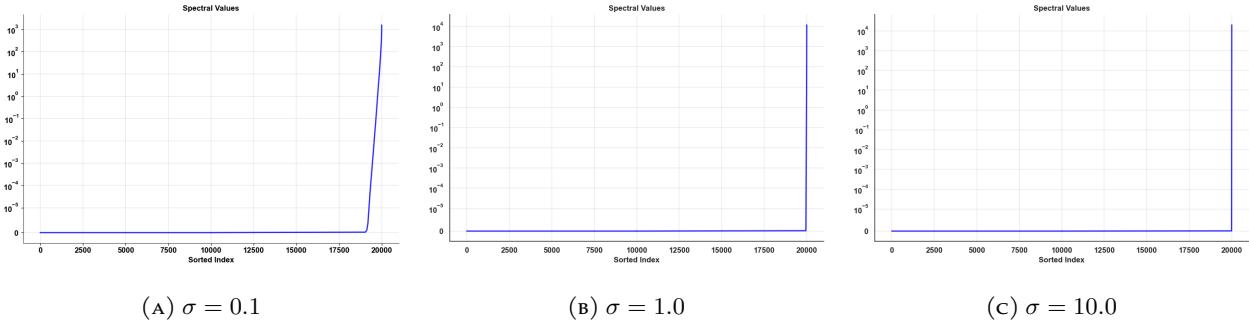


FIGURE 3. Spectral values for 3D-spatial network data.

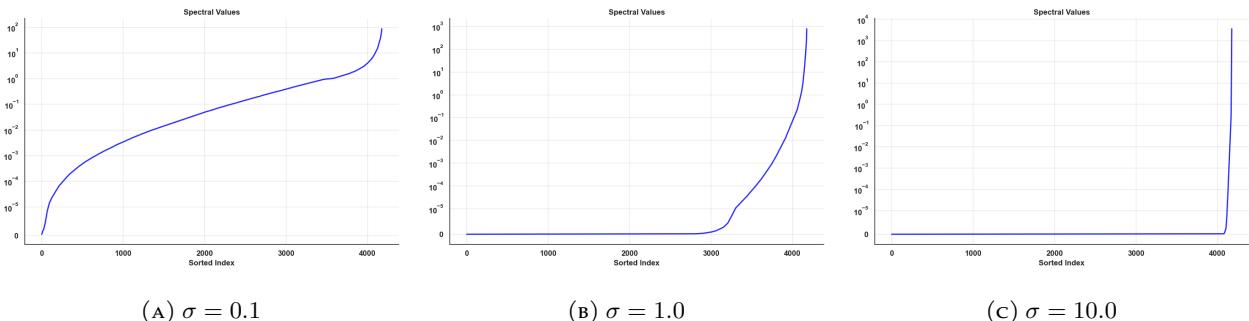


FIGURE 4. Spectral values for abalone data.

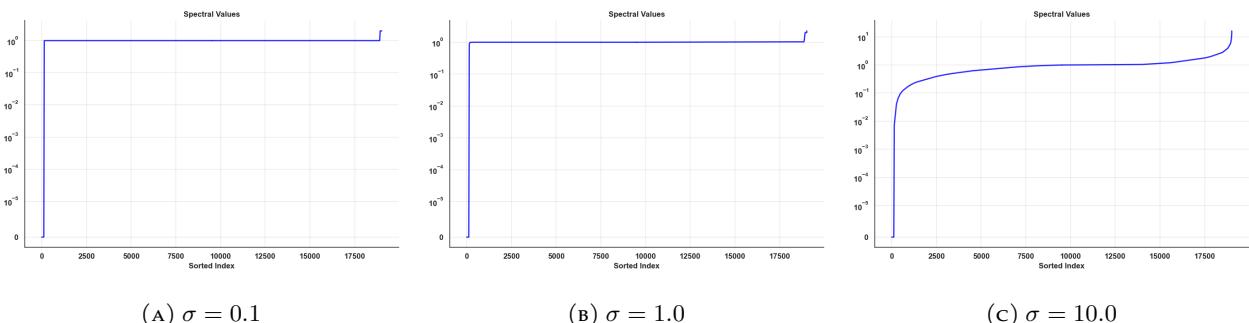


FIGURE 5. Spectral values for magic04 data.

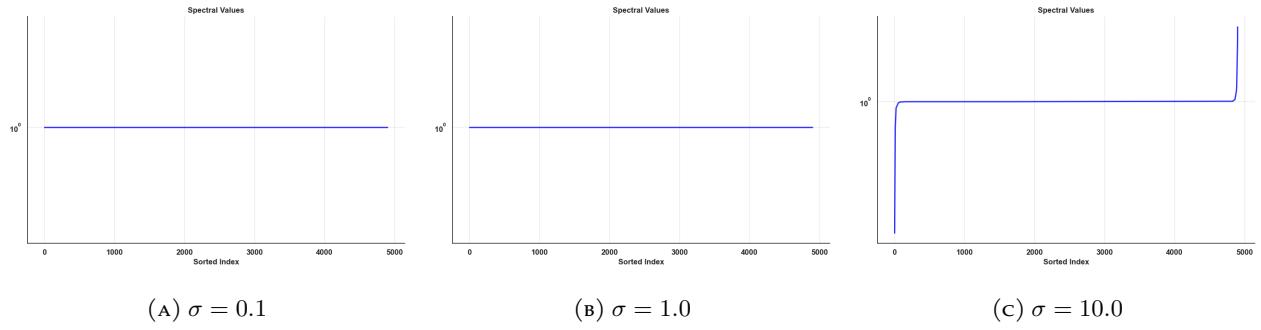


FIGURE 6. Spectral values for stock market data.

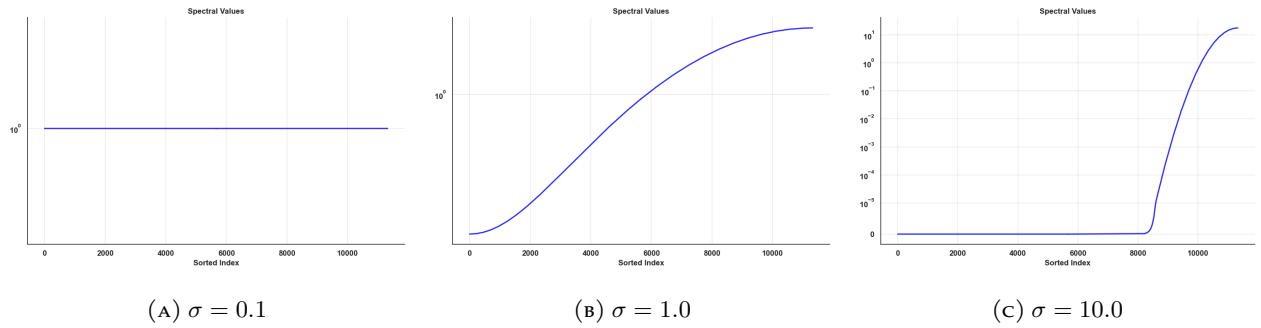


FIGURE 7. Spectral values for temperature data.

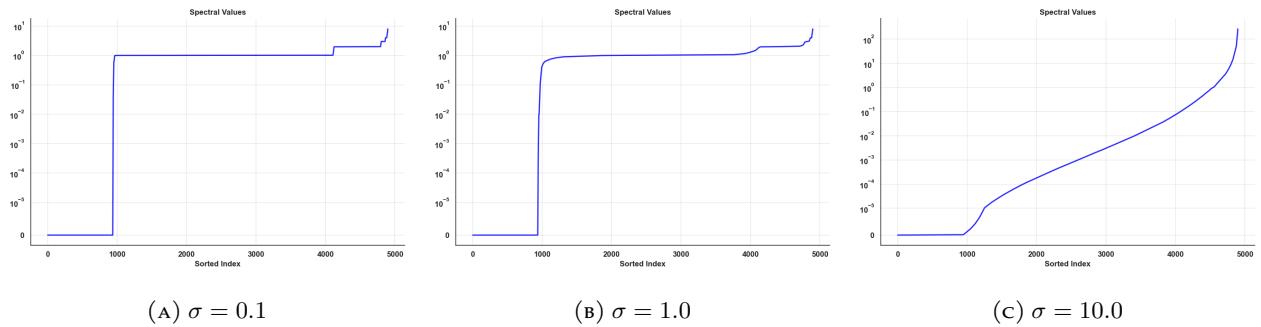


FIGURE 8. Spectral values for wine data.

A.2. Nystrom Scores.

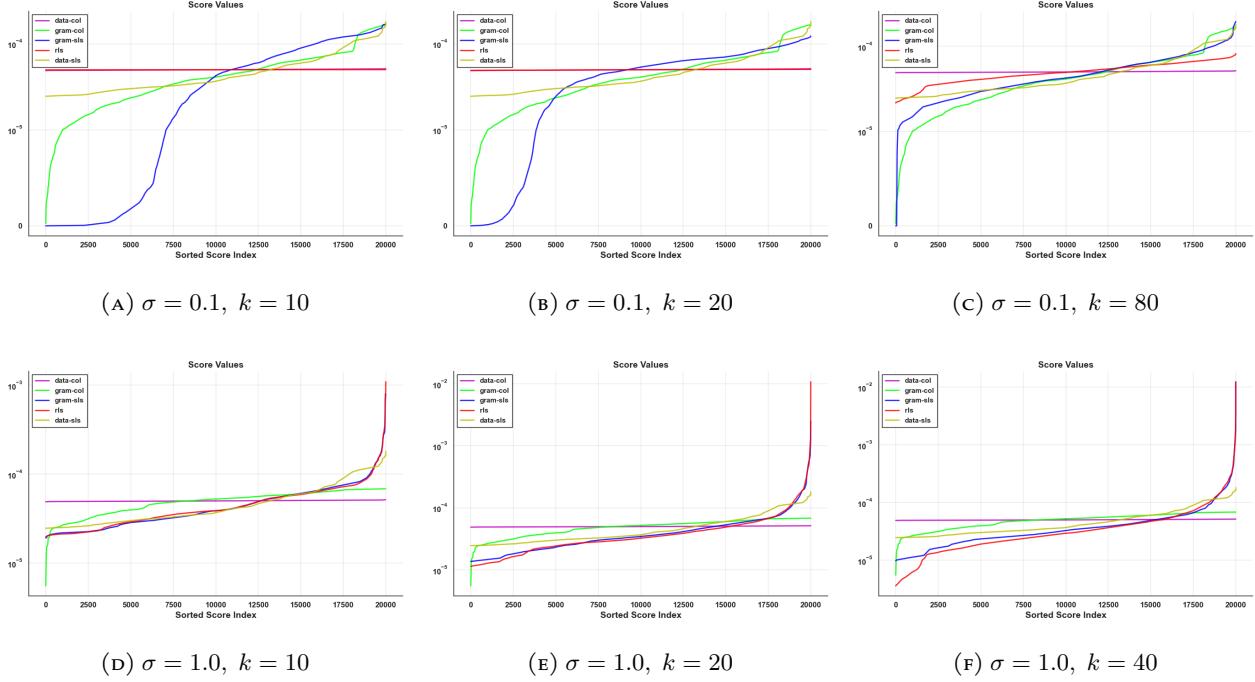


FIGURE 9. Nystrom scores for the 3D-spatial network data set.

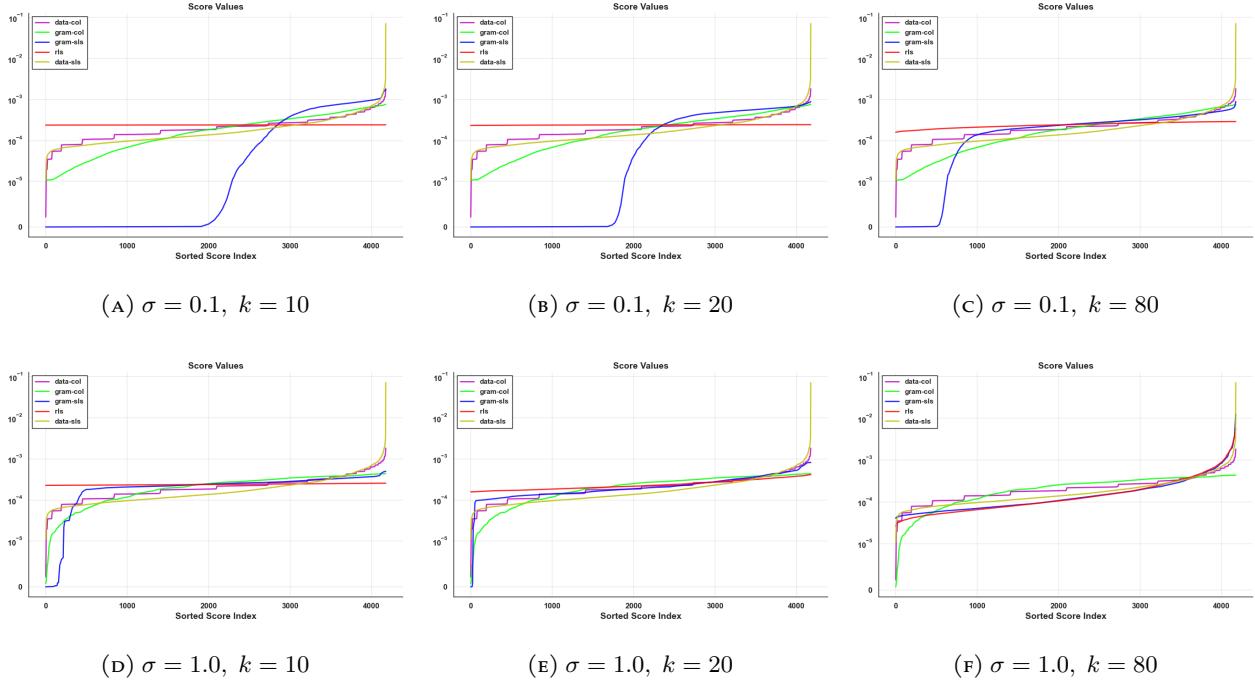


FIGURE 10. Nystrom scores for the Abalone data set.

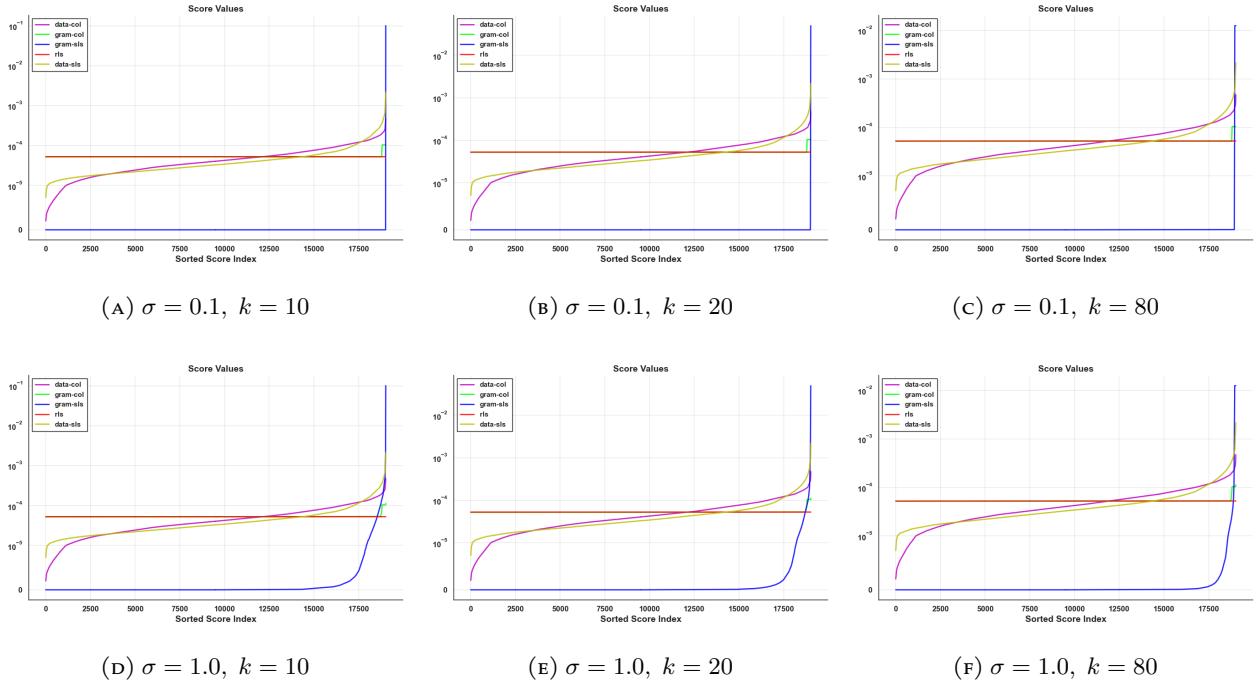


FIGURE 11. Nystrom scores for the Magic data set.

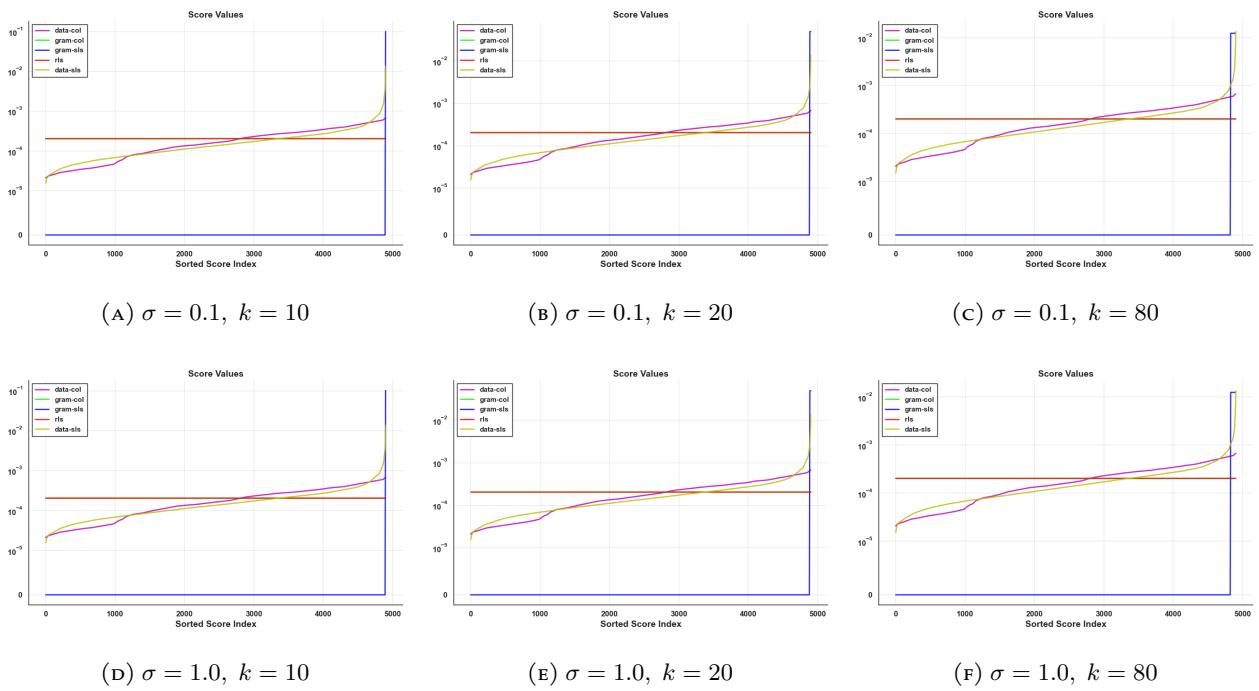


FIGURE 12. Nystrom scores for the stock market data set.

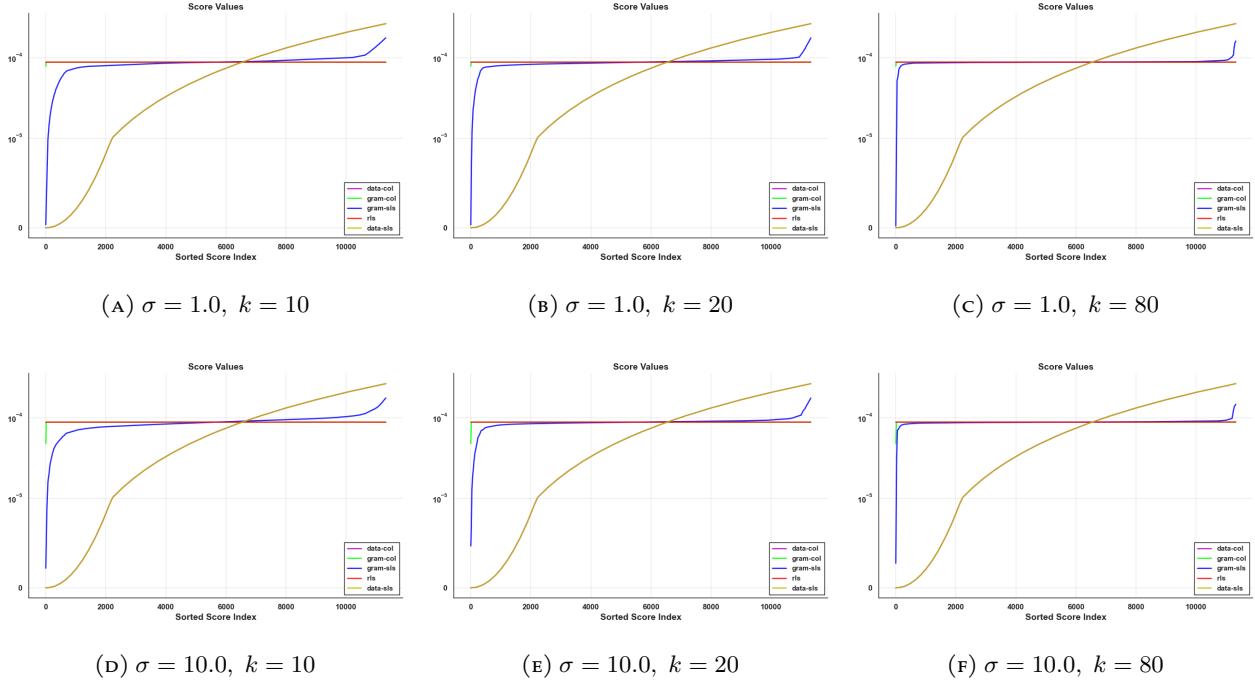


FIGURE 13. Nystrom scores for the temperature data set.

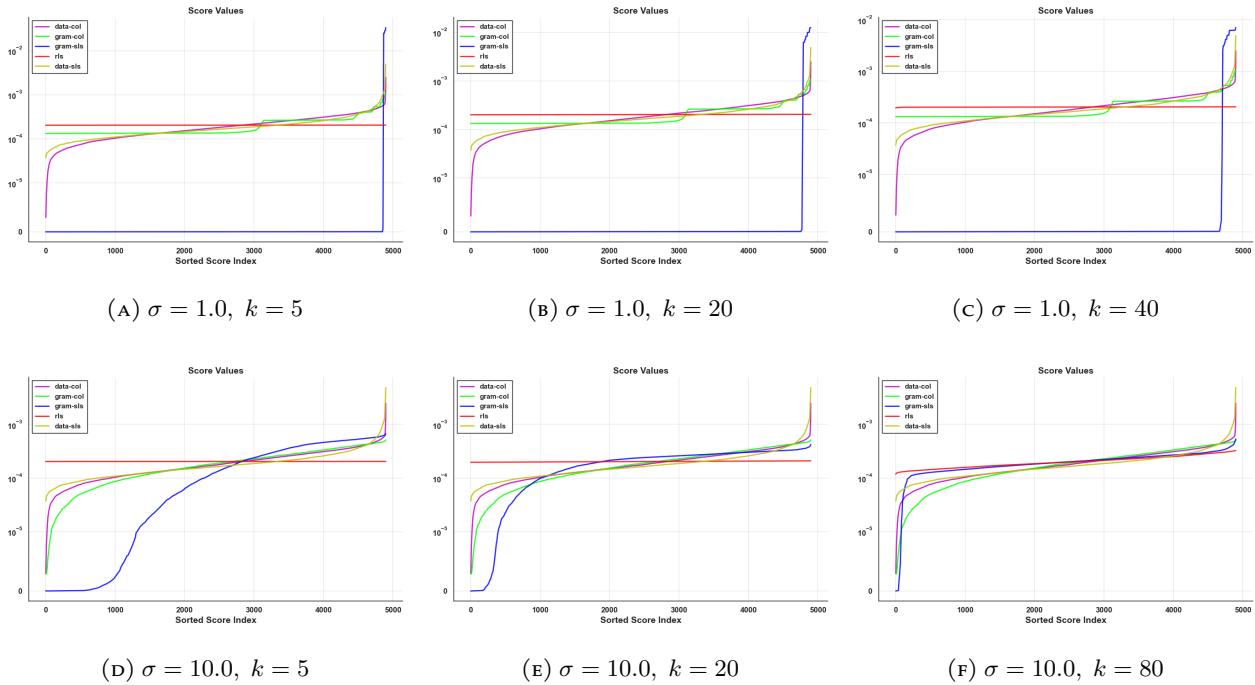


FIGURE 14. Nystrom scores for the wine data set.

A.3. Ridge Leverage Scores.

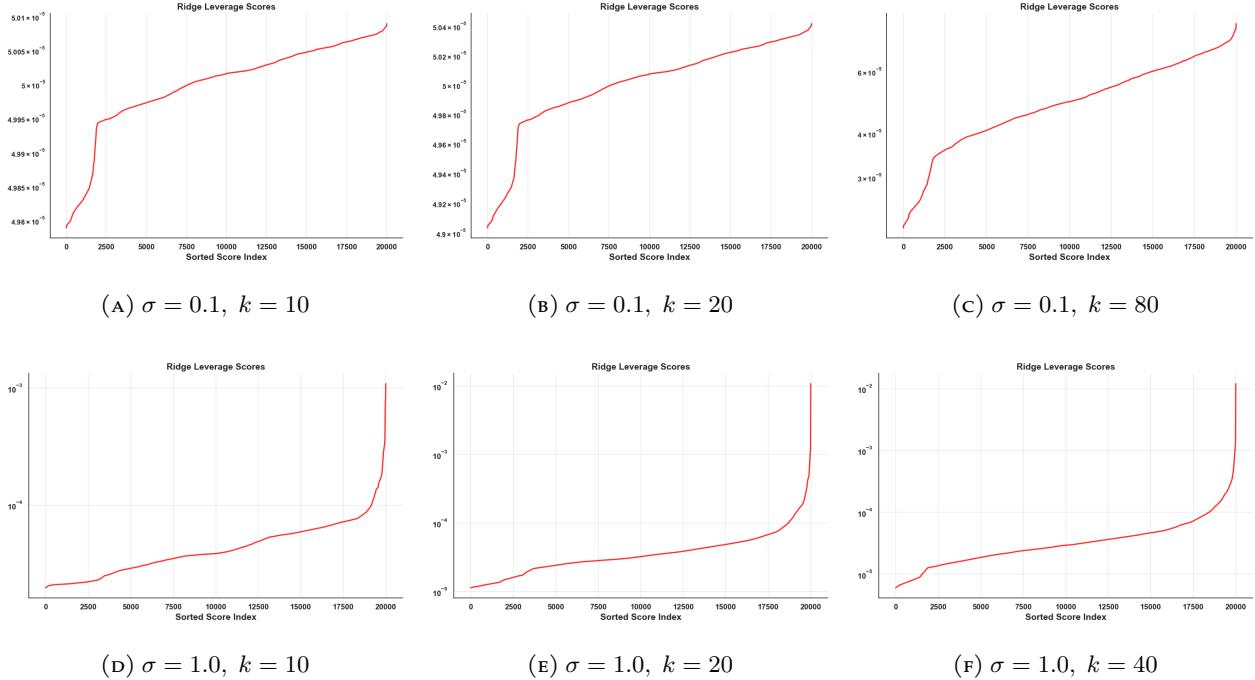


FIGURE 15. Ridge Leverage scores for the 3D-spatial network data set.

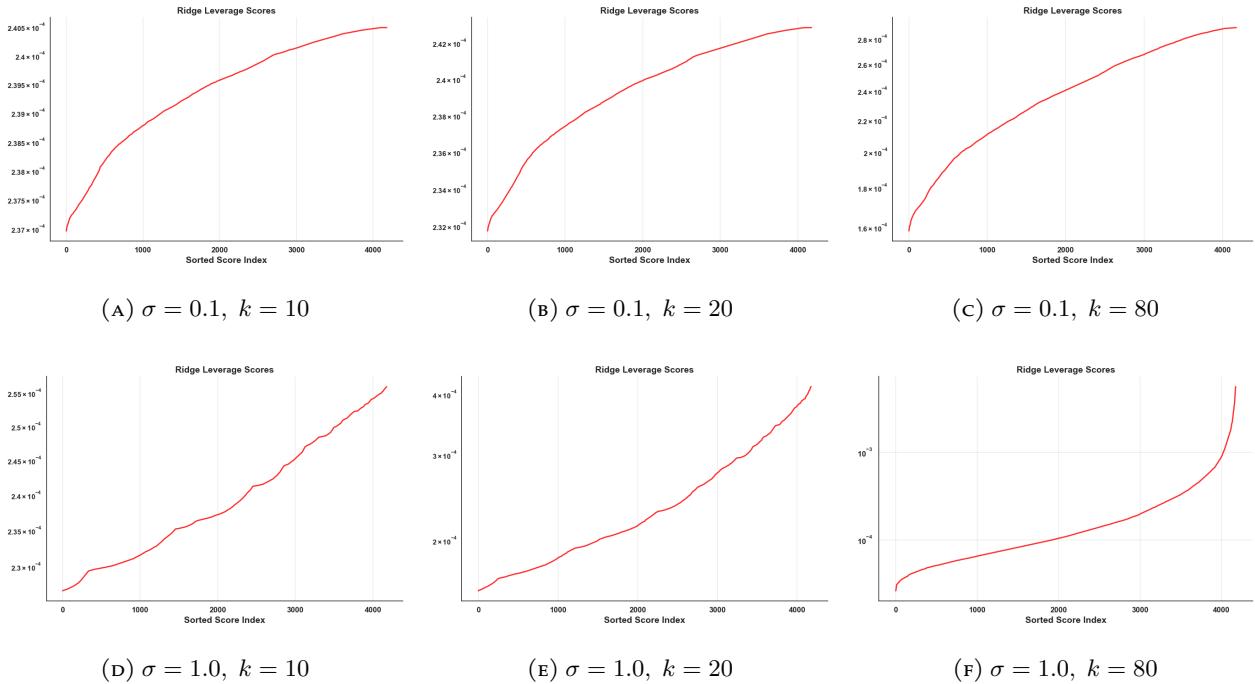


FIGURE 16. Ridge Leverage scores for the Abalone data set.

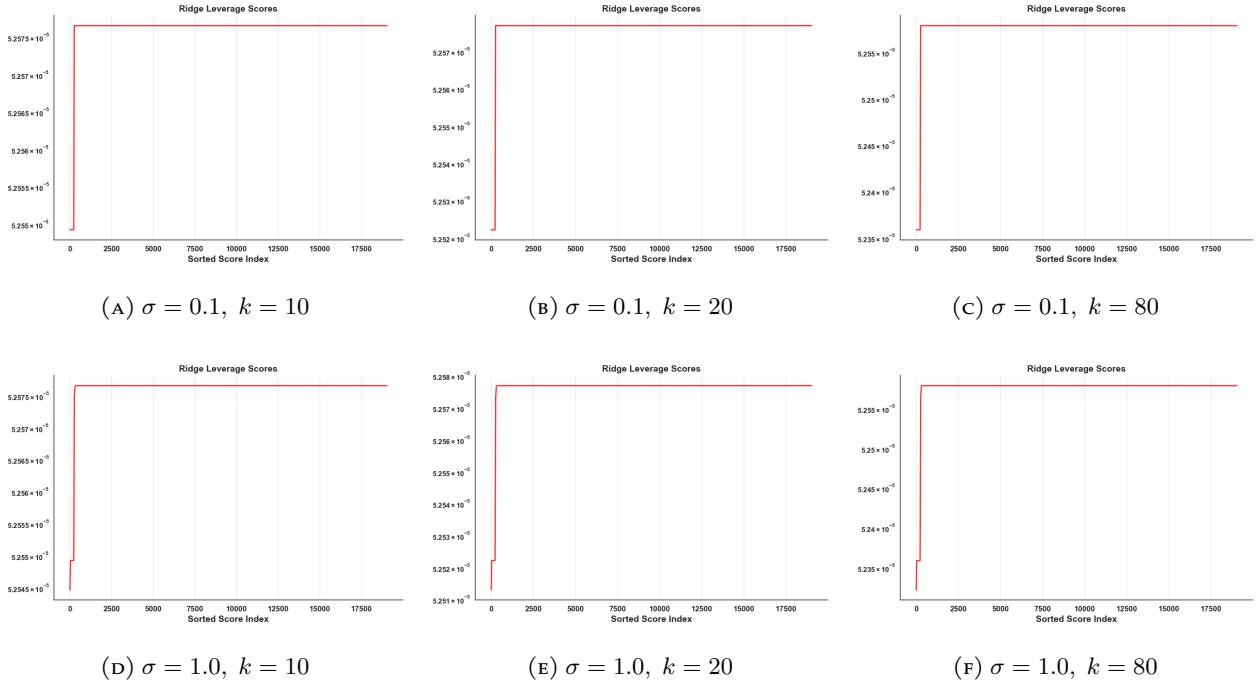


FIGURE 17. Ridge Leverage scores for the Magic data set.

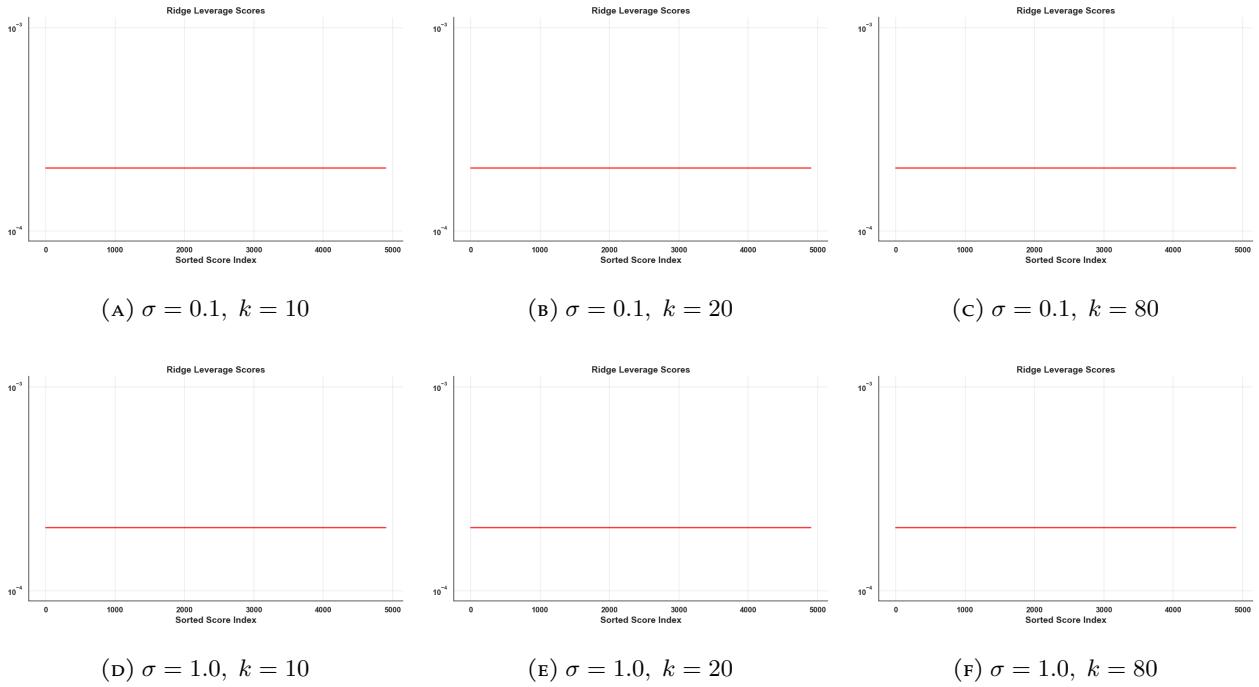


FIGURE 18. Ridge Leverage scores for the stock market data set.

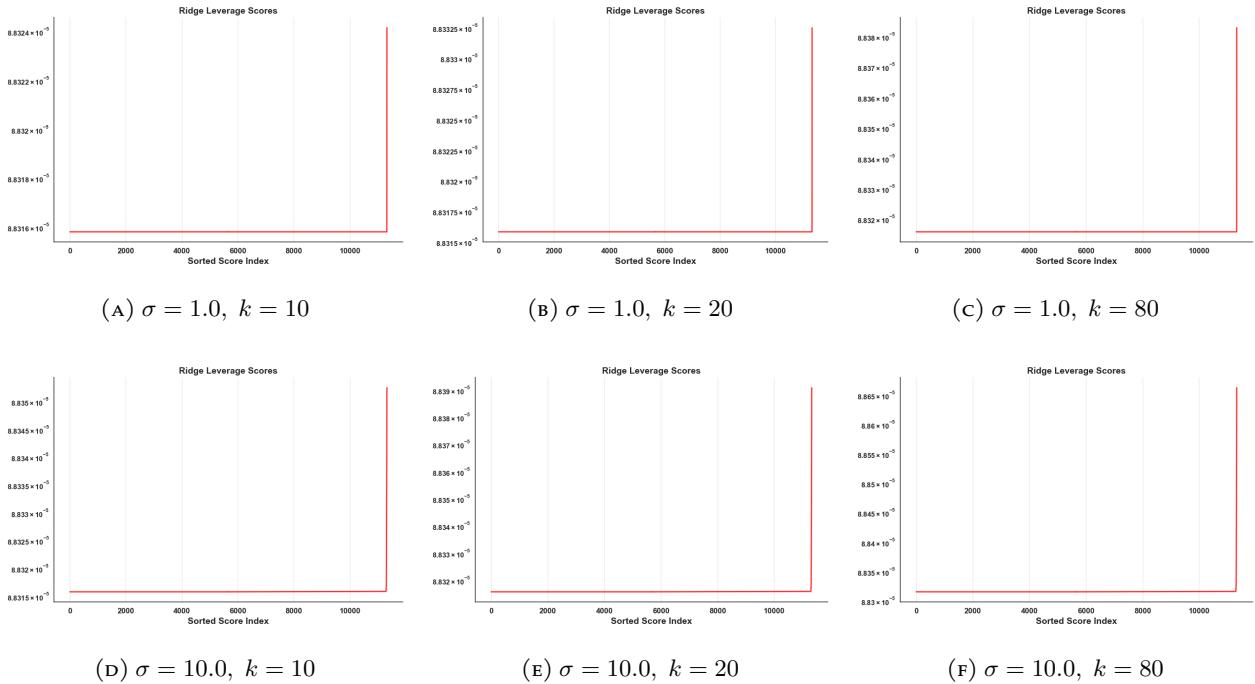


FIGURE 19. Ridge Leverage scores for the temperature data set.

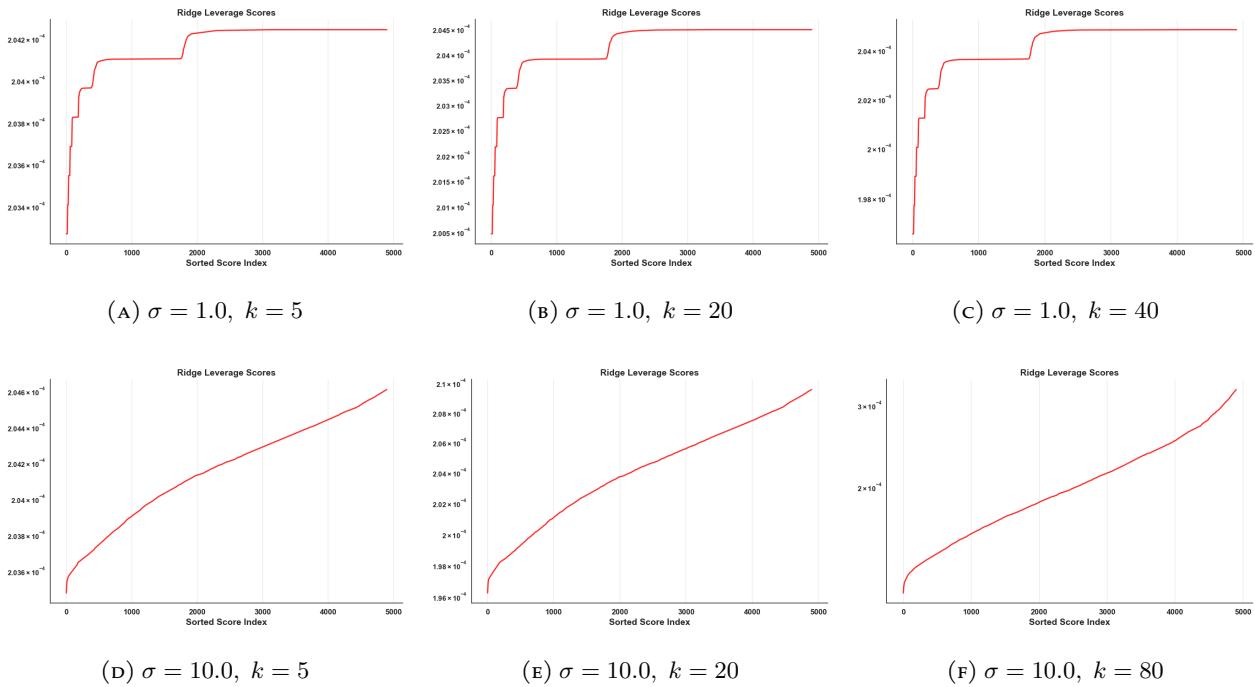


FIGURE 20. Ridge Leverage scores for the wine data set.

A.4. Nystrom Errors.

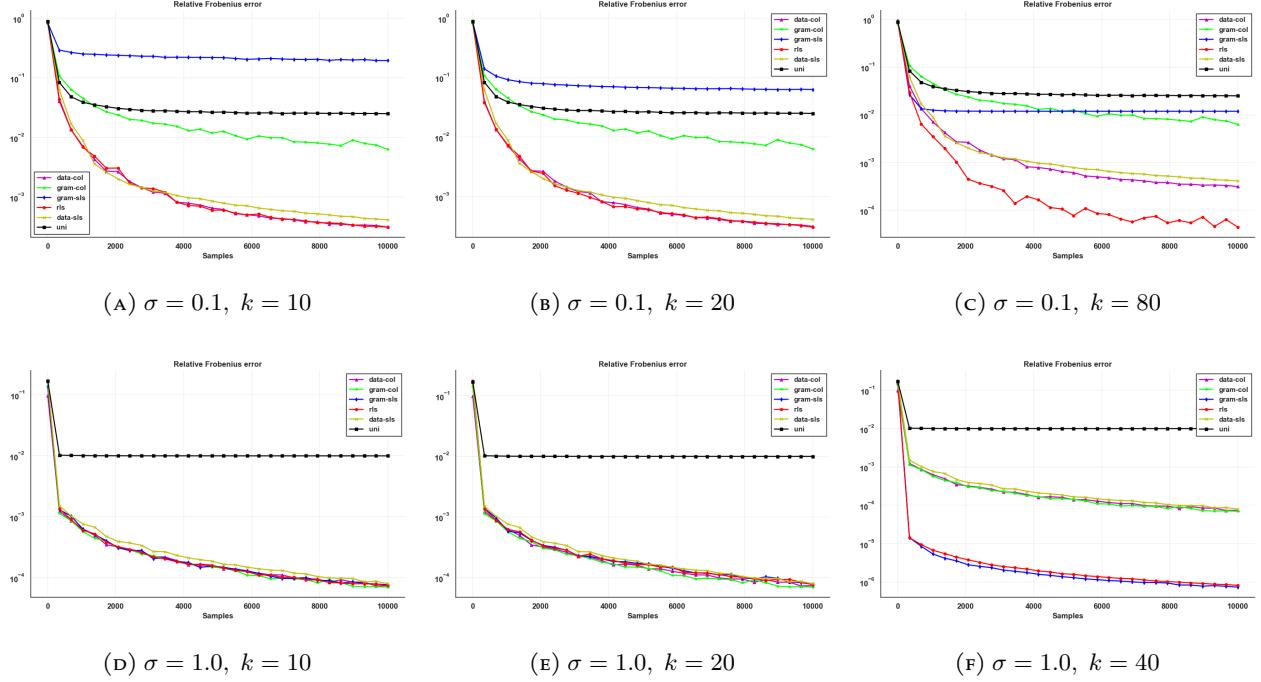


FIGURE 21. Nystrom Frobenius errors for the 3DSN data set.

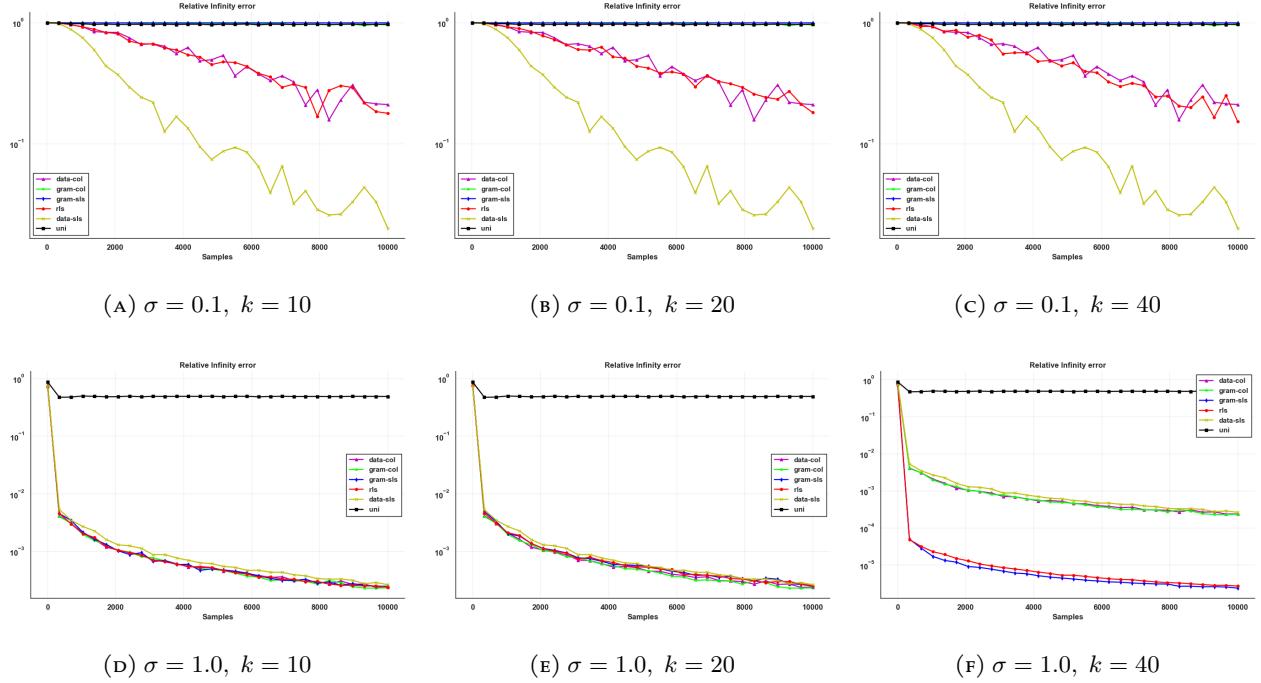


FIGURE 22. Nystrom infinity errors for the 3DSN data set.

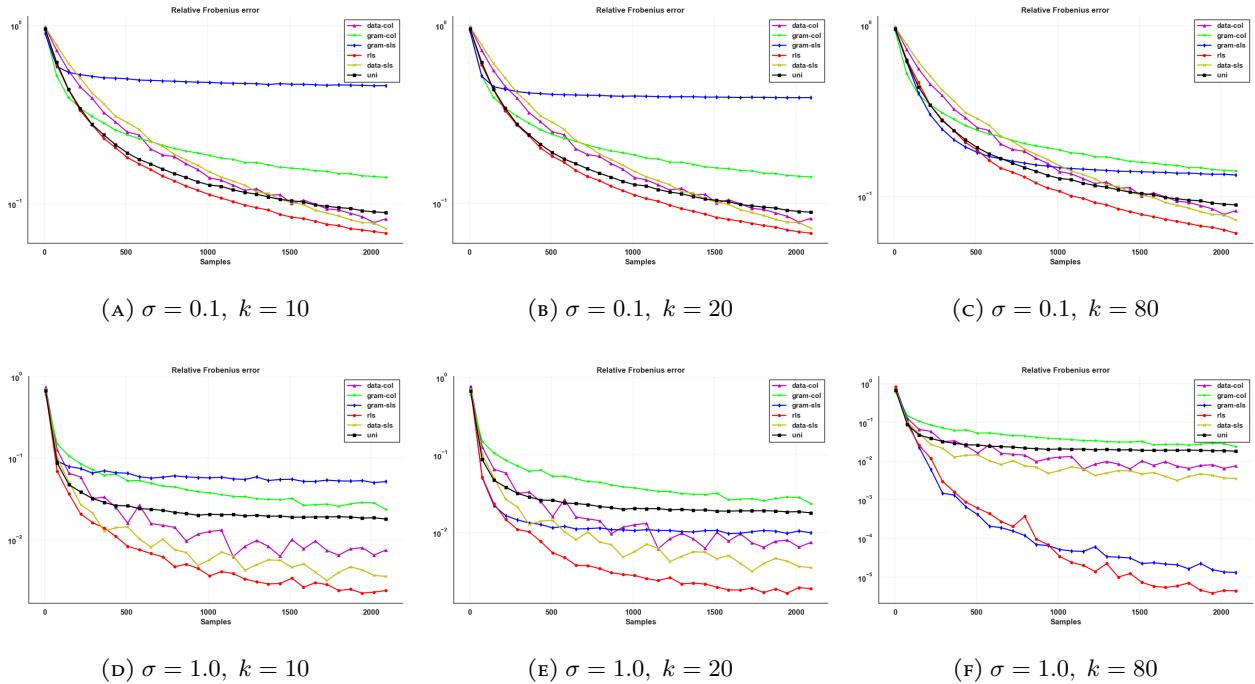


FIGURE 23. Nystrom Frobenius errors for the Abalone data set.

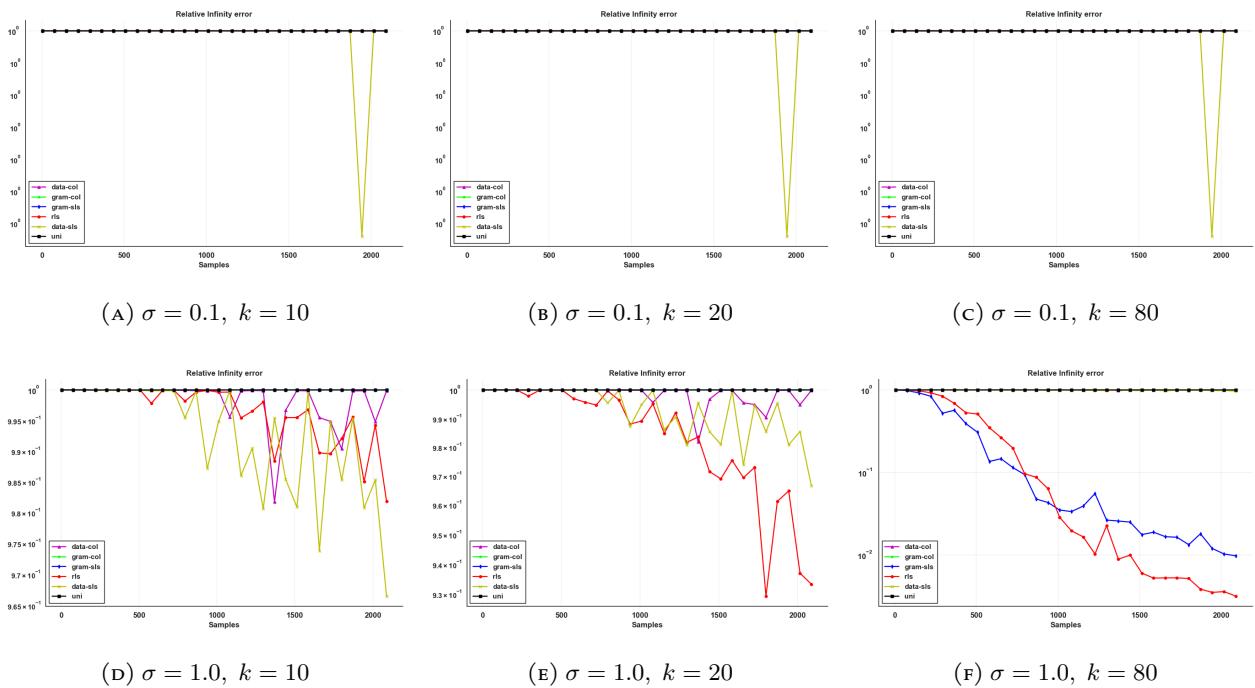


FIGURE 24. Nystrom infinity errors for the Abalone data set.

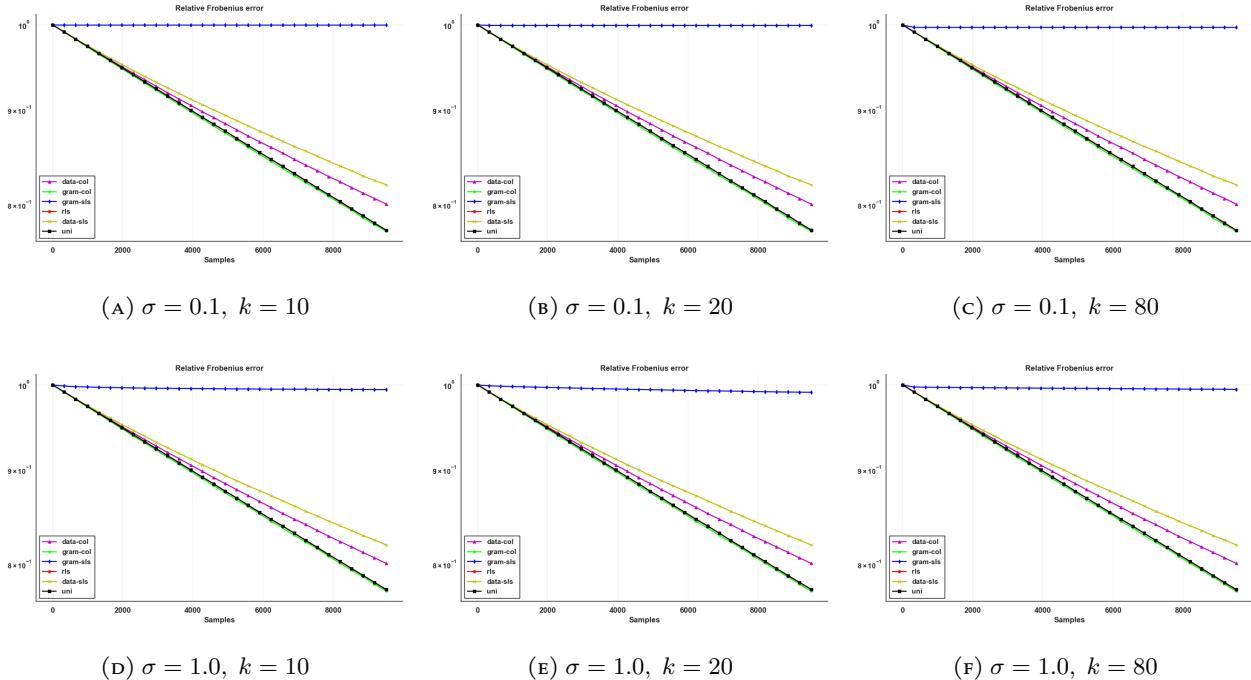


FIGURE 25. Nystrom Frobenius errors for the Magic data set.

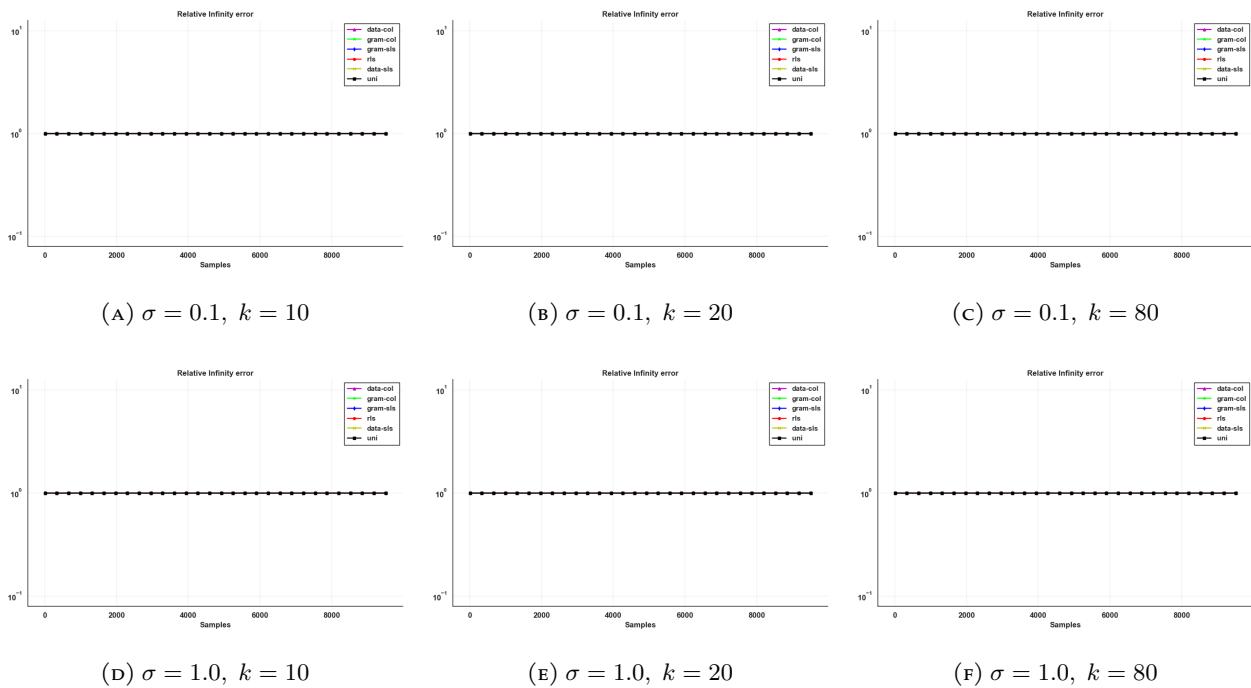


FIGURE 26. Nystrom infinity errors for the Magic data set.

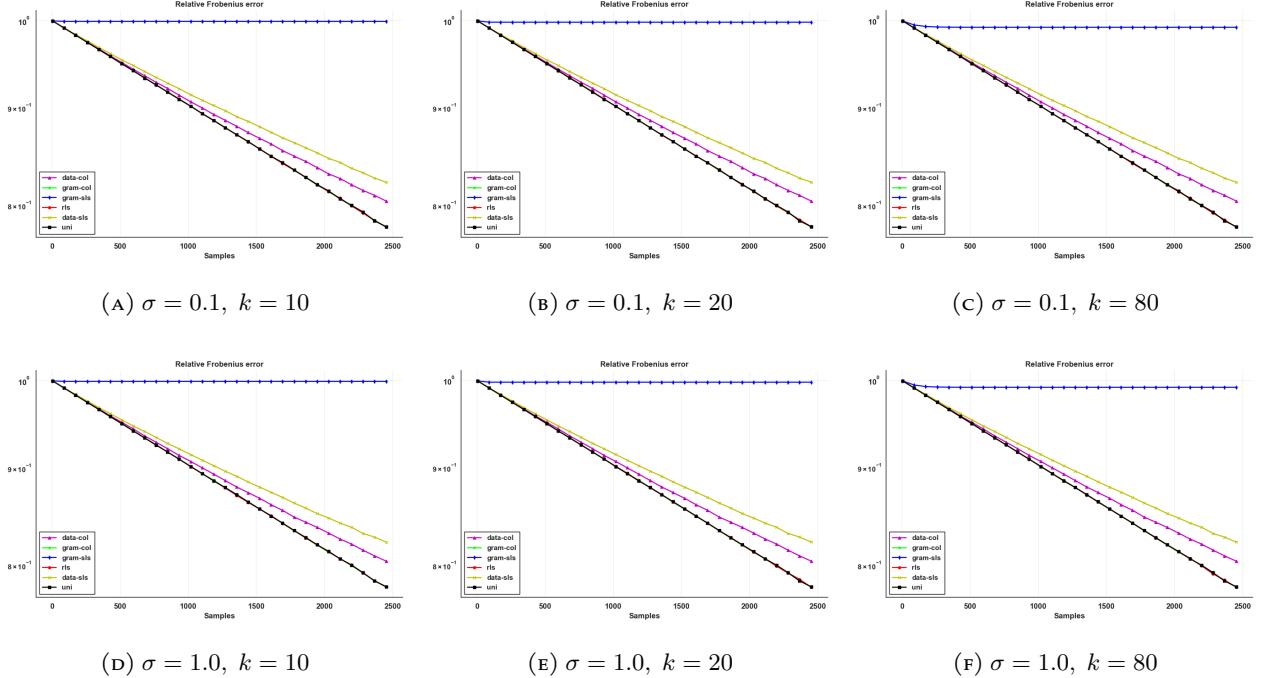


FIGURE 27. Nystrom Frobenius errors for the Stocks data set.

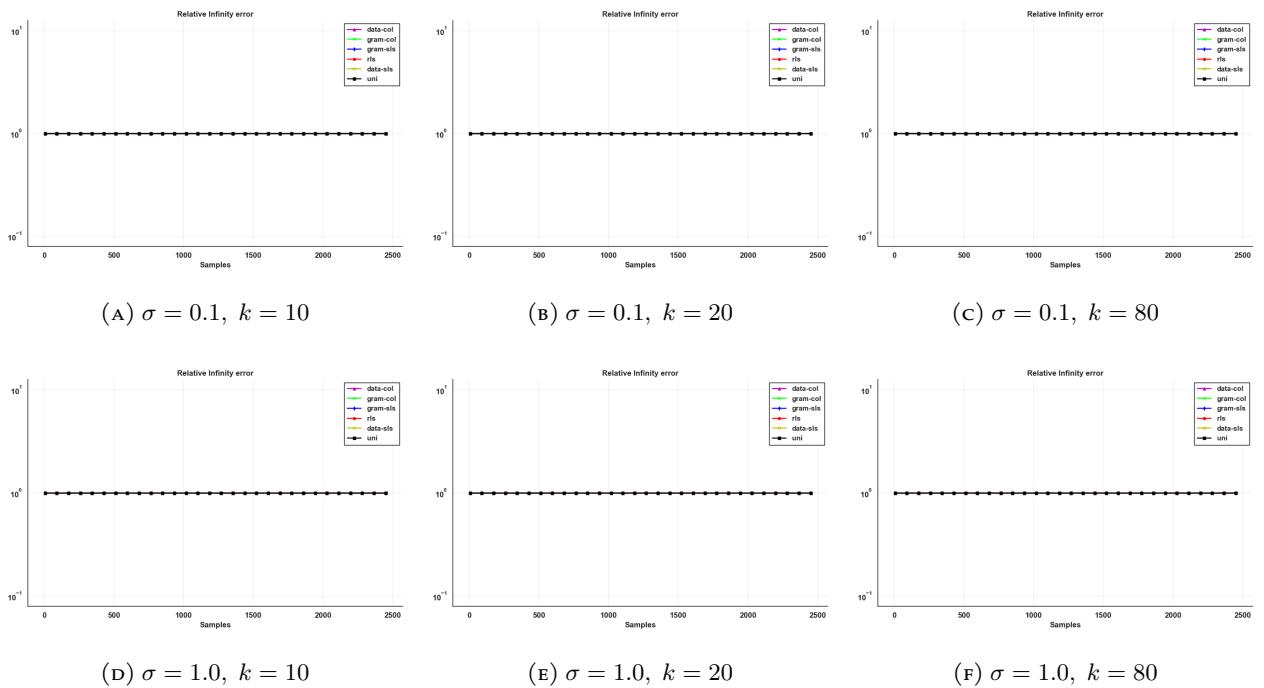


FIGURE 28. Nystrom infinity errors for the Stocks data set.

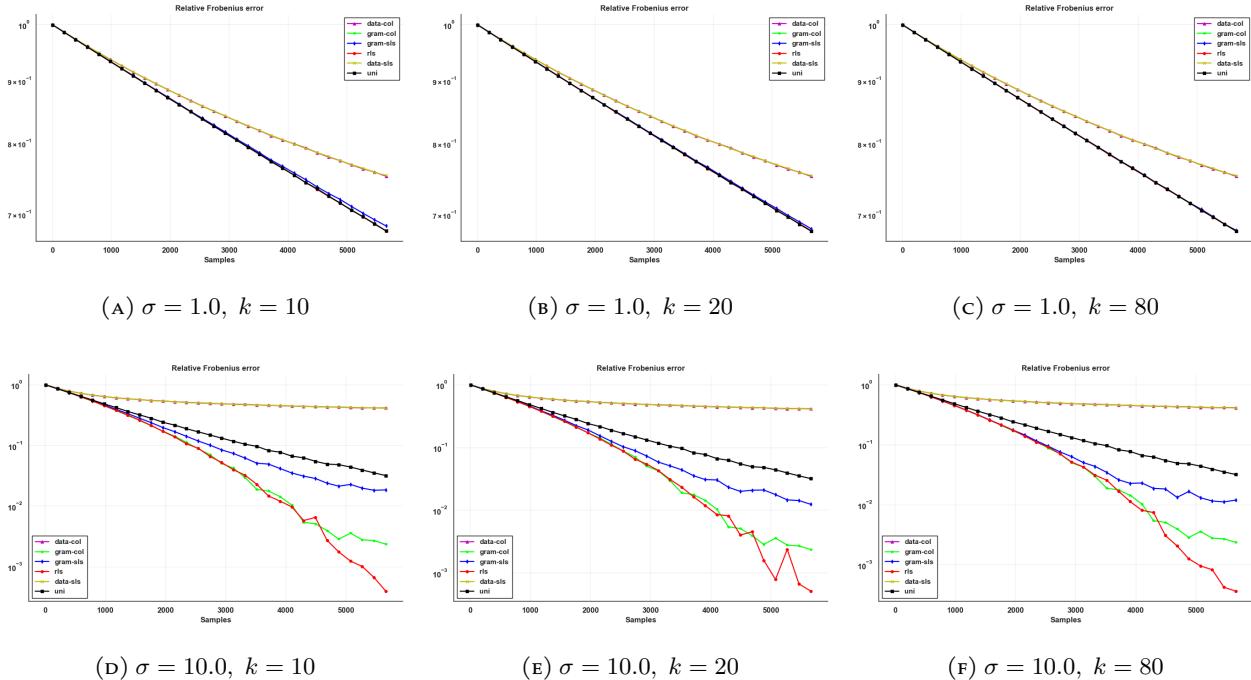


FIGURE 29. Nystrom Frobenius errors for the Temp data set.

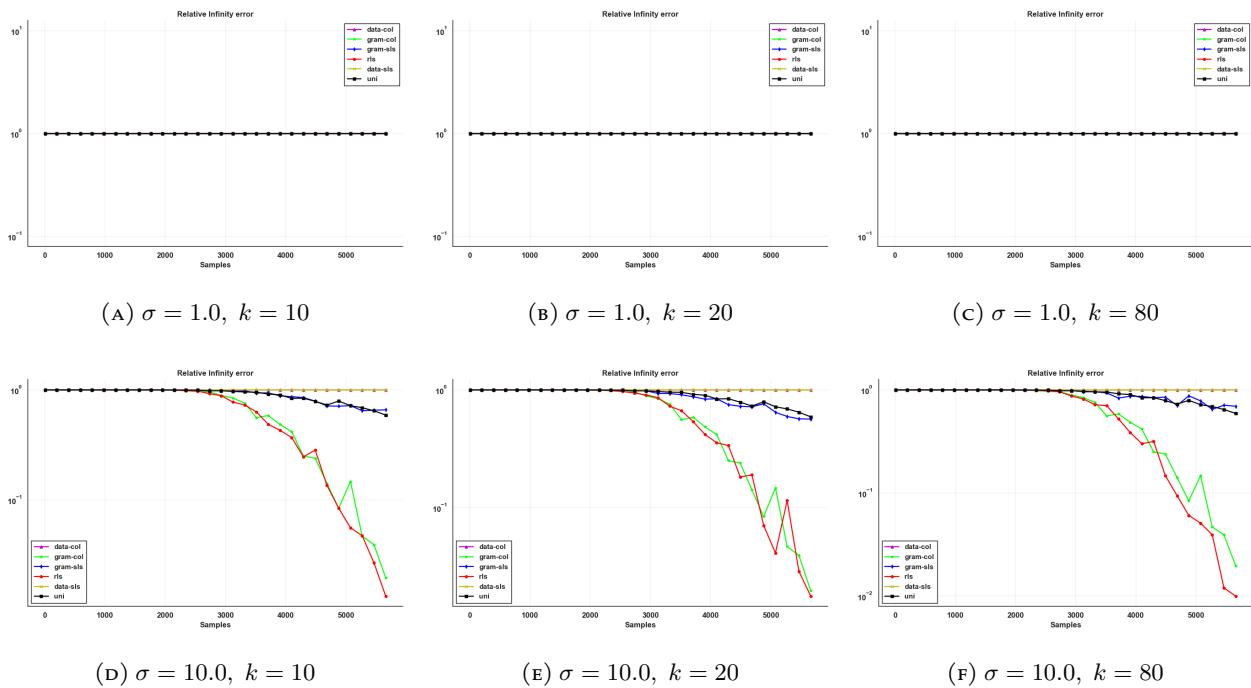


FIGURE 30. Nystrom infinity errors for the Temp data set.

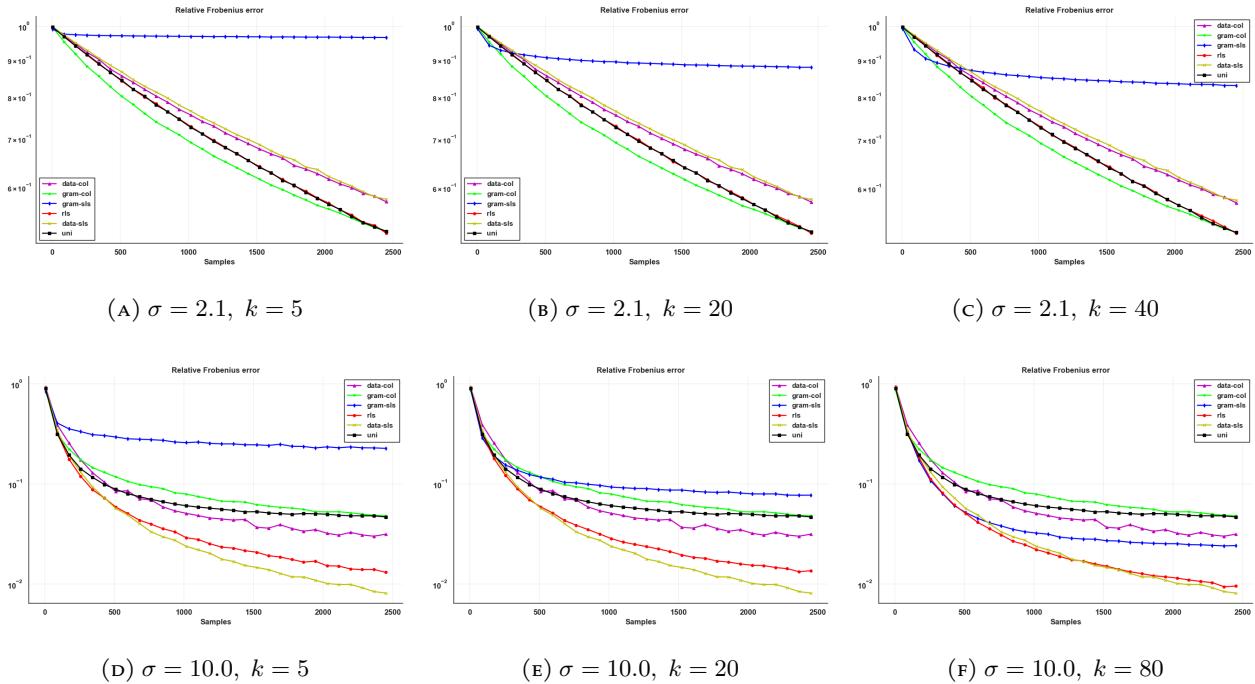


FIGURE 31. Nystrom Frobenius errors for the wine data set.

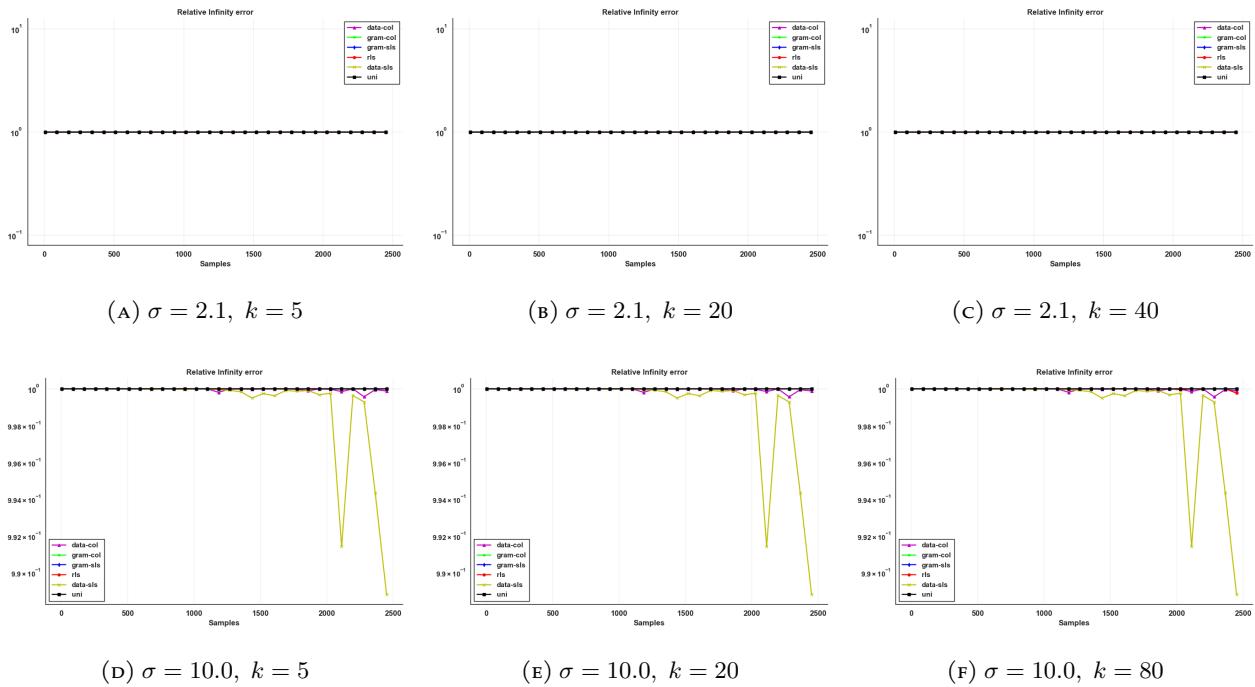


FIGURE 32. Nystrom infinity errors for the wine data set.

A.5. Nyström Errors.

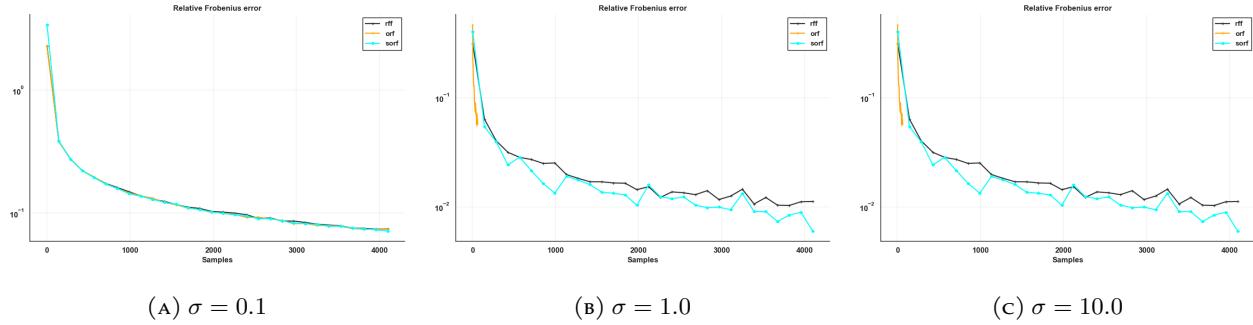


FIGURE 33. RFF Frobenius errors for the 3DSN data set.

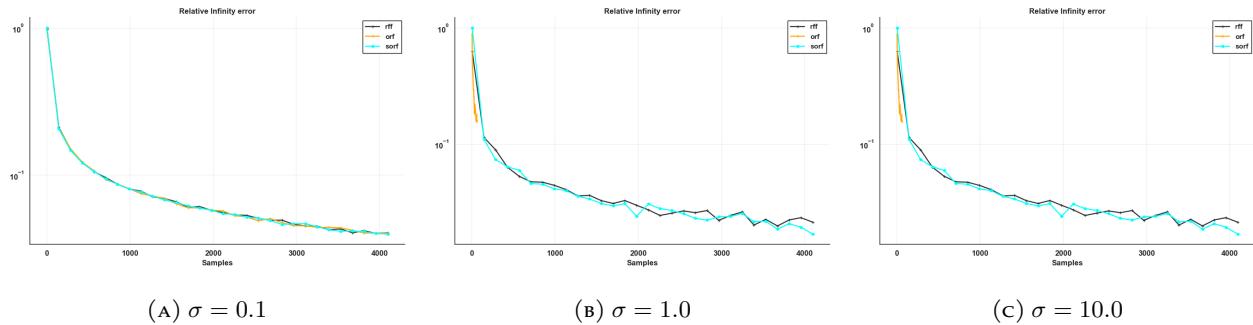


FIGURE 34. RFF absolute errors for the 3DSN data set.

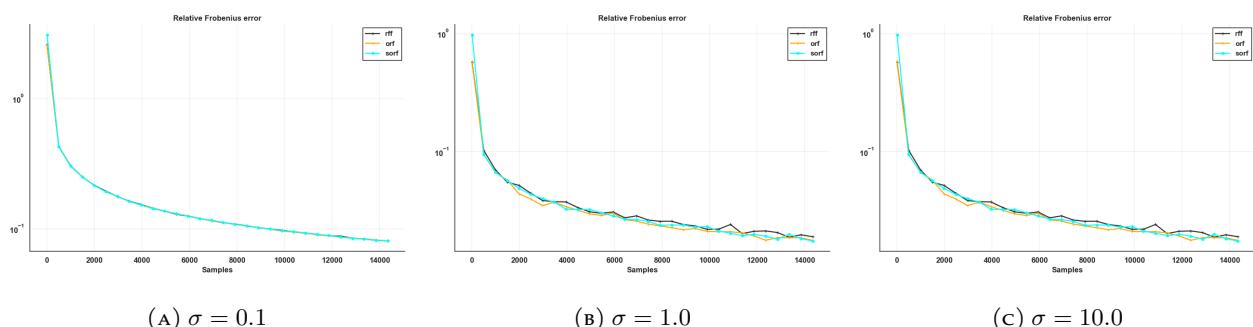


FIGURE 35. RFF Frobenius errors for the Abalone data set.

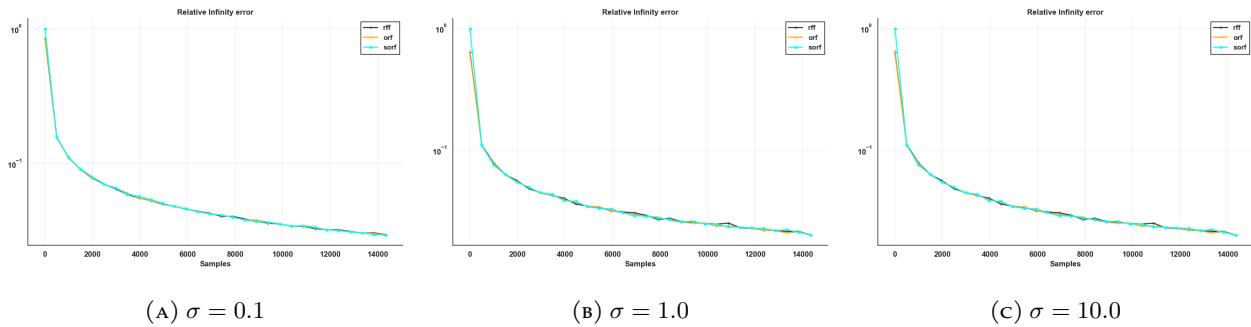


FIGURE 36. RFF infinity errors for the Abalone data set.

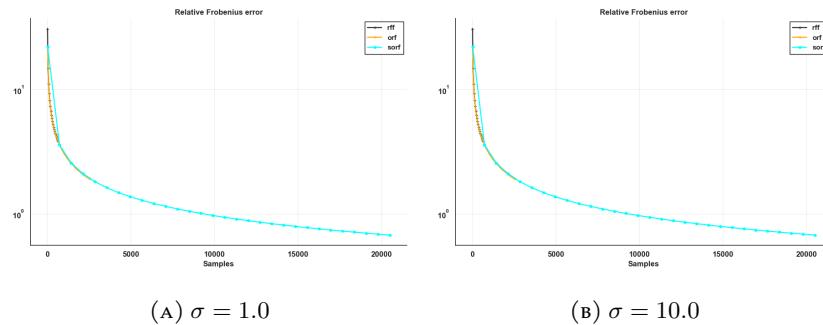


FIGURE 37. RFF Frobenius errors for the magic04 data set.

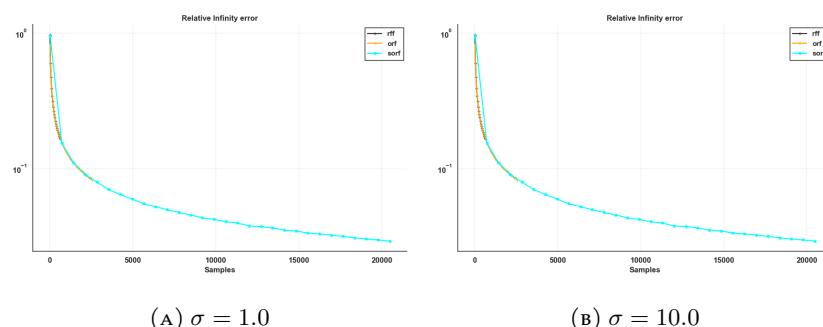


FIGURE 38. RFF infinity errors for the magic04 data set.

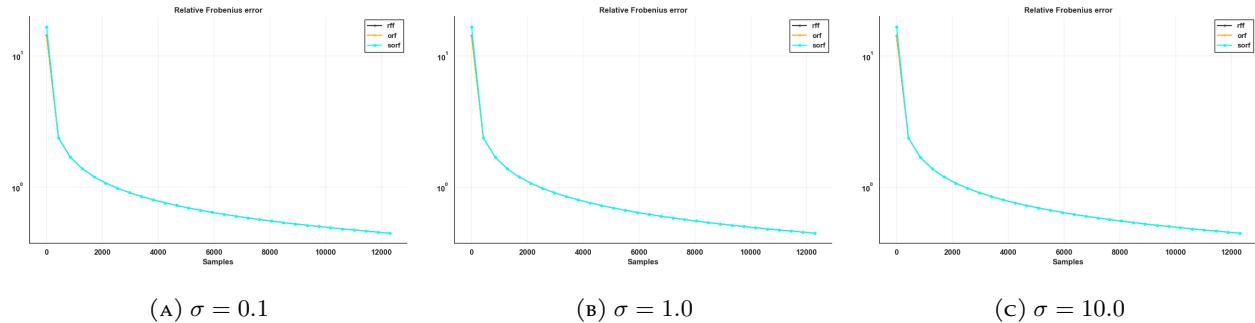


FIGURE 39. RFF Frobenius errors for the Stocks data set.

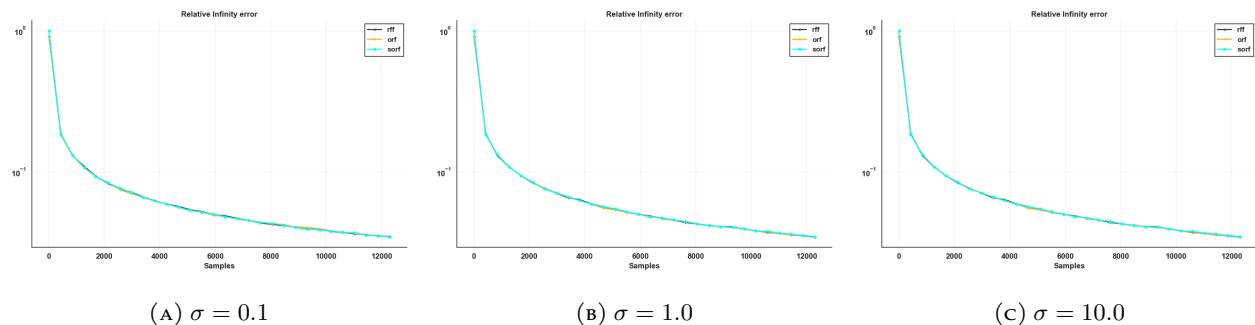


FIGURE 40. RFF infinity errors for the Stocks data set.

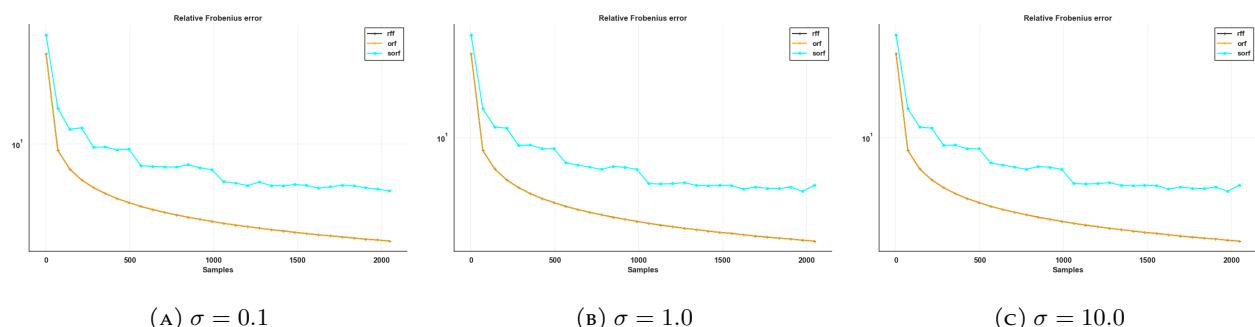


FIGURE 41. RFF Frobenius errors for the Temp data set.

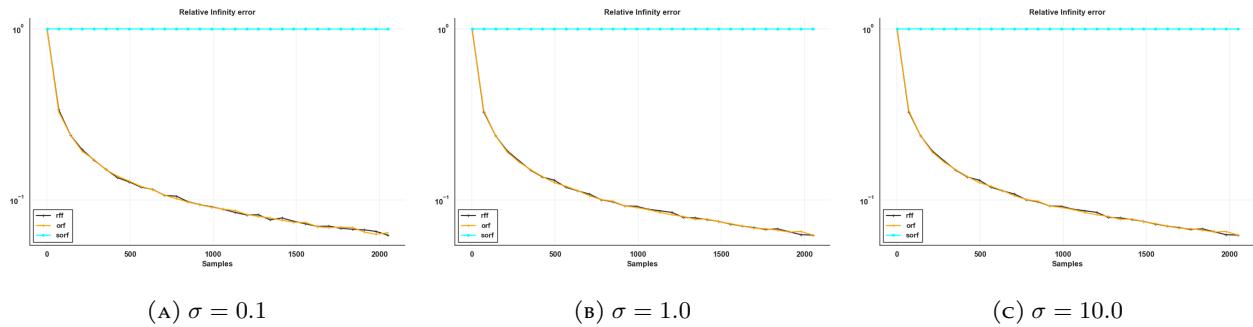


FIGURE 42. RFF infinity errors for the Temp data set.

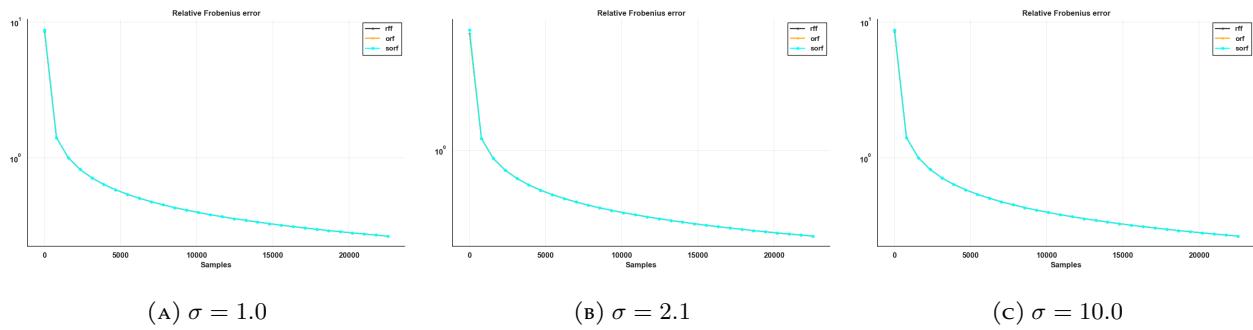


FIGURE 43. RFF Frobenius errors for the Wine data set.

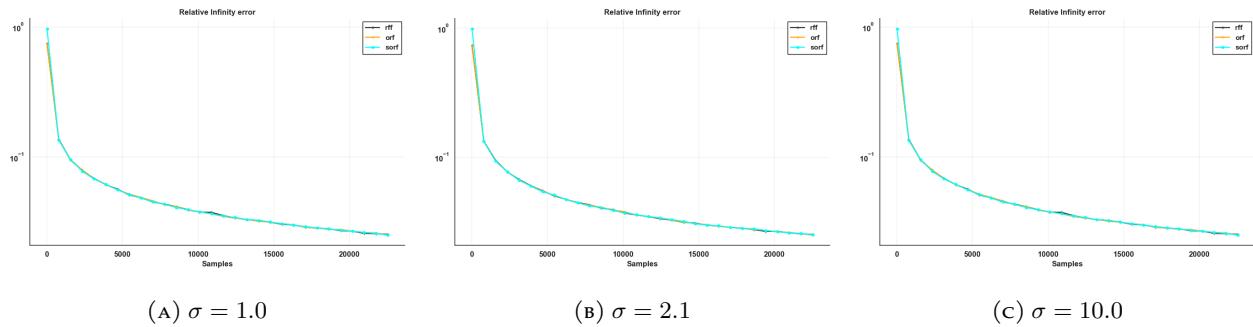


FIGURE 44. RFF infinity errors for the Wine data set.

A.6. Kernel Comparisons.

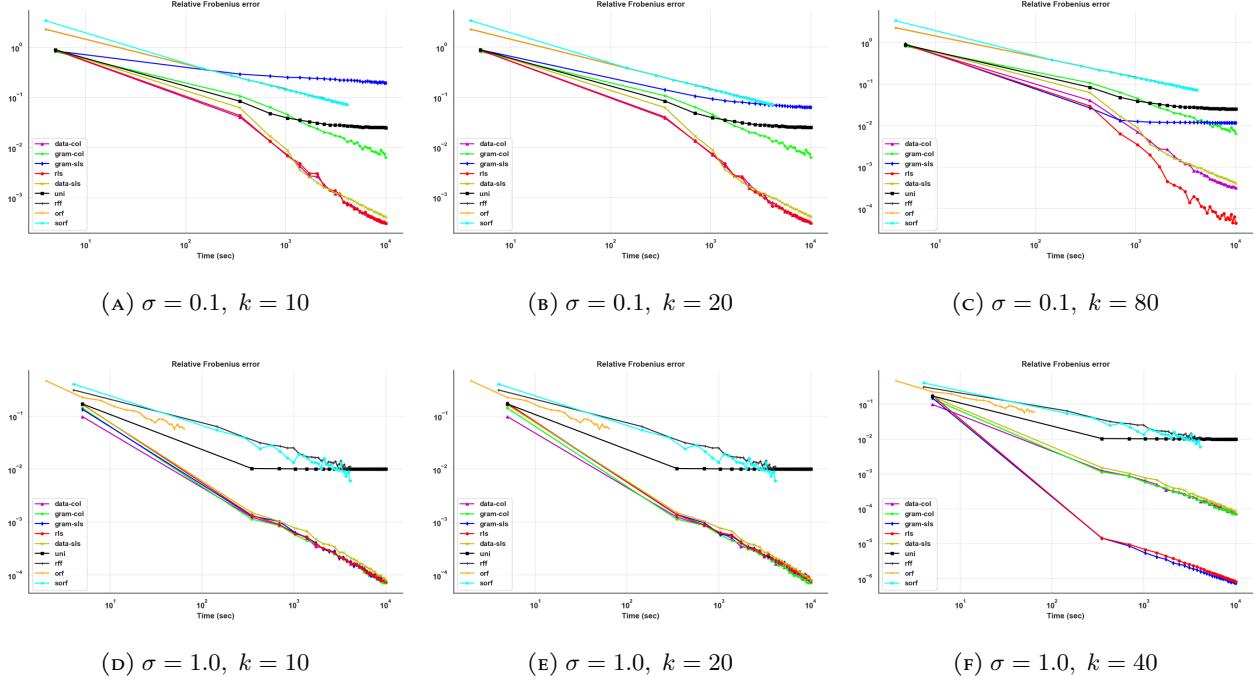


FIGURE 45. Comparing Frobenius errors of kernel methods for the 3DSN data set.

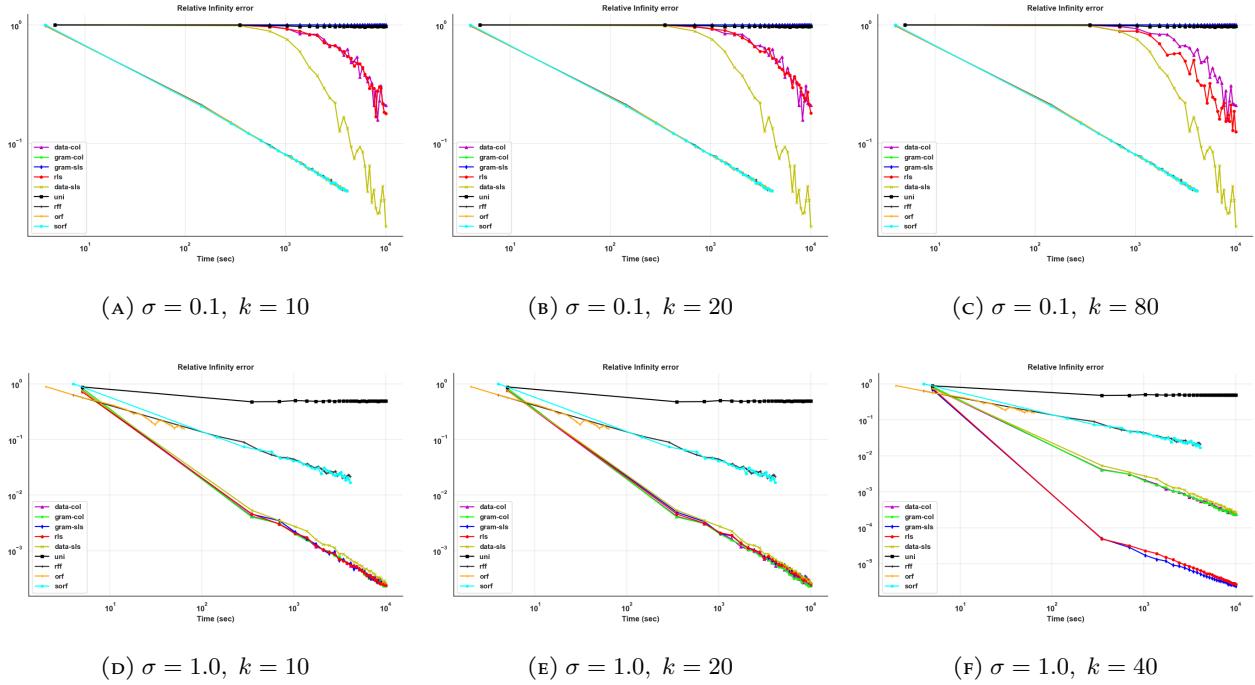


FIGURE 46. Comparing infinity errors of kernel methods for the 3DSN data set.

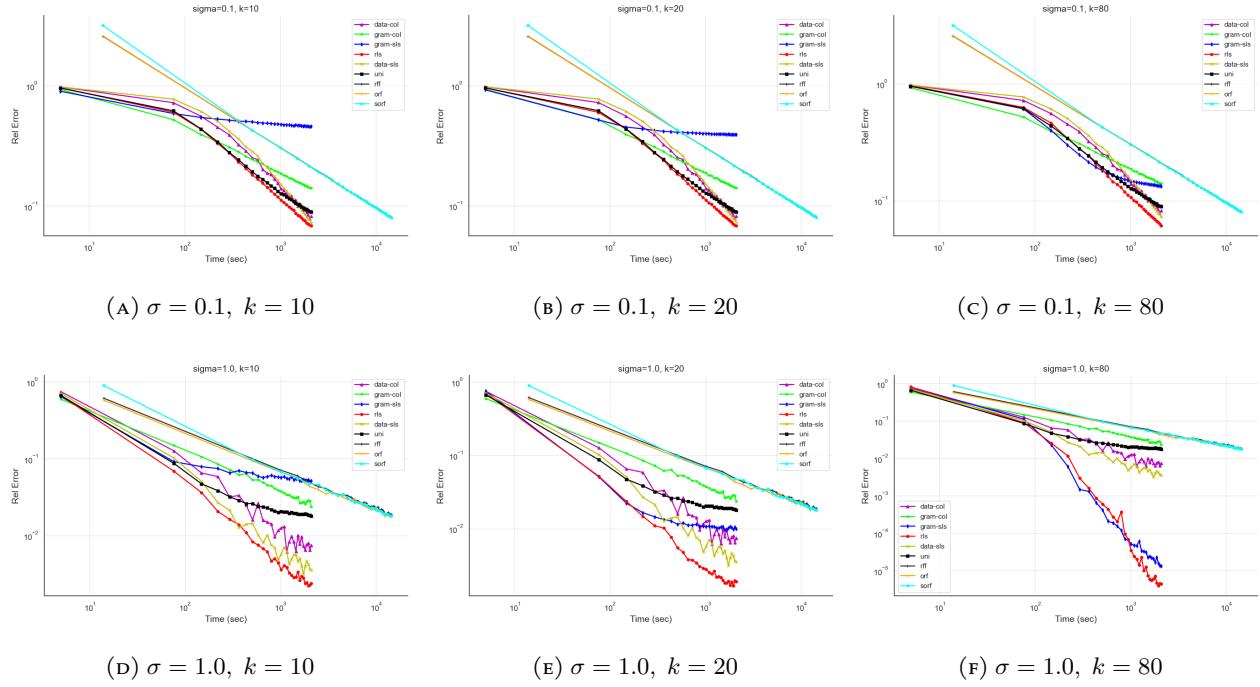


FIGURE 47. Comparing Frobenius errors of kernel methods for the Abalone data set.

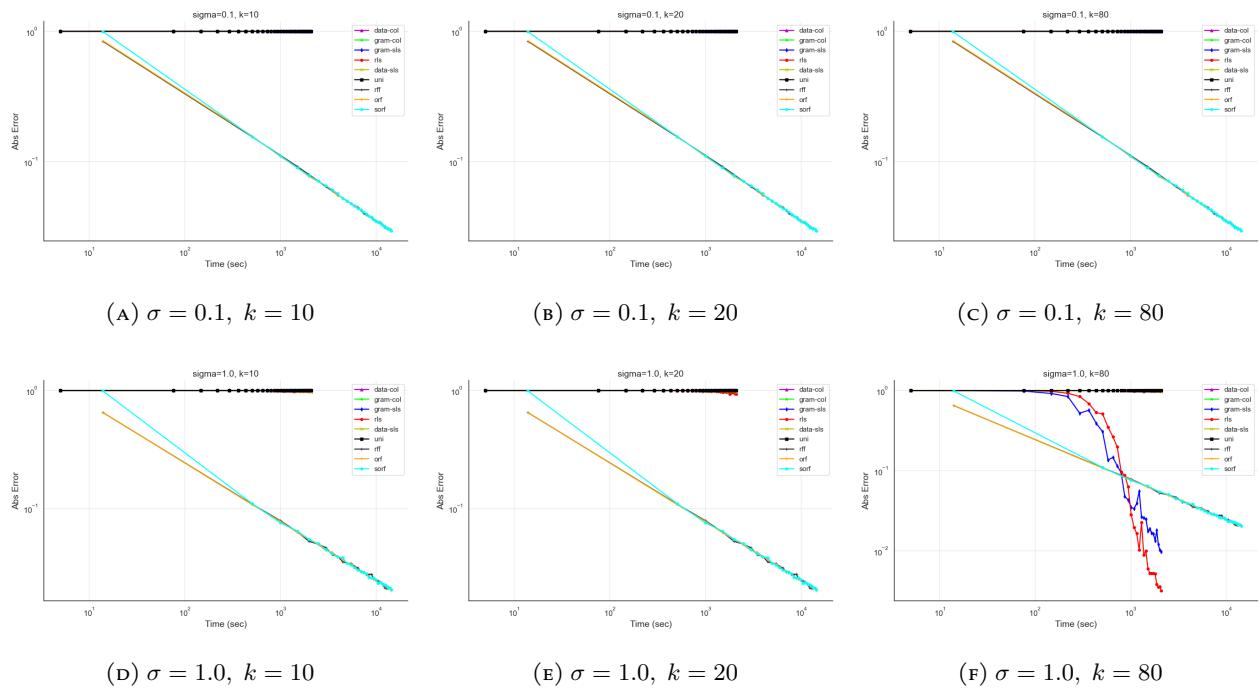


FIGURE 48. Comparing infinity errors of kernel methods for the Magic data set.

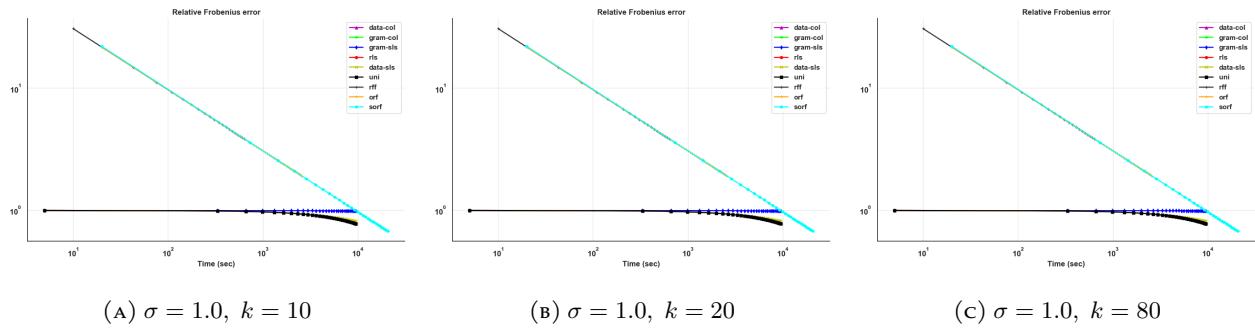


FIGURE 49. Comparing Frobenius errors of kernel methods for the Magic data set.

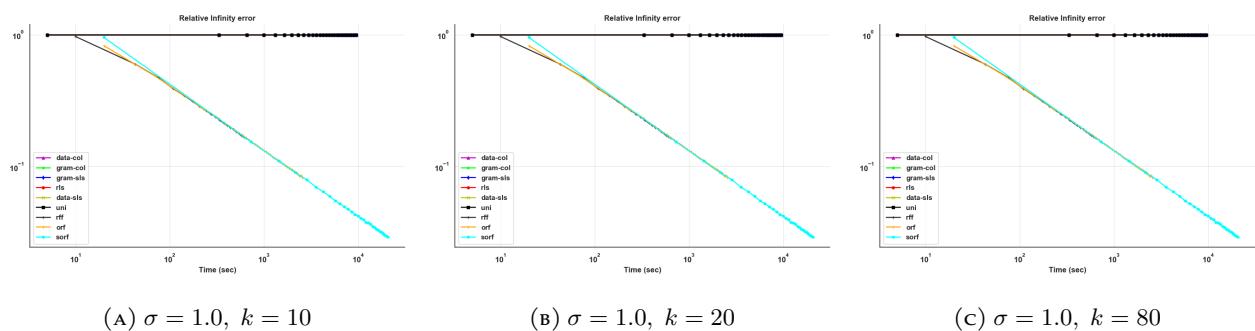


FIGURE 50. Comparing infinity errors of kernel methods for the Magic data set.

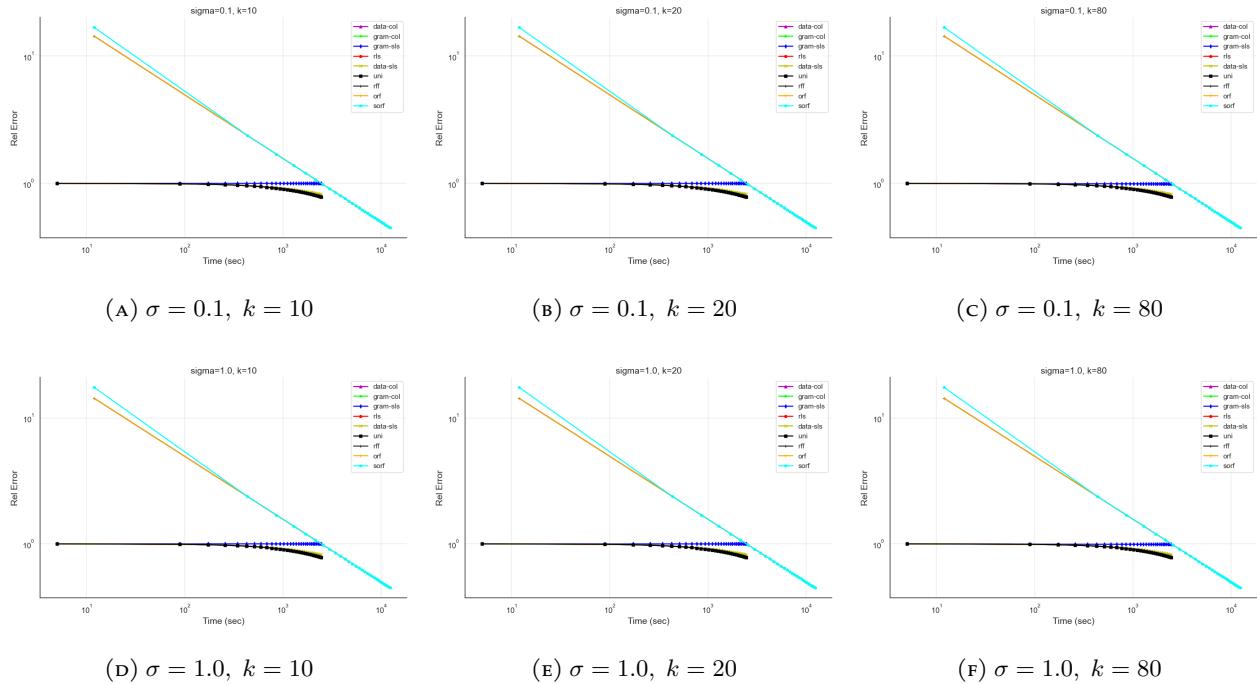


FIGURE 51. Comparing Frobenius errors of kernel methods for the Stocks data set.

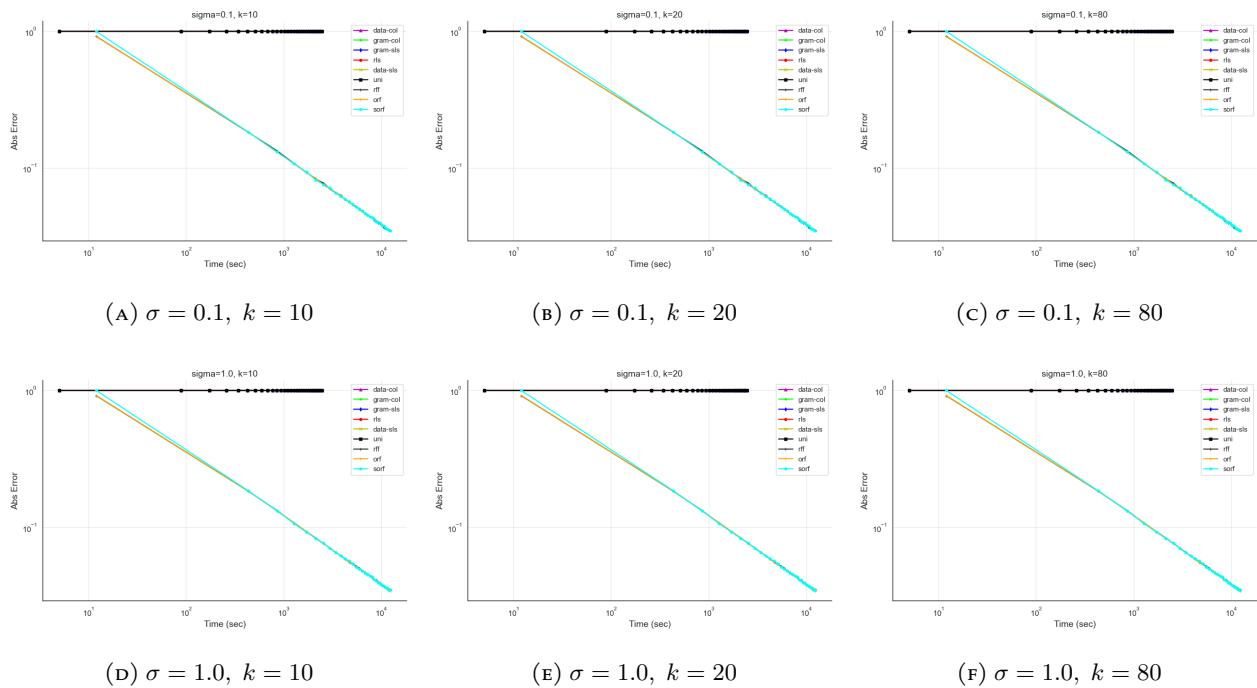


FIGURE 52. Comparing infinity errors of kernel methods for the Stocks data set.

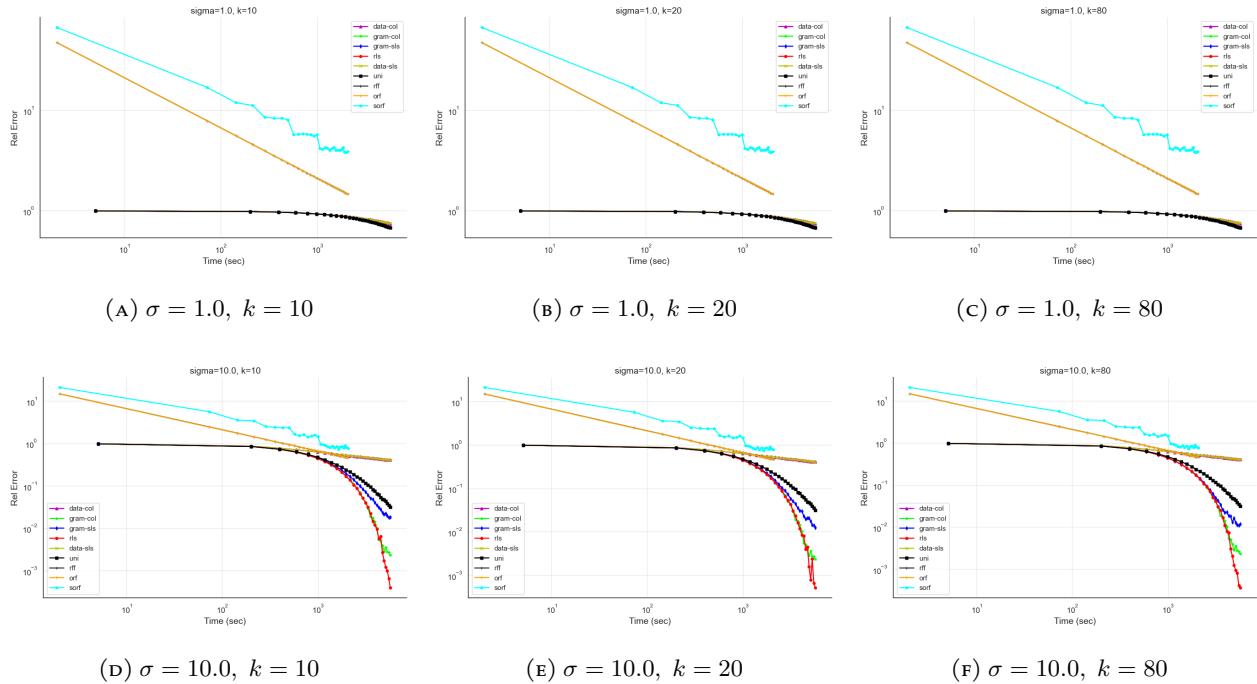


FIGURE 53. Comparing Frobenius errors of kernel methods for the Temp data set.

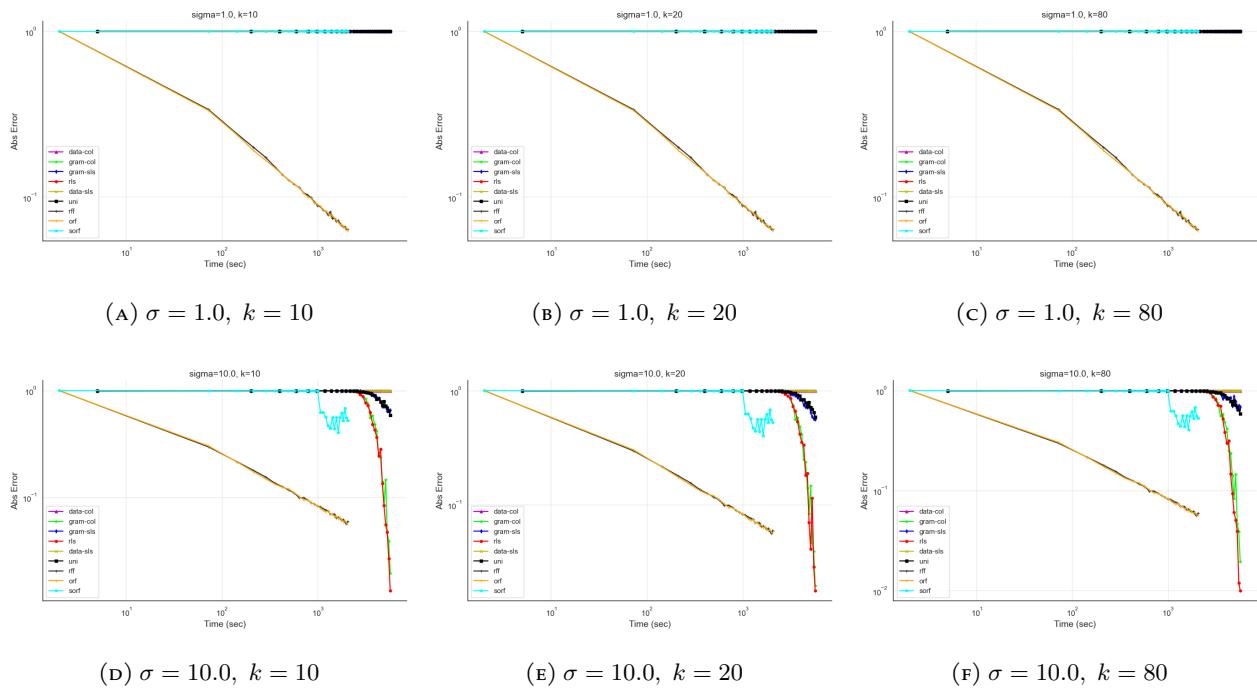


FIGURE 54. Comparing infinity errors of kernel methods for the Temp data set.

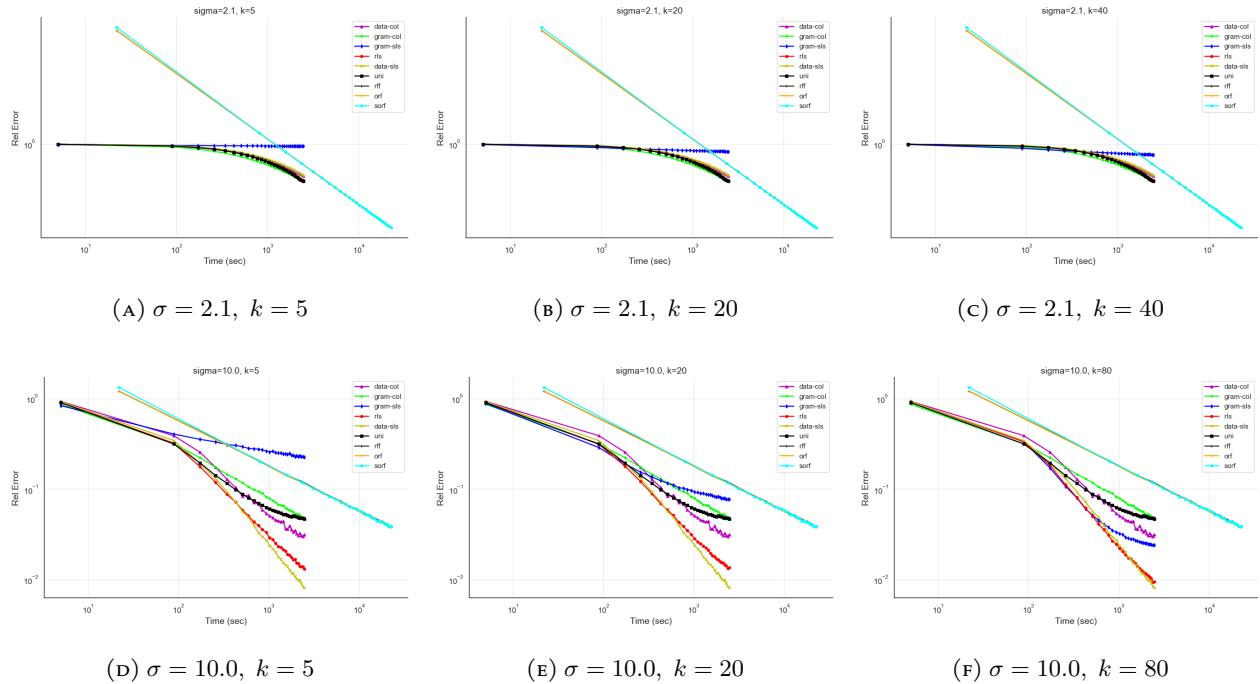


FIGURE 55. Comparing Frobenius errors of kernel methods for the wine data set.

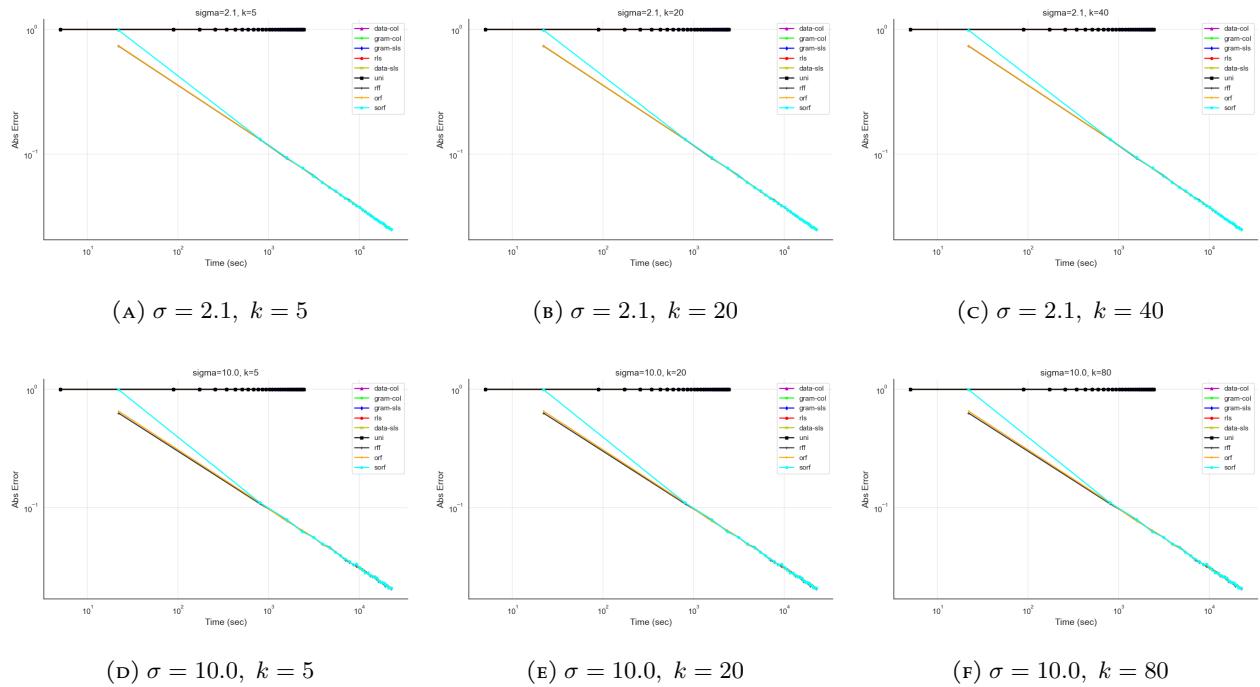


FIGURE 56. Comparing infinity errors of kernel methods for the wine data set.