

IJCNN 2021

Towards Unbiased Random Features with Lower Variance For Stationary Indefinite Kernels

Qin Luo, Kun Fang, Jie Yang, Xiaolin Huang

Institute of Image Processing and Pattern Recognition
Shanghai Jiao Tong University



INNS - International
Neural Network Society



IEEE - Computational
Intelligence Society

Contents

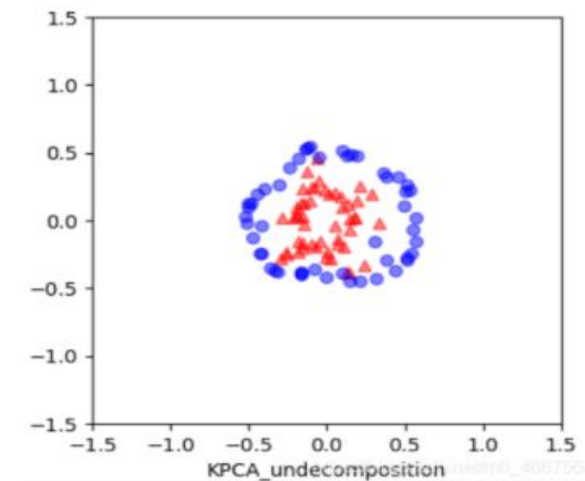
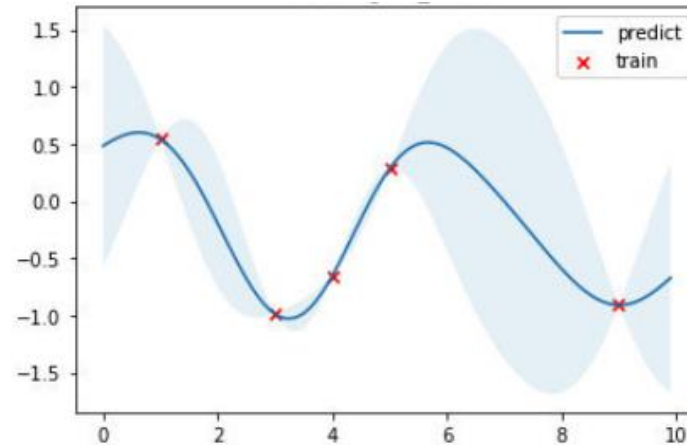
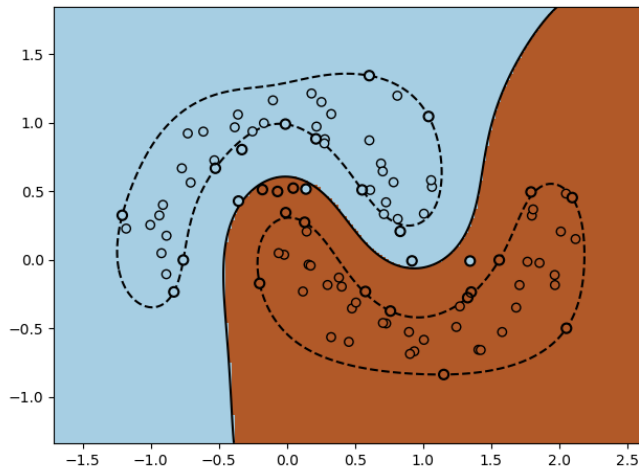
- Background
- Motivation
- Methods
- Experiments
- Conclusion

Contents

- Background
- Motivation
- Methods
- Experiments
- Conclusion

Background

- Kernel methods are extensively used in classification [1], regression [2] and dimension reduction [3].
- Kernel methods scale poorly to large datasets because of $\mathcal{O}(N^3)$ time complexity and $\mathcal{O}(N^2)$ space complexity
- Random Fourier Features reduce the computation cost and storage space to $\mathcal{O}(Ns^2)$ and $\mathcal{O}(Ns)$ ($s \ll N$)



- [1] Schölkopf B, Smola A J, Bach F. Learning with kernels: support vector machines, regularization, optimization, and beyond[M]. MIT press, 2002.
- [2] Wilson A, Adams R. Gaussian process kernels for pattern discovery and extrapolation[C]//International conference on machine learning. PMLR, 2013: 1067-1075.
- [3] Schölkopf B, Smola A, Müller K R. Kernel principal component analysis[C]//International conference on artificial neural networks. Springer, Berlin, Heidelberg, 1997: 583-588.

Background

- Random Fourier Features are restricted to the kernels:

1) **shift-invariant (stationary)**

$$k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y})$$

2) **positive definite (PD)**

$$\boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \geq 0 \text{ for all } \boldsymbol{\alpha} \neq \mathbf{0} \text{ and } K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

- Not satisfy the requirement

1) **non-stationary kernels**: polynomial kernel, neural tangent kernel (NTK)

When data is restricted on the sphere \rightarrow stationary but indefinite kernel

2) **non-PD kernel**: linear combination of Gaussian kernel (Delta-Gaussian kernel), TL1-kernel

(Bohner's Theorem) A continuous and **stationary** function $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is **positive definite** if and only if it can be represented as

$$k(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^d} \exp(i\mathbf{w}^T(\mathbf{x} - \mathbf{y})) p(d\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim p(\mathbf{w})} [\exp(i\mathbf{w}^T(\mathbf{x} - \mathbf{y}))]$$

where $p(\mathbf{w})$ is the positive finite measure over \mathbf{w} , i is imaginary unit.

Background

Related Work

Methods	Kernel Types	Unbiasedness	Variance
Random Maclaurin (RM) [1]	Polynomial	✓	$\mathcal{O}\left(\left(\frac{32RL}{\epsilon}\right)^{2d} \exp\left(-\frac{D\epsilon^2}{8C_\Omega^2}\right)\right)$
Tensor Sketch (TS) [2]	Polynomial	✓	$\mathcal{O}\left(\exp\left(-\frac{t\epsilon^2}{2R^{4p}}\right)\right)$
Spherical Random Features (SRF) [3]	Stationary Indefinite	✗	_____
Double Variation Random Features (DIGMM) [4]	Stationary Indefinite	✗	_____
Generalized Random Fourier Features (GRFF) [5]	Stationary Indefinite	✓	$\mathcal{O}\left(\left(\frac{2\sigma R}{\epsilon}\right)^{2d} \exp\left(-\frac{s\epsilon^2}{32(d+2)}\right)\right)$

[1] Kar P, Karnick H. Random feature maps for dot product kernels[C]//Artificial intelligence and statistics. PMLR, 2012: 583-591.

[2] Pham N, Pagh R. Fast and scalable polynomial kernels via explicit feature maps[C]//Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining. 2013: 239-247.

[3] Pennington J, Felix X Y, Kumar S. Spherical Random Features for Polynomial Kernels[C]//NIPS. 2015.

[4] Liu F, Huang X, Shi L, et al. A double-variational bayesian framework in random fourier features for indefinite kernels[J]. IEEE transactions on neural networks and learning systems, 2019, 31(8): 2965-2979.

[5] Liu F, Huang X, Chen Y, et al. Fast Learning in Reproducing Kernel Krein Spaces via Signed Measures[C]//International Conference on Artificial Intelligence and Statistics. PMLR, 2021: 388-396.

Contents

- Background
- **Motivation**
- Methods
- Experiments
- Conclusion

Motivation

Objective: Unbiased random Fourier approximation with lower variance for stationary indefinite kernels

Contribution:

- Unbiased approximation and lower the variance utilizing orthogonal sampling
- Theoretical analysis of the unbiasedness and variance reduction
- Experimental validation of the approximation error and classification or regression performance compared with the existing approximation method

Contents

- Background
- Motivation
- **Methods**
- Experiments
- Conclusion

Methods

Preliminaries

(Signed Measure) Let Ω be some set, and \mathcal{A} be a σ -algebra of subsets on Ω . A signed measure is a function $\mu: \mathcal{A} \rightarrow [-\infty, +\infty)$ or $(-\infty, +\infty]$ satisfying σ -additivity.

(Jordan Decomposition) Let μ be a signed measure defined on the σ -algebra \mathcal{A} . There exists two nonnegative measures μ_+ and μ_- (one of the measure) such that $\mu = \mu_+ - \mu_-$. The total mass is defined as: $||\mu|| = ||\mu_+|| + ||\mu_-||$

Methods

$p(\mathbf{w})$ is not a probability measure, viewed as a signed measure.

$$\begin{aligned} k(\mathbf{x} - \mathbf{y}) = k(\mathbf{z}) &= \int_{\mathbb{R}^d} \exp(i\mathbf{w}^T \mathbf{z}) p(d\mathbf{w}) = \int_{\mathbb{R}^d} \exp(i\mathbf{w}^T \mathbf{z}) p_+(d\mathbf{w}) - \int_{\mathbb{R}^d} \exp(i\mathbf{w}^T \mathbf{z}) p_-(d\mathbf{w}) \\ &= \|p_+\| \mathbb{E}_{\mathbf{w} \sim \widetilde{p}_+}(\exp(i\mathbf{w}^T \mathbf{z})) - \|p_-\| \mathbb{E}_{\mathbf{w} \sim \widetilde{p}_-}(\exp(i\mathbf{w}^T \mathbf{z})) \end{aligned}$$

where $\widetilde{p}_+ = \frac{p_+(\mathbf{w})}{\|p_+\|}$ and $\widetilde{p}_- = \frac{p_-(\mathbf{w})}{\|p_-\|}$

Define $\phi(\mathbf{x}) = \frac{1}{\sqrt{s}} [\psi_1(\mathbf{x}), \psi_2(\mathbf{x}), \psi_3(\mathbf{x}), \dots, \psi_s(\mathbf{x})]^T$ with $\psi_i(\mathbf{x})$:

$$\psi_i(\mathbf{x}) = \left[\sqrt{\|p_+\|} \cos(\mathbf{w}_i^T \mathbf{x}), \sqrt{\|p_+\|} \sin(\mathbf{w}_i^T \mathbf{x}), i\sqrt{\|p_-\|} \cos(\mathbf{v}_i^T \mathbf{x}), i\sqrt{\|p_-\|} \sin(\mathbf{v}_i^T \mathbf{x}) \right]^T$$

Generalized Random
Fourier Features

$$k(\mathbf{x} - \mathbf{y}) \approx \frac{1}{s} \sum_{i=1}^s \langle \psi_i(\mathbf{x}), \psi_i(\mathbf{y}) \rangle = \phi(\mathbf{x})^T \phi(\mathbf{y})$$

Methods

Unbiased Approximation

$$\mathbb{E}(K_{GRFF}(\mathbf{z})) = k(\mathbf{z})$$

Variance determines the whole approximation error. Orthogonal sampling could reduce the variance.

QR decomposition

- 1) Amplitude sampling $||\mathbf{w}_i||_2 \sim \widetilde{p_+(\mathbf{w})} \quad ||\mathbf{v}_i||_2 \sim \widetilde{p_-(\mathbf{w})}$
- 2) Orthogonal direction $\mathbf{a}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2m}) \quad \mathbf{b}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2m}) \quad \mathbf{M} = [\mathbf{a}_1, \dots, \mathbf{a}_m, \mathbf{b}_1, \dots, \mathbf{b}_m]$
 $M^{orth} = QR(\mathbf{M})$
- 3) Composition $\mathbf{w}_i = ||\mathbf{w}_i||_2 M_i^{orthn} \quad \mathbf{v}_i = ||\mathbf{v}_i||_2 M_{s+i}^{orthn}$

Methods

Variance Reduction

$$\begin{aligned}
 & \text{Orthogonality on } W_{pos} \quad \text{Mutual orthogonality on } W_{pos} \text{ and } W_{neg} \\
 & \text{Orthogonality on } W_{neg} \\
 & \text{Var}(K_{GORF}(\mathbf{z})) - \text{Var}(K_{GRFF}(\mathbf{z})) = \boxed{||p_+||^2 G_{\widetilde{k}_+}(\mathbf{z})} + \boxed{||p_-||^2 G_{\widetilde{k}_-}(\mathbf{z})} + \boxed{H(\mathbf{z})}
 \end{aligned}$$

Theorem 5 [1] For a PD radial kernel k on \mathbb{R}^d with Fourier measure $p(\mathbf{w})$ and $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, writing $\mathbf{z} = \mathbf{x} - \mathbf{y}$, we have:

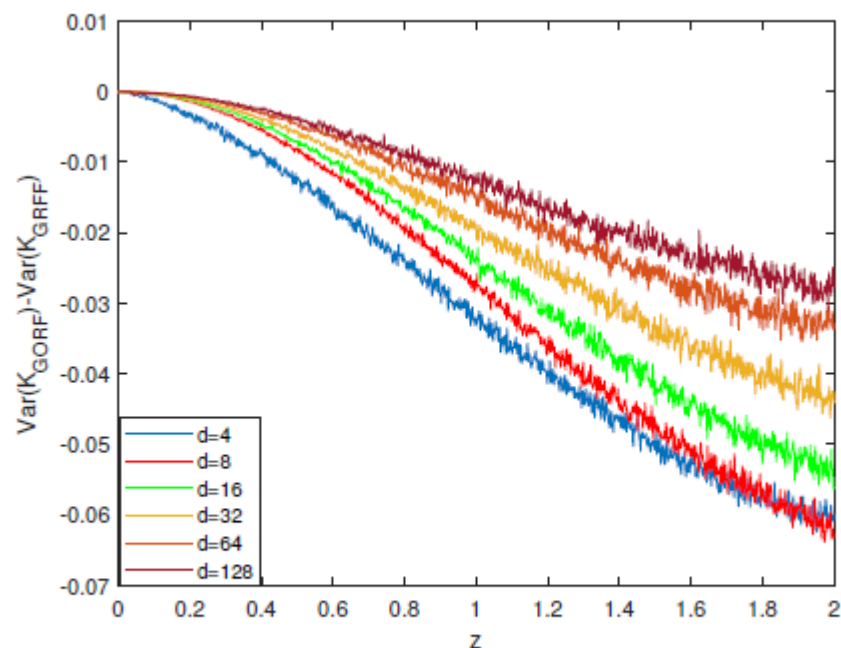
$$G_k(\mathbf{z}) = \text{Var}(K_{ORF}(\mathbf{z})) - \text{Var}(K_{RFF}(\mathbf{z})) = \frac{s-1}{s} \mathbb{E}_{R_1} \left[\frac{J_{\frac{d}{2}-1}(R_1 ||\mathbf{z}||) \Gamma\left(\frac{d}{2}\right)}{\left(\frac{R_1 ||\mathbf{z}||}{2}\right)^{\frac{d}{2}-1}} \right]^2 - \frac{s-1}{s} \mathbb{E}_{R_1, R_2} \left[\frac{J_{\frac{d}{2}-1}(\sqrt{R_1^2 + R_2^2} ||\mathbf{z}||) \Gamma\left(\frac{d}{2}\right)}{\left(\frac{\sqrt{R_1^2 + R_2^2} ||\mathbf{z}||}{2}\right)^{\frac{d}{2}-1}} \right]^2$$

where $R_1, R_2 \sim p(\mathbf{w})$, and J_α is the Bessel function of the first kind of degree α .

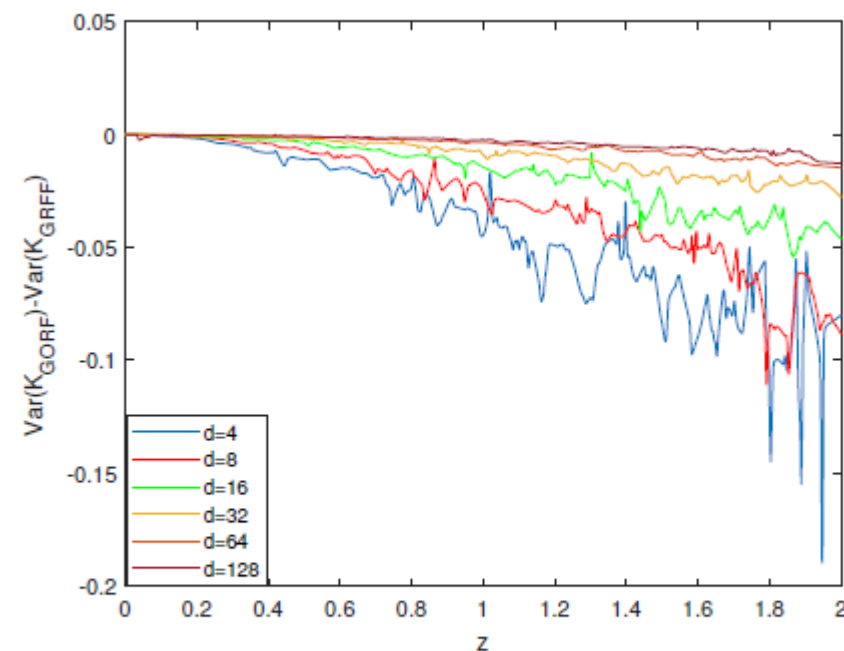
$$H(\mathbf{z}) = 2||p_+|| ||p_-|| [\mathbb{E}(a_1) \mathbb{E}(b_1) - \mathbb{E}(a_1 b_1)], \quad a_1 = \cos(\mathbf{w}_1^T \mathbf{z}), \quad b_1 = \cos(\mathbf{v}_1^T \mathbf{z})$$

Methods

Variance Reduction



(a) polynomial kernel on the unit sphere ($a = 3, m = 1$)



(b) delta-gaussian kernel ($a_1 = 1, a_2 = -1, \sigma_1 = 1, \sigma_2 = 10$)

Contents

- Background
- Motivation
- Methods
- Experiments
- Conclusion

Experiments

Setup

1) Kernels

Polynomial kernel

$$k(\mathbf{x}, \mathbf{y}) = \alpha(q + \langle \mathbf{x}, \mathbf{y} \rangle)^m = \left(1 - \frac{\|\mathbf{x} - \mathbf{y}\|^2}{a^2}\right)^m$$

where $q = a^2/2 - 1$, $\alpha = (2/a^2)^m$

Delta-Gaussian kernel

$$k(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \alpha_i e^{-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma_i^2}}$$

2) Datasets

Datasets	d	training	testing
<i>letter</i>	16	12000	6000
<i>ijcnn1</i>	22	49990	91701
<i>usps</i>	256	7291	2007

housing: d=13, training=405, testing=101

Experiments

Approximation error

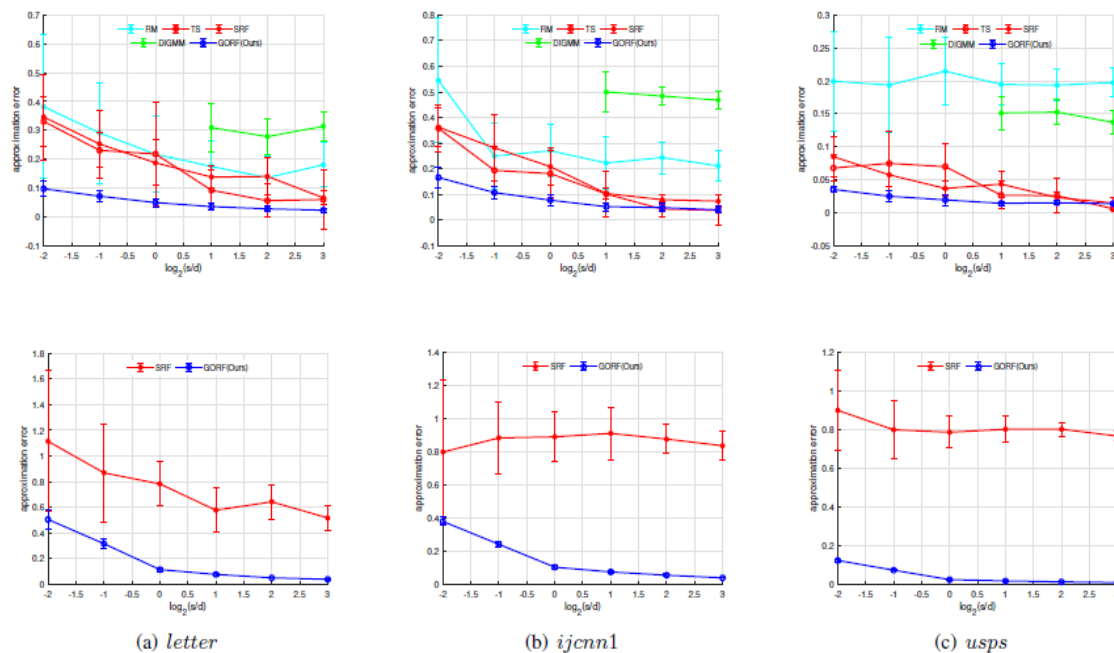


Fig. 2. Comparisons of various algorithms for kernel approximation in terms of approximation error across two typical stationary indefinite kernels and three datasets with different dimensions. Top: polynomial kernel on the unit sphere. Below: delta-gaussian kernel

TABLE II
COMPARISON RESULTS BETWEEN APPLYING ORTHOGONAL SAMPLING AND I.I.D SAMPLING ON STATIONARY INDEFINITE KERNELS IN TERMS OF APPROXIMATION ERROR (MEAN \pm STD.). THE LOWEST ERROR IS HIGHLIGHTED IN **BOLDFACE**

Kernel	DataSet	Method	s=1/2d	s=d	s=2d	s=8d
polynomial	letter	GRFF	0.0859 \pm 0.0309	0.0547 \pm 0.0078	0.0469 \pm 0.0109	0.0261 \pm 0.0059
		GORF	0.0716\pm0.0175	0.0495\pm0.0139	0.0360\pm0.0110	0.0231\pm0.0078
	ijcnn1	GRFF	0.1159 \pm 0.0158	0.0907 \pm 0.0194	0.0794 \pm 0.0142	0.0433 \pm 0.0059
		GORF	0.1072\pm0.0228	0.0775\pm0.0204	0.0487\pm0.0155	0.0397\pm0.0135
	usps	GRFF	0.0270 \pm 0.0056	0.0213 \pm 0.0063	0.0160 \pm 0.0029	0.0137 \pm 0.0020
		GORF	0.0251\pm0.0078	0.0194\pm0.0087	0.0143\pm0.0030	0.0137\pm0.0016
delta-gaussian	letter	GRFF	0.3918 \pm 0.0428	0.2736 \pm 0.0345	0.1887 \pm 0.0201	0.1017 \pm 0.0088
		GORF	0.3154\pm0.0424	0.1133\pm0.0181	0.0760\pm0.0090	0.0376\pm0.0039
	ijcnn1	GRFF	0.2924 \pm 0.0188	0.2171 \pm 0.0222	0.1504 \pm 0.0134	0.0757 \pm 0.0081
		GORF	0.2415\pm0.0190	0.1026\pm0.0129	0.0739\pm0.0065	0.0383\pm0.0022
	usps	GRFF	0.1005 \pm 0.0061	0.0690 \pm 0.0050	0.0500 \pm 0.0024	0.0253 \pm 0.0023
		GORF	0.0724\pm0.0049	0.0235\pm0.0009	0.0166\pm0.0008	0.0083\pm0.0003

Unbiased estimation + Lower variance = Lower approximation error

Experiments

SVM classification problem

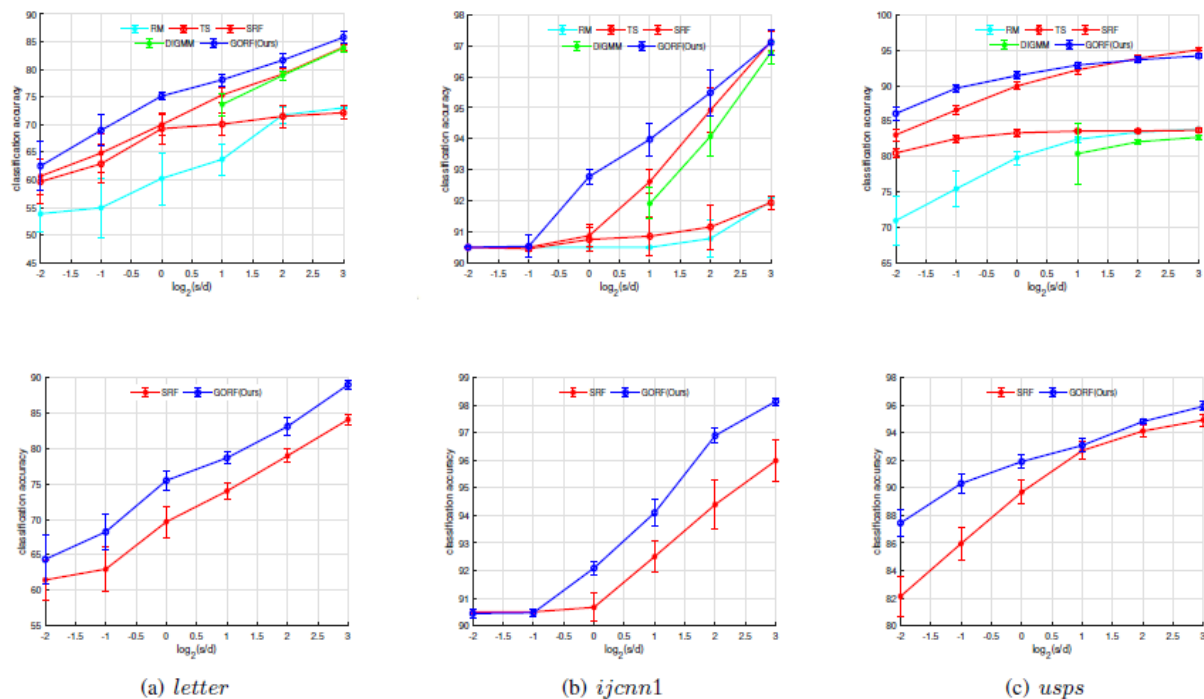


Fig. 3. Comparisons of various algorithms for SVM classification task in terms of accuracy across two typical stationary indefinite kernels and three datasets with different dimensions. Top: polynomial kernel on the unit sphere. Below: delta-gaussian kernel

SVR regression problem

TABLE III
REGRESSION ERROR FOR KERNEL APPROXIMATION METHODS ON POLYNOMIAL KERNEL AND *HOUSING* DATASET (RMSE: MEAN \pm STD). THE LOWEST ERROR IS HIGHLIGHTED IN **BOLDFACE**

Methods	s=2d	s=4d	s=8d
RM	7.153 \pm 1.772	5.436 \pm 0.917	4.491 \pm 0.008
TS	5.414 \pm 0.879	4.772 \pm 0.177	4.657 \pm 0.316
SRF	4.391 \pm 0.368	3.906 \pm 0.219	3.555 \pm 0.130
DIGMM	4.897 \pm 0.368	4.130 \pm 0.324	4.000 \pm 0.475
GORF(OURS)	4.079\pm0.233	3.817\pm0.204	3.472\pm0.137

TABLE IV
REGRESSION ERROR FOR KERNEL APPROXIMATION METHODS ON DELTA-GAUSSIAN KERNEL AND *HOUSING* DATASET (RMSE: MEAN \pm STD). THE LOWEST ERROR IS HIGHLIGHTED IN **BOLDFACE**

Methods	s=2d	s=4d	s=8d
SRF	5.432 \pm 0.729	3.845 \pm 0.379	3.321 \pm 0.274
GORF(OURS)	3.739\pm0.360	3.474\pm0.330	3.164\pm0.452

Contents

- Background
- Motivation
- Methods
- Experiments
- **Conclusion**

Conclusion

- 1) Propose an unbiased random feature approximation with lower variance for stationary indefinite kernels
- 2) Verify the unbiasedness and numerically calculate the reduced variance.
- 3) Experimentally demonstrate the approximation error and performance in classification and regression task compared with other methods.