

Learning with invariances in random features and kernel models

Song Mei¹ Theodor Misiakiewicz² Andrea Montanari^{2,3}

¹Department of Statistics, UC Berkeley

²Department of Statistics, Stanford University

³Department of Electrical Engineering, Stanford University

Introduction

- In many learning tasks, the **data present some natural symmetries** (e.g., labels are invariant under translation of the images in image recognition tasks).
- One important goal of machine learning has been to design **predictive models that take advantage of these symmetries** to make a more efficient use of data.
- For instance, **convolutional networks** are believed to owe their success to their ability to encode translation invariance.
- Empirically, models that exploit invariances perform better than models that do not.

Focus of this work:

Quantifying the performance gain of using invariant architectures over non-invariant ones in random features and kernel models.

Setting and models

- Data:** $\mathbf{x} \sim \text{Unif}(\mathcal{A}_d)$, $\mathcal{A}_d = \mathbb{S}^{d-1}(\sqrt{d})$ or $\mathcal{A}_d = \{-1, +1\}^d$.
- Invariance group:** \mathcal{G}_d subgroup of orthogonal group $\mathcal{O}(d)$ (that preserves the hypercube if $\mathcal{A}_d = \{-1, +1\}^d$).
- Goal:** learn a \mathcal{G}_d -invariant function f_*

$$\text{i.e., } f_*(g \cdot \mathbf{x}) = f_*(\mathbf{x}) \text{ for all } g \in \mathcal{G}_d,$$

given iid samples $\{y_i, \mathbf{x}_i\}_{i \leq n}$ with $\mathbf{x}_i \sim_{iid} \text{Unif}(\mathcal{A}_d)$ and

$$y_i = f_*(\mathbf{x}_i) + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i] = 0, \quad \mathbb{E}[\varepsilon_i^2] \leq \tau^2.$$

Models:

- Random Features models:** $(\sqrt{d}\mathbf{w}_i) \sim_{iid} \text{Unif}(\mathcal{A}_d)$ fixed,

$$\hat{f}_{\text{RF}}(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$$

$$\rightarrow \hat{f}_{\text{RF}}^{\text{inv}}(\mathbf{x}; \mathbf{a}) = \sum_{j=1}^N a_j \int_{\mathcal{G}_d} \sigma(\langle \mathbf{w}_j, g \cdot \mathbf{x} \rangle) \pi_d(dg),$$

where π_d is the Haar measure on the group \mathcal{G}_d .

Fit the coefficients with Ridge Regression (RFRR):

$$\hat{\mathbf{a}}^{\text{inv}}(\lambda) = \arg \min_{\mathbf{a} \in \mathbb{R}^N} \left\{ \sum_{i=1}^n (y_i - \hat{f}_{\text{RF}}^{\text{inv}}(\mathbf{x}_i; \mathbf{a}))^2 + N\lambda \|\mathbf{a}\|_2^2 \right\}.$$

- Kernel models:** inner-prod. kernel $H(\mathbf{x}, \mathbf{z}) = h\left(\frac{\langle \mathbf{x}, \mathbf{z} \rangle}{d}\right)$,

$$\rightarrow H^{\text{inv}}(\mathbf{x}, \mathbf{z}) = \int_{\mathcal{G}_d} h(\langle \mathbf{x}, g \cdot \mathbf{z} \rangle / d) \pi_d(dg).$$

Fit the function with Kernel Ridge Regression (KRR):

$$\hat{f}_\lambda^{\text{inv}} = \arg \min_{\hat{f} \in \mathcal{H}^{\text{inv}}} \left\{ \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2 + \lambda \|\hat{f}\|_{\mathcal{H}^{\text{inv}}}^2 \right\}.$$

Example: 2-layer CNN

- The cyclic group $\mathcal{G}_d = \{g_0, g_1, \dots, g_{d-1}\}$:
- $$g_i \cdot \mathbf{x} = (x_{d-i+1}, x_{d-i+2}, \dots, x_d, x_1, x_2, \dots, x_{d-i}).$$

- Cyclic invariant RF model:

$$f_{\text{RF}}^{\text{inv}}(\mathbf{x}; \mathbf{a}) = \frac{1}{d} \sum_{j=1}^N a_j \sum_{k=1}^d \sigma(\langle \mathbf{w}_j, g_k \cdot \mathbf{x} \rangle).$$

2-layers CNN with global average pooling & filters $\mathbf{w}_j \in \mathbb{R}^d$.

- Inner-prod. kernel: NTK of fully-connected NNs; vs Cyclic invariant kernel: NTK of 2-layer CNN with global pooling. \rightarrow performance gap FC-NN vs. CNN in kernel regime.

Degeneracy of the invariance group

- Identify **degeneracy** of \mathcal{G}_d as the measure of the approx. and generalization power gain of using invariant models.
- $V_{d,k}$: subspace of degree- k polynomials orthogonal to degree- $(k-1)$ polynomials in $L^2(\mathcal{A}_d)$.
- $V_{d,k}(\mathcal{G}_d)$: subspace of $V_{d,k}$ of \mathcal{G}_d -invariant polynomials.

Groups of degeneracy $\alpha \in \mathbb{R}_{>0}$

\mathcal{G}_d has degeneracy α if for any $k \geq \alpha$, we have $\dim(V_{d,k}) / \dim(V_{d,k}(\mathcal{G}_d)) \asymp d^\alpha$.

- d^α : ‘**effective dimension**’ of the action of the group.
 - $\alpha = 1$ for cyclic group.
 - Not necessarily equal to the size of \mathcal{G}_d : e.g., translation invariance on band-limited signals $\text{Sft}_d = \{g_u, u \in [0, 1]\}$
- $$g_u \cdot \mathbf{x} = (x_1, \cos(2\pi u)x_2 + \sin(2\pi u)x_3, \dots).$$
- Sft_d has degeneracy $\alpha = 1$.

Counting invariant polynomials

- $\{Y_{ks}\}_{s \leq B_{d,k}}$ orthonormal basis of $V_{d,k}$ ($B_{d,k} = \dim(V_{d,k})$).
- $\{\bar{Y}_{ks}\}_{s \leq D_{d,k}}$ orth. basis of $V_{d,k}(\mathcal{G}_d)$ ($D_{d,k} = \dim(V_{d,k}(\mathcal{G}_d))$).
- Gegenbauer polynomial on \mathcal{A}_d of degree- k :

$$Q_k(\langle \mathbf{x}, \mathbf{z} \rangle) = \frac{1}{B_{d,k}} \sum_{s \leq B_{d,k}} Y_{ks}(\mathbf{x}) Y_{ks}(\mathbf{z}).$$

Representation lemma

Lemma 1 ([3]) We have

$$\frac{1}{D_{d,k}} \sum_{s \leq D_{d,k}} \bar{Y}_{ks}(\mathbf{x}) \bar{Y}_{ks}(\mathbf{z}) = \frac{B_{d,k}}{D_{d,k}} \int_{\mathcal{G}_d} Q_k(\langle \mathbf{x}, g \cdot \mathbf{z} \rangle) \pi_d(dg).$$

- To compute degeneracy, it is sufficient to show for all $k \geq \alpha$:

$$\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathcal{A}_d)} \left[\int_{\mathcal{G}_d} Q_k(\langle \mathbf{x}, g \cdot \mathbf{z} \rangle) \pi_d(dg) \right] = \frac{D_{d,k}}{B_{d,k}} = \Theta_d(d^{-\alpha}).$$

Test error of invariant models

- Let f_* be \mathcal{G}_d -invariant with \mathcal{G}_d group of degeneracy α .
- Test error with square loss:

$$R(f_*, \mathbf{X}, \mathbf{W}, \lambda) = \mathbb{E}_{\mathbf{x}} \left\{ (f_*(\mathbf{x}) - \hat{f}_{\text{RF}}(\mathbf{x}; \hat{\mathbf{a}}(\lambda)))^2 \right\}.$$

Test error of RFRR

Theorem 1 ([3]) Assume $\max(N/n, n/N) \geq d^\delta$ and $\lambda = O_d(1 \vee (N/n))$, σ satisfies some conditions, then:

- (Standard RF) If $d^{\ell+\delta} \leq \min(N, n) \leq d^{\ell+1-\delta}$,
- $$R(f_*, \mathbf{X}, \mathbf{W}, \lambda) = \|\mathbb{P}_{>\ell} f_*\|_{L^2}^2 + o_d(\mathbb{P}(\cdot)).$$
- (Invariant RF) If $d^{\ell+\delta}/d^\alpha \leq \min(N, n) \leq d^{\ell+1-\delta}/d^\alpha$,
- $$R^{\text{inv}}(f_*, \mathbf{X}, \mathbf{W}, \lambda/d^\alpha) = \|\mathbb{P}_{>\ell} f_*\|_{L^2}^2 + o_d(\mathbb{P}(\cdot)).$$

$\mathbb{P}_{>\ell}$: projection orthogonal to the space of degree- ℓ polynomials.

- RFRR learns the best degree- ℓ polynomial approx. to f_* .
- Same result for KRR as above with $N = \infty$: invariant KRR saves a factor d^α in sample size compared to standard KRR.

Invariant RF saves a factor d^α in sample size and number of hidden units to achieve same test error as std. RF.

Assumptions on σ

- Results: **consequence of a general framework** in [2].
- Technical general conditions of [2] checked for
 - Cyclic group and σ assumed $(\ell+1)$ -differentiable.
 - General groups of degeneracy α and σ polynomial.
- Deferred weaker conditions to future work (if σ diff., our proof techniques generalize to subgroups of permutations).

Symmetrization and data augmentation

Compare 4 approaches to learning invariant models:

- Standard KRR:** with inner-prod. kernel.
- Invariant KRR:** with invariant kernel (‘intrinsic approach’: invariance directly enforced in the model).
- Output symmetrization:** take \hat{f}_λ solution of standard KRR and symmetrize it:

$$\hat{f}_\lambda^{\text{inv}}(\mathbf{x}) := \int_{\mathcal{G}_d} \hat{f}_\lambda(g \cdot \mathbf{x}) \pi_d(dg).$$

Does not significantly improve on standard KRR.

- Full data augmentation:** add $\{(y_i, g \cdot \mathbf{x}_i)\}_{i \leq n, g \in \mathcal{G}_d}$ to the training set with standard KRR. \Rightarrow this is equivalent to invariant KRR [1].

Test errors: (b) = (d) \ll (c) \approx (a).

Sketch of the proof for KRR

- Inner-prod. kernels have eigenspaces $V_{d,k}$:

$$H_d(\mathbf{x}, \mathbf{z}) = \sum_{k=0}^{\infty} \xi_{d,k}^2 \sum_{s \leq B_{d,k}} Y_{ks}(\mathbf{x}) Y_{ks}(\mathbf{z}).$$

- Space $V_{d,k}$ is **preserved under the action of \mathcal{G}_d** :

$$H_d^{\text{inv}}(\mathbf{x}, \mathbf{z}) = \sum_{k=0}^{\infty} \xi_{d,k}^2 \sum_{s \leq D_{d,k}} \bar{Y}_{ks}(\mathbf{x}) \bar{Y}_{ks}(\mathbf{z}).$$

- H_d^{inv} has the same eigenvalues $\xi_{d,k}^2 = \Theta_d(d^{-k})$ as H_d but with degeneracy lower by a factor $B_{d,k}/D_{d,k} = \Theta_d(d^\alpha)$.

- Theorem [2]:** kernel eigenval. $\{\lambda_{d,k}\}$ in decreas. order, eigenvect. $\{\psi_k\}$ + technical conditions. Let $m \in \mathbb{N}$ s.t.

$$\lambda_{d,m+1} \cdot n^{1+\delta} \leq \sum_{k \geq m+1} \lambda_{d,k}, \quad m \leq n^{1-\delta}.$$

Define $s^{\text{eff}} = \lambda + \sum_{k \geq m+1} \lambda_k$.

Then **KRR acts as shrinkage operator**, i.e.,

$$\hat{f}_\lambda(\mathbf{x}) \approx \sum_{k \geq 1} \frac{\lambda_{d,k}}{\lambda_{d,k} + s^{\text{eff}}/n} \cdot \langle f_*, \psi_k \rangle_{L^2} \cdot \psi_k(\mathbf{x}).$$

Learn eigendirection if $\lambda_{d,k} \gg \frac{s^{\text{eff}}}{n}$, not at all if $\lambda_{d,k} \ll \frac{s^{\text{eff}}}{n}$.

- Std KRR:** $m = \#\{Y_{ks}\}_{k \leq \ell} = \Theta_d(d^\ell)$, $s^{\text{eff}} = \Theta_d(1)$. Hence if $d^{\ell+\delta} \leq n \leq d^{\ell+1-\delta}$, learns degree- ℓ polynomial approx.
- Inv KRR:** $m = \#\{\bar{Y}_{ks}\}_{k \leq \ell} = \Theta_d(d^{\ell-\alpha})$, $s^{\text{eff}} = \Theta_d(d^{-\alpha})$. Hence if $d^{\ell-\alpha+\delta} \leq n \leq d^{\ell+1-\alpha-\delta}$, learns degree- ℓ approx.

Numerical illustration

- Data $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ with $d = 30$.
- Target functions invariant w.r.t. cyclic group.
- Degeneracy $\alpha = 1$, hence save factor d in sample size.
- $f_{\text{lin}} = \sum_{i \leq d} x_i$, $f_{\text{quad}} = \sum_{i \leq d} x_i x_{i+1}$, $f_{\text{cube}} = \sum_{i \leq d} x_i x_{i+1} x_{i+2}$.

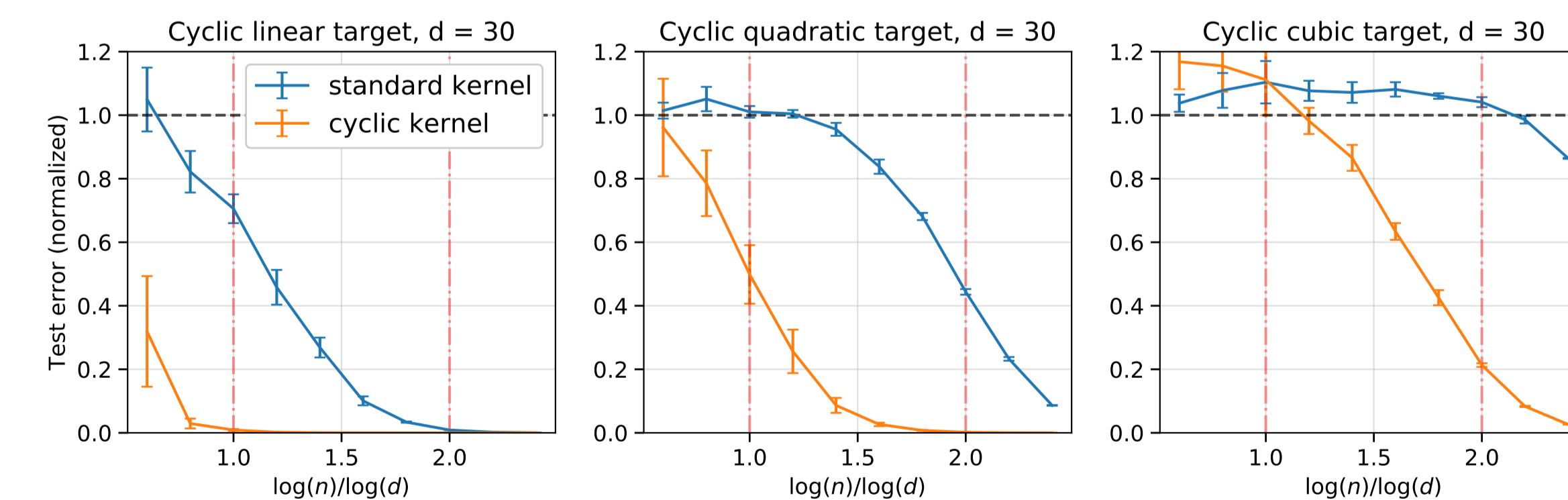


Figure 1: Normalized test error of KRR with cyclic vs standard kernels.

Bibliography

- Z. Li, R. Wang, D. Yu, S. S. Du, W. Hu, R. Salakhutdinov, and S. Arora. Enhanced convolutional neural tangent kernels. *arXiv preprint arXiv:1911.00809*, 2019.
- S. Mei, T. Misiakiewicz, and A. Montanari. Generalization error of random features and kernel methods: hypercontractivity and kernel matrix concentration. *arXiv preprint arXiv:2101.10588*, 2021.
- S. Mei, T. Misiakiewicz, and A. Montanari. Learning with invariances in random features and kernel models. *arXiv preprint arXiv:2102.13219*, 2021.