

1. Motivation

Question: How do we quantify statistical difficulty or target-kernel alignment after the kernel has been learned?

Classical rates $n^{-s\beta/(1+s\beta)}$ for spectral algorithms with kernel $\mathbf{k}(x, x') = \sum_j \lambda_j \phi_j(x) \phi_j(x')$ usually assume the polynomial eigen-decay and source conditions:

$$\lambda_j \asymp j^{-\beta}, \quad \sum_j \frac{\langle f^*, \phi_j \rangle^2}{\lambda_j^s} \leq R_s.$$

Classical assumptions may fail when the kernel evolves over time: the learned spectrum may not decay polynomially and the index s may be unclear.

Goal: A complexity measure that stays valid and informative if the kernel is learned from data.

2. Effective Span Dimension

A sequence model with zero mean & uncorrelated noise:

$$Z_j = \underbrace{\theta_j^*}_{\text{obs. signal}} + \underbrace{\xi_j}_{\text{noise}}, \quad \text{Var}(\xi_j) = \sigma^2, \quad j = 1, \dots, d \quad (d \text{ can be } \infty).$$

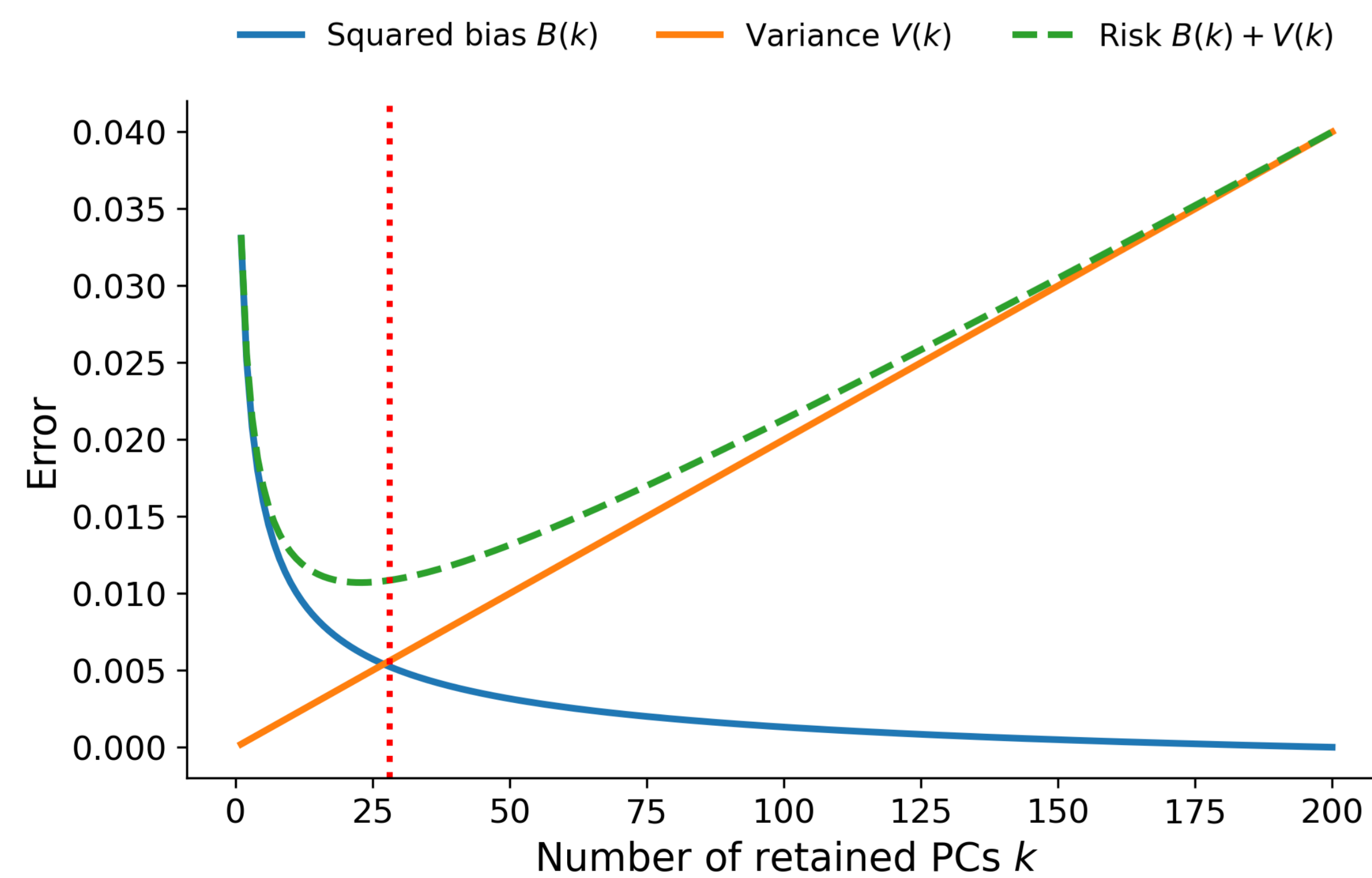
Given a spectrum $\lambda = (\lambda_1, \lambda_2, \dots)$ and a filter $\psi_\nu(\cdot)$, spectral estimator: $\hat{\theta}_j = (1 - \psi_\nu(\lambda_j))Z_j$, $j \in [d]$; e.g., a spectral-cutoff estimator keeps coordinates with larger λ_j .

Definition

Let $\{\pi_j\}_1^d$ order the eigenvalues decreasingly. Define the Effective Span Dimension (ESD) as

$$d^\dagger(\sigma^2; \theta^*, \lambda) = \min \left\{ k : \frac{1}{k} \sum_{i=k+1}^d (\theta_{\pi_i}^*)^2 \leq \sigma^2 \right\}.$$

d^\dagger is the cutoff where the squared bias of the spectral-cutoff estimator $B(k) = \sum_{i>k} (\theta_{\pi_i}^*)^2$ is no larger than the accumulated variance $V(k) = k\sigma^2$, indicated by the dotted line in the bias-variance plot below.



Unlike signal-agnostic measures such as effective dimension, ESD captures target-kernel alignment.

3. Minimax Risk

Given a quota $K \in [d]$, define the ESD-bounded class

$$\mathcal{F}_{K, \lambda}^{(\sigma^2)} = \{ \theta : d^\dagger(\sigma^2; \theta, \lambda) \leq K \}.$$

Minimax characterization

Consider risk $\mathcal{R}(\hat{\theta}; \theta^*) = \mathbb{E}_{\theta^*} \|\hat{\theta} - \theta^*\|_2^2$. The minimax risk over this class grows linearly with K :

$$\inf_{\hat{\theta}} \sup_{\theta^* \in \mathcal{F}_{K, \lambda}^{(\sigma^2)}} \mathcal{R}(\hat{\theta}; \theta^*) \asymp K\sigma^2.$$

$d^\dagger\sigma^2$ captures the statistical difficulty of the problem.

4. Examples

- Classical source ellipsoid.** If $\lambda_i = i^{-\beta}$, the source condition is the ellipsoid

$$\Theta_s(R) = \left\{ \theta : \sum_{i=1}^d \lambda_i^{-s} \theta_i^2 \leq R \right\} = \left\{ \theta : \sum_{i=1}^d i^{s\beta} \theta_i^2 \leq R \right\}.$$

It embeds into an ESD class with

$$K_{\text{src}} = \min \left\{ d, \left\lceil (R/\sigma^2)^{1/(1+s\beta)} \right\rceil \right\}, \quad \Theta_s(R) \subseteq \mathcal{F}_{K_{\text{src}}, \lambda}^{(\sigma^2)}.$$

Taking $\sigma^2 = \sigma_0^2/n$ gives the rate

$$\sigma^2 K_{\text{src}} \asymp \sigma_0^2 \min \{ n^{-s\beta/(1+s\beta)}, d/n \}.$$

- Near-parametric rates beyond classical assumptions.** Let λ be decreasing. For $b \geq 1$, set $g(x) = \sigma_0^2 x e^{-x^b}$, $\theta_1^* = 0$, and $(\theta_{j+1}^*)^2 = g(j) - g(j+1)$. Then with $K_n = \lceil (\log n)^{1/b} \rceil$,

$$\inf_{\hat{\theta}} \sup_{\theta^* \in \mathcal{F}_{K_n, \lambda}^{(\sigma_0^2/n)}} \mathcal{R}(\hat{\theta}; \theta^*) \asymp \sigma_0^2 K_n/n \asymp \sigma_0^2 (\log n)^{1/b}/n.$$

The classical rate does not yield the $(\log n)^{1/b}/n$ term.

- Same eigenvalues, different difficulty.** Let $S = \text{supp}(\theta^*)$ and suppose $|S| = s \ll d$. Compare two spectra with the same eigenvalue multiset.
 - [Well-aligned] If the largest s eigenvalues of $\lambda^{(1)}$ are on S , then $d^\dagger(\sigma^2; \theta^*, \lambda^{(1)}) \leq s$.
 - [Misaligned] If the largest $d-s$ eigenvalues of $\lambda^{(2)}$ are on S^c , then $d^\dagger(\sigma^2; \theta^*, \lambda^{(2)}) \geq \min\{d-s, \|\theta^*\|_2^2/\sigma^2\}$. Spectrum-only measures cannot distinguish (1) & (2).

5. ESD in Regression Models with Noise Variance σ_0^2

- Fixed-design linear model** $\mathbf{Y} = \mathbf{X}\beta^* + \epsilon$. Given SVD $n^{-1/2}\mathbf{X} = \mathbf{U}_{n \times r} \mathbf{S}_{r \times r} \mathbf{V}_{r \times p}^\top$, transform

$$\mathbf{Z} = n^{-1/2} \mathbf{U}^\top \mathbf{Y}, \quad \theta^* = \mathbf{S} \mathbf{V}^\top \beta^*, \quad \lambda_j = s_j^2, \quad \sigma^2 = \sigma_0^2/n.$$

Define $d^\dagger(\sigma_0^2/n; \beta^*, \mathbf{X}) := d^\dagger(\sigma_0^2/n; \theta^*, \lambda)$.

- RKHS** $y = f^*(x) + \epsilon$. For $\mathbf{k}(x, x')$, set

$$\theta_j^* = \langle f^*, \phi_j \rangle_{L^2(\mu)}, \quad \bar{\sigma}^2 = \frac{\sigma_0^2 + \|f^*\|_\infty^2}{n}.$$

Define $d^\dagger(\bar{\sigma}^2; f^*, \mathbf{k}) := d^\dagger(\bar{\sigma}^2; \theta^*, \lambda)$.

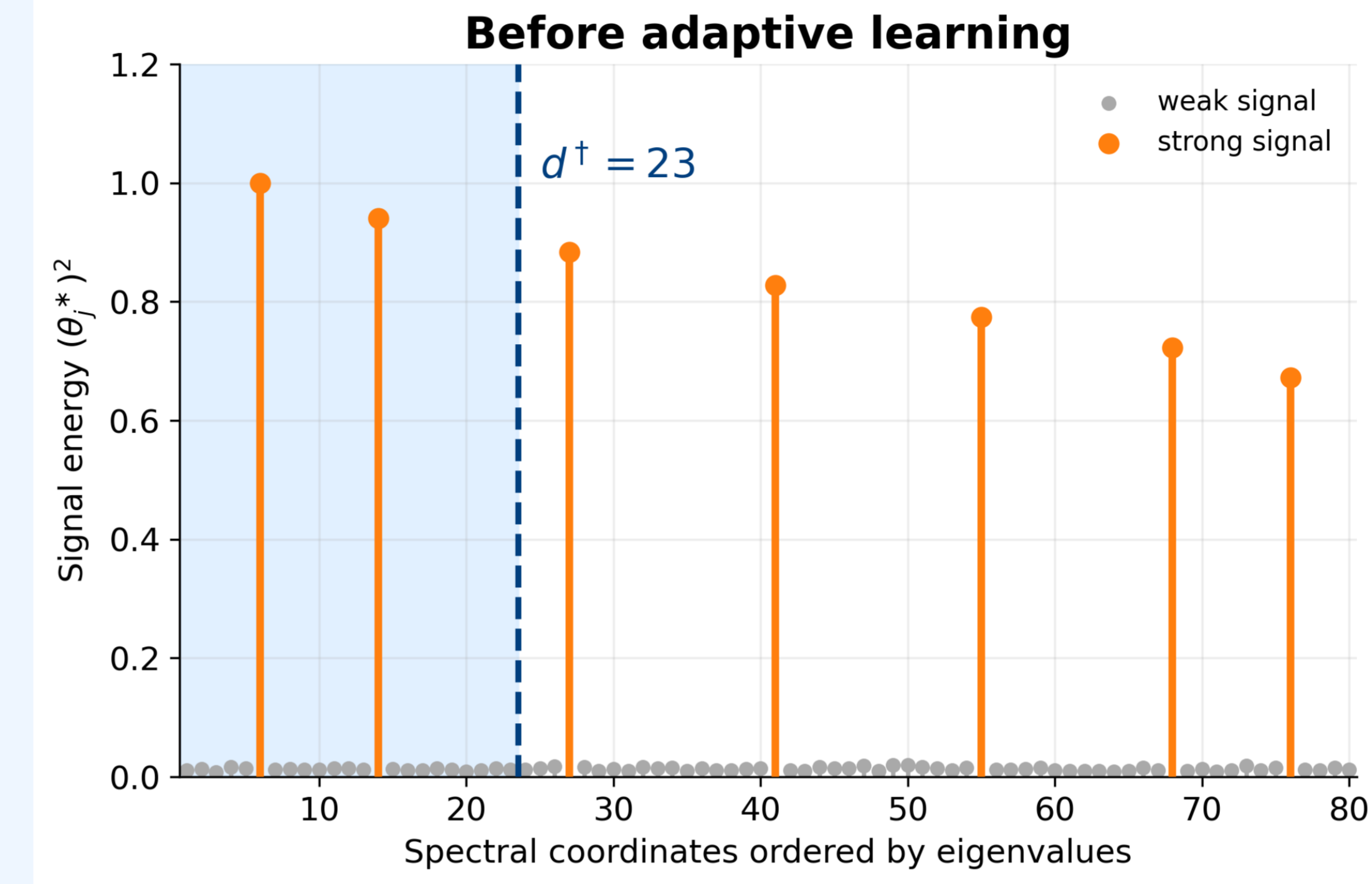
Minimax characterization can also be established.

6. Feature Learning

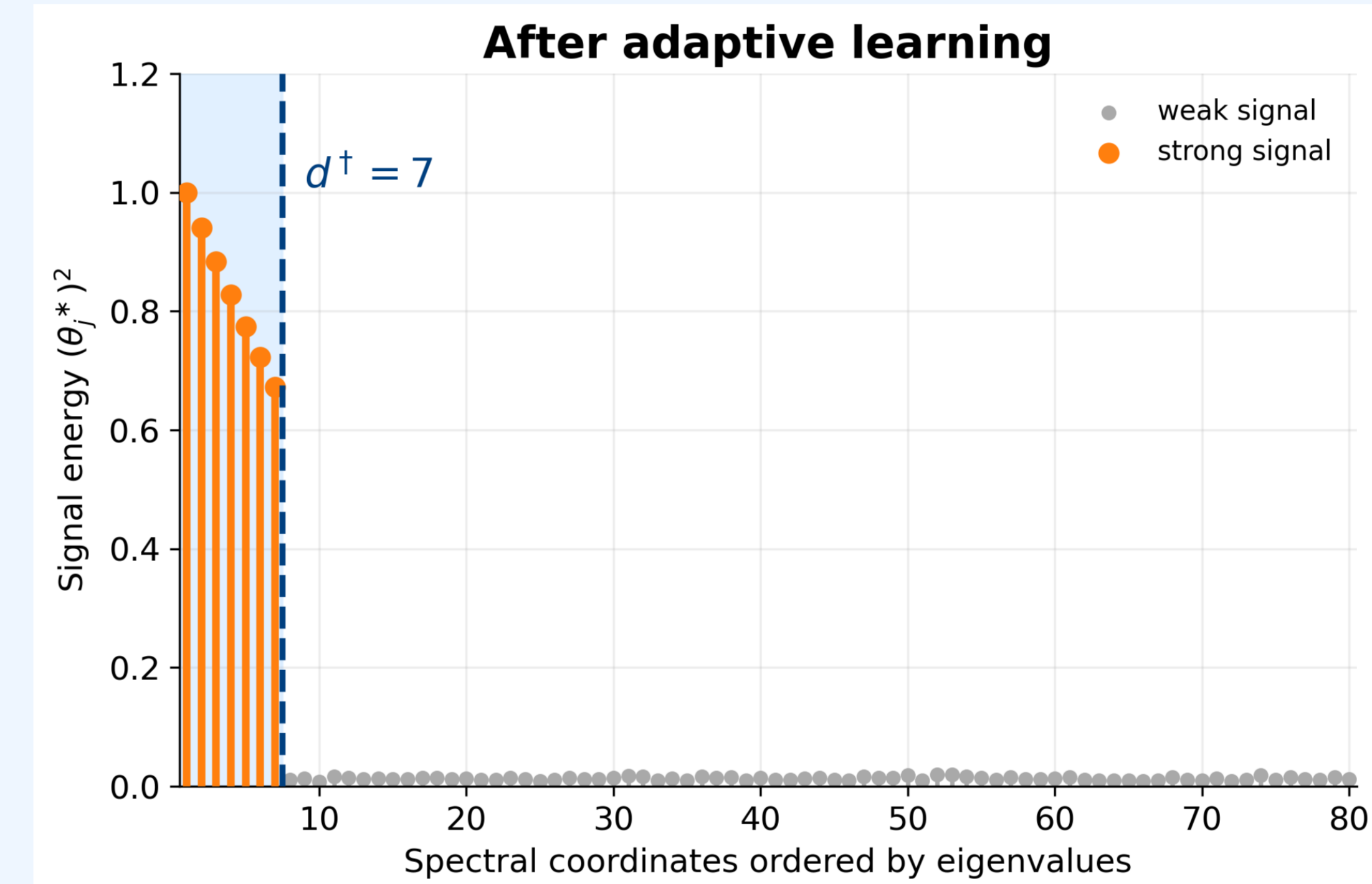
If the kernel is changed during training, one may compare the ESD before and after feature learning.

Example: one signal, two spectral orderings

Before adaptation: a badly-aligned kernel, strong signal coordinates appear late in the spectral order, $d^\dagger = 23$



After adaptation: a better-aligned kernel, strong signal is concentrated on leading directions, $d^\dagger = 7$



Risk interpretation

If adaptation changes $\lambda^{(0)}$ to $\lambda^{(a)}$ and

$$d^\dagger(\sigma^2; \theta^*, \lambda^{(a)}) < d^\dagger(\sigma^2; \theta^*, \lambda^{(0)}),$$

the same signal moves from a **minimax-hard class** with a larger ESD quota to an **easier class** with a smaller quota.

General mechanism. Feature learning reshapes the kernel so the same target has a smaller ESD at the relevant noise level, which means it moves into an easier class.

Trade-off function and span profile. Define

$$\mathbf{H}_{\theta^*, \lambda}(k) = \frac{1}{k} \sum_{i=k+1}^d (\theta_{\pi_i}^*)^2, \quad \mathbf{D}_{\theta^*, \lambda}(\tau) = \min\{k : \mathbf{H}_{\theta^*, \lambda}(k) \leq \tau\}.$$

$\sigma^{-2} \mathbf{H}_{\theta^*, \lambda}(k)$ is the spectral-cutoff bias-variance ratio, and

$\mathbf{D}_{\theta^*, \lambda}(\tau) = d^\dagger(\tau; \theta^*, \lambda)$ is the span profile.

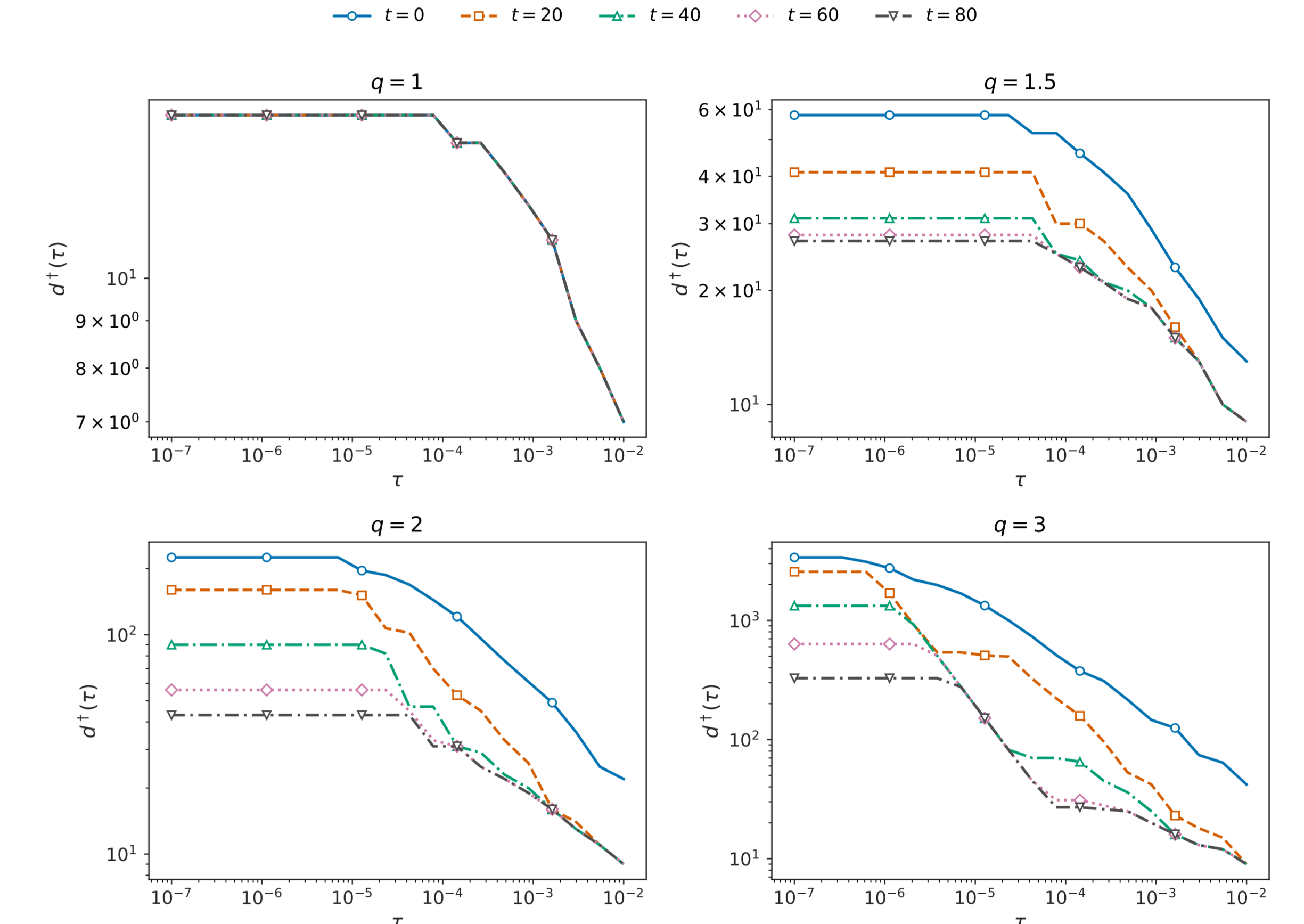
We expect $\mathbf{D}_{\theta^*, \lambda}(\cdot)$ moves down during feature learning.

7. ESD reduction in over-parameterized sequence models

Estimate: $\hat{\theta}_j(t) = \tilde{\lambda}_j(t)^{1/2} \beta_j(t)$, where $\tilde{\lambda}_j(t) = (a_j(t)b_j(t)^D)^2$ and (a_j, b_j, β_j) are learned using gradient flow. This learns the eigenvalues while fixing the eigencoordinates. We analyze

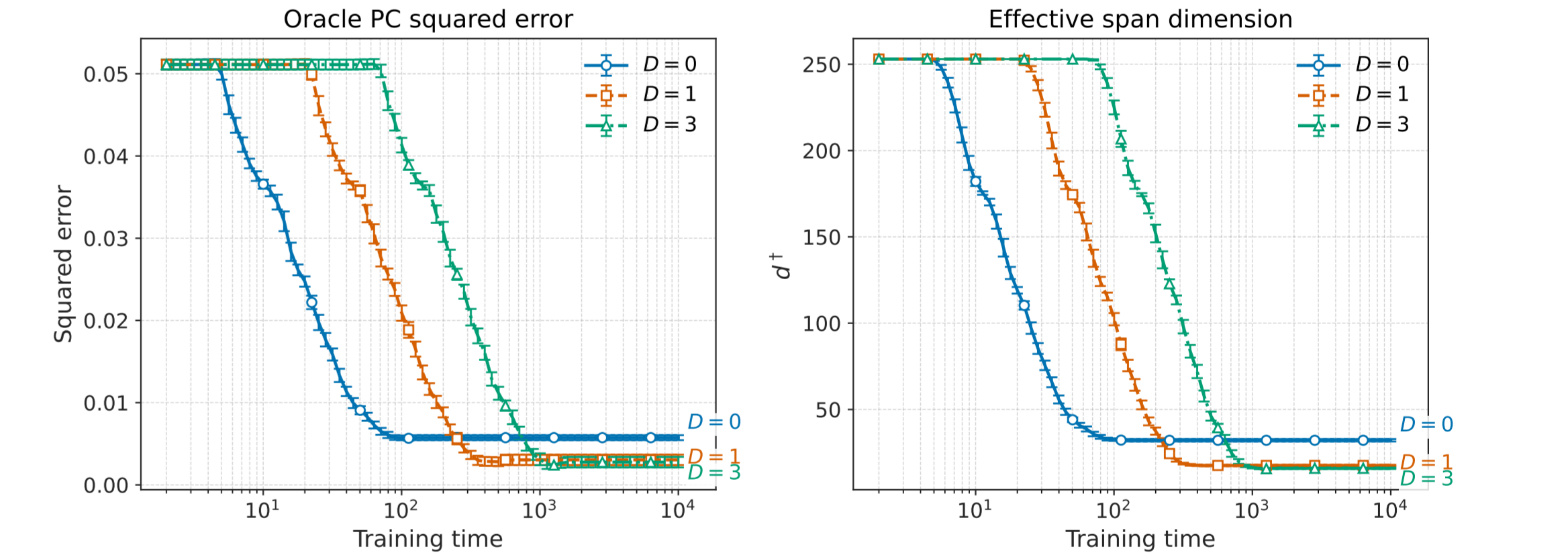
$$d^\dagger(t) = d^\dagger(\sigma^2; \theta^*, \tilde{\lambda}(t)).$$

Span-profile evolution under OP-GF



Setup: $\lambda_j = j^{-\gamma}$; nonzeros at $\ell(j) = \lfloor j^\eta \rfloor$, $\theta_{\ell(j)}^* = Cj^{-(\eta+1)/2}$; others zero. Here $q \in \{1, 1.5, 2, 3\}$, $n = 10^4$, $\sigma_0 = 1$, $d = 5000$, $J = 15$, $p = 2.5$, $\gamma = 1$ and depth $D = 0$. Larger q : worse alignment.

Risk and ESD along training



Here $q = 2.5$, with the same $n, \sigma_0, d, J, p, \gamma$ as above; compare $D \in \{0, 1, 3\}$. Mean over 20 Monte Carlo runs; bars: one standard error. The plotted risk is the squared error of the spectral-cutoff estimator tuned at the current ESD.

8. Beyond fixed eigenbasis: pathwise ESD

For a learned kernel \mathbf{k}_t at time t , we track the pathwise ESD

$$d^\dagger(t) = d^\dagger(\bar{\sigma}^2; f^*, \mathbf{k}_t).$$

Example: training a 4-layer linear neural network

- $\mathbf{A}(t)$ is the learned representation matrix in $f_i(x) = w(t)^\top \mathbf{A}(t)x$
- $\mathbf{k}_t(x, x') = x^\top \mathbf{A}(t)^\top \mathbf{A}(t)x'$

Both $d^\dagger(t)$ and risk decrease during training. $p = 900$, $n = 1000$, $\beta_j^* = j^{-1}$ for $j \leq 200$. Risk is $\|\mathbf{A}(t)^\top w(t) - \beta^*\|_2^2$.

9. Take-Home Message

- ESD depends jointly on the target, kernel, and noise level.
- For an ESD quota K , minimax risk $\asymp K\sigma^2$.
- Successful adaptation should reduce the ESD of the same signal.

Scan for paper



arXiv:2509.20294