



# No Double Descent in Principal Component Regression: A High-Dimensional Analysis



UPPSALA  
UNIVERSITET

Daniel Gedon, Antônio H. Ribeiro, Thomas B. Schön  
Uppsala University (Sweden) contact: daniel.gedon@it.uu.se

## Background

**Spiked covariance model.** Population covariance  $\mathbf{C}$  of  $\mathbf{x}$ , eigenvectors  $\mathbf{v}$ , eigenvalues  $\lambda$ ,

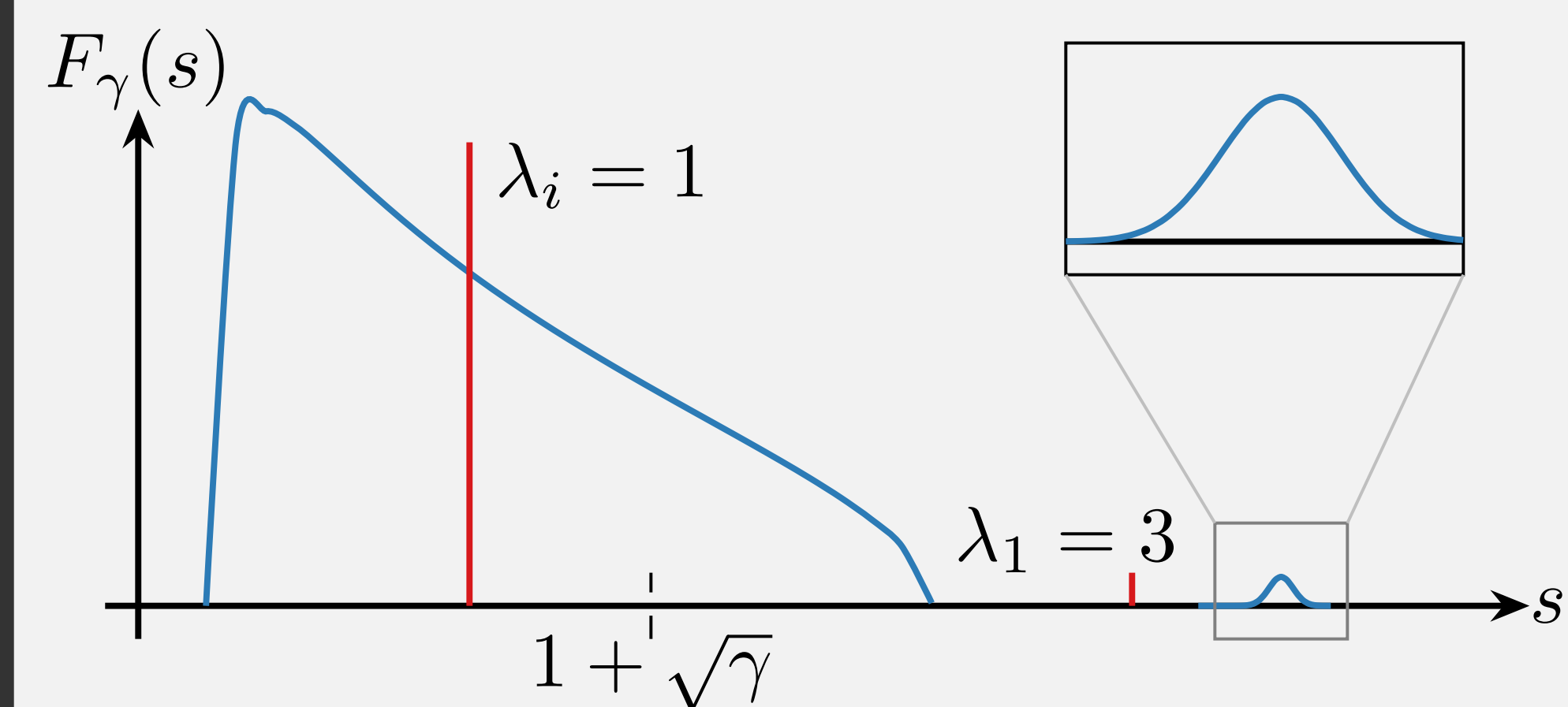
$$\mathbf{C} = \mathbf{I} + \sum_{i=1}^d \mathbf{v}_i \lambda_i \mathbf{v}_i^\top.$$

**Eigenvalue shift.** Distribution of spike sample eigenvalue  $\hat{\lambda}_1$ , depending on  $\lambda_1$ :

- $\lambda_1 \in [1, 1 + \sqrt{\gamma}]$ : spike at  $\mu(\gamma) = (1 + \sqrt{\gamma})^2$ .
- $\lambda_1 > 1 + \sqrt{\gamma}$ : spike Normal with  $\mu(\lambda, \gamma) = \gamma \frac{\lambda}{\lambda-1} + \lambda$ .

**Eigenvector shift.** Sample eigenvectors  $\hat{\mathbf{v}}$ . As  $p/n \rightarrow \gamma$ ,

$$(\mathbf{v}_i^\top \hat{\mathbf{v}}_i)^2 \rightarrow \begin{cases} \frac{1-\gamma/(\lambda_i-1)^2}{1+\gamma/(\lambda_i+1)} & \text{for } \lambda_i > 1 + \sqrt{\gamma} \\ 0 & \text{for } \lambda_i \in [1, 1 + \sqrt{\gamma}]. \end{cases}$$



[1] Johnstone, I. M. On the distribution of the largest eigenvalue in principal components analysis. *Annals of Statistics*, 2001.

[2] Johnstone, I. M. and Paul, D. PCA in high dimensions: An orientation. *Proceedings of the IEEE*, 2018.

## Theory – preliminary

**Risk.**  $R(\hat{\theta}) = \mathbb{E}_{(\mathbf{x}_0, y_0) \sim \mathcal{D}} [(y_0 - \hat{y}_0(\mathbf{x}_0))^2]$

**Eigenvector shift.** Generalisation to all eigenvectors

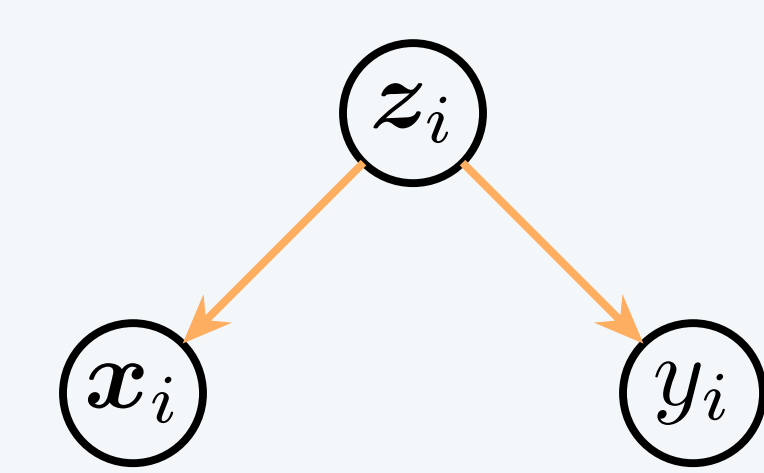
$$\mathbf{P}_k = \begin{cases} \text{diag}((\mathbf{v}_1^\top \hat{\mathbf{v}}_1)^2, \dots, (\mathbf{v}_k^\top \hat{\mathbf{v}}_k)^2, 0, \dots, 0), & k < d, \\ \text{diag}((\mathbf{v}_1^\top \hat{\mathbf{v}}_1)^2, \dots, (\mathbf{v}_d^\top \hat{\mathbf{v}}_d)^2), & k = d, \\ \text{diag}((\mathbf{v}_1^\top \hat{\mathbf{v}}_1)^2 + c_1^2, \dots, (\mathbf{v}_d^\top \hat{\mathbf{v}}_d)^2 + c_d^2), & k > d. \end{cases}$$

## Motivation

- Overparameterized models,  $\gamma = \frac{p}{n} > 1 \rightarrow$  double descent
- Real-world data often on a low-dimensional manifold
- PCR (= PCA + linear regression) widely adopted in practice

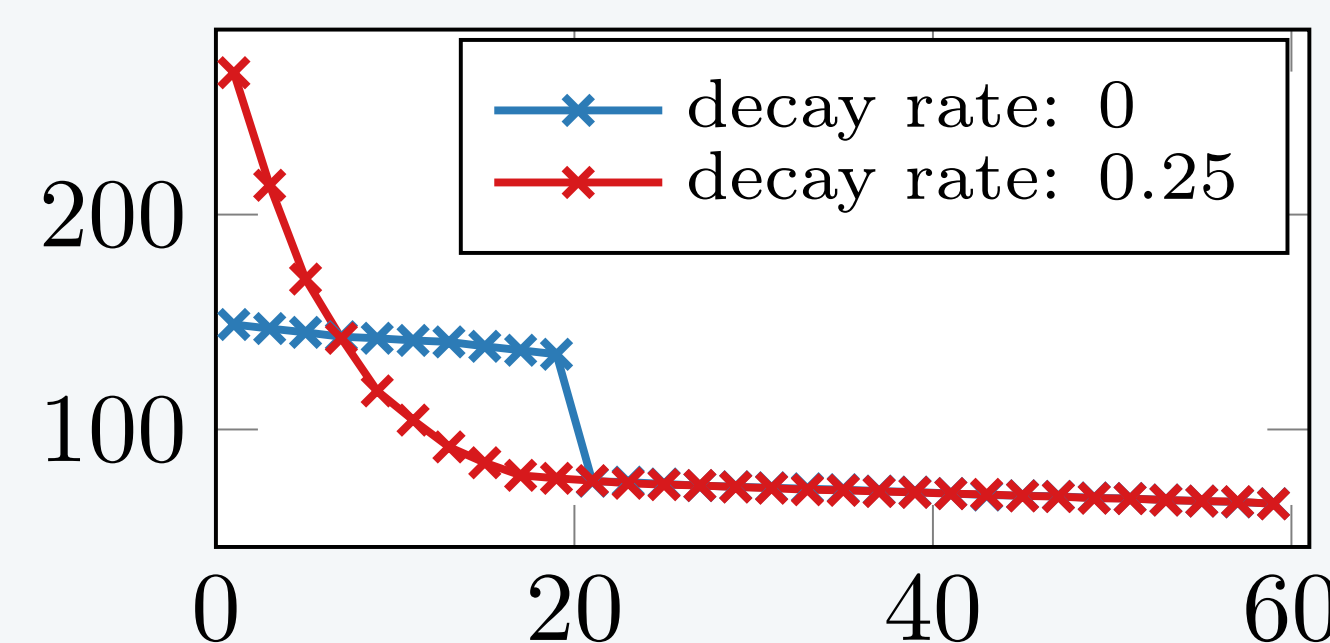
[3] Belkin, M., Hsu, D. J., Ma, S., and Mandal, S. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *PNAS*, 2019

## Data generator



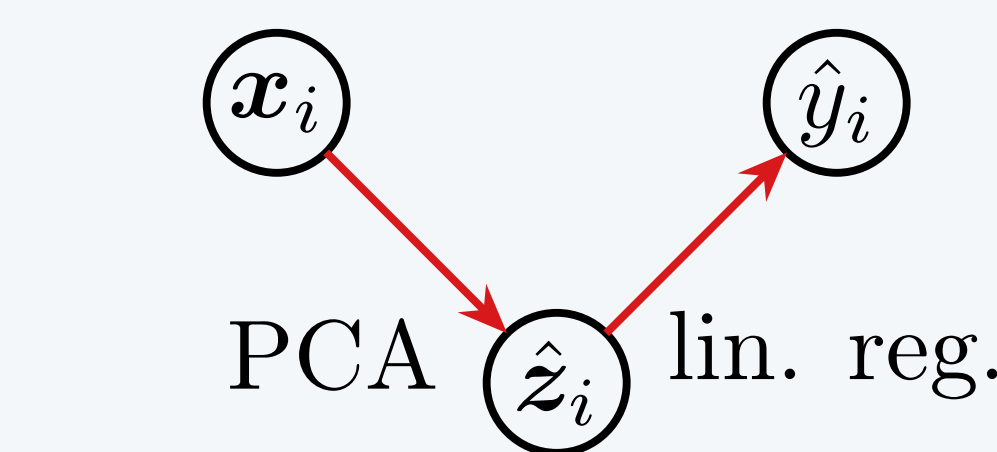
$$\mathbf{x}_i = \mathbf{W} r_w \mathbf{z}_i + \mathbf{e}_i, \\ \mathbf{y}_i = \boldsymbol{\theta}^\top \mathbf{z}_i + \varepsilon_i.$$

Eigenvalues of  $\mathbf{x}$ .



Dimensions:  
 $\mathbf{x} \in \mathbb{R}^p, \mathbf{z} \in \mathbb{R}^d, \hat{\mathbf{z}} \in \mathbb{R}^k.$

## Model



$$\text{SVD: } \mathbf{X} \approx \hat{\mathbf{U}} \hat{\boldsymbol{\Sigma}} \hat{\mathbf{V}}^\top, \\ \text{PCA: } \hat{\mathbf{z}}_i = \hat{\mathbf{V}}^\top \mathbf{x}_i, \\ \text{lin. reg. } \hat{y}_i = \hat{\boldsymbol{\theta}}^\top \hat{\mathbf{z}}_i.$$

## Theory

**Theorem.** Let  $n, p \rightarrow \infty$  with  $\frac{p}{n} \rightarrow \gamma \in (0, \infty)$ , we have a.s.

$$\mathbb{E}_\nu [R(\hat{\theta})] \rightarrow \text{Bias}_\gamma(\hat{\theta})^2 + \text{Var}_\gamma(\hat{\theta}) + \sigma_\nu^2,$$

$$\text{Bias}_\gamma(\hat{\theta})^2 = \bar{\boldsymbol{\beta}}^\top (\boldsymbol{\Lambda}_d - \boldsymbol{\Lambda}_d \mathbf{P}_k - \mathbf{P}_k \boldsymbol{\Lambda}_d + \mathbf{P}_k + \mathbf{P}_k r_w^2 \mathbf{C}_z \mathbf{P}_k) \bar{\boldsymbol{\beta}},$$

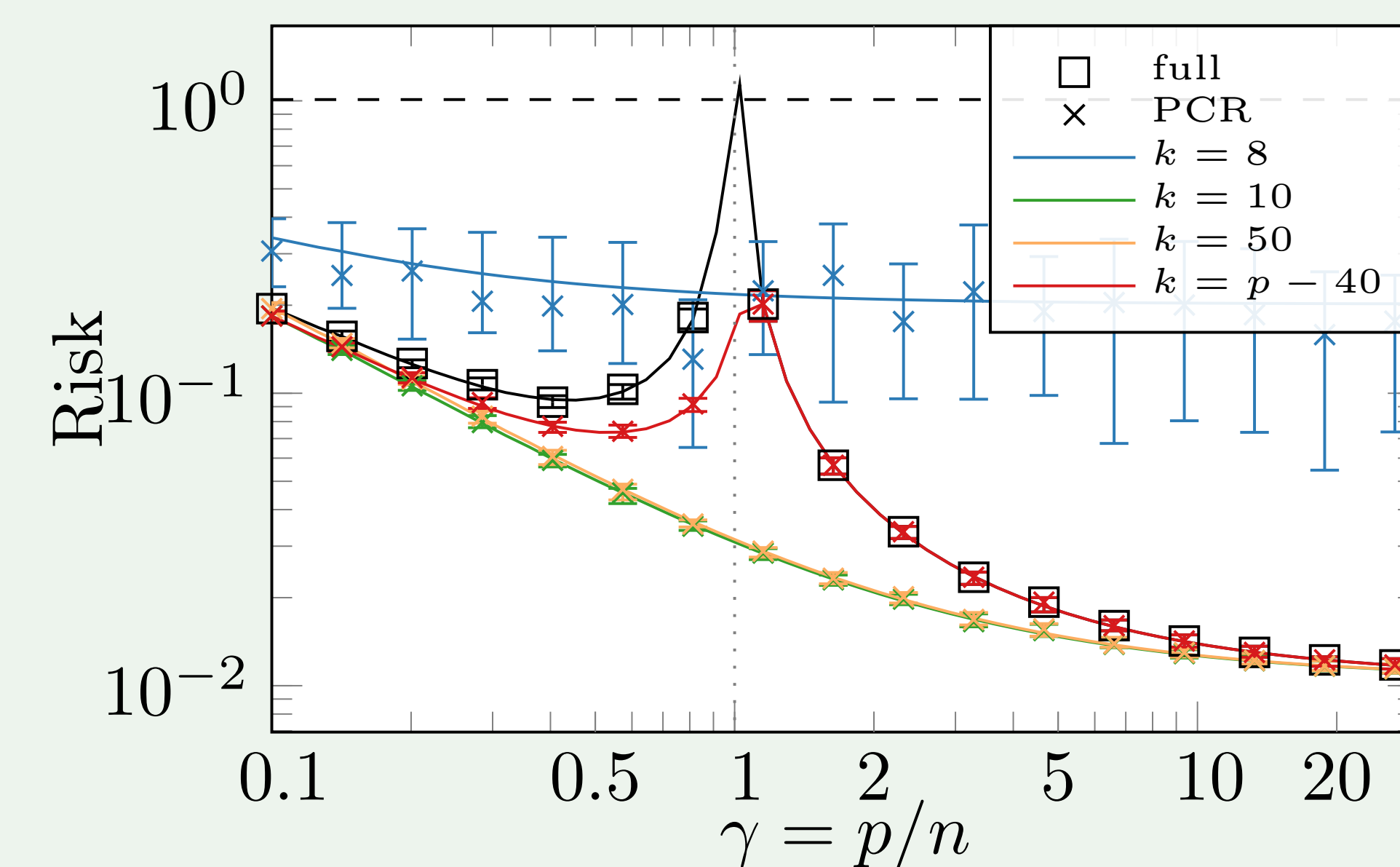
$$\text{Var}_\gamma(\hat{\theta}) = \frac{\sigma_\nu^2}{n} \left( \text{Tr} \left[ (\mathbf{P}_k r_w^2 \mathbf{C}_z + \mathbf{I}_k) \frac{1}{\mu(\boldsymbol{\Lambda}, \gamma)} \right] + (p-d) \int_{s_c}^{(1+\sqrt{\gamma})^2} \frac{1}{s} dF_\gamma(s) \right).$$

## Interpretation.

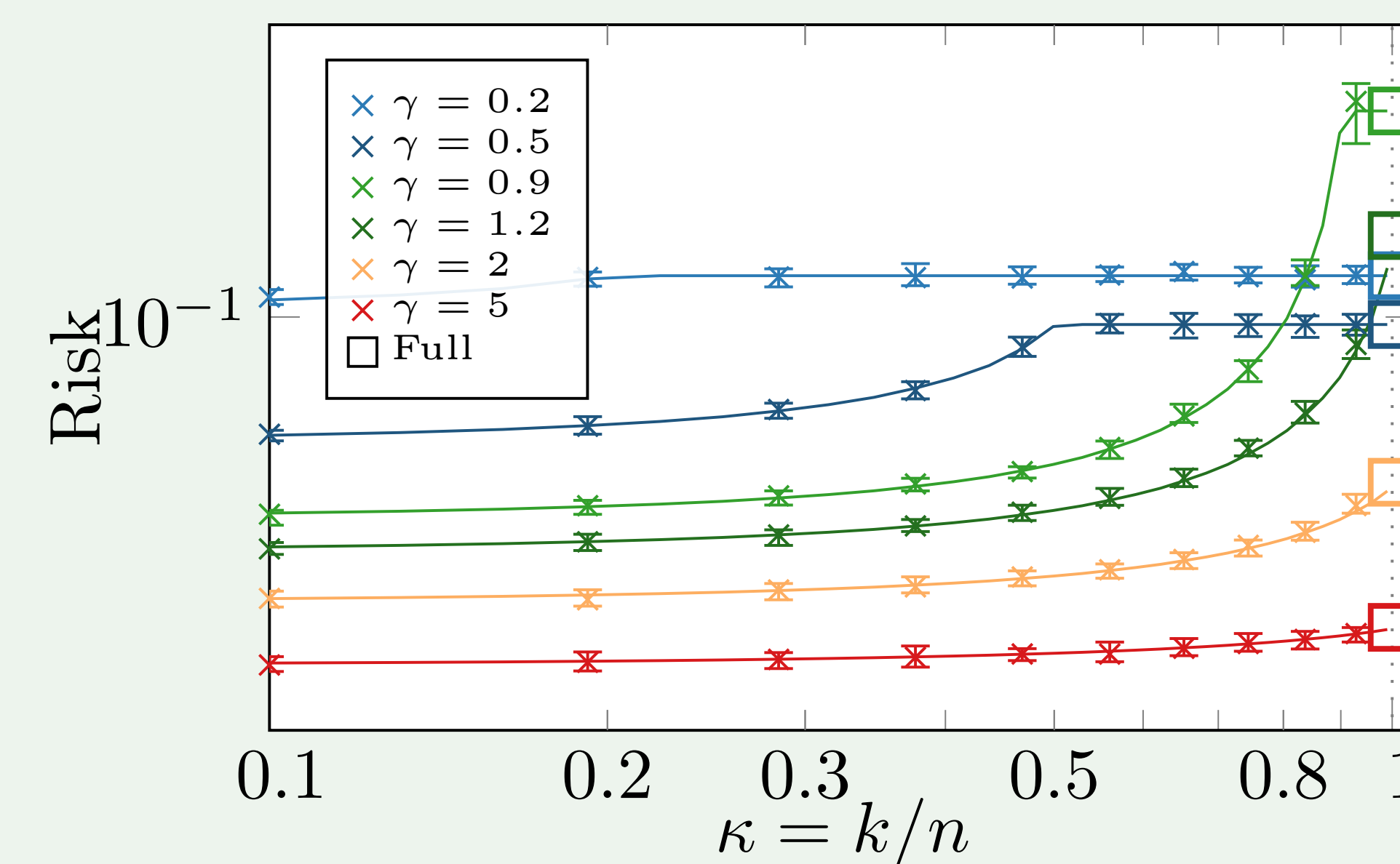
1. *Bias* as scaled  $d$ -dimensional subspace of eigenvalues  $\boldsymbol{\Lambda}$ .
2. *Variance* with  $k \leq d \rightarrow s_c = (1 + \sqrt{\gamma})^2 \rightarrow$  integral term=0.
3. Principal components  $k \rightarrow$  cut-off for considered data distribution.

## Results

**Main result.** Number of spikes  $d = 10$ .



## PCA-projection space.



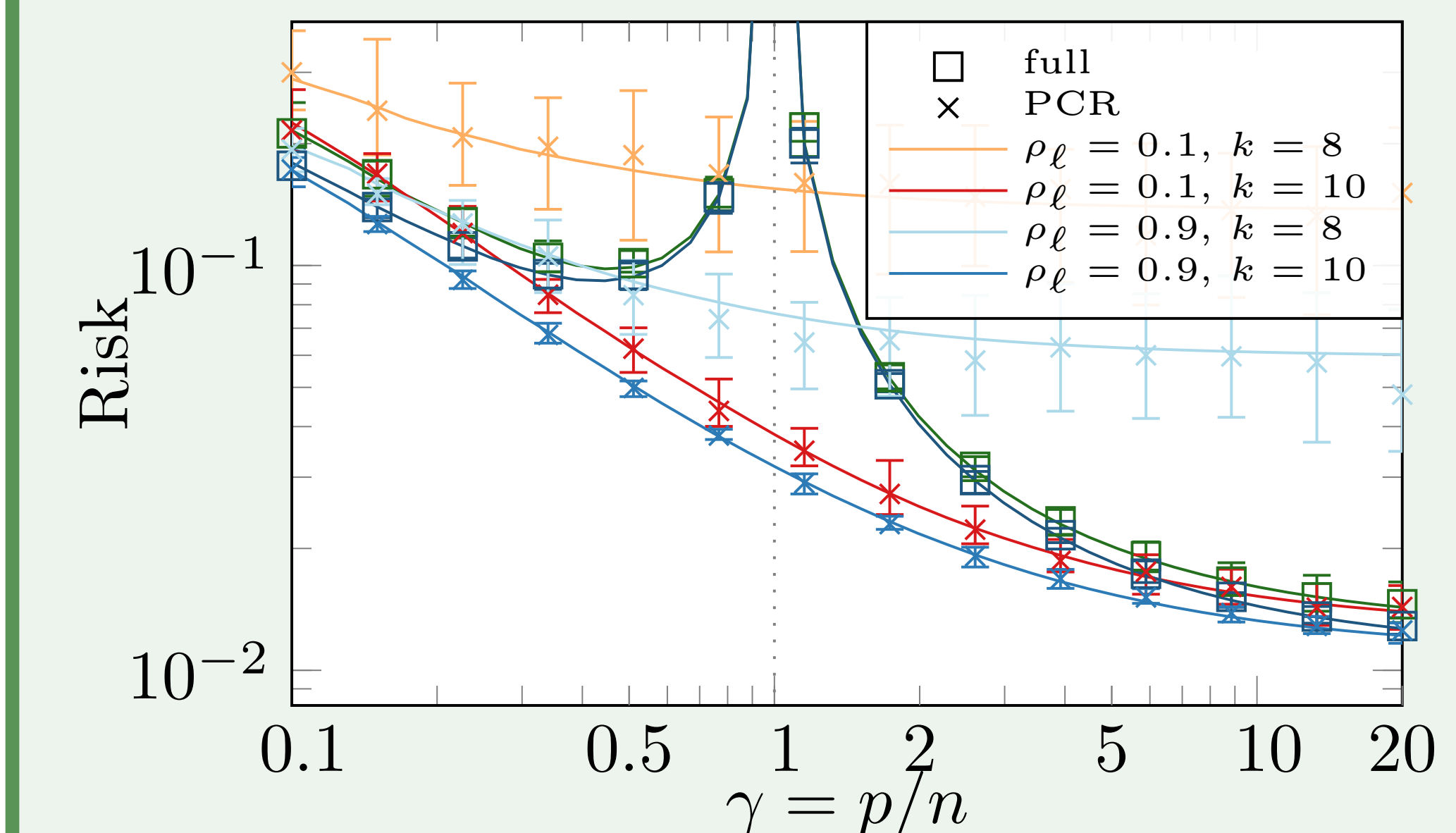
## Covariate-shift

**Assumption.** Train/source  $\mathbf{C}_S = \mathbf{V} \boldsymbol{\Lambda}_S \mathbf{V}^\top$  and test covariance  $\mathbf{C}_T = \mathbf{V} \boldsymbol{\Lambda}_T \mathbf{V}^\top$ .

**Theorem.** Similar to in-distribution risk but scaled. But combinations of  $\mathbf{C}_S, \mathbf{C}_T$  here.

## Results.

Data generator with  $\text{Cov}(\mathbf{v}_S, \mathbf{v}_T) = \sigma_\ell^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ .



## Observation.

1. Same behaviour as for in-distribution data.
2. Higher correlation  $\rho \rightarrow$  lower risk.

## Conclusion

### Summary.

- Asymptotic risk of PCR under spiked covariance model. Tool: random matrix theory.
- Guarantee: widely used model & real-world data structures

### Guide to choosing $k$ .

- High  $k$  increases variance contribution.
- For  $0.5 < \gamma < 2$  a suitable  $k$  has little effect.

### Limitations.

- Linear supervised setting.
- Asymptotic results — no finite sample guarantees.