



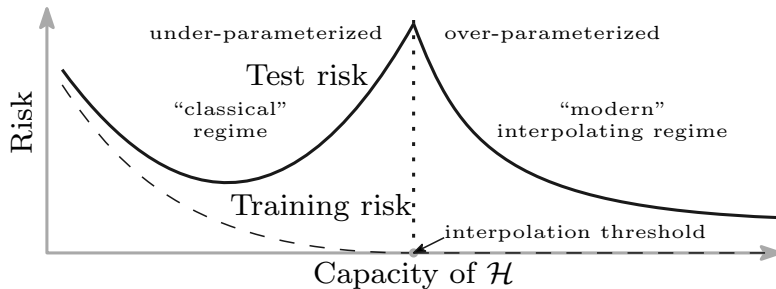
UPPSALA  
UNIVERSITET

# No double descent in PCA: Training and pre-training in high dimensions

---

Daniel Gedon  
Uppsala University

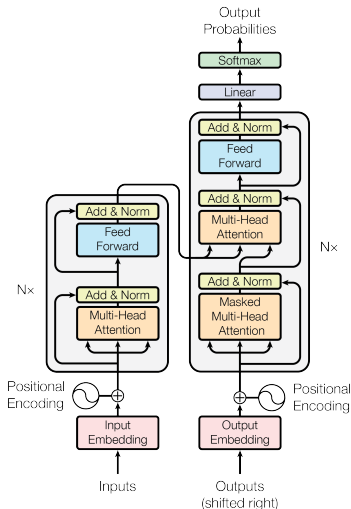
Belkin Lab weekly meeting  
San Diego, March 06, 2023



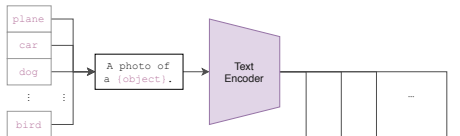
Found/analysed in: linear/logistic regression, random forests, adversarial training, neural networks, ...

---

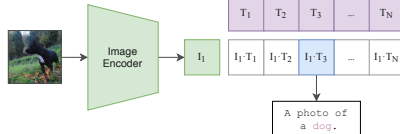
Belkin et al., "Reconciling modern machine-learning practice and the classical bias–variance trade-off".



(2) Create dataset classifier from label text



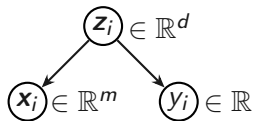
(3) Use for zero-shot prediction



Vaswani et al., "Attention is all you need".

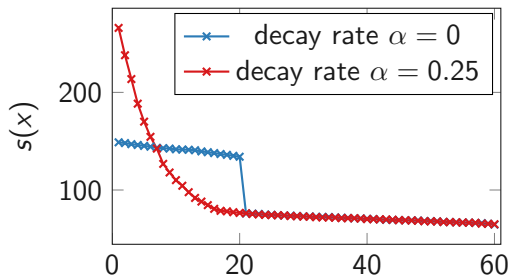
Radford et al., "Learning transferable visual models from natural language supervision".

## Data generator

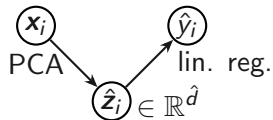


$$\mathbf{x}_i = \mathbf{D}\mathbf{z}_i + \mathbf{e}_i$$

$$y_i = \boldsymbol{\theta}\mathbf{z}_i + \varepsilon_i$$



## Model



$$\text{SVD: } \mathbf{X} \approx \hat{\mathbf{U}}\hat{\boldsymbol{\Sigma}}\hat{\mathbf{V}}^T,$$

$$\text{PCA: } \hat{\mathbf{z}}_i = \hat{\mathbf{V}}^T \mathbf{x}_i,$$

$$\text{lin. reg. } \hat{y}_i = \hat{\boldsymbol{\theta}}^T \hat{\mathbf{z}}_i.$$

Interested in risk:  $R(\hat{\theta}) = \mathbb{E}_{y_0} [(y_0 - \hat{y}_0)^2]$ .

Write data generator directly from features to outputs as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  with  $\boldsymbol{\beta} \in \mathbb{R}^m$ .

## Lemma

Sample covariance  $\hat{\mathbf{C}} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$  and the true covariance  $\mathbf{C}$ . Orthogonal projectors  $\boldsymbol{\Pi} = \mathbf{I}_m - \hat{\mathbf{V}}\hat{\mathbf{V}}^\top$ . Then,

$$\mathbb{E}_\epsilon [R(\hat{\theta})] = \boldsymbol{\beta}^\top \boldsymbol{\Pi} \mathbf{C} \boldsymbol{\Pi} \boldsymbol{\beta} + \frac{\sigma_\epsilon^2}{n} \text{Tr}(\hat{\mathbf{V}}^\top \mathbf{C} \hat{\mathbf{V}} \hat{\mathbf{V}}^\top \hat{\mathbf{C}} + \hat{\mathbf{V}}) + \sigma_\epsilon^2.$$

Compare with Hastie et al. for direct linear regression:

$$\mathbb{E}_\epsilon [R(\hat{\theta})] = \boldsymbol{\beta}^\top \boldsymbol{\Pi} \mathbf{C} \boldsymbol{\Pi} \boldsymbol{\beta} + \frac{\sigma_\epsilon^2}{n} \text{Tr}(\mathbf{C} \hat{\mathbf{C}}^+) + \sigma_\epsilon^2.$$

---

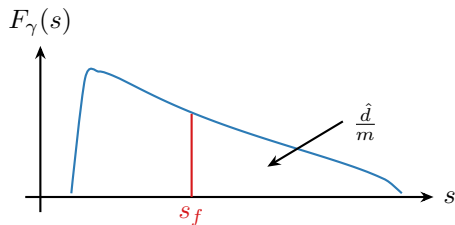
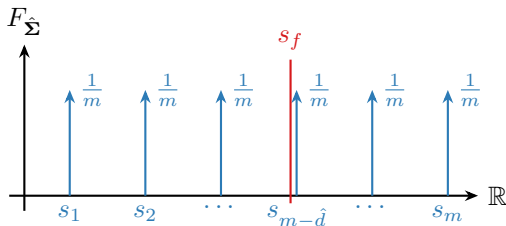
Hastie et al., “Surprises in high-dimensional ridgeless least squares interpolation”.

## Theorem

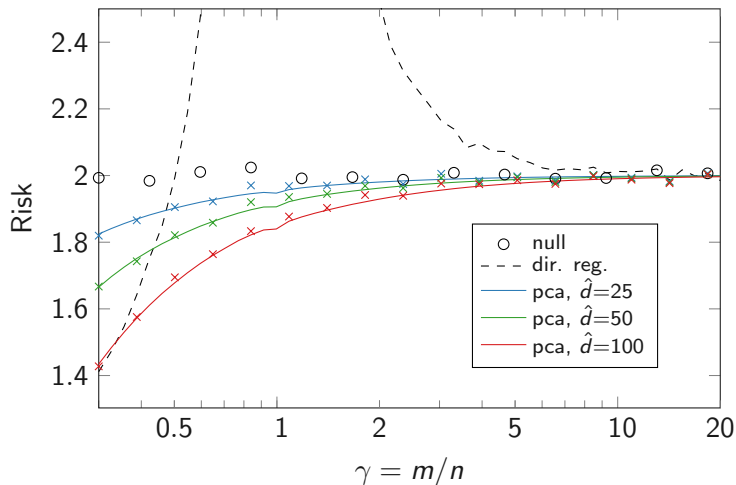
Assume isotropic features  $\mathbf{C} = \mathbf{I}_m$ , which implies  $d = m$  and choose constant  $\hat{d}$ . Then, as  $m, n \rightarrow \infty$ , such that  $\frac{m}{n} \rightarrow \gamma$ , the expected risk satisfies almost surely

$$\mathbb{E}_\epsilon \left[ R(\hat{\theta}) \right] \rightarrow \sigma_\epsilon^2 \frac{m}{n} \int_{s_f}^\infty \frac{1}{s} dF_\gamma(s) + \sigma_\epsilon^2 + \begin{cases} \beta^\top \beta \left( 1 - \min(\hat{d}, m)/m \right) & \text{for } \gamma < 1 \\ \beta^\top \beta \left( 1 - \min(\hat{d}, n)/m \right) & \text{for } \gamma > 1 \end{cases}$$

with  $F_\gamma$  the Marčenko-Pastur law and  $s_f$  the value in  $\mathbb{R}$  that satisfies  $\frac{\hat{d}}{m} = \int_{s_f}^\infty dF_\gamma$ .

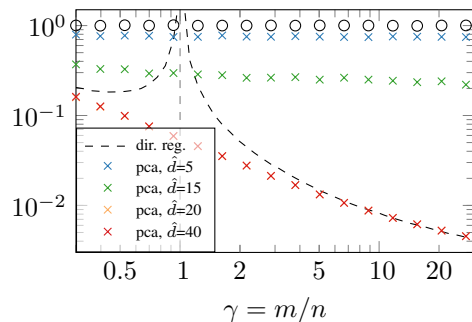


Isotropic data: Theorem from last slide.

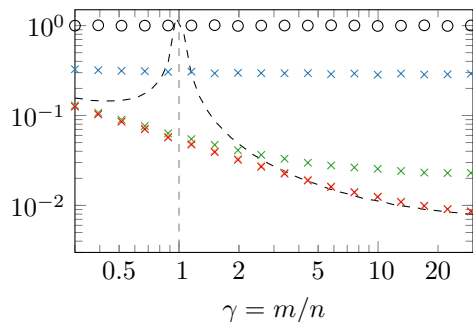


**Latent variable data:** No theorem but empirical results.

Risk,  $\alpha = 0$



Risk,  $\alpha = 0.25$

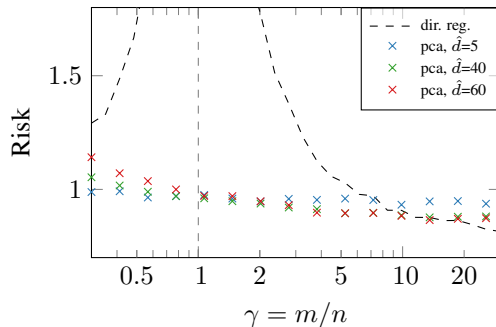
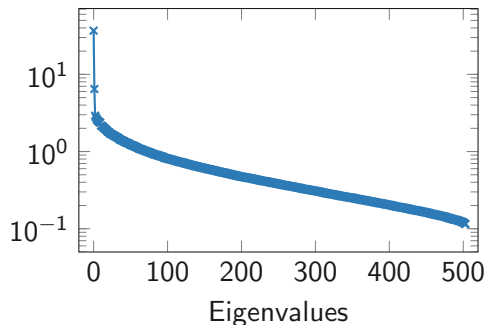




**Real world example:** Diverse MAGIC wheat genetics data set.

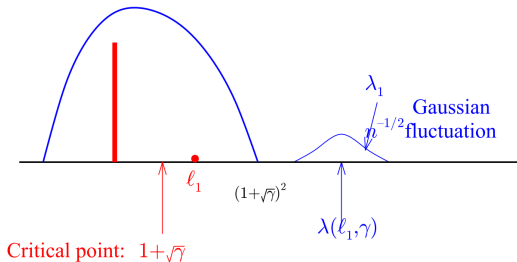
Input: genome sequence (=1.1M nucleotides) of 504 wheat lines.

Outcome: real-valued phenotypes.



**Aim:** Asymptotic risk for latent variable data generator.

Data generator = spiked covariance model  $\mathbf{C} = \sigma_0 \mathbf{I} + \sum_{i=0}^K s_i \mathbf{v}_i \mathbf{v}_i^\top$



Need eigenvector product:  $\mathbb{E}_\epsilon [R(\hat{\theta})] = \beta^\top \mathbf{\Pi} \mathbf{C} \mathbf{\Pi} \beta + \frac{\sigma_\epsilon^2}{n} \text{Tr}(\hat{\mathbf{V}}^\top \mathbf{C} \hat{\mathbf{V}} \hat{\mathbf{V}}^\top \hat{\mathbf{C}} + \hat{\mathbf{V}}) + \sigma_\epsilon^2$ .

Results:

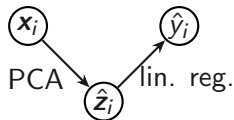
$$(\hat{\mathbf{v}}_i^\top \mathbf{v}_i)^2 \rightarrow \begin{cases} \frac{1-\gamma/(\ell_i-1)^2}{1-\gamma/(\ell_i-1)} & \text{for } \ell_i > 1 + \sqrt{\gamma} \\ 0 & \text{for } \ell_i \in [1, 1 + \sqrt{\gamma}] \end{cases}$$

Johnstone and Paul, "PCA in High Dimensions: An Orientation".

Paul, "Asymptotics of sample eigenstructure for a large dimensional spiked covariance model".

Use two data sets:

- Pre-training data set  $\{\mathbf{x}_i\}_{i=1}^{n_p}$
- Training data set  $\{\mathbf{x}_i, y_i\}_{i=1}^n$



Two step training procedure:

1. Unsupervised pre-training of PCA.
2. Train linear regression on the PCA features  $\hat{\mathbf{z}}_i$ .

For technical reasons: orthogonalize features and noise  $\mathbf{x}_i = \mathbf{D}\mathbf{z}_i + \mathbf{D}_\perp \mathbf{e}_i$ . Then:

$$\text{Model: } \hat{\mathbf{z}}_i = \hat{\mathbf{V}}^\top \mathbf{x}_i,$$

$$\text{Data generator: } \mathbf{z}_i = \mathbf{D}^+(\mathbf{x}_i - \mathbf{D}_\perp \mathbf{e}_i) = \mathbf{D}^+ \mathbf{x}_i.$$

Define projection loss:

$$\mathcal{L}(\mathbf{D}) = \mathbb{E} [\|\mathbf{x}\|_2^2 - \|\mathbf{D}^+ \mathbf{x}\|_2^2]; \quad \mathcal{L}(\hat{\mathbf{V}}) = \mathbb{E} [\|\mathbf{x}\|_2^2 - \|\hat{\mathbf{V}}^\top \mathbf{x}\|_2^2]$$

**Lemma**

$$\mathcal{L}(\hat{\mathbf{V}}) - \mathcal{L}(\mathbf{D}) = \sum_{i=1}^{\min(d, \hat{d})} \sum_{j=1}^m (\hat{\mathbf{v}}_i^\top \mathbf{v}_j)^2 (s_i - s_j) + \underbrace{\sum_{i=\hat{d}}^d s_i}_{=0 \text{ for } \hat{d} \geq d} + \underbrace{\sum_{i=d}^{\hat{d}} \sum_{j=1}^m (\hat{\mathbf{v}}_i^\top \mathbf{v}_j)^2 s_j}_{=0 \text{ for } \hat{d} \leq d}.$$

→ Correct estimation of eigenvectors  $\hat{\mathbf{V}}$  crucial for small loss difference.

## Theorem

Take  $t > 0$ ,  $k_j^2 = s_j(s_j + \text{Tr}(\mathbf{C}))$ , then

$$P\left(\mathcal{L}(\hat{\mathbf{V}}) - \mathcal{L}(\mathbf{D}) > t\right) \leq \frac{4}{t n_p} \left( \sum_{i=1}^{\min(d, \hat{d})} \sum_{j=i+1}^m \frac{k_j^2}{|s_i - s_j|} + \sum_{i=\hat{d}}^d \sum_{j=1}^m \frac{k_j^2 s_i}{(s_i - s_j)^2} + \sum_{i=d}^{\hat{d}} \sum_{j=1}^m \frac{k_j^2 s_j}{(s_i - s_j)^2} \right).$$

Tighter bound by:

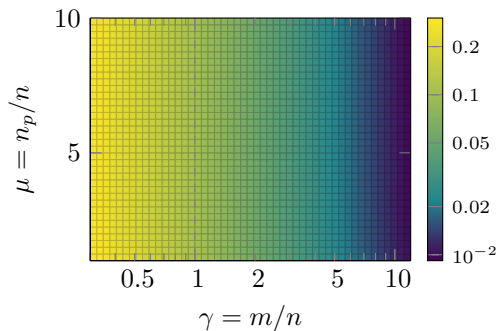
- rapidly decaying eigenvalues  $\rightarrow$  large  $|s_i - s_j| \geq 0$ ,
- more pre-training samples  $n_p$ .

---

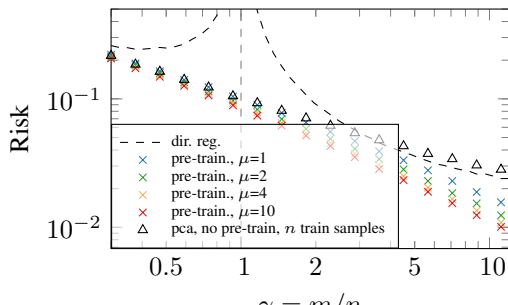
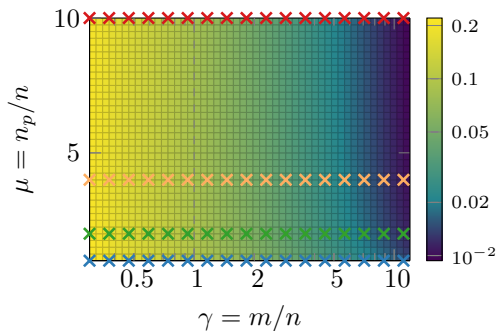
Loukas, “How close are the eigenvectors of the sample and actual covariance matrices?”

# Pre-training the PCA – Numerical results

Risk with pre-training,  $\alpha = 0$



Risk with pre-training,  $\alpha = 0.25$



**Aim:** Closed form for the risk with pre-training.

Move away from sample complexity towards asymptotic results.

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

to quantify asymptotic risk for spiked covariance model

## Extensions:

- Use real-world nonlinear data
- Neural networks: multi-layer approaches
- Nonlinear setting: replace PCA with kernel PCA
- Distribution shift between pre-training and training / testing

---

Paul, "Asymptotics of sample eigenstructure for a large dimensional spiked covariance model".

Supervised case:

- Generalize results from Hastie et al. for direct regression
- Selecting sufficiently large latent dimension  $\hat{d}$  is crucial for low risk

→ formal guarantees for performance of PCA-regression on real-world data structures

Pre-training:

- more pre-training data only help to improve eigenvector estimates
- certain decay rate  $\alpha$  is necessary such that more pre-training data are helpful



# Thank you!

**Daniel Gedon, Uppsala University**

*E-mail:* `daniel.gedon@it.uu.se`

*Web:* `dgedon.github.io`

*Twitter:* @danigedon

Supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by Knut and Alice Wallenberg Foundation.