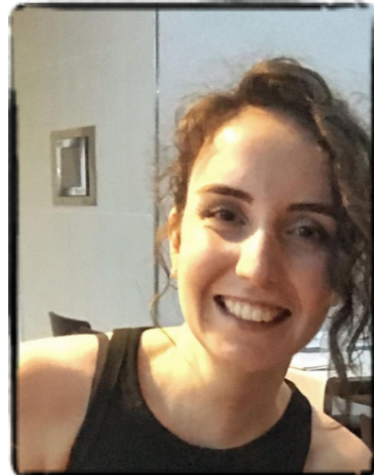


# Generalisation error in learning with random features and the hidden manifold model

---



**B. Loureiro**  
(IPhT)



F. Gerace  
(IPhT)



M. Mézard  
(ENS)



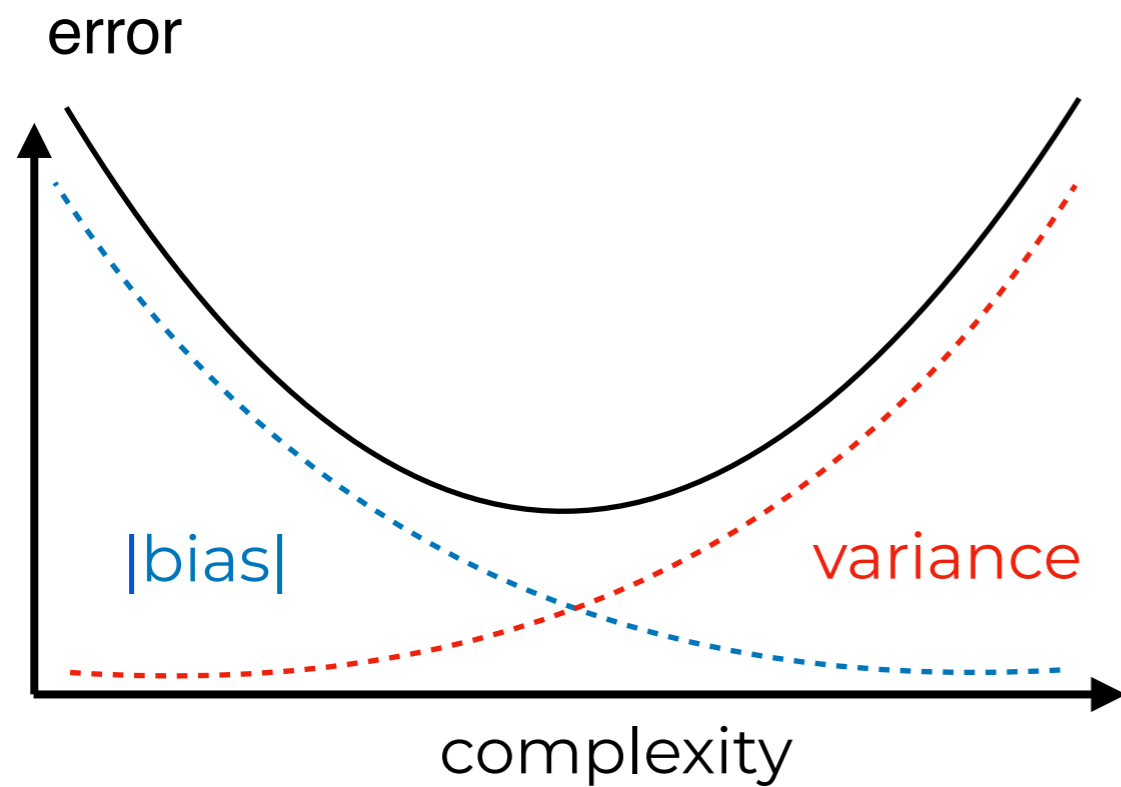
L. Zdeborová  
(IPhT)



F. Krzakala  
(ENS)

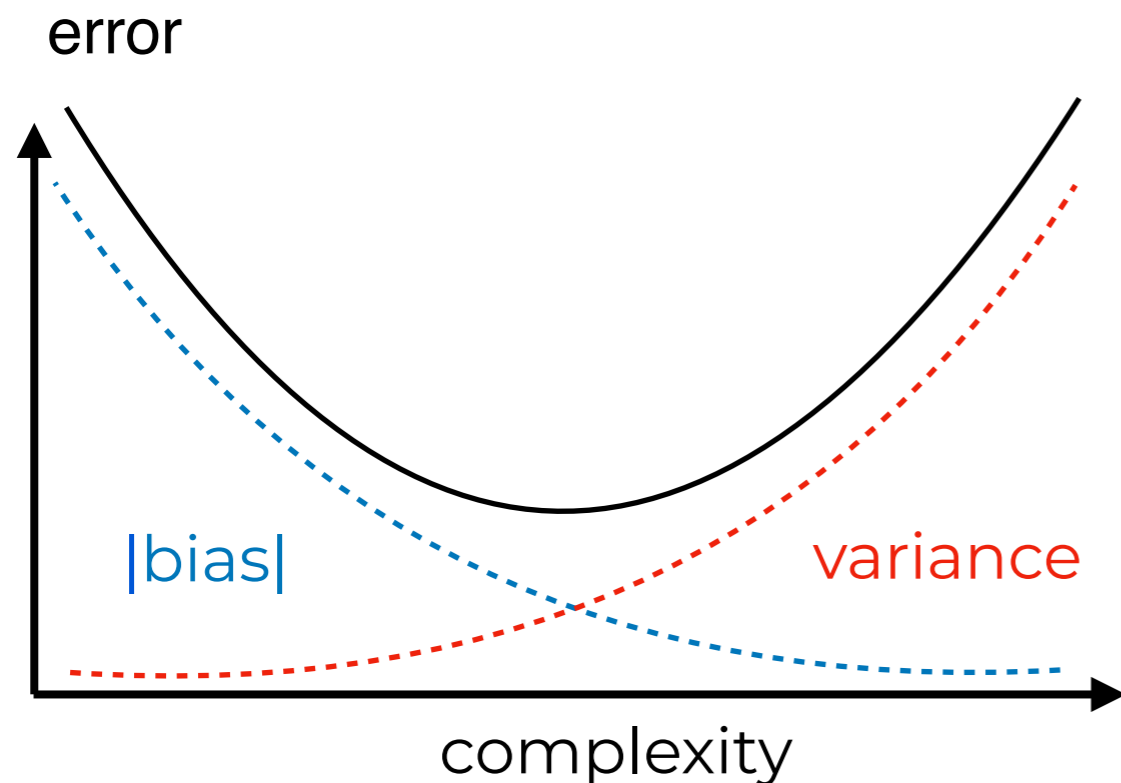
# TRAINING A NEURAL NET

## EXPECTATIONS

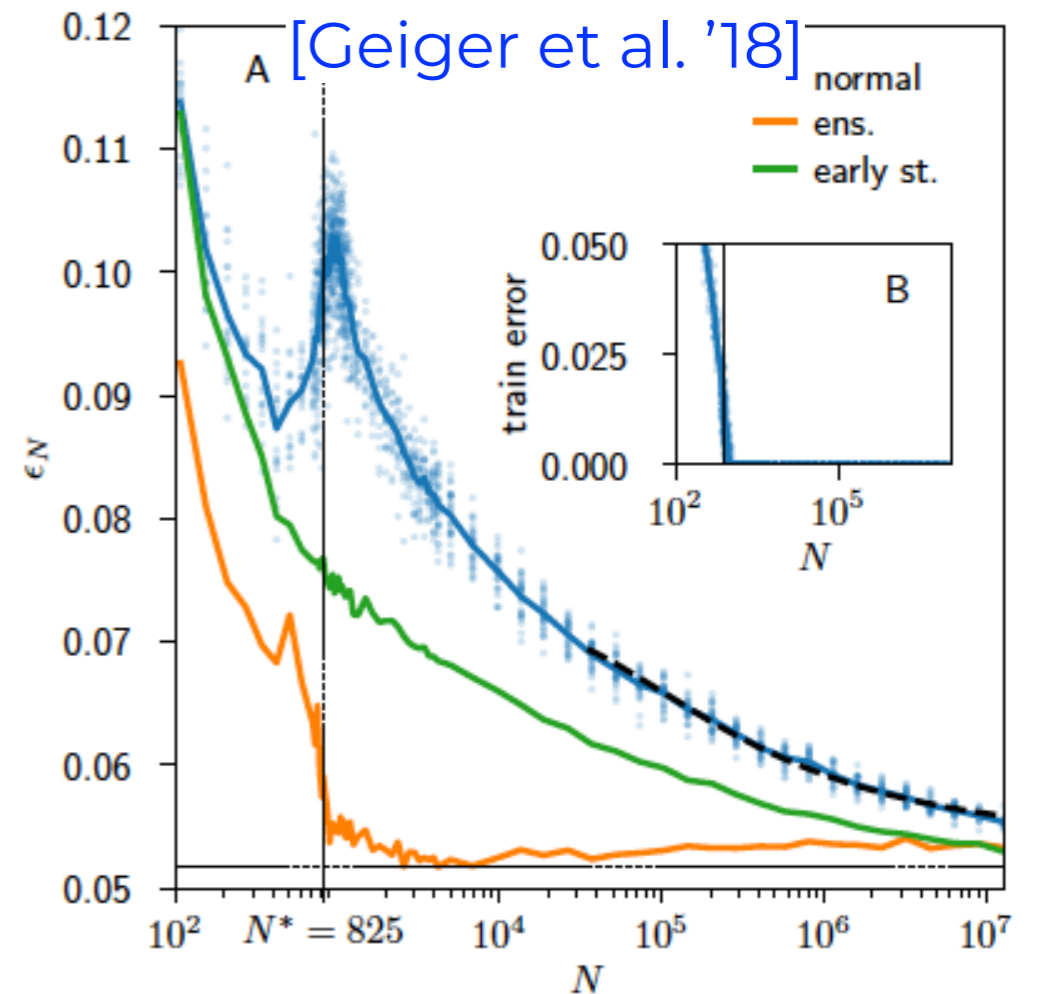


# TRAINING A NEURAL NET

## EXPECTATIONS



## REALITY



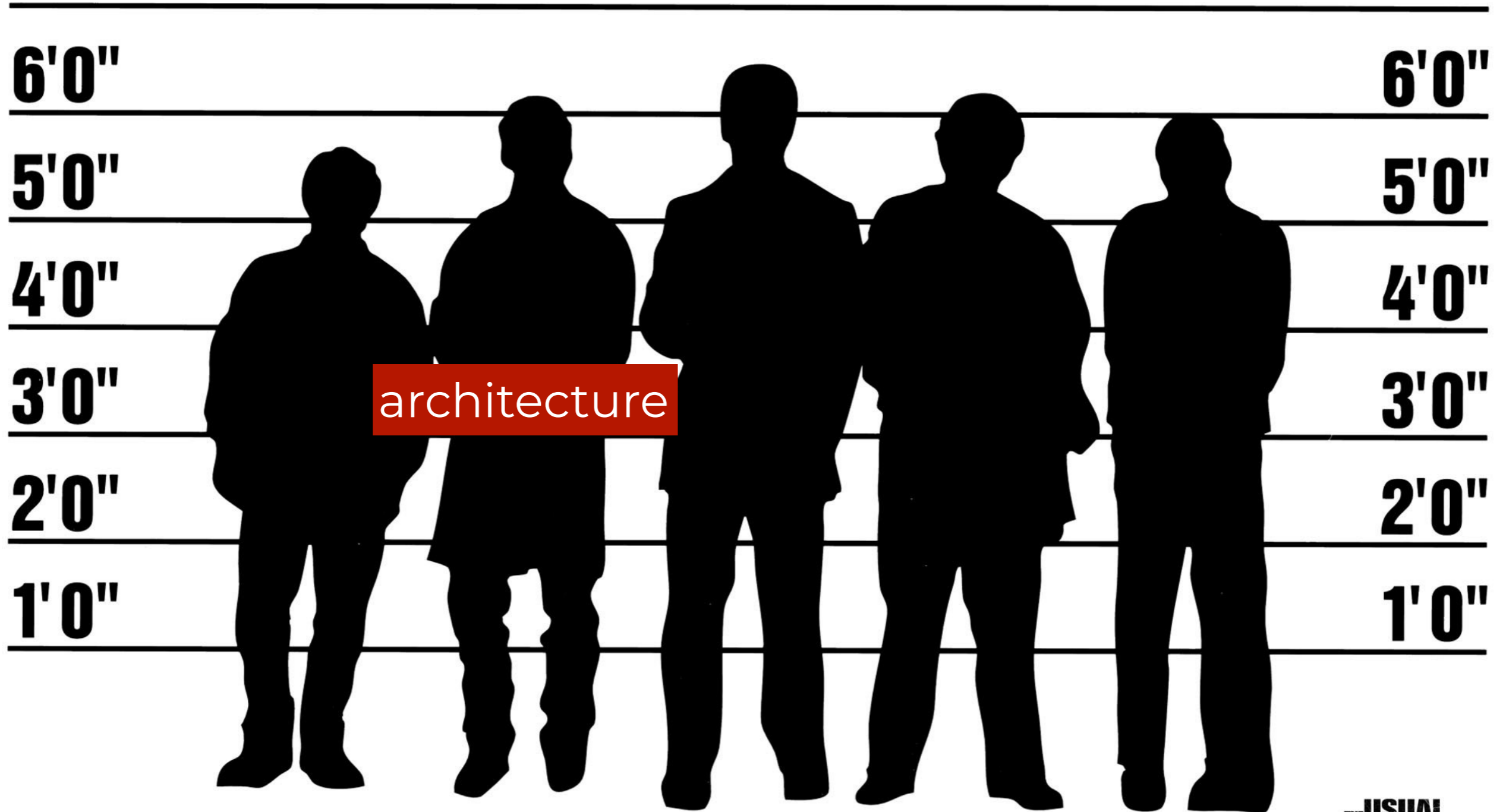
See also [Geman et al. '92; Opper '95; Neyshabur, Tomiyoka, Srebro, 2015; Advani-Saxe 2017; Belkin, Hsu, Ma, Soumik, Mandal 2019; Nakkiran et al. 2019]

# The usual suspects

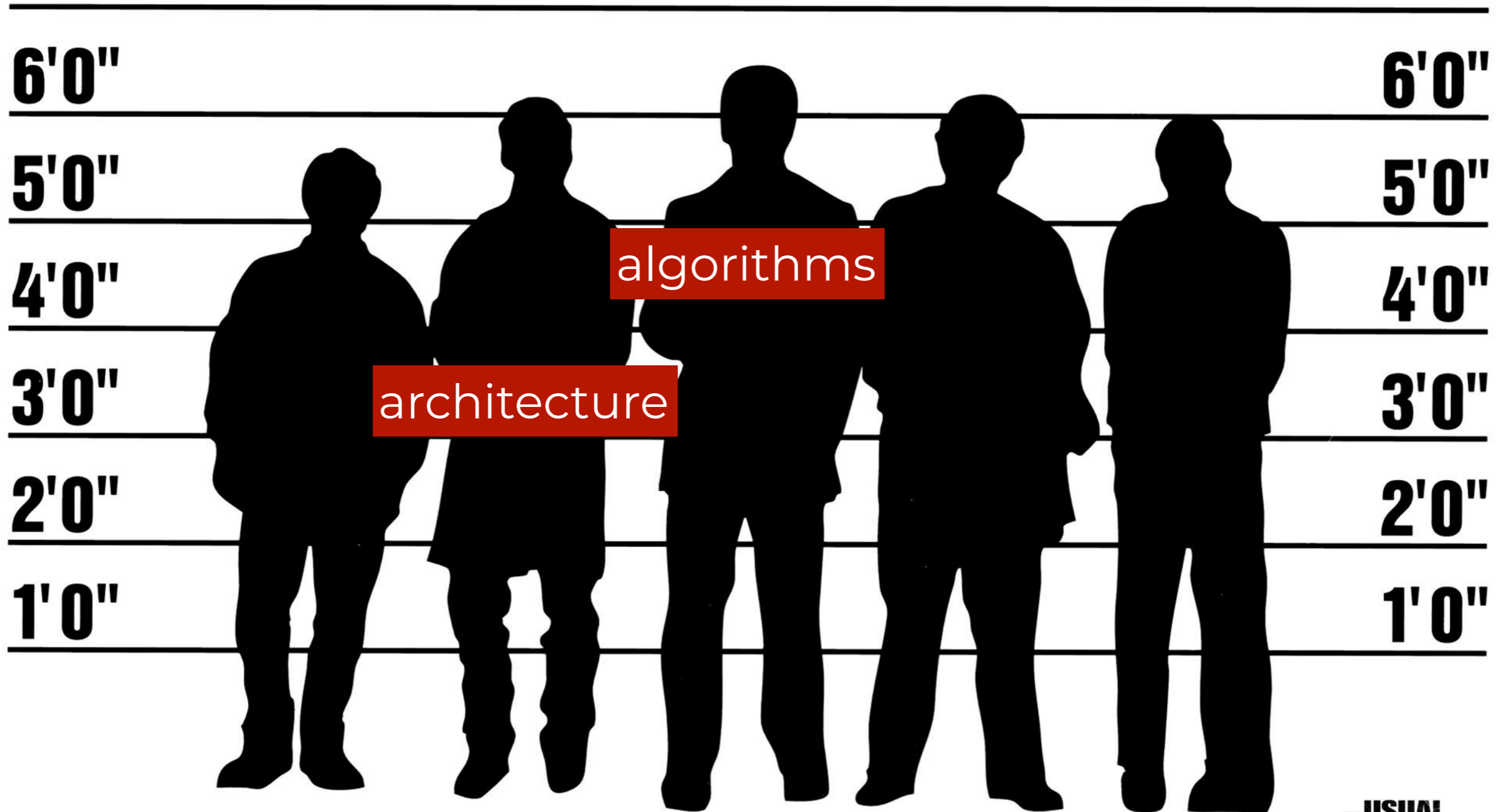


**THE USUAL SUSPECTS**

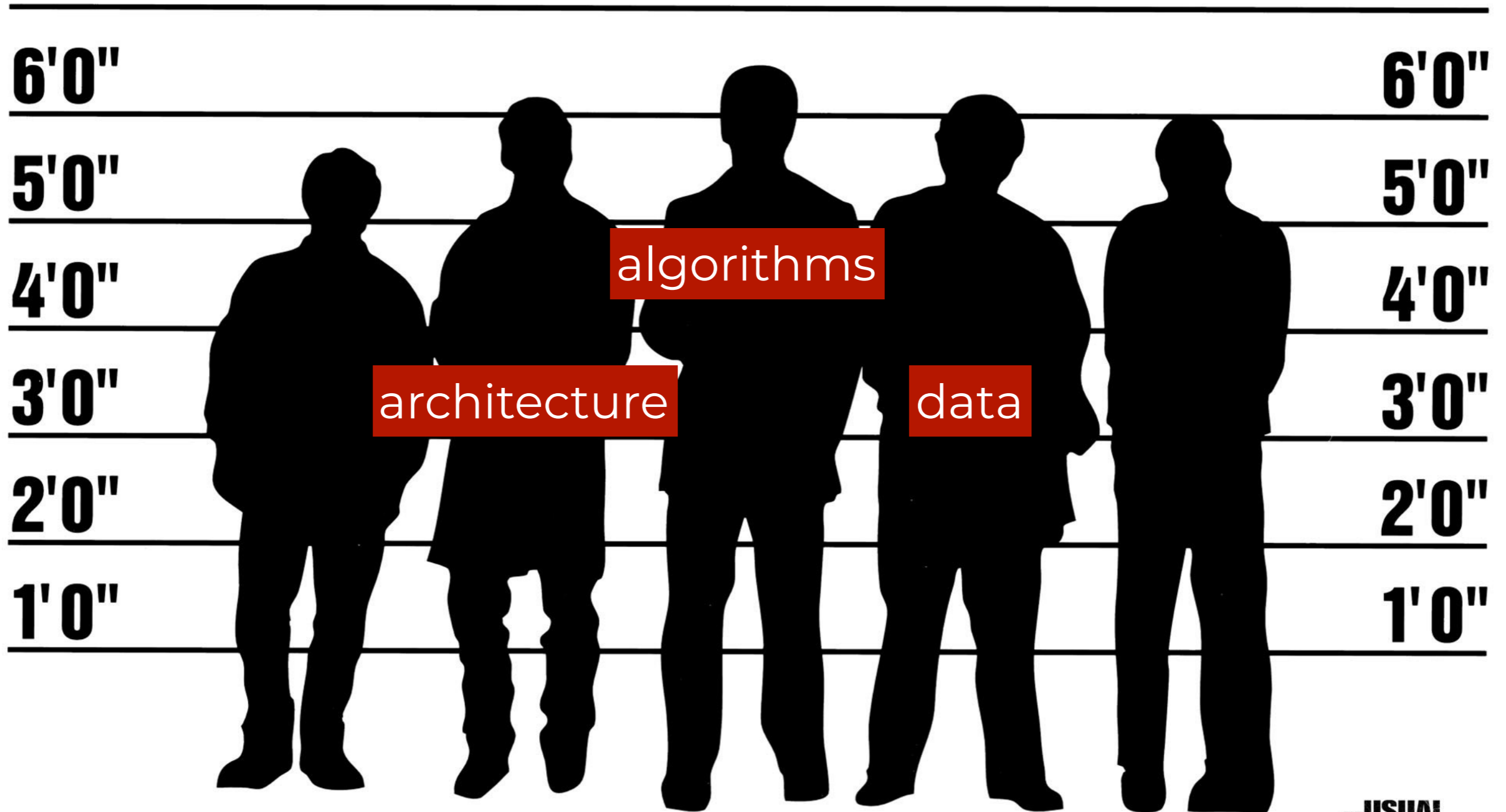
# The usual suspects



# The usual suspects



# The usual suspects

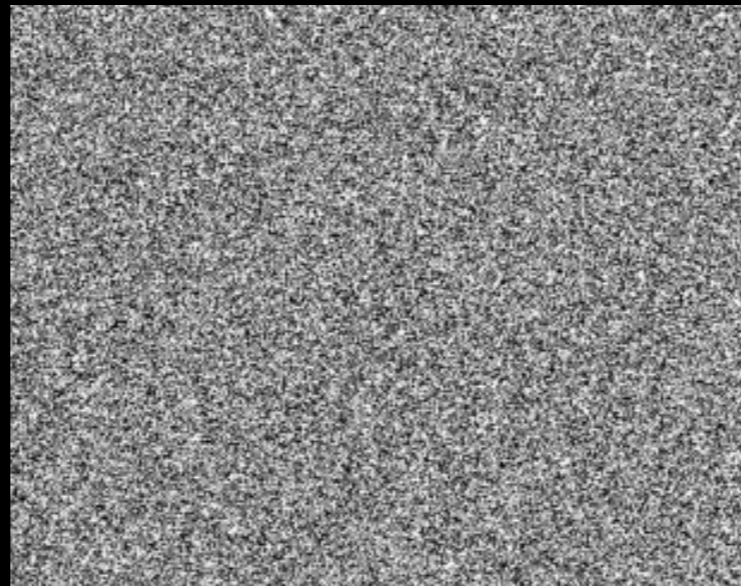


# The two theory cultures

## DATA



**What worst-case  
analysis  
think it looks like**



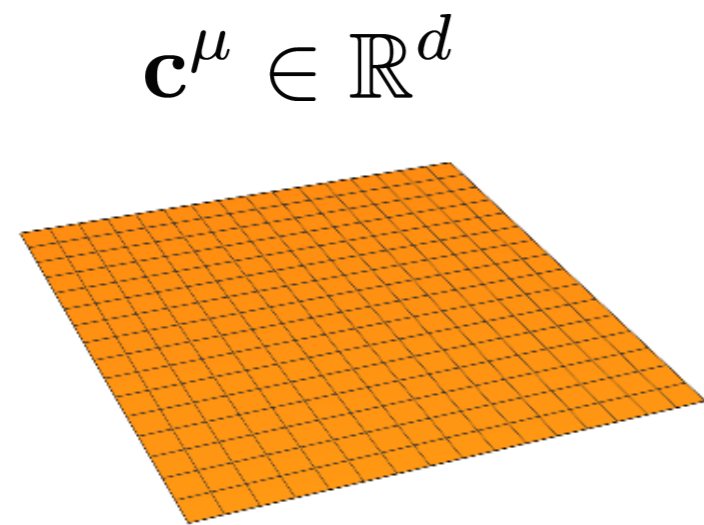
**What typical-case  
analysis  
think it looks like**



**What it really  
looks like**

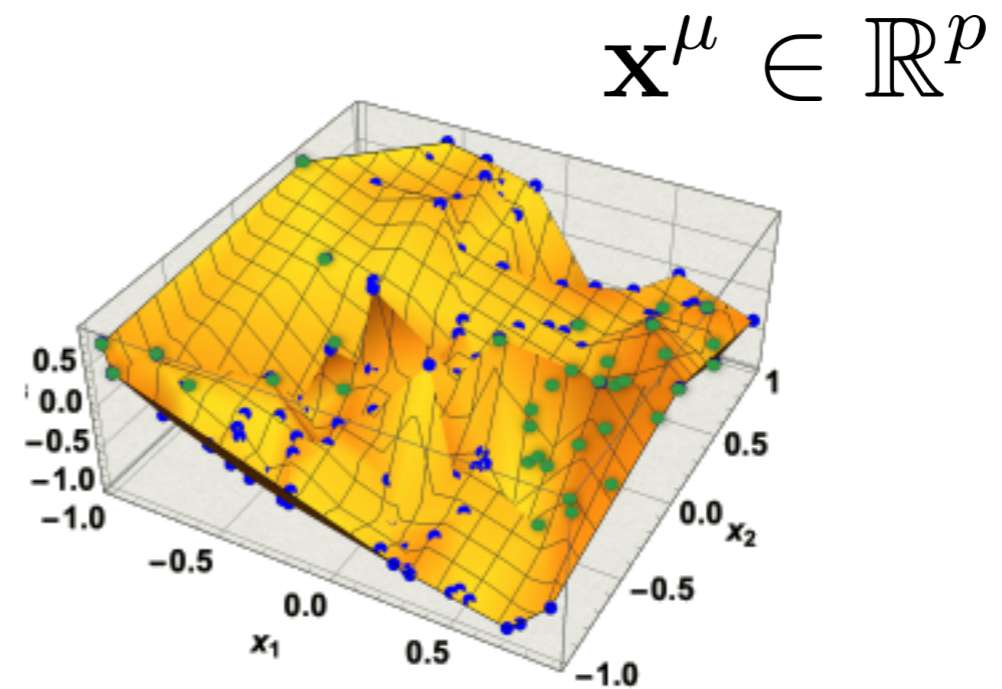


# Spoiler



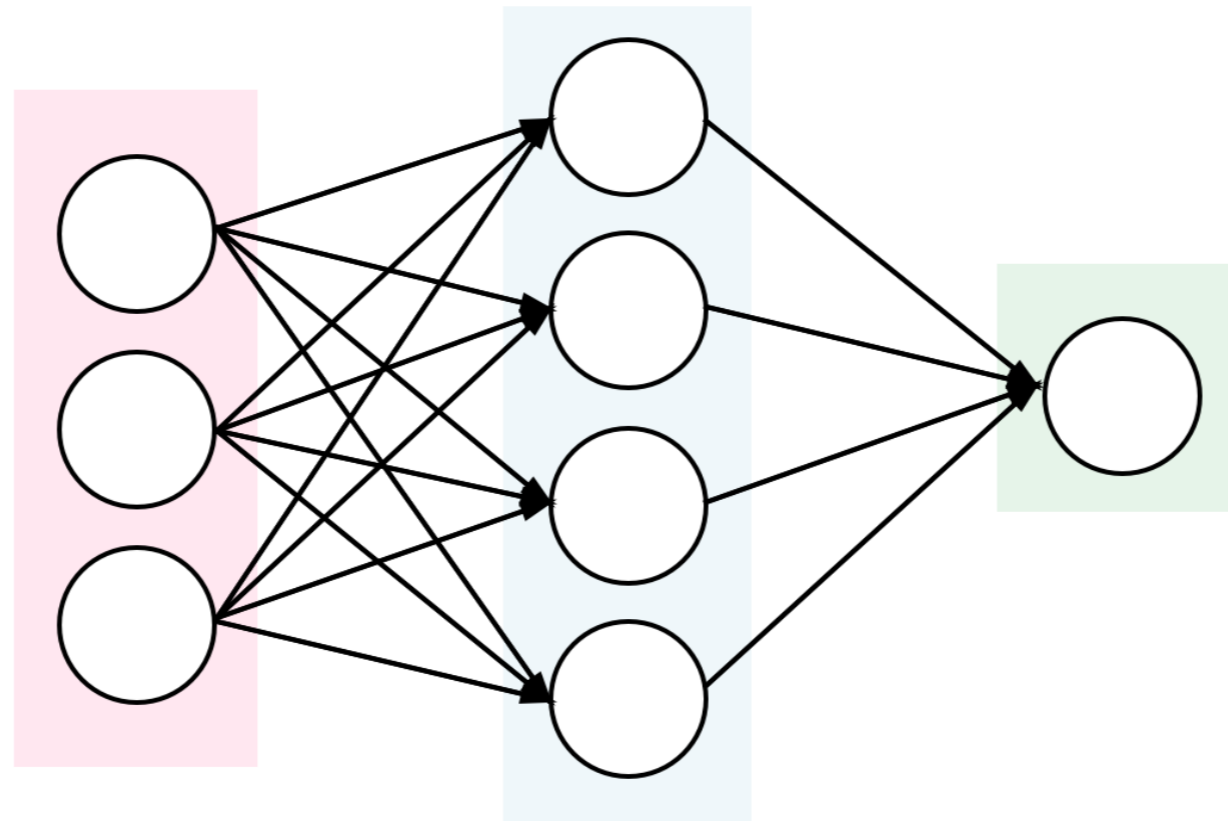
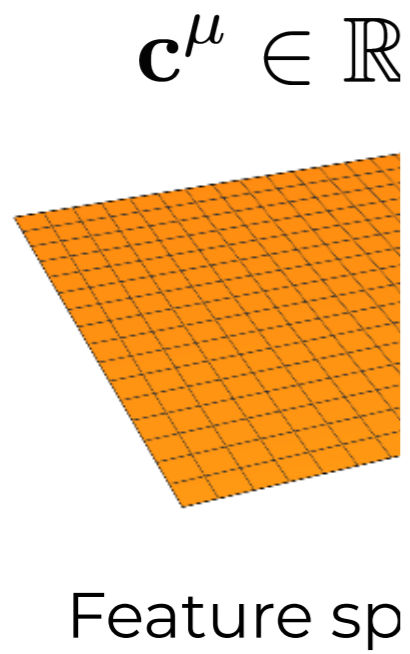
Feature space

$\sigma(\mathbf{F}^\top \cdot)$

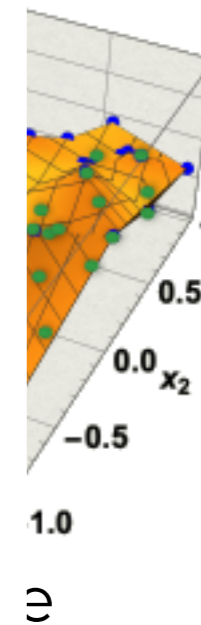


Input space

# Spoiler

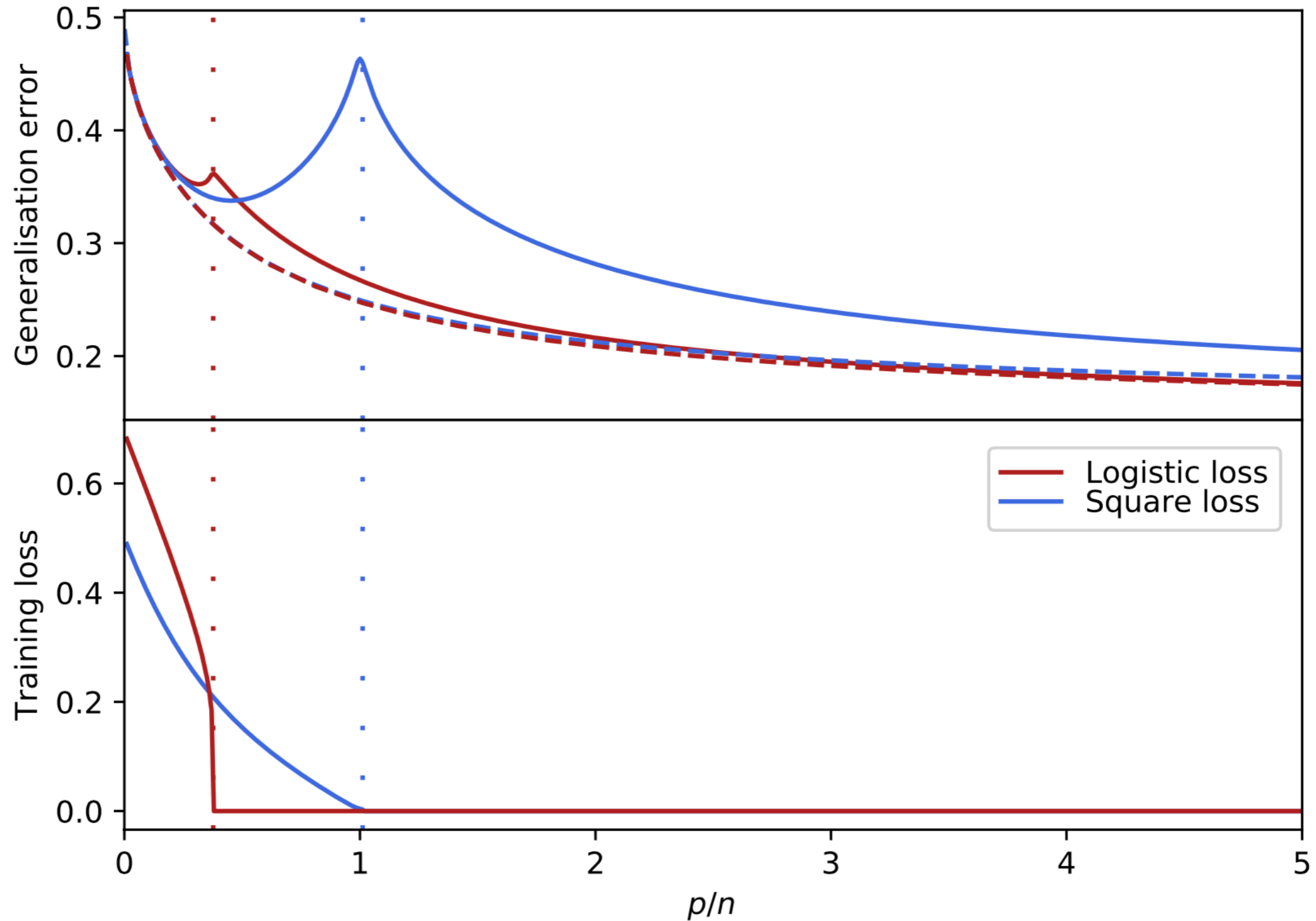


$\mathbf{x}^\mu \in \mathbb{R}^p$

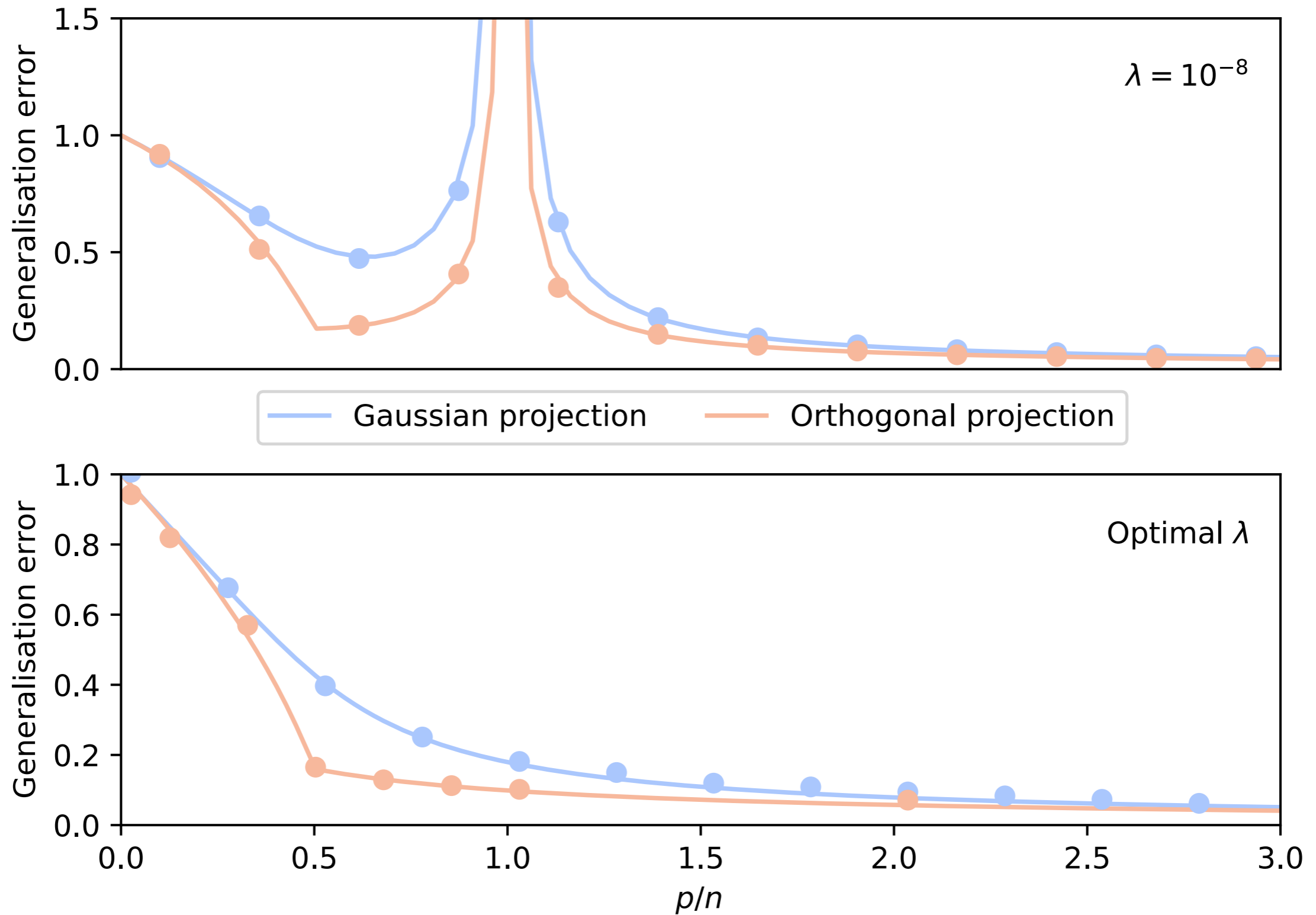


# Spoiler

**SPOILER  
ALERT**

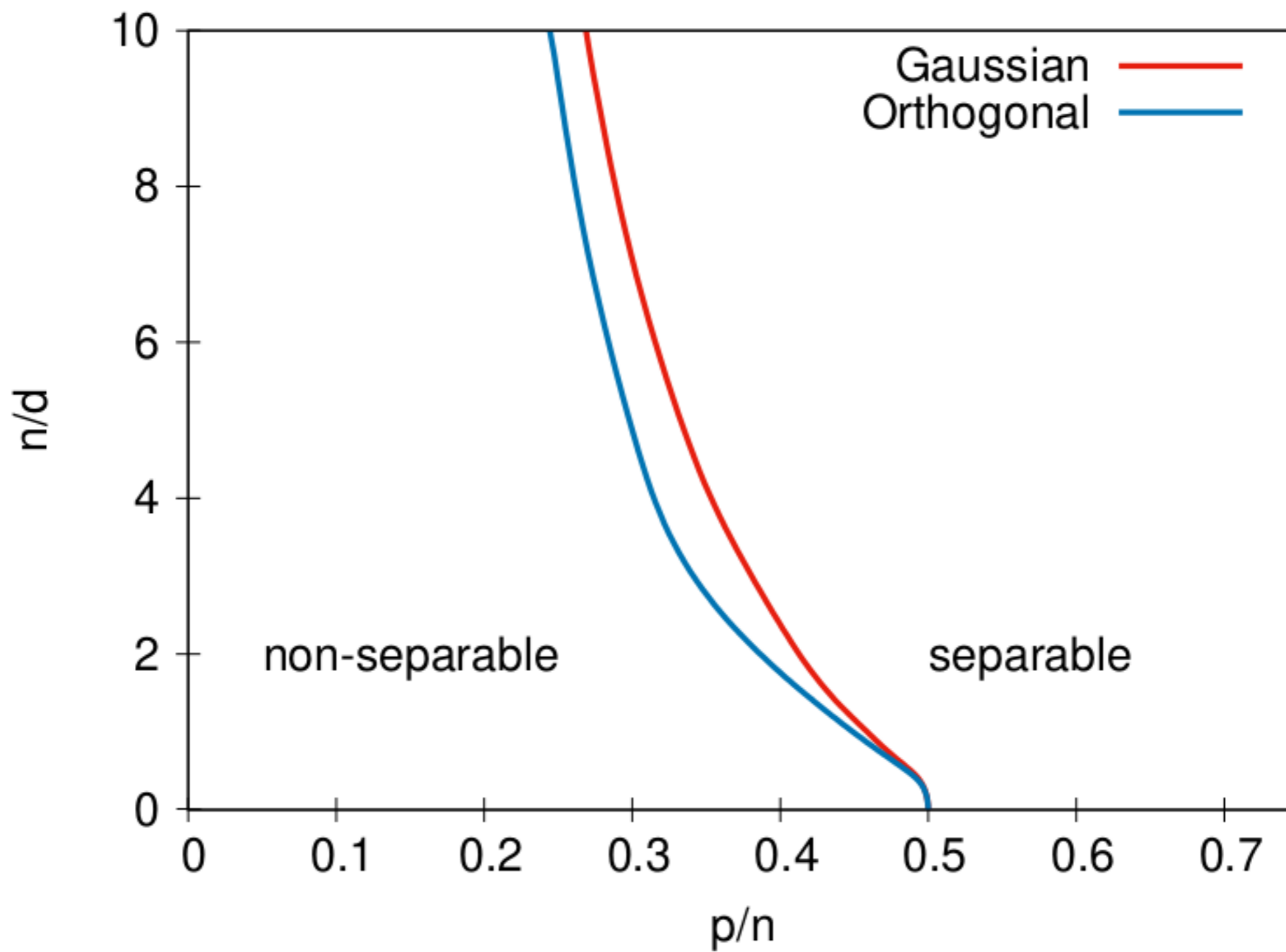


# Spoiler



# Spoiler

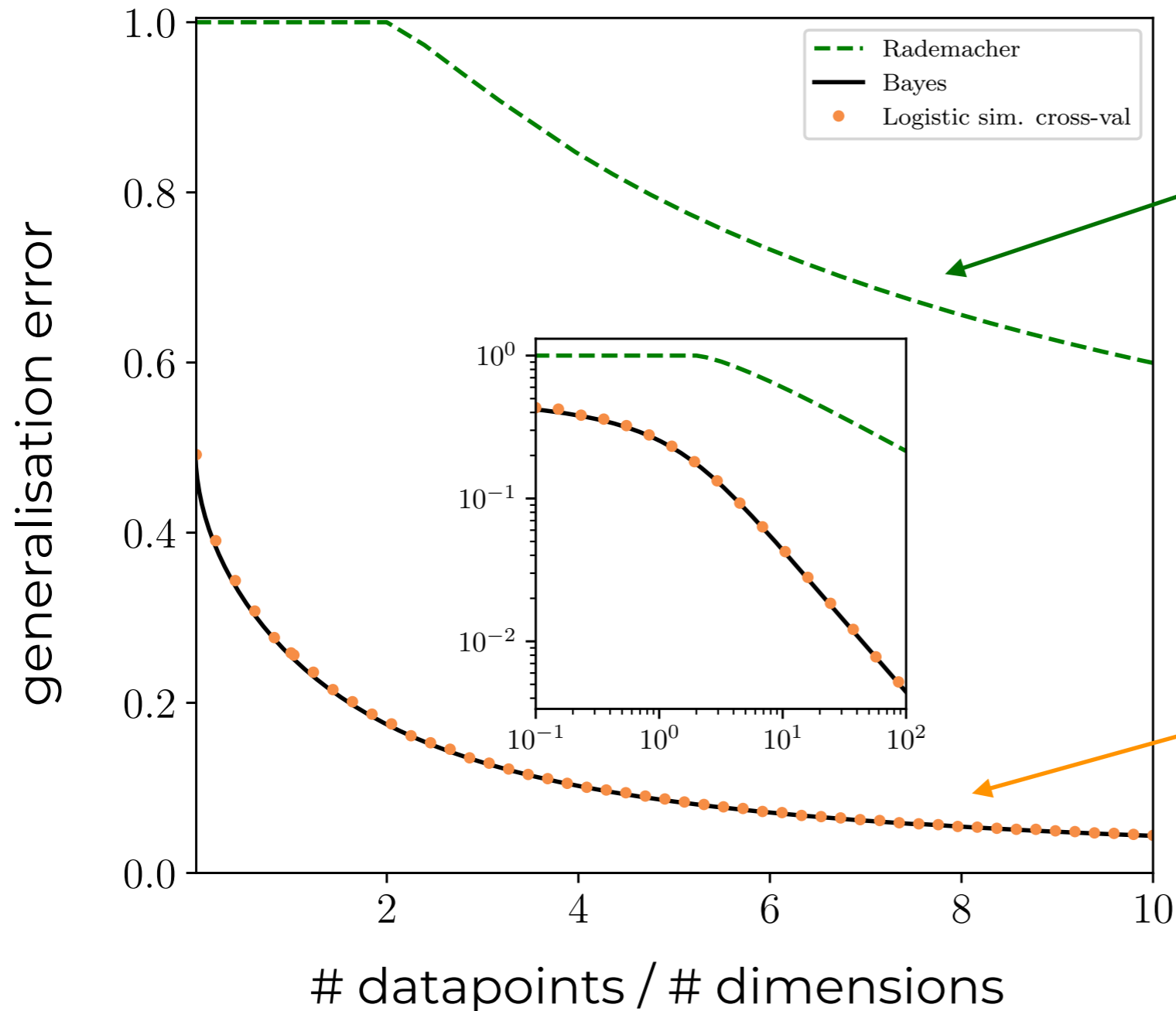
**SPOILER  
ALERT**



Worst-case vs. typical-case:  
A concrete example

# Concrete example

Dataset  $D = \{\mathbf{x}^\mu, y^\mu\}_{\mu=1}^n$  with  $\mathbf{x}^\mu \sim \mathcal{N}(0, \mathbf{I}_d)$  and labels  $y^\mu = \text{sign}(\mathbf{x}^\mu \cdot \boldsymbol{\theta}^0)$



Radamacher bound  
for function class  
 $f_{\boldsymbol{\theta}}(\mathbf{x}) = \text{sign}(\mathbf{x} \cdot \boldsymbol{\theta})$

Out-of-the-box  
Logistic regression  
(Sklearn)

Can we do better?



# Hidden Manifold Model

[Goldt, Mézard, Krzakala, Zdeborová '19]

**Idea:** dataset where both data points and labels only depend on a subset of latent variables.

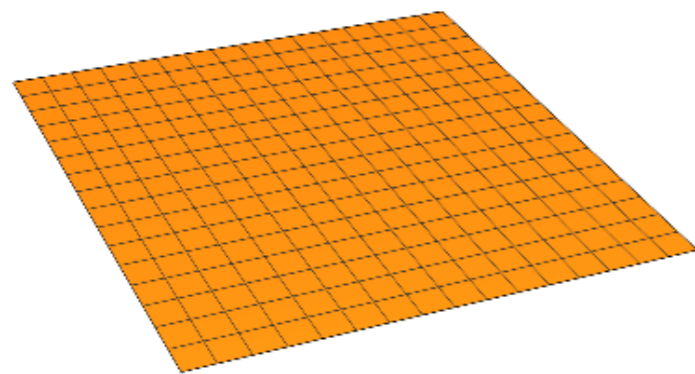
$$D = \{ \mathbf{x}^\mu, y^\mu \}_{\mu=1}^n$$

$$y^\mu = f^0 \left( \frac{\mathbf{c}^\mu \cdot \boldsymbol{\theta}^0}{\sqrt{d}} \right)$$

$$\mathbf{x}^\mu = \sigma \left( \frac{\mathbf{F}^\top \mathbf{c}^\mu}{\sqrt{d}} \right)$$

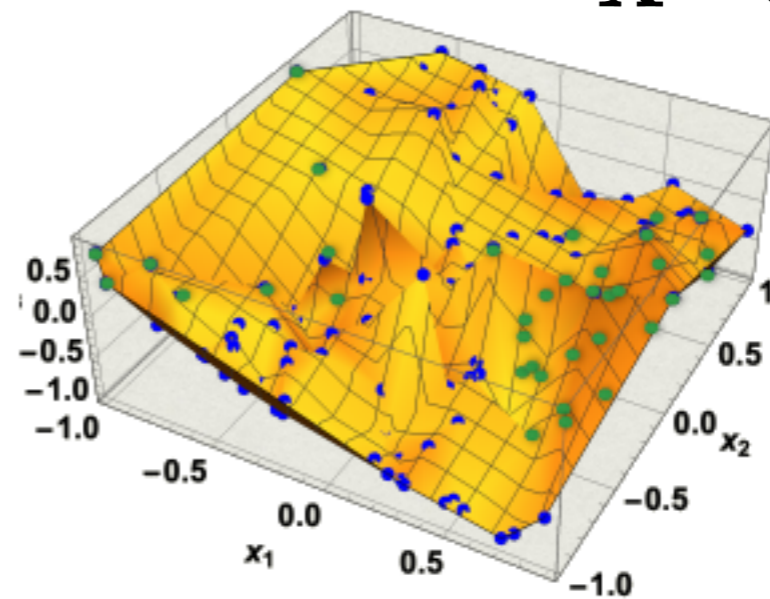
$$\mathbf{x}^\mu \in \mathbb{R}^p$$

$$\mathbf{c}^\mu \in \mathbb{R}^d$$



Feature space

$$\sigma(\mathbf{F}^\top \cdot)$$



Input space

Aim: study classification and regression tasks on this dataset

# The task

Learn the labels using a linear model with empirical risk minimisation

$$\hat{y}^\mu = \hat{f}(\mathbf{x}^\mu \cdot \hat{\mathbf{w}})$$

where:

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \left[ \frac{1}{n} \sum_{\mu=1}^n \ell(y^\mu, \mathbf{x}^\mu \cdot \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right]$$

loss function                      ridge penalty

examples:

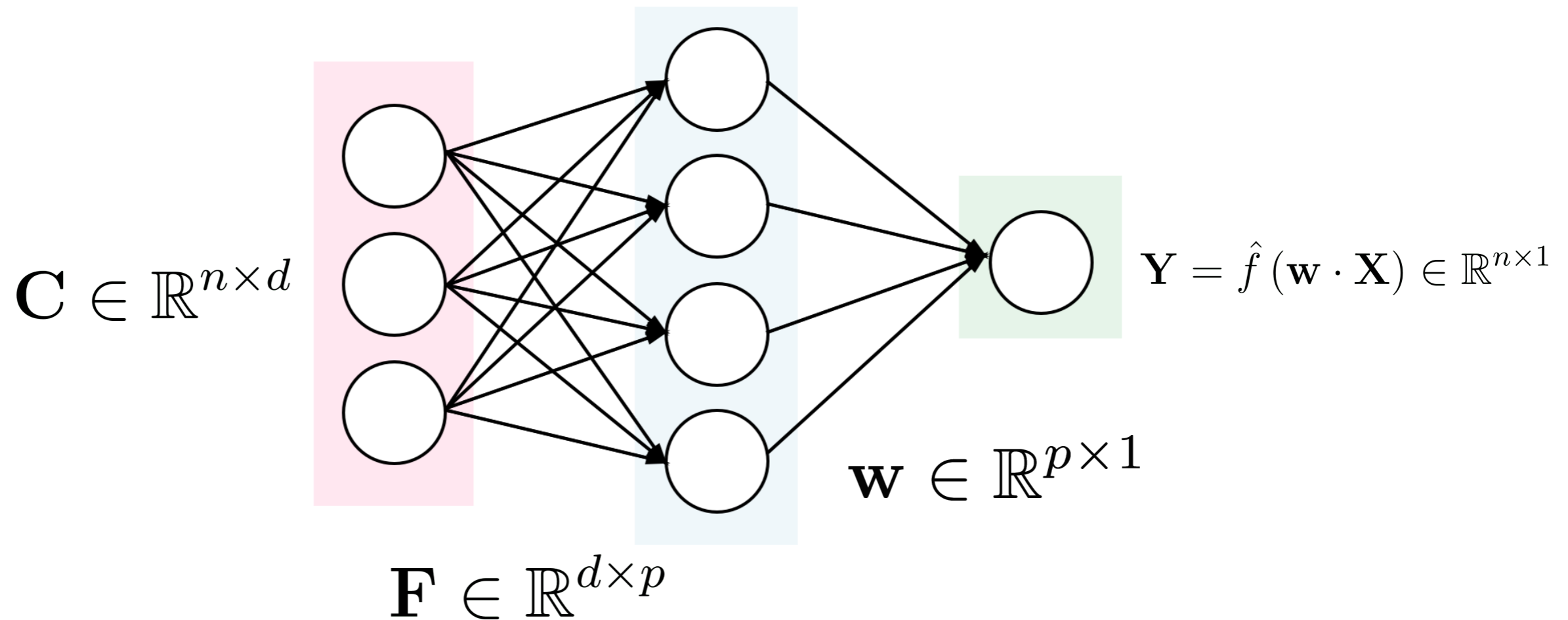
- Ridge regression:  $f^0(x) = \hat{f}(x) = x$                        $\ell(x, y) = \frac{1}{2} (y - x)^2$
- Logistic regression:  $f^0(x) = \hat{f}(x) = \operatorname{sign}(x)$                        $\ell(x, y) = \log(1 + e^{-xy})$

# Two alternative points of view

[Williams '98,'07; Retch, Raimi '07; Montanari, Mei 19']

Dataset  $D = \{\mathbf{c}^\mu, y^\mu\}_{\mu=1}^n$

$$\mathbf{X} = \Phi_{\mathbf{F}}(\mathbf{C}) = \sigma\left(\frac{\mathbf{C}\mathbf{F}}{\sqrt{d}}\right) \in \mathbb{R}^{n \times p}$$

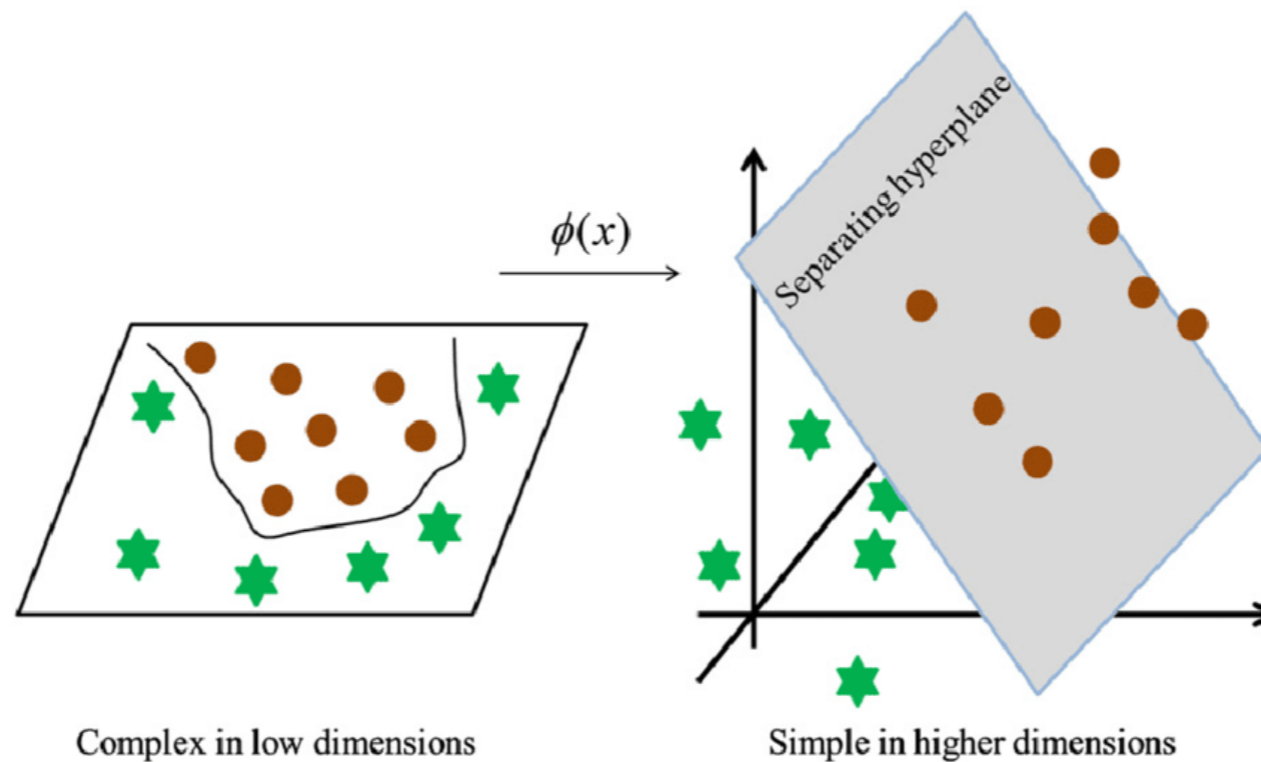


# Two alternative points of view

[Williams '98,'07; Retch, Raimi '07; Montanari, Mei 19']

Dataset  $D = \{\mathbf{c}^\mu, y^\mu\}_{\mu=1}^n$

Feature map  $\Phi_F(\mathbf{c}) = \sigma(\mathbf{F}^\top \mathbf{c})$



$$\Phi_F(\mathbf{c})\Phi_F(\mathbf{c}') \xrightarrow{p \rightarrow \infty} K(\mathbf{c}, \mathbf{c}')$$

Mercer's theorem

Main result:  
Asymptotic generalisation error  
for arbitrary  
loss  $\ell$  and projection  $F$

## Definitions:

Consider the unique fixed point of the following system of equations

$$\left\{ \begin{array}{l} \hat{V}_s = \frac{\alpha}{\gamma} \kappa_1^2 \mathbb{E}_{\xi,y} \left[ \mathcal{L}(y, \omega_0) \frac{\partial_\omega \eta(y, \omega_1)}{V} \right], \\ \hat{q}_s = \frac{\alpha}{\gamma} \kappa_1^2 \mathbb{E}_{\xi,y} \left[ \mathcal{L}(y, \omega_0) \frac{(\eta(y, \omega_1) - \omega_1)^2}{V^2} \right], \\ \hat{m}_s = \frac{\alpha}{\gamma} \kappa_1 \mathbb{E}_{\xi,y} \left[ \partial_\omega \mathcal{L}(y, \omega_0) \frac{(\eta(y, \omega_1) - \omega_1)}{V} \right], \\ \hat{V}_w = \alpha \kappa_\star^2 \mathbb{E}_{\xi,y} \left[ \mathcal{L}(y, \omega_0) \frac{\partial_\omega \eta(y, \omega_1)}{V} \right], \\ \hat{q}_w = \alpha \kappa_\star^2 \mathbb{E}_{\xi,y} \left[ \mathcal{L}(y, \omega_0) \frac{(\eta(y, \omega_1) - \omega_1)^2}{V^2} \right], \end{array} \right. \quad \left\{ \begin{array}{l} V_s = \frac{1}{\hat{V}_s} \left( 1 - z g_\mu(-z) \right), \\ q_s = \frac{\hat{m}_s^2 + \hat{q}_s}{\hat{V}_s} \left[ 1 - 2z g_\mu(-z) + z^2 g'_\mu(-z) \right] \\ \quad - \frac{\hat{q}_w}{(\lambda + \hat{V}_w) \hat{V}_s} \left[ -z g_\mu(-z) + z^2 g'_\mu(-z) \right], \\ m_s = \frac{\hat{m}_s}{\hat{V}_s} \left( 1 - z g_\mu(-z) \right), \\ V_w = \frac{\gamma}{\lambda + \hat{V}_w} \left[ \frac{1}{\gamma} - 1 + z g_\mu(-z) \right], \\ q_w = \gamma \frac{\hat{q}_w}{(\lambda + \hat{V}_w)^2} \left[ \frac{1}{\gamma} - 1 + z^2 g'_\mu(-z) \right], \\ \quad + \frac{\hat{m}_s^2 + \hat{q}_s}{(\lambda + \hat{V}_w) \hat{V}_s} \left[ -z g_\mu(-z) + z^2 g'_\mu(-z) \right], \end{array} \right. \quad \left\{ \begin{array}{l} \eta(y, \omega) = \operatorname{argmin}_{x \in \mathbb{R}} \left[ \frac{(x - \omega)^2}{2V} + \ell(y, x) \right] \\ \mathcal{L}(y, \omega) = \int \frac{dx}{\sqrt{2\pi V^0}} e^{-\frac{1}{2V^0}(x - \omega)^2} \delta(y - f^0(x)) \end{array} \right.$$

where  $V = \kappa_1^2 V_s + \kappa_\star^2 V_w$ ,  $V^0 = \rho - \frac{M^2}{Q}$ ,  $Q = \kappa_1^2 q_s + \kappa_\star^2 q_w$ ,  $M = \kappa_1 m_s$ ,  $\omega_0 = M/\sqrt{Q}\xi$ ,  $\omega_1 = \sqrt{Q}\xi$  and  $g_\mu$  is the Stieltjes transform of  $FF^T$

$$\kappa_0 = \mathbb{E}[\sigma(z)], \kappa_1 \equiv \mathbb{E}[z\sigma(z)], \kappa_\star \equiv \mathbb{E}[\sigma(z)^2] - \kappa_0^2 - \kappa_1^2 \quad \text{and} \quad \vec{z}^\mu \sim \mathcal{N}(\vec{0}, I_p)$$

In the high-dimensional limit:

$$\epsilon_{gen} = \mathbb{E}_{\lambda, \nu} \left[ (f^0(\nu) - \hat{f}(\lambda))^2 \right]$$

$$\text{with } (\nu, \lambda) \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \rho & M^\star \\ M^\star & Q^\star \end{pmatrix} \right)$$

$$\mathcal{L}_{\text{training}} = \frac{\lambda}{2\alpha} q_w^\star + \mathbb{E}_{\xi,y} \left[ \mathcal{L}(y, \omega_0^\star) \ell(y, \eta(y, \omega_1^\star)) \right]$$

$$\text{with } \omega_0^\star = M^\star/\sqrt{Q^\star}\xi, \omega_1^\star = \sqrt{Q^\star}\xi$$

Agrees with [\[Mei-Montanari '19\]](#) who solved a particular case using random matrix theory:  
linear function  $f^0$ ,  $\ell(x, y) = \|x - y\|_2^2$  & Gaussian random weights  $\mathbf{F}$

## Technical note: replicated Gaussian Equivalence

An important step in the derivation of this result is the observation that the generalisation and training properties of the dataset  $\{\mathbf{x}^\mu, y^\mu\}_{\mu=1}^n$  are *statistically equivalent* to the following dataset  $\{\tilde{\mathbf{x}}^\mu, y^\mu\}_{\mu=1}^n$  with the same labels but:

$$\tilde{\mathbf{x}}^\mu = \kappa_1 \frac{1}{\sqrt{d}} \mathbf{F}^\top \mathbf{c}^\mu + \kappa_\star \mathbf{z}^\mu \quad \mathbf{z}^\mu \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$$

where the coefficients  $\kappa_1, \kappa_\star$  are chosen to match

$$\kappa_1 = \mathbb{E}_\xi [\xi \sigma(\xi)] \quad \kappa_\star^2 = \mathbb{E}_\xi [\sigma(\xi)^2] - \kappa_1^2 \quad \xi \sim \mathcal{N}(0, 1)$$

Generalisation of an observation in  
[Mei, Montanari 19'; Goldt, Mézard, Krzakala, Zdeborová '19]

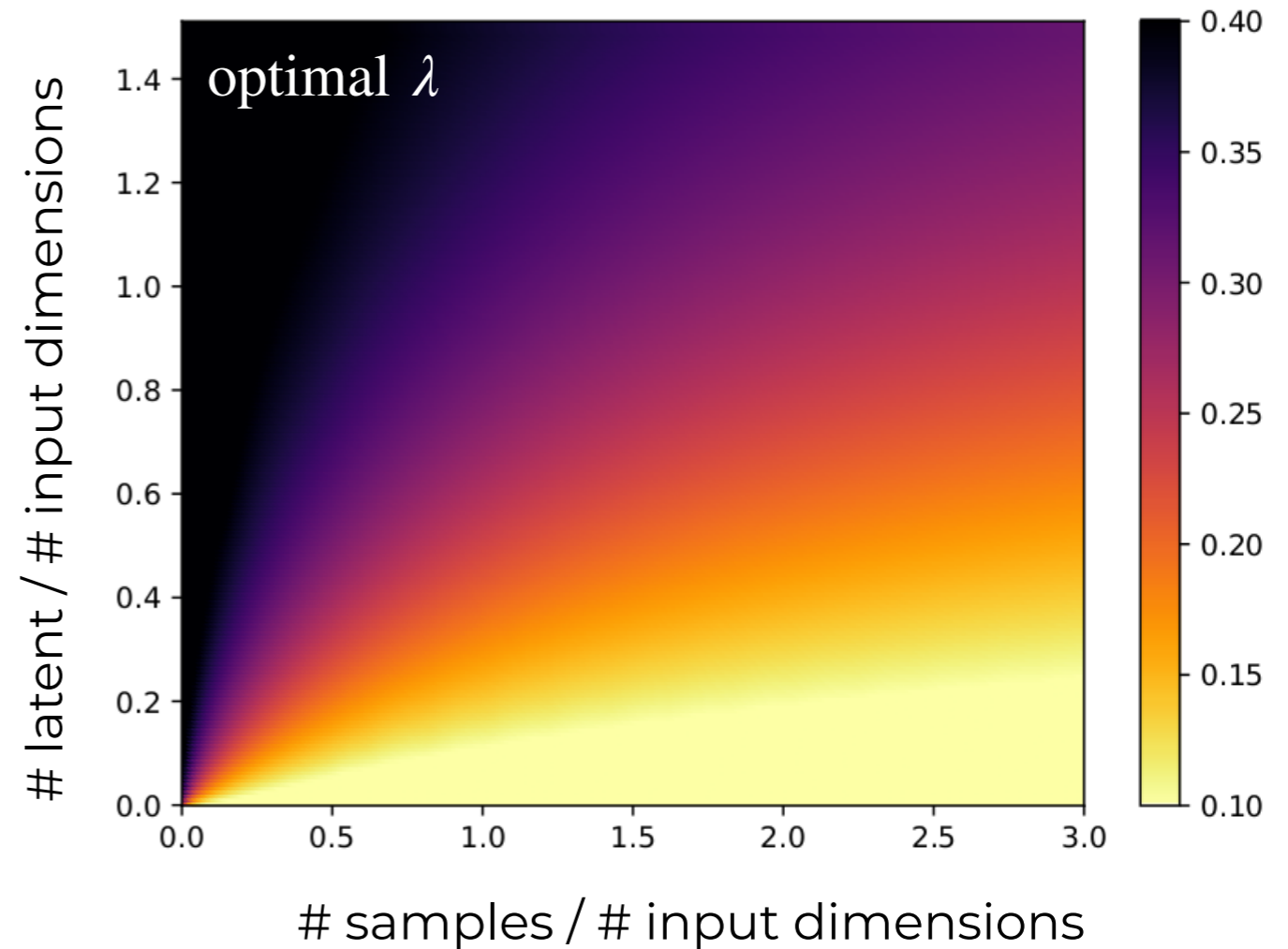
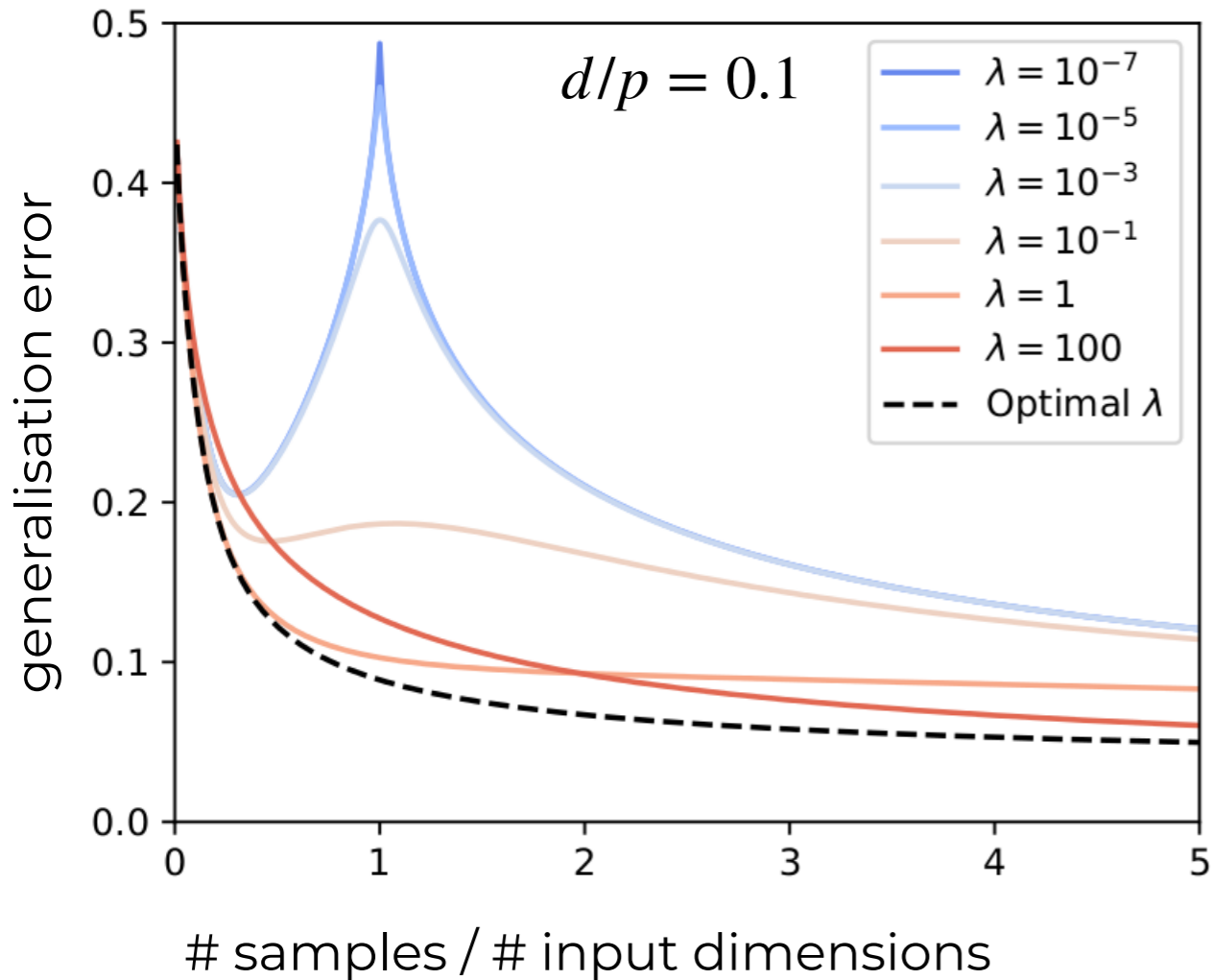


Drawing the consequences  
of our formula

# Learning in the HMM

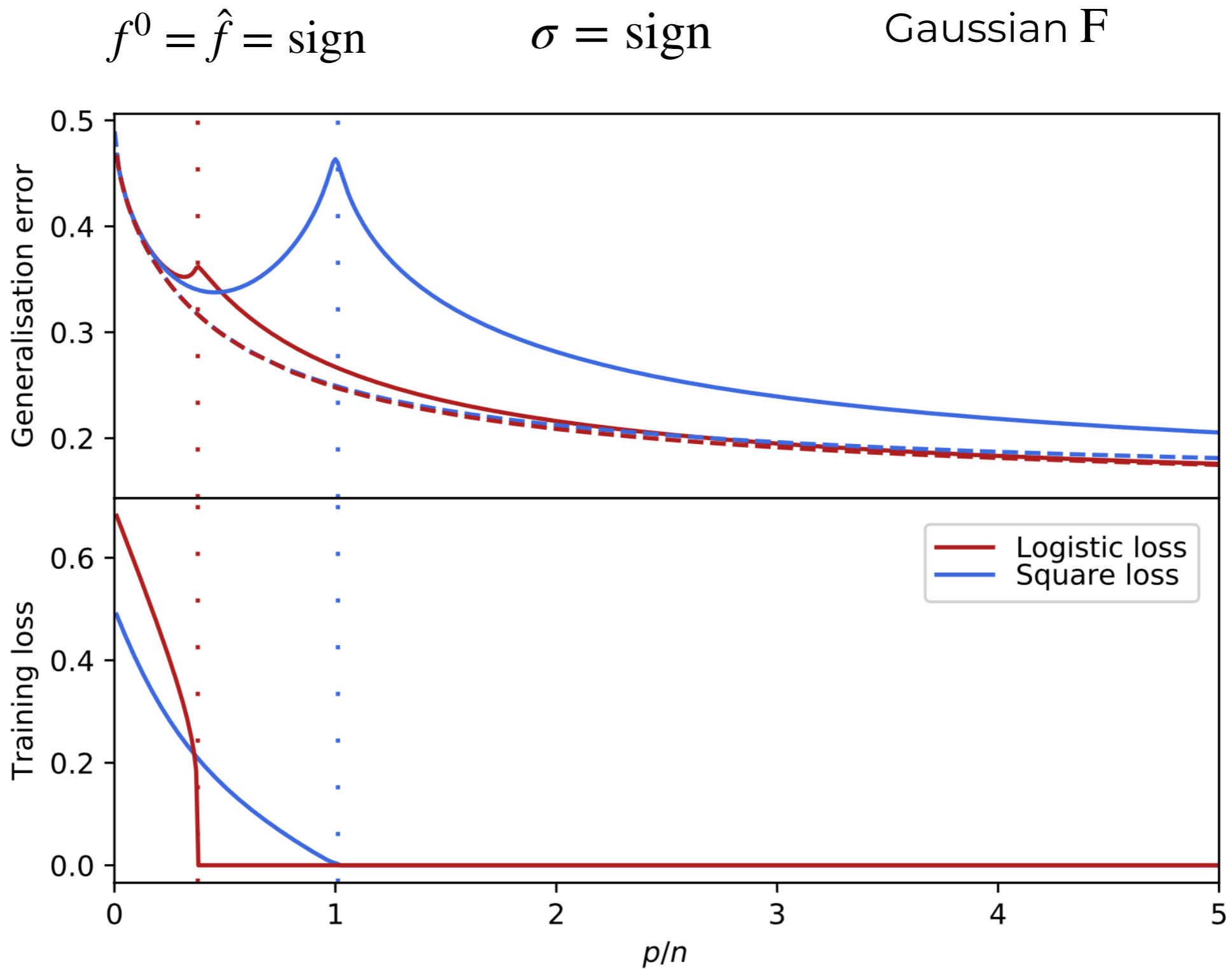
$$f^0 = \hat{f} = \text{sign} \quad \sigma = \text{erf}$$

$$l(x, y) = \frac{1}{2}(x - y)^2 \quad \text{Gaussian F}$$



Good generalisation performance for small latent space,  
even for small sample complexities

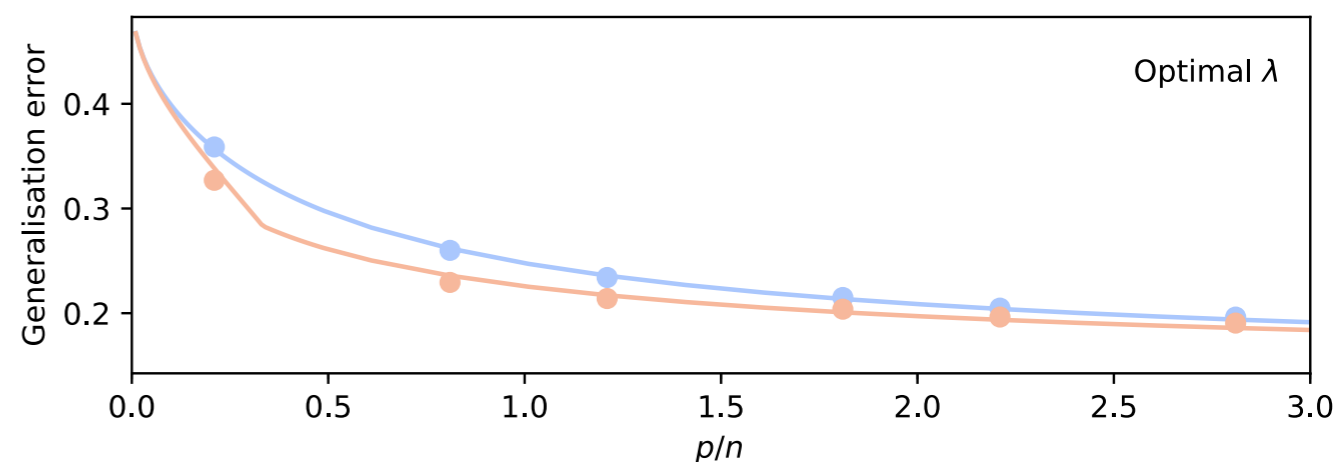
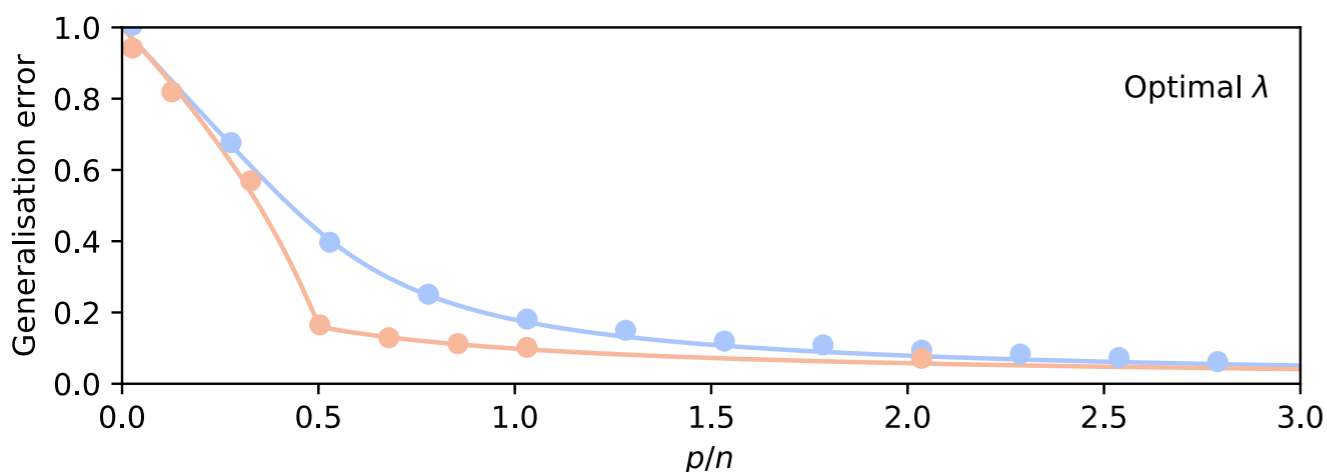
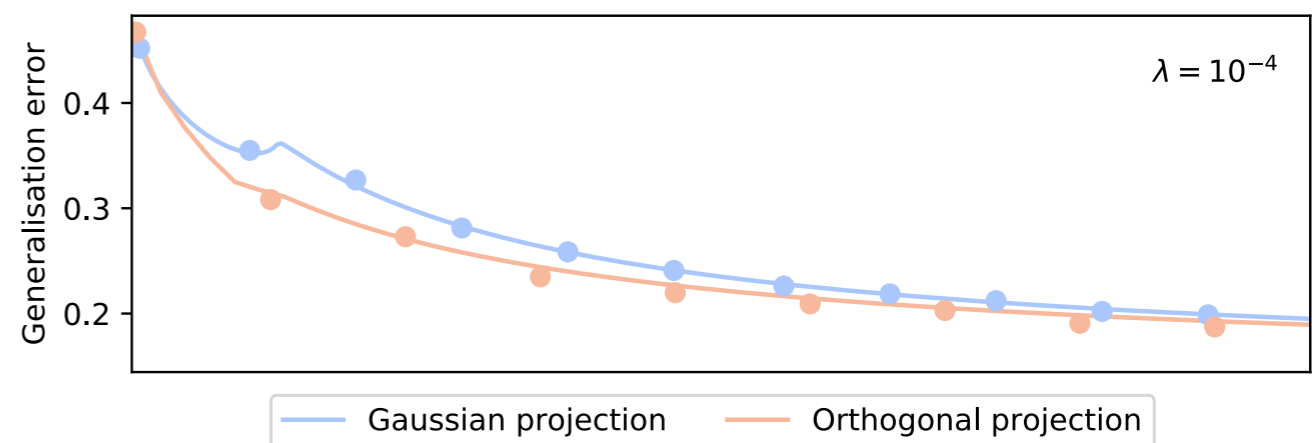
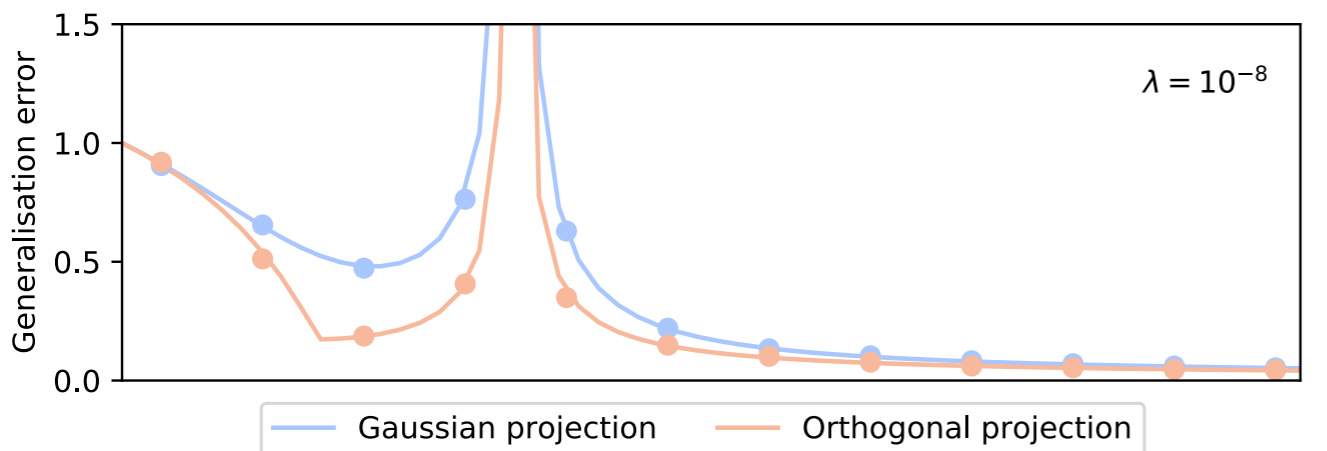
# Classification tasks



# Random vs. orthogonal projections

## Ridge regression

## Logistic regression



— First layer: random i.i.d. Gaussian Matrix

$$F_{\rho}^i \sim \mathcal{N}(0, 1/d)$$

— First layer: subsampled Fourier matrix

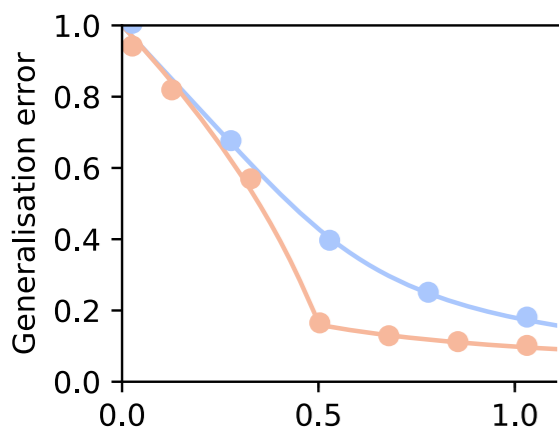
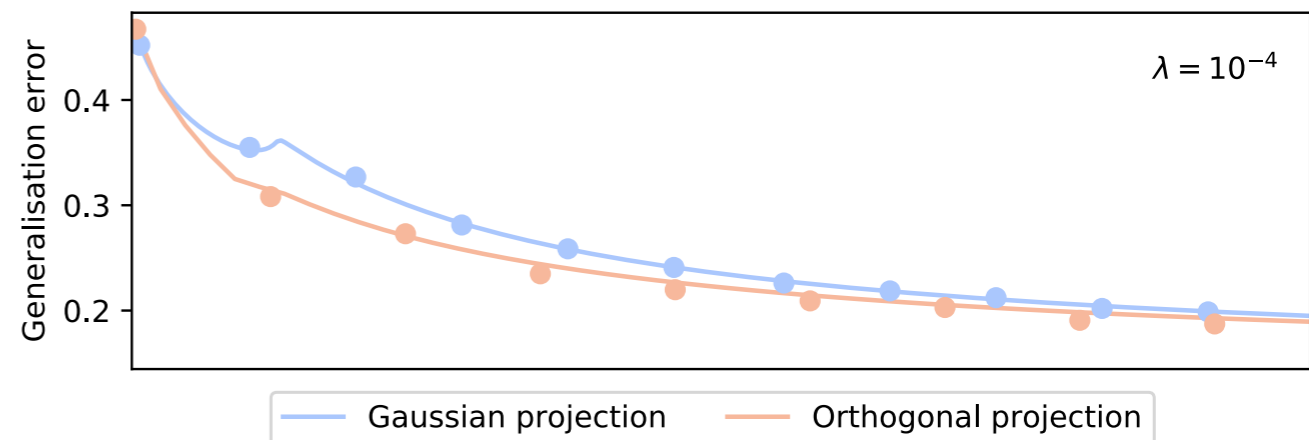
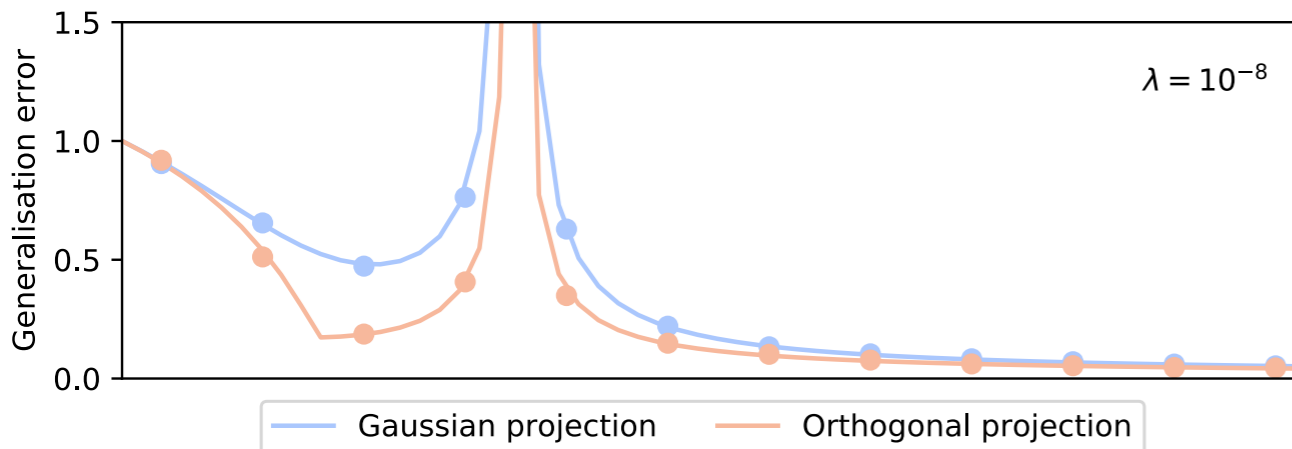
$$F = U^T D V$$

$$U, V \sim \text{Haar}$$

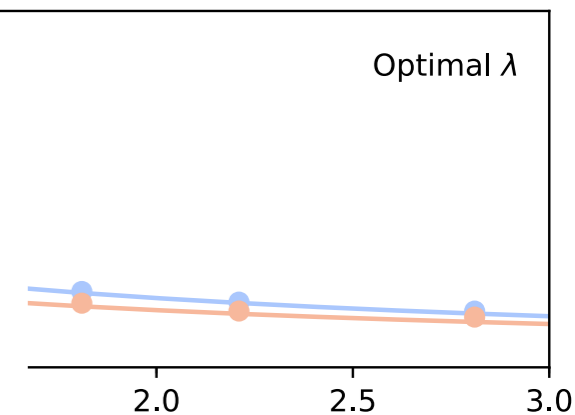
# Random vs. orthogonal projections

Ridge regression

Logistic regression



**The Unreasonable Effectiveness of Structured  
Random Orthogonal Embeddings**



**First**  
**First**

**Krzysztof Choromanski \***  
Google Brain Robotics  
kchoro@google.com

**Mark Rowland \***  
University of Cambridge  
mr504@cam.ac.uk

**Adrian Weller**  
University of Cambridge and Alan Turing Institute  
aw665@cam.ac.uk

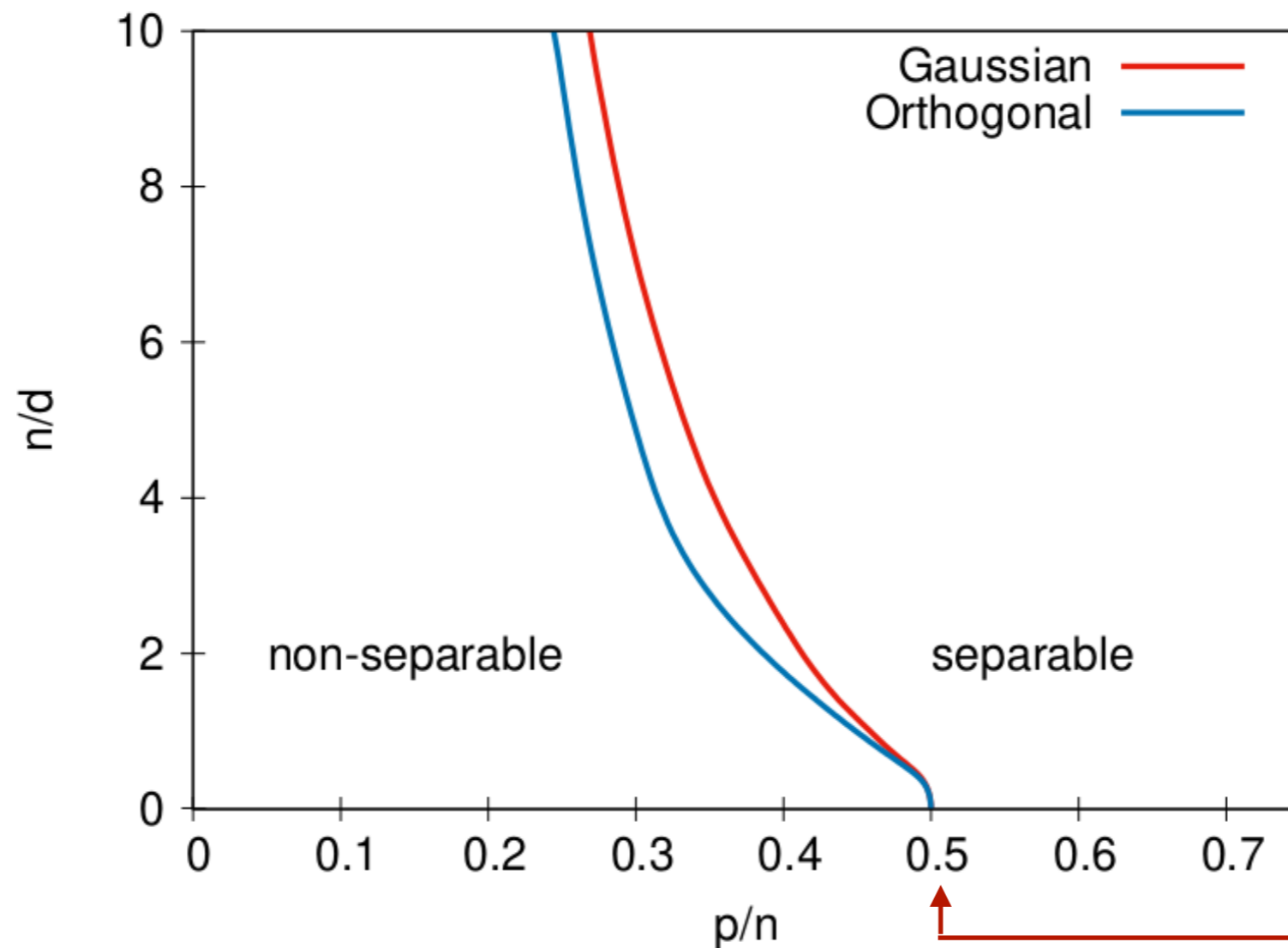
$1/d$ )  
[NIPS, '17]  
 $\mathbb{R}$

# Separability transition in logistic regression

$$f^0 = \hat{f} = \text{sign}$$

$$\sigma = \text{erf}$$

$$l(x, y) = \log(1 + e^{-xy})$$

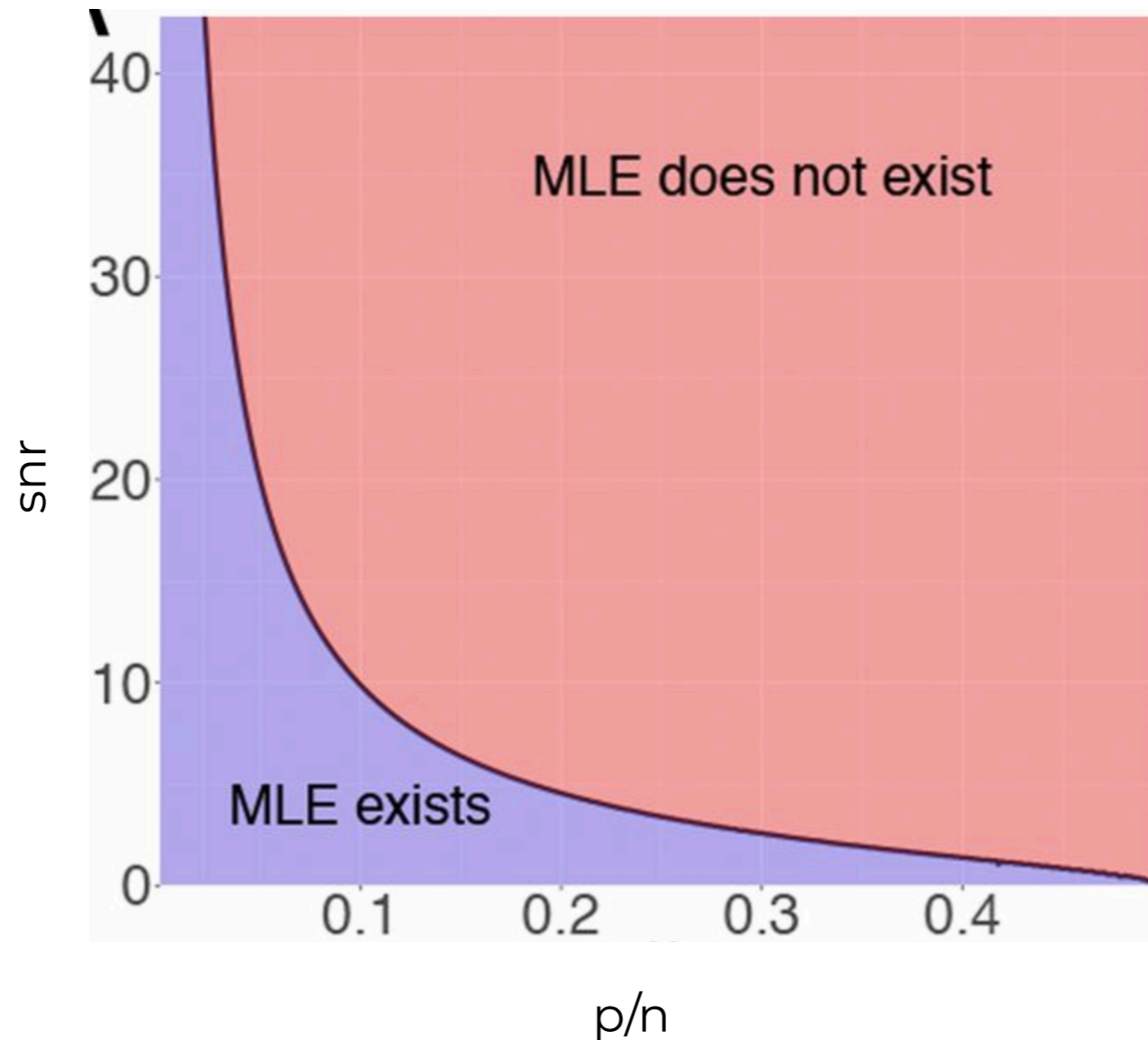


Cover theory '65

# Separability transition in logistic regression

## A modern maximum-likelihood theory for high-dimensional logistic regression

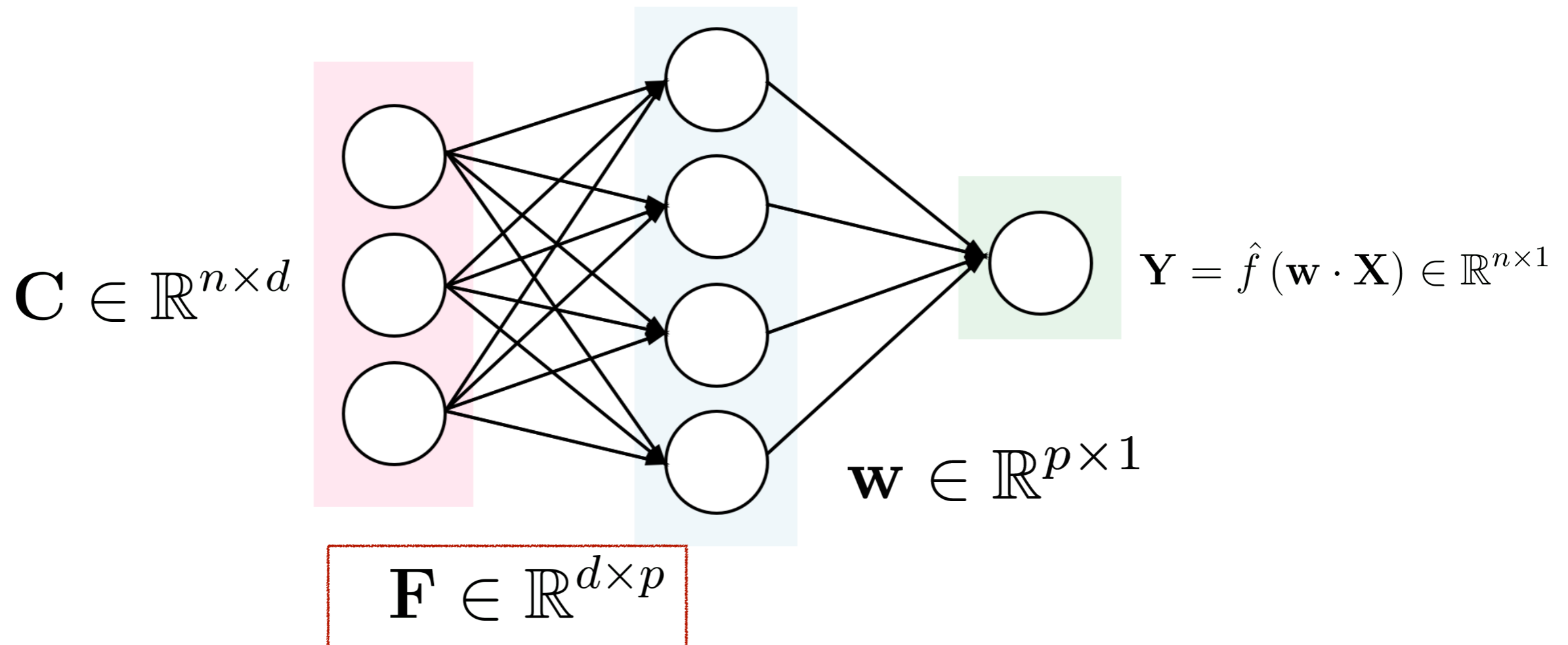
Pragya Sur and Emmanuel J. Candès



[Sur & Candès, '18]

# Next steps

$$\mathbf{X} = \Phi_{\mathbf{F}}(\mathbf{C}) = \sigma\left(\frac{\mathbf{C}\mathbf{F}}{\sqrt{d}}\right) \in \mathbb{R}^{n \times p}$$



Learning F?



Thank you for your attention!

---

Check our paper @  
[arXiv: 2002.09339 \[mat.ST\]](https://arxiv.org/abs/2002.09339)

contact: [brloureiro@gmail.com](mailto:brloureiro@gmail.com)

# References in this talk

F. Gerace, B. Loureiro, F. Krzakala, M. Mézard, L. Zdeborobá, *“Generalisation error in learning with random features and the hidden manifold model”*, arXiv: 2002.09339

M. Geiger, S. Spigler, S. d’Ascoli, L. Sagun, M. Baity-Jesi, G. Biroli and M. Wyart, *“Jamming transition as a paradigm to understand the loss landscape of deep neural networks”*, Physical Review E, 100(1):012115

S. Goldt, M. Mézard, F. Krzakala, L. Zdeborobá, *“Modelling the influence of data structure on learning in neural networks: the hidden manifold model”*, arXiv: 1909.11500

A. Abbara, B. Aubin, F. Krzakala, L. Zdeborobá, *“Rademacher complexity and spin glasses: A link between the replica and statistical theories of learning”*, arXiv: 1912.02729

C. Williams , *“Computing with infinite networks”*, NIPS 98’

A. Rahimi, B. Recht, *“Random Features for Large-Scale Kernel Machines”*, NIPS 07’

S. Mei, A. Montanari, *“The generalization error of random features regression: Precise asymptotics and double descent curve”*, arXiv: 1908.05355

K. Choromanski, M. Rowland, A. Weller, *“The Unreasonable Effectiveness of Structured Random Orthogonal Embeddings”*, NIPS 07’

P. Sur and E.J. Candès, *“A modern maximum-likelihood theory for high-dimensional logistic regression”*, PNAS 19’