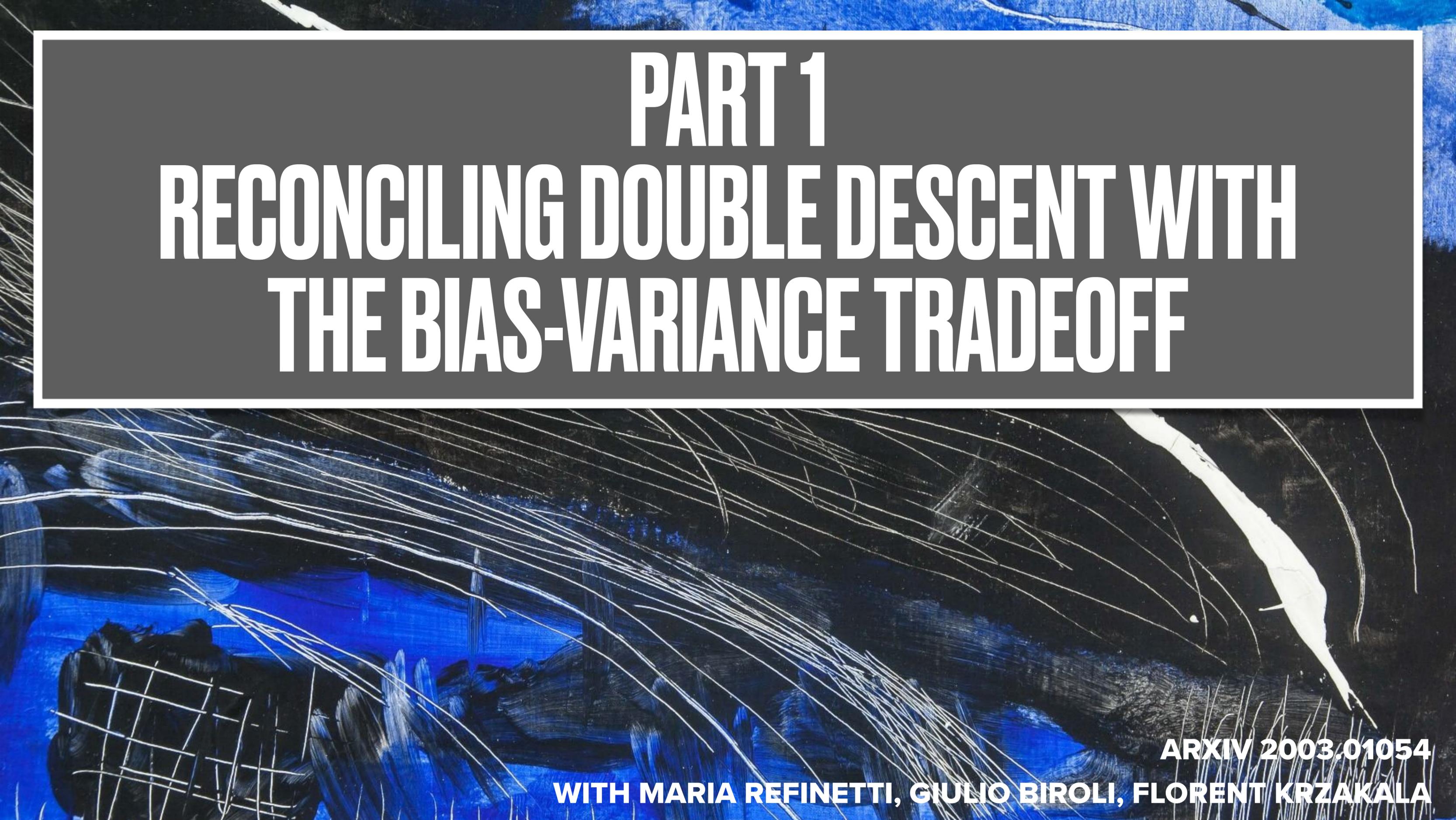


# RECONCILING DOUBLE DESCENT WITH OLDER IDEAS

**STÉPHANE D'ASCOLI**

ÉCOLE NORMALE SUPÉRIEURE & FACEBOOK AI RESEARCH, PARIS

JOINT WORK WITH MARIA REFINETTI, LEVENT SAGUN, GIULIO BIROLI, FLORENT KRZAKALA



# **PART 1**

# **RECONCILING DOUBLE DESCENT WITH THE BIAS-VARIANCE TRADEOFF**

**ARXIV 2003.01054**

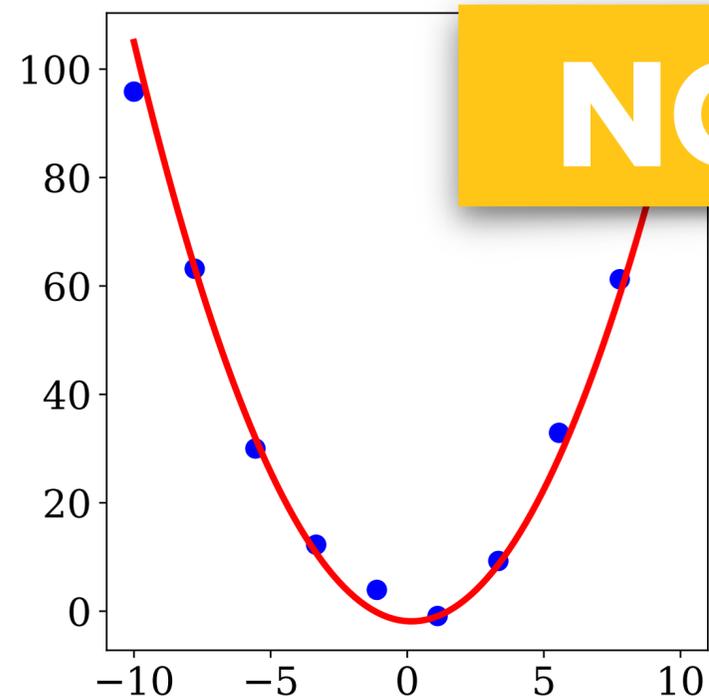
**WITH MARIA REFINETTI, GIULIO BIROLI, FLORENT KRZAKALA**

# FROM CLASSICAL THEORY

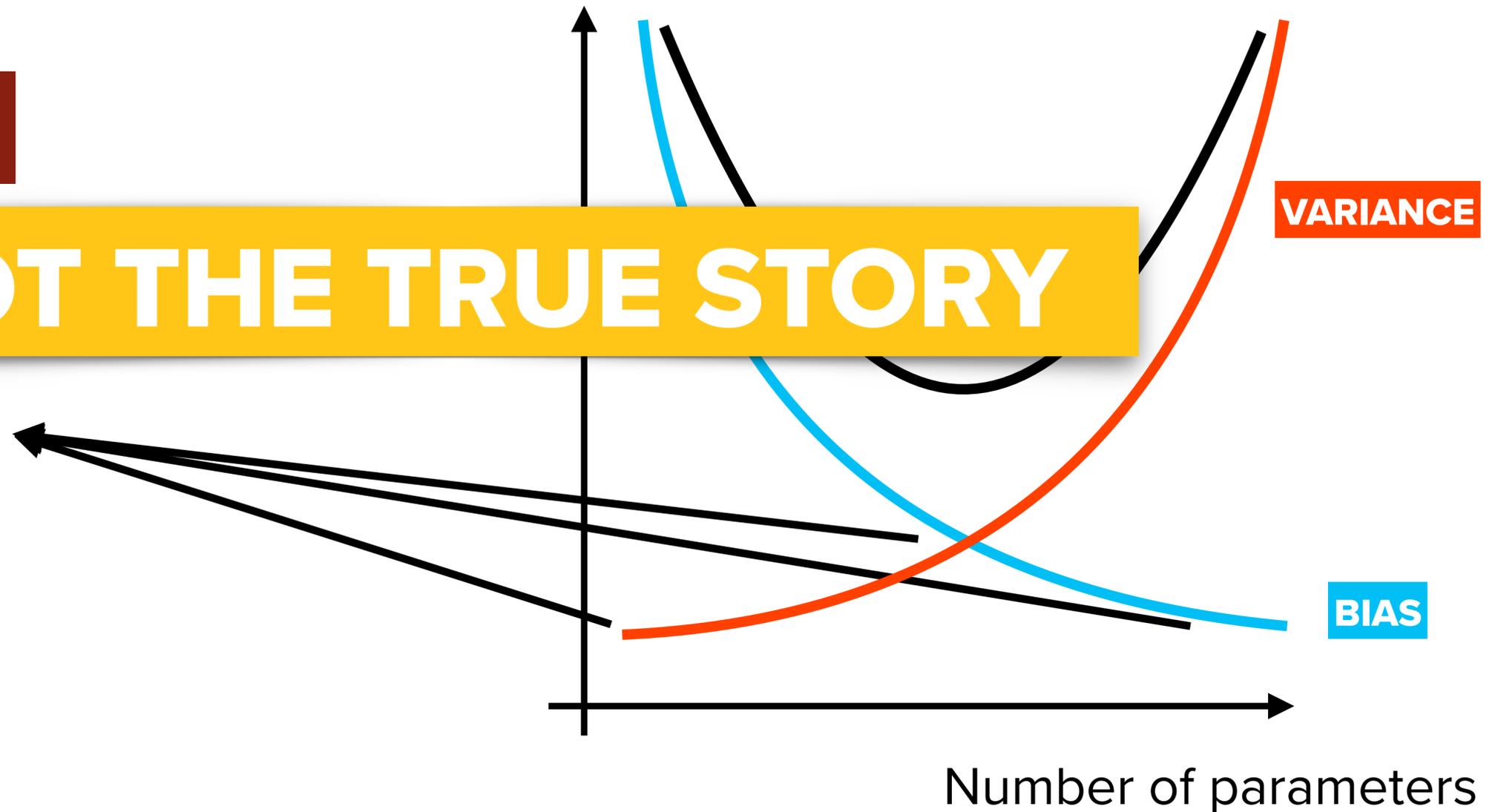
## BIAS-VARIANCE TRADEOFF

*"The price to pay for achieving low bias is high variance"*  
(Geman et al., 1992)

**TOO MANY PARAMETERS.  
BALANCED  
LOW BIAS, HIGH VARIANCE**



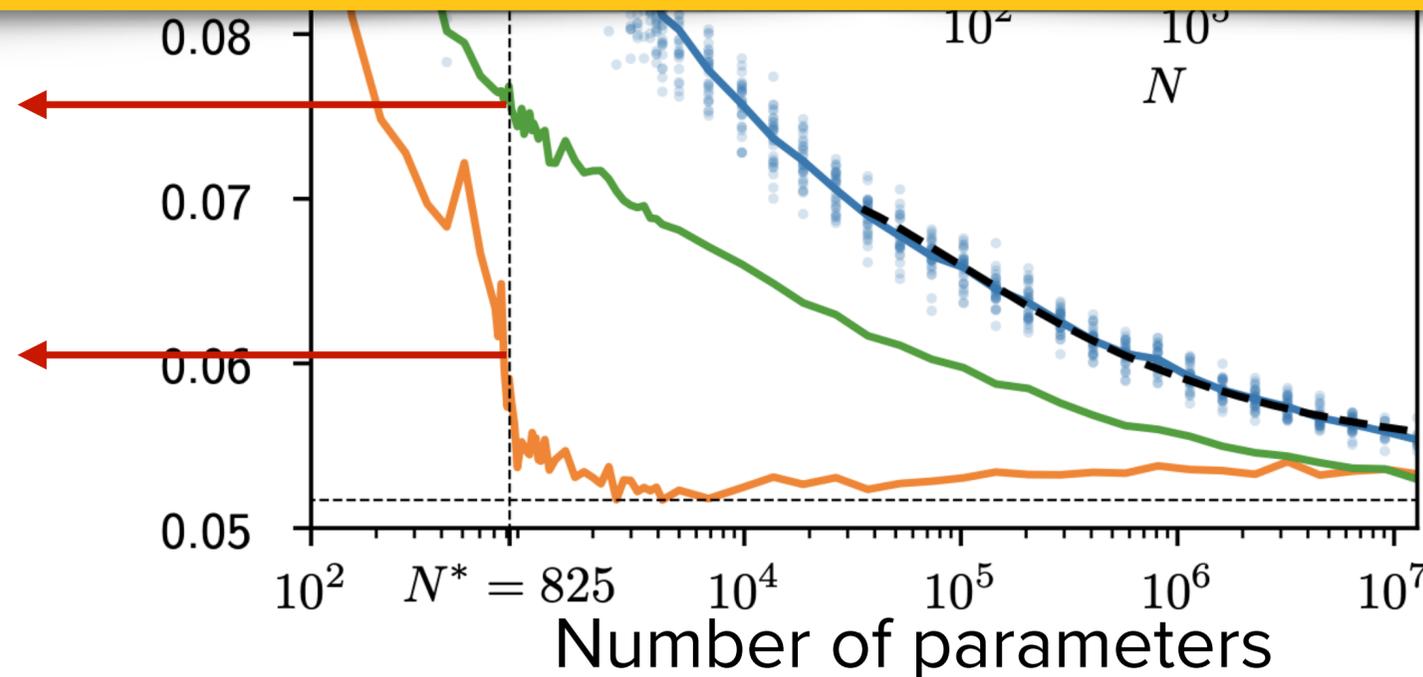
**NOT THE TRUE STORY**



# TO DOUBLE DESCENT



## WHAT IS HAPPENING TO THE BIAS AND VARIANCES ?



### See also

Belkin et al. 2019

Nakkiran et al. 2019

Hastie et al. 2019

...nari 2019

20

2020

Ba et al. 2020

....

Geiger, M., Jacot, A., Spigler, S., Gabriel, F., Sagun, L., d'Ascoli, S., Biroli, G., Hongler, C., and Wyart, M. Scaling description of generalization with number of parameters in deep learning. *arXiv preprint arXiv:1901.01608*, 2019a.

Peak in Test Error

At interpolation t

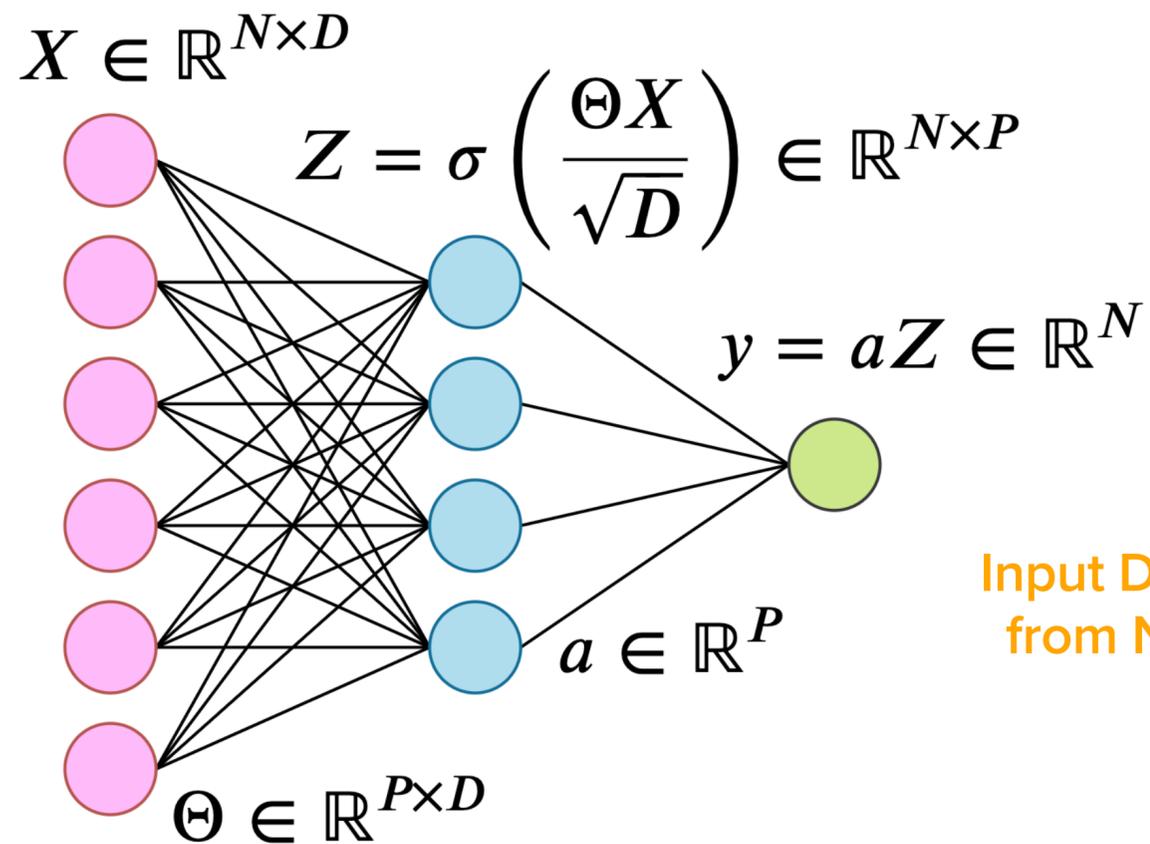
“U” shaped curve  
underparametrized

Monotonous decrease  
overparametrized re

Peak suppressed by  
**regularisation**

Peak suppressed by  
**ensembling**

# RANDOM FEATURES: WHY?



First introduced as **approximation** for **kernel methods** [1]

## ANALYTICALLY TRACTABLE MODEL

Input Dimension disentangled  
from Number of Parameters

Double descent  
curve [2]

Effect of Initialisation:  
Study Ensembling

Relevant to the lazy  
regime of neural  
networks [3,4]

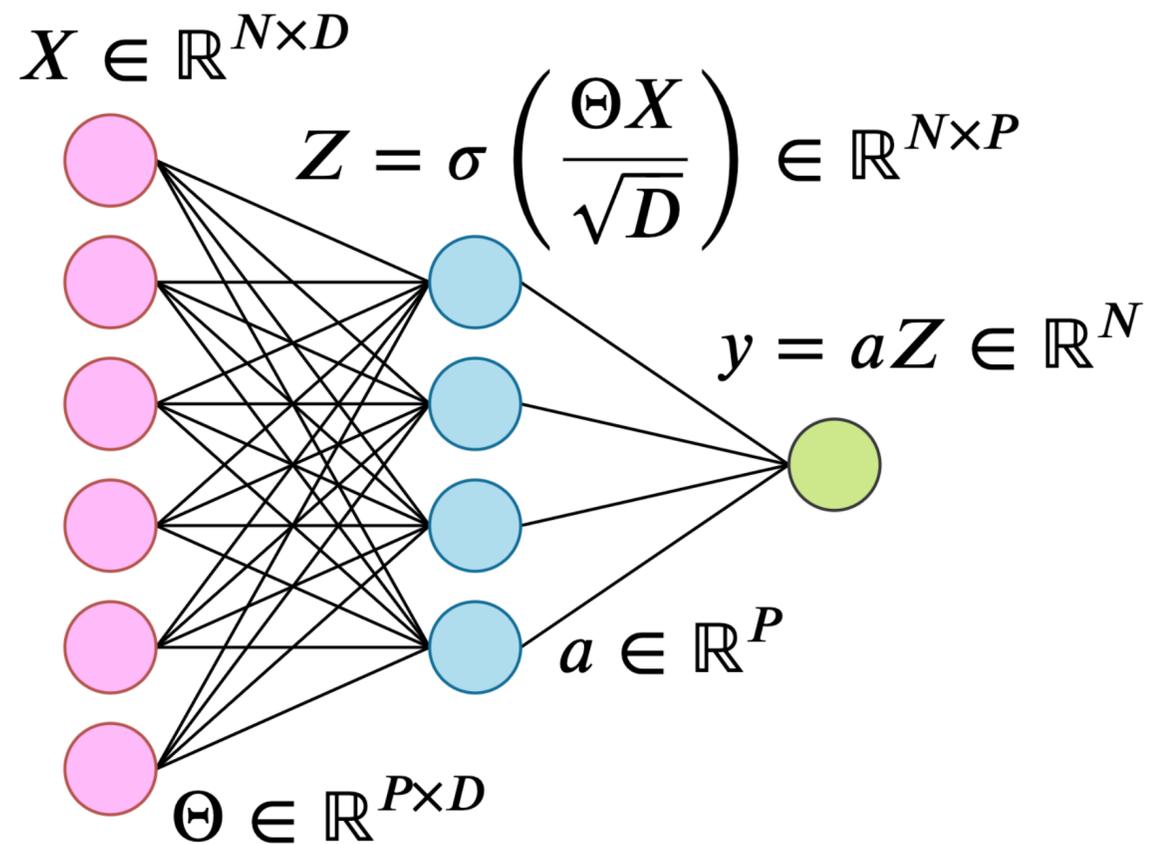
[1] Rahimi, Ali, and Benjamin Recht. "Random features for large-scale kernel machines." *Advances in neural information processing systems*. 2008.

[2] Mei, Song, and Andrea Montanari. "The generalization error of random features regression: Precise asymptotics and double descent curve." arXiv preprint arXiv:1908.05355 (2019).

[3] Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural tangent kernel: Convergence and generalization in neural networks." In *Advances in neural information processing systems*, pages 8571–8580, 2018.

[4] Chizat, Lenaic, and Francis Bach. "A note on lazy training in supervised differentiable programming." arXiv preprint arXiv:1812.07956 8 (2018).

# THE RANDOM FEATURES MODEL



$\sigma = \text{ReLU}$

$X, \Theta, \beta \sim \mathcal{N}(0,1)$

LEARNER: ONE HIDDEN LAYER

$$\hat{f}(X_\mu) = \sum_{i=1}^P a_i \sigma \left( \frac{\langle \Theta_i, X_\mu \rangle}{\sqrt{D}} \right)$$

GROUND TRUTH: LINEAR WITH NOISE

$$y_\mu = f(X_\mu) = \langle \beta, X_\mu \rangle + \epsilon_\mu$$

$$\|\beta\| = F, \quad \epsilon_\mu \sim \mathcal{N}(0, \tau)$$

$$\text{SNR} = F/\tau$$

TRAIN ERROR

$$\mathcal{L}_{\text{RF}}(\mathbf{a}) \equiv \frac{1}{N} \sum_{\mu=1}^N \left( y_\mu - \hat{f}(X_\mu) \right)^2 + \frac{P\lambda}{D} \|\mathbf{a}\|_2^2$$

$$\hat{\mathbf{a}} \equiv \arg \min_{\mathbf{a} \in \mathbb{R}^P} \mathcal{L}_{\text{RF}}(\mathbf{a})$$

TEST ERROR

$$\mathcal{R}_{\text{RF}} = \mathbb{E}_x \left[ \left( f(x) - \hat{f}(x) \right)^2 \right]$$

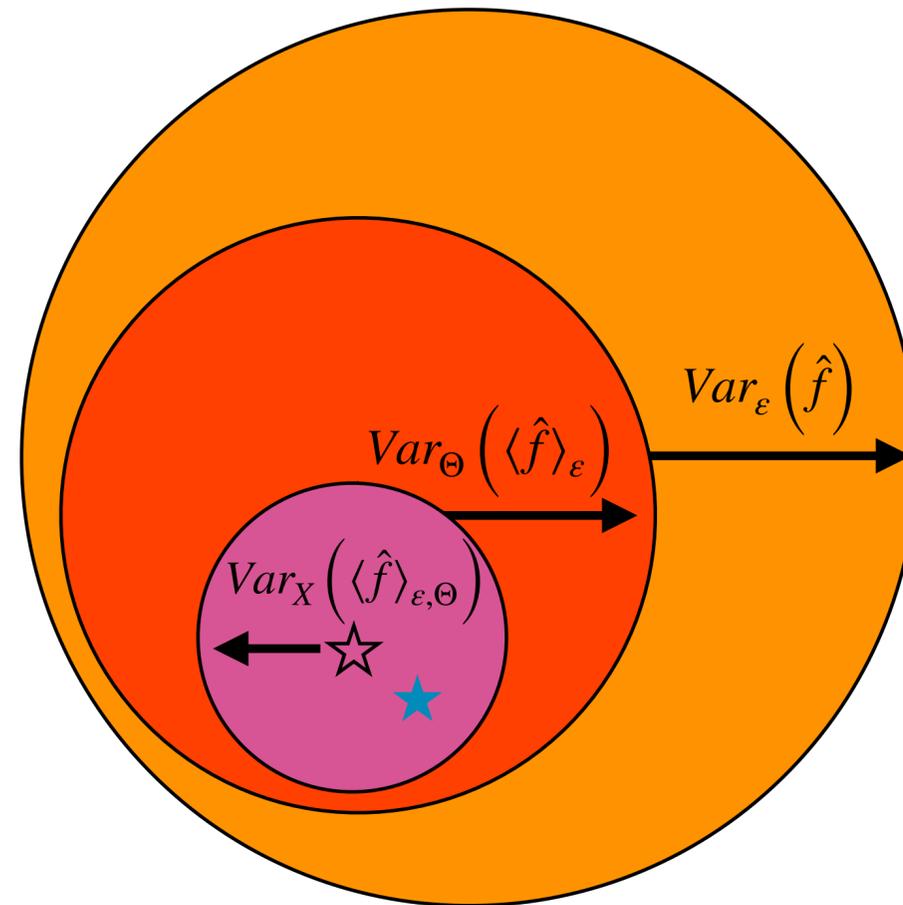
# BIAS AND VARIANCE(S)

$$\begin{aligned}
 \mathcal{R}_{\text{RF}} &= \mathbb{E} \left[ (f - \hat{f})^2 \right] \\
 &= F^2 + 2\mathbb{E} [f\hat{f}] + \mathbb{E} [\hat{f}^2] \\
 &= F^2 + 2\mathbb{E} [f\hat{f}] + \mathbb{E} \left[ \text{Var}_\varepsilon(\hat{f}) + \langle \hat{f} \rangle_\varepsilon^2 \right] \\
 &= F^2 + 2\mathbb{E} [f\hat{f}] + \mathbb{E} \left[ \text{Var}_\varepsilon(\hat{f}) + \text{Var}_\Theta(\langle \hat{f} \rangle_\varepsilon) + \langle \hat{f} \rangle_{\varepsilon, \Theta}^2 \right] \\
 &= F^2 + 2\mathbb{E} [f\hat{f}] + \mathbb{E} \left[ \text{Var}_\varepsilon(\hat{f}) + \text{Var}_\Theta(\langle \hat{f} \rangle_\varepsilon) + \text{Var}_X(\langle \hat{f} \rangle_{\varepsilon, \Theta}) + \langle \hat{f} \rangle_{\varepsilon, \Theta, X}^2 \right]
 \end{aligned}$$

Noise

Initialization

Sampling

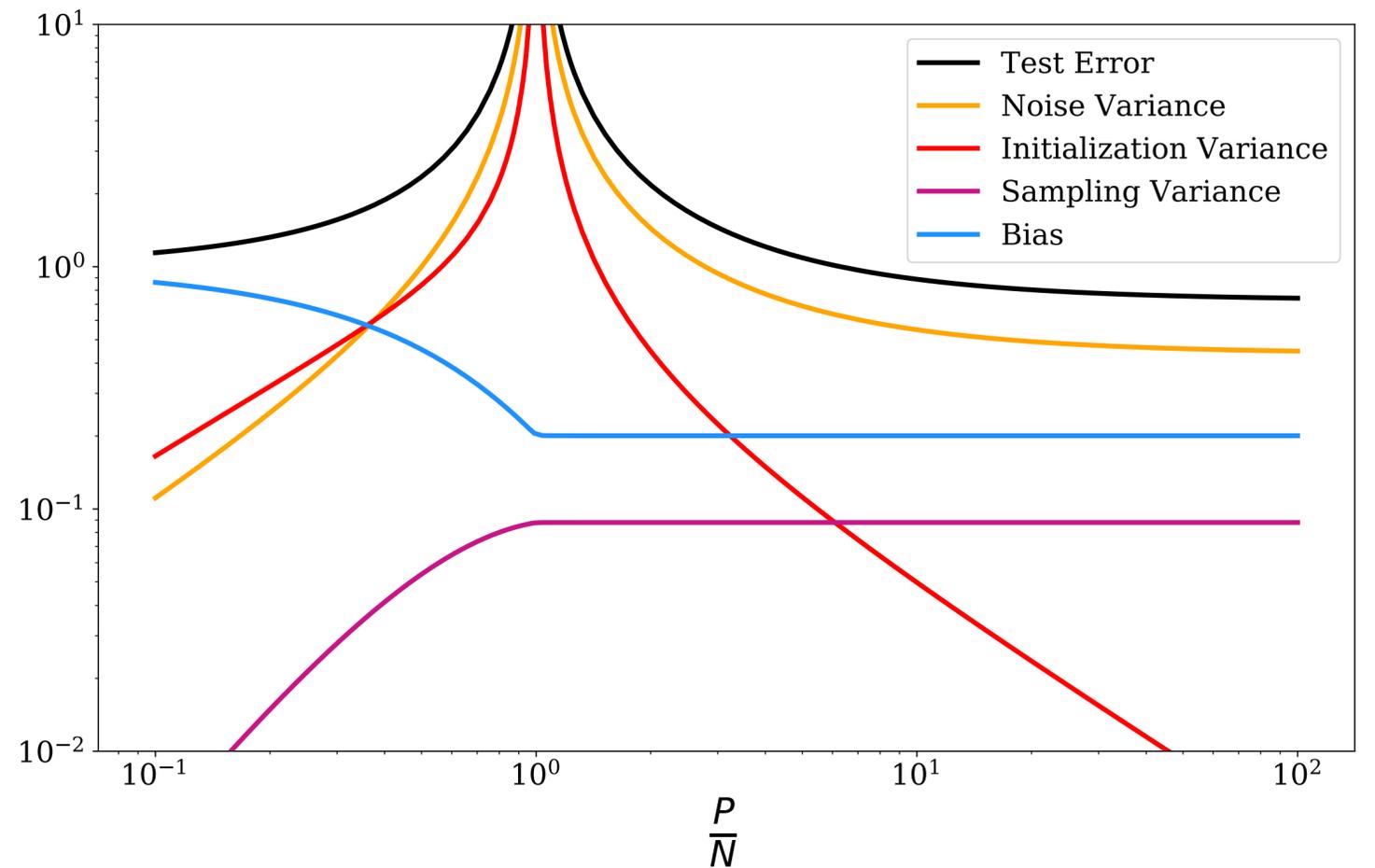
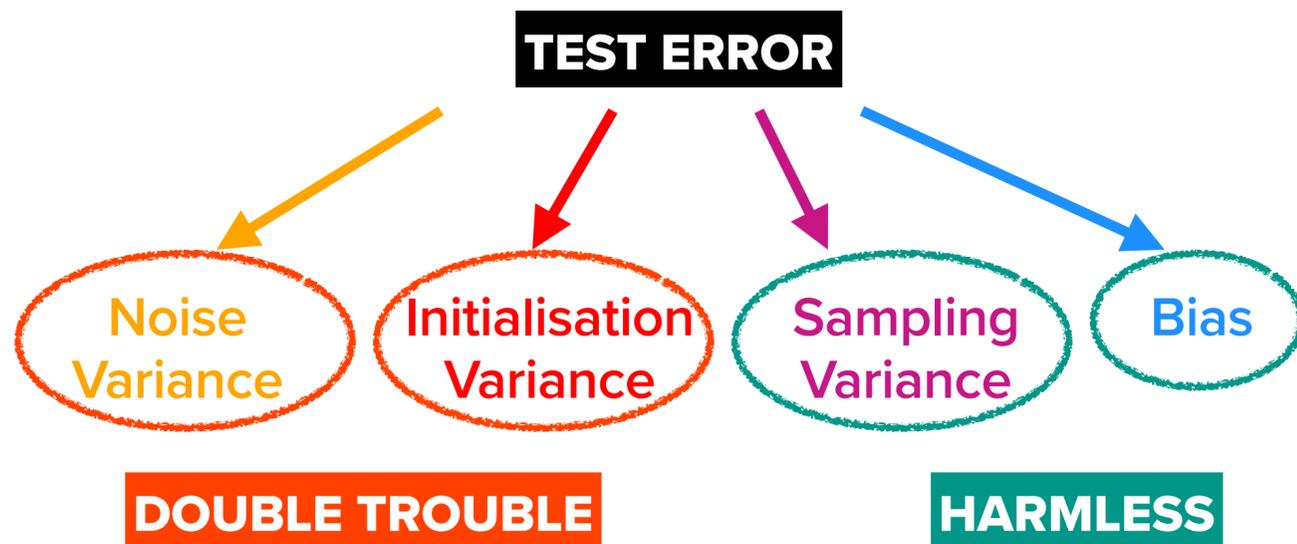


- ★  $\langle \hat{f} \rangle_\varepsilon$
- ★  $\langle \hat{f} \rangle_{\varepsilon, \Theta}$
- ★  $\langle \hat{f} \rangle_{\varepsilon, \Theta, X}$
- ★  $f$

# BIAS AND VARIANCE(S)

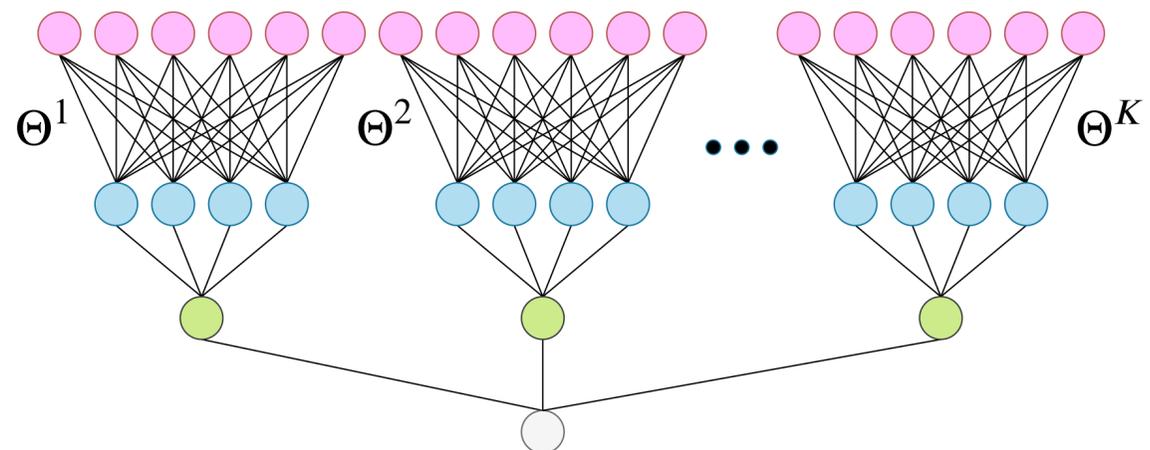
## HIGH DIMENSIONAL LIMIT

$$N, D, P \rightarrow \infty, \quad \frac{D}{P} = \mathcal{O}(1), \quad \frac{D}{N} = \mathcal{O}(1)$$



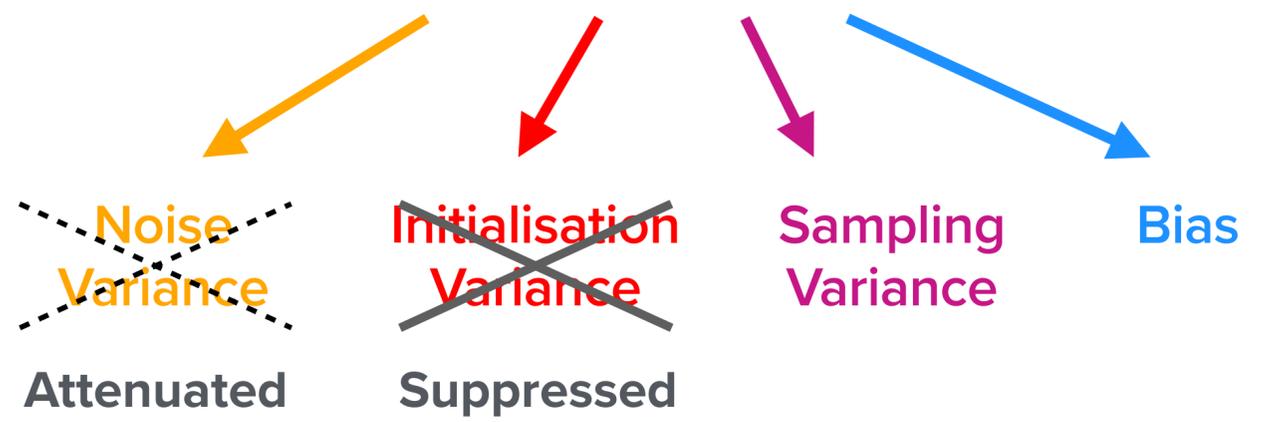
Analytical results for Random Feature Networks

# ENSEMBLING

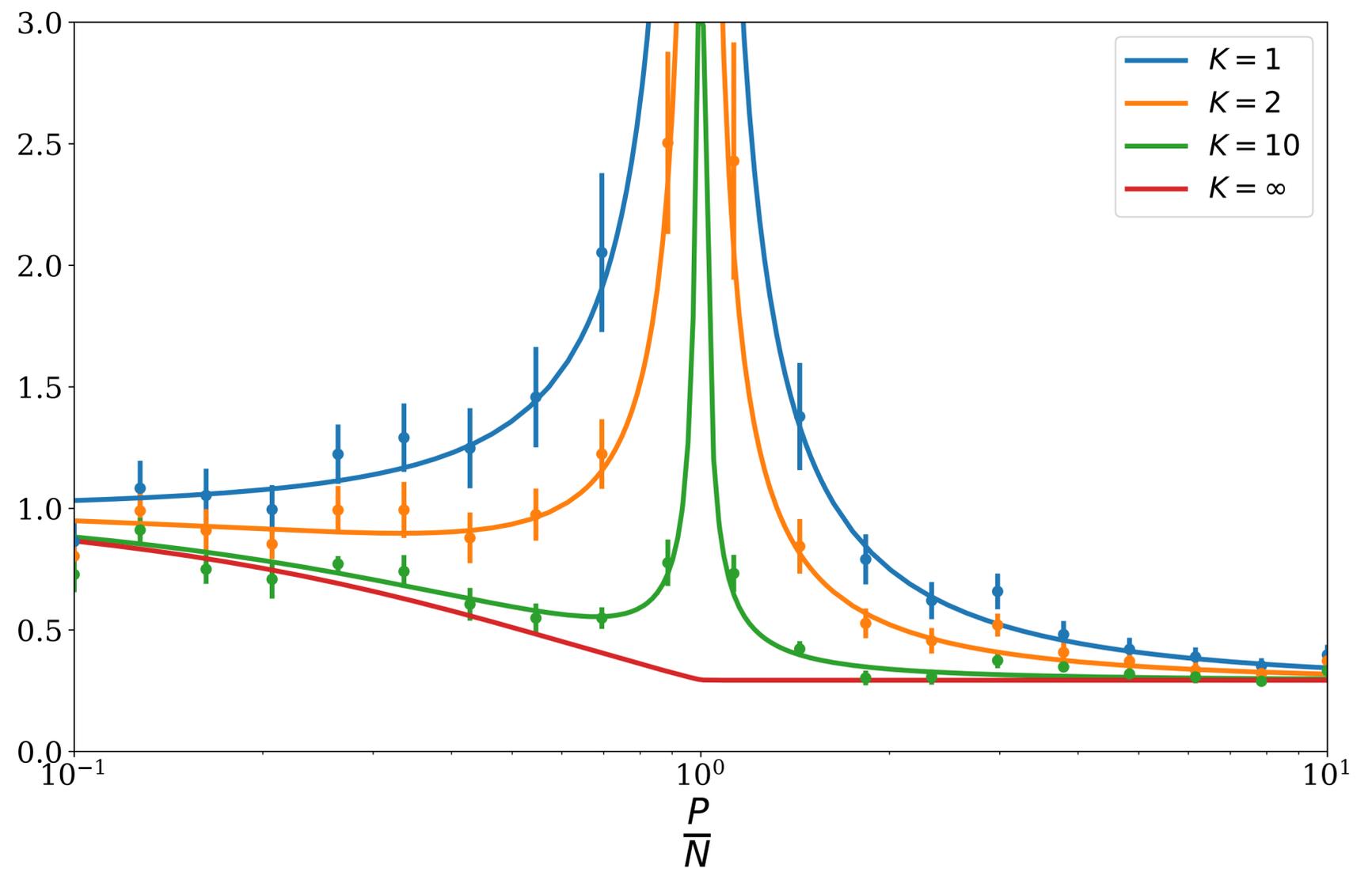


$$f(x) = \frac{1}{K} \sum_{k=1}^K f_{\Theta^k}(x)$$

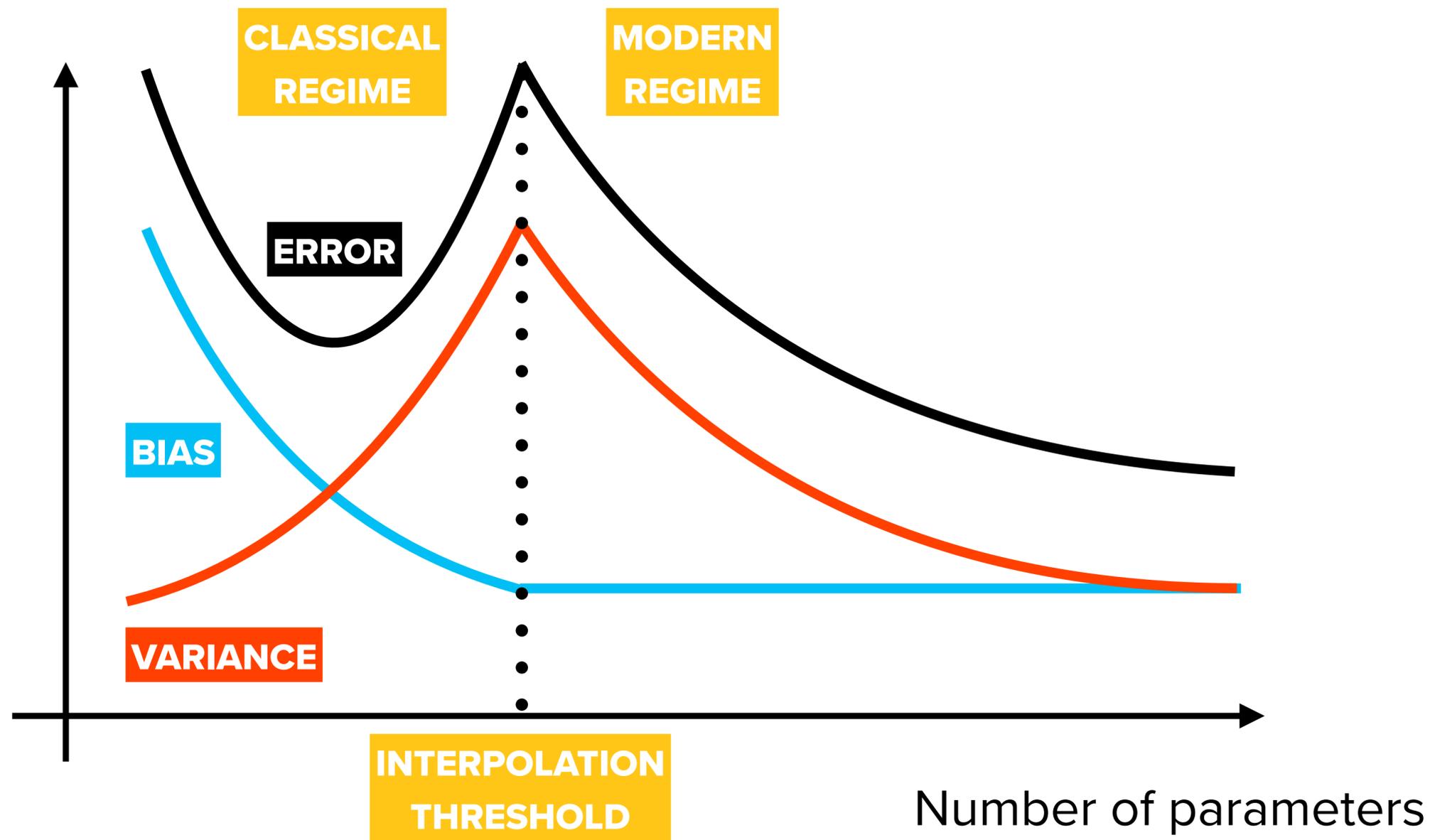
**TEST ERROR**

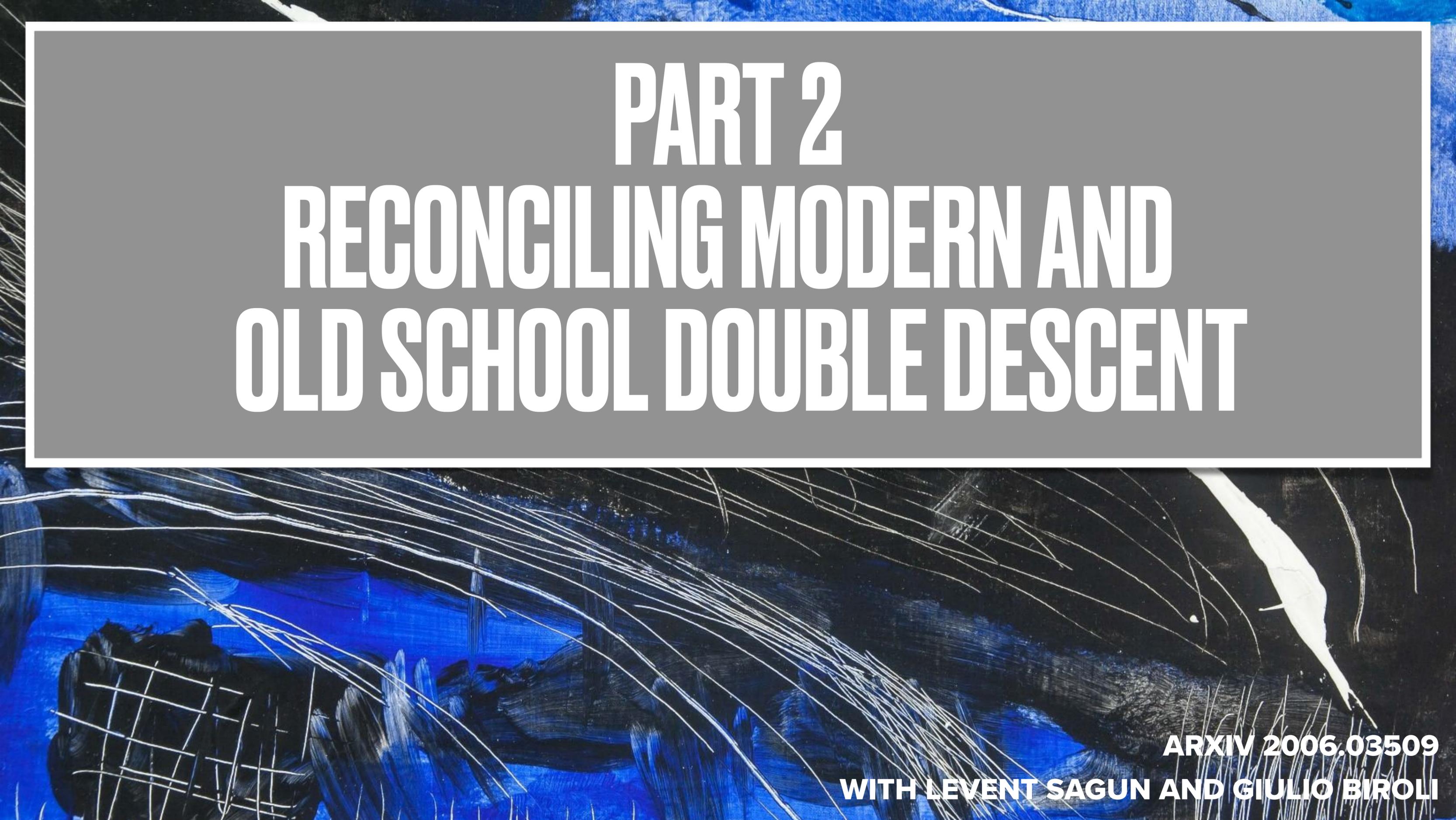


$K \rightarrow \infty$



# TAKEAWAY





# **PART 2**

# **RECONCILING MODERN AND OLD SCHOOL DOUBLE DESCENT**

**ARXIV 2006.03509**

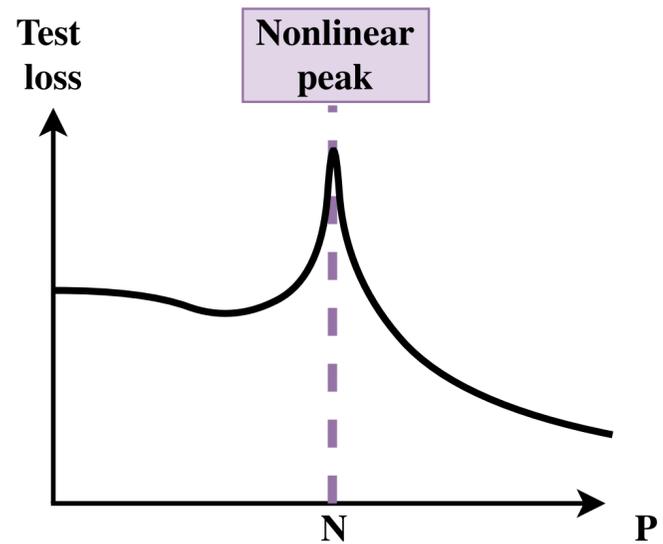
**WITH LEVENT SAGUN AND GIULIO BIROLI**

# LINEAR VS NONLINEAR MODELS

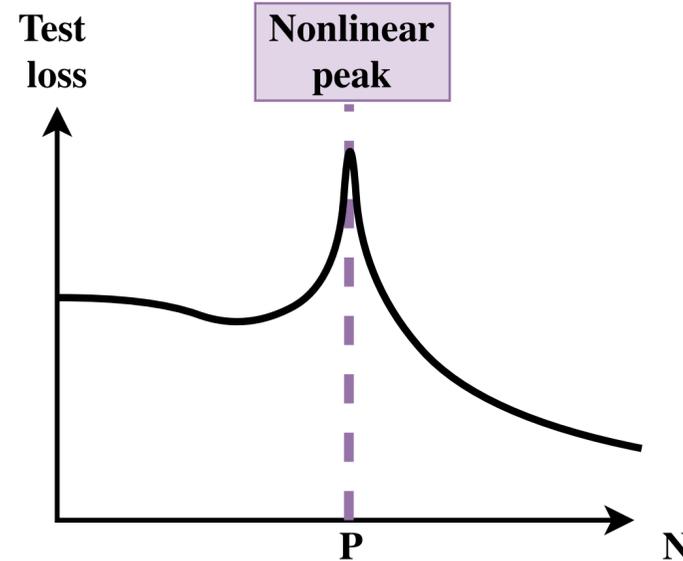
NONLINEAR NETWORKS

“MODERN” DOUBLE DESCENT [GEIGER '19]

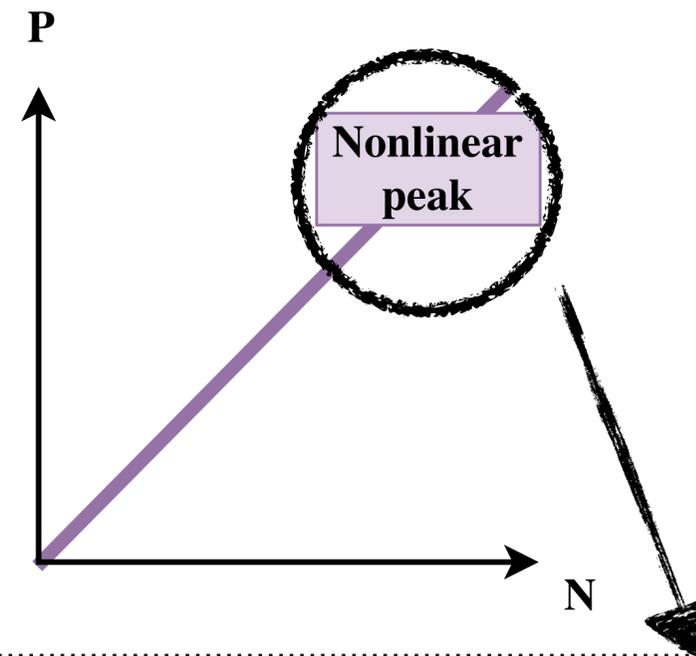
PARAMETER-WISE



SAMPLE-WISE

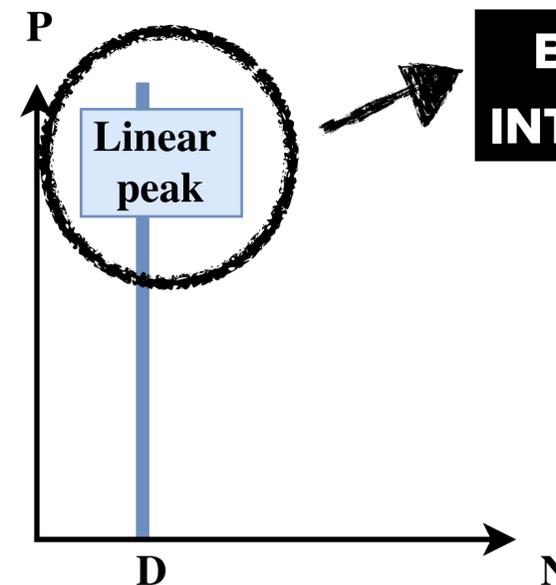
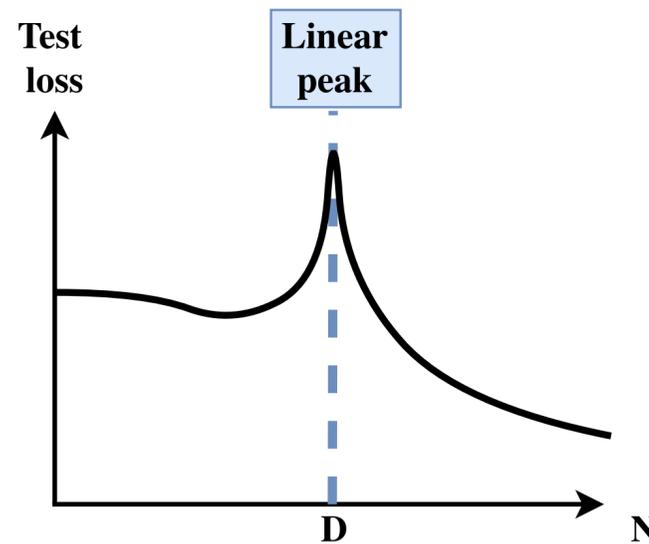
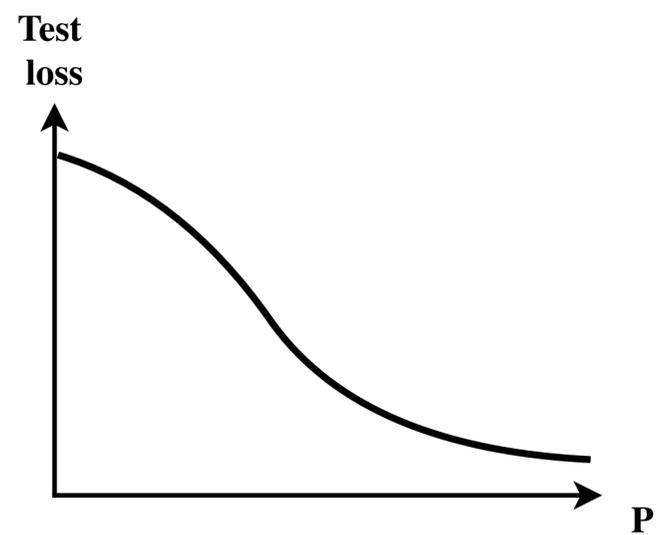


PHASE-SPACE



LINEAR NETWORKS

“OLD SCHOOL” DOUBLE DESCENT [OPPER '97]



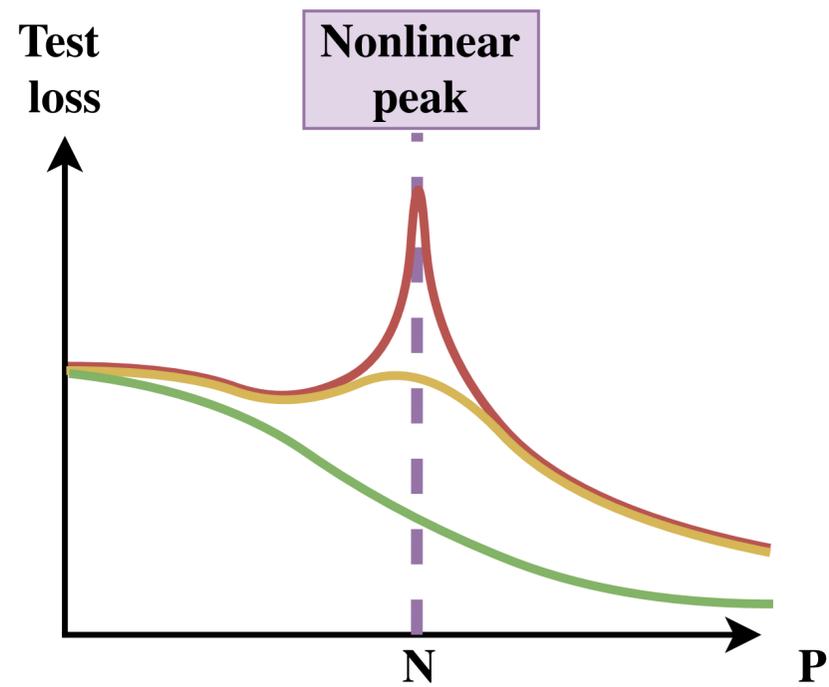
BOTH OCCUR AT INTERP THRESHOLD

ARE THEY THE SAME ?

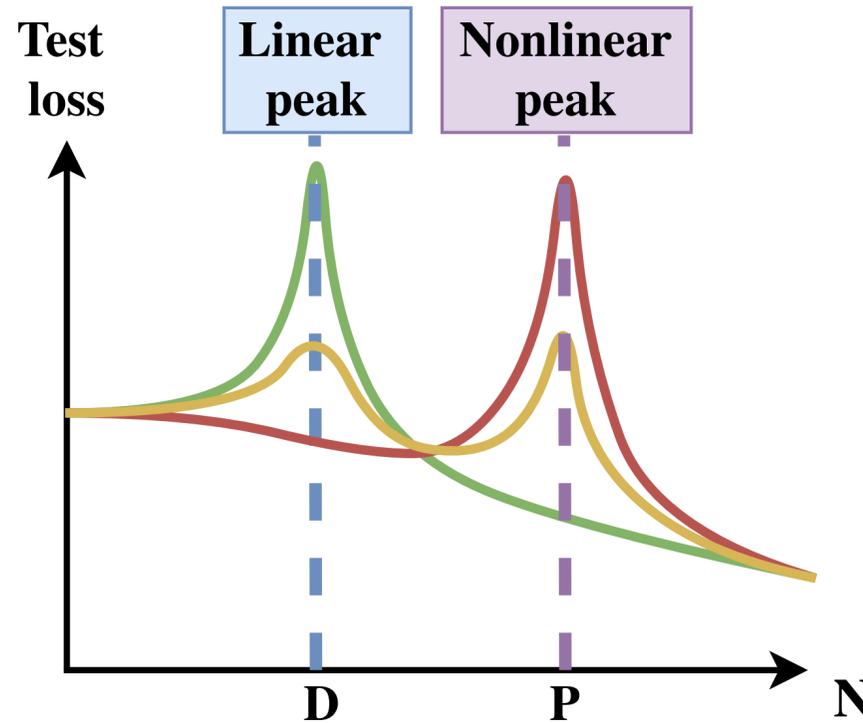
NO !

# FROM LINEAR TO LINEAR

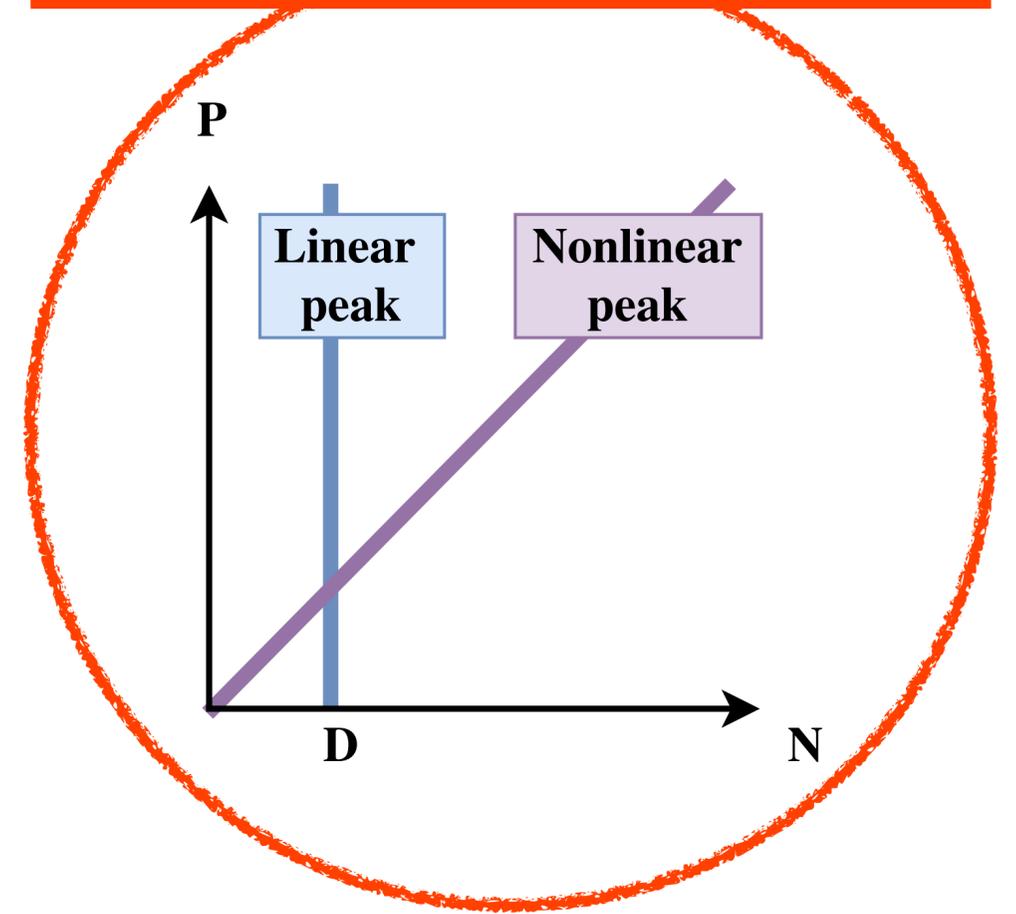
**PARAMETER-WISE**



**SAMPLE-WISE**



**BOTH PEAKS COEXIST AT HIGH NOISE**



**Activation function**

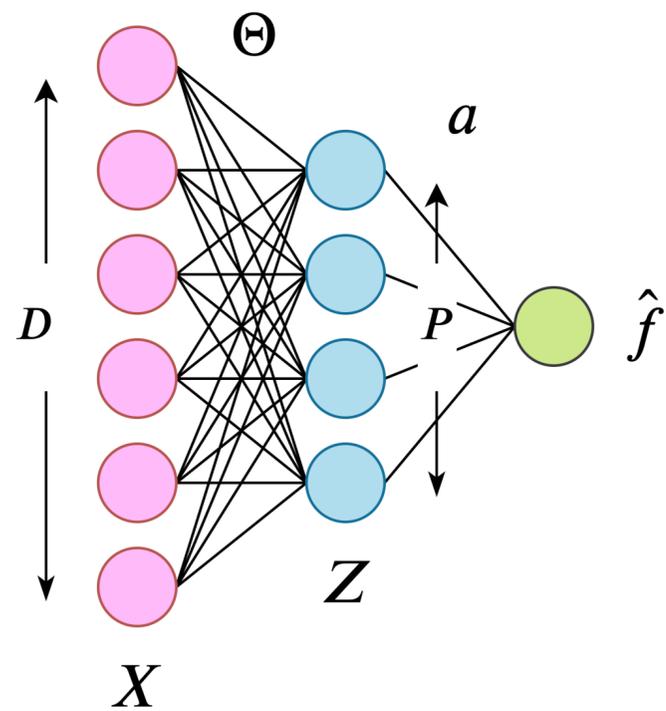
- Strongly nonlinear
- Weakly nonlinear
- Linear

**WHAT MECHANISMS UNDERLIE THESE PEAKS ?  
HOW ARE THEY DIFFERENT ?**

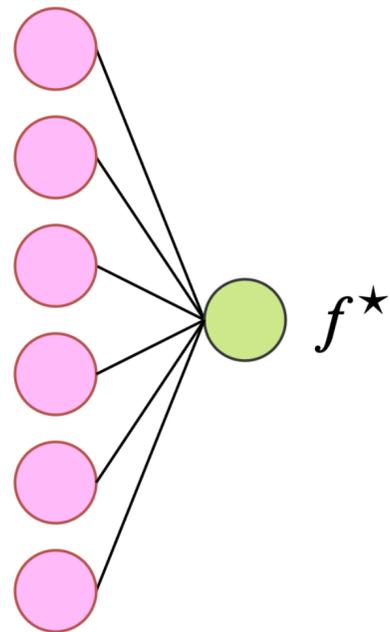
# THE TWO MODELS

RF MODEL

RF STUDENT



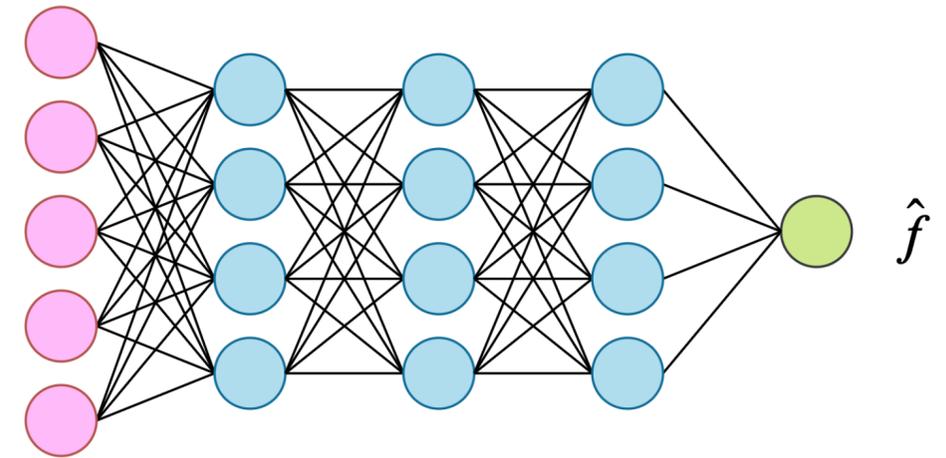
LINEAR TEACHER



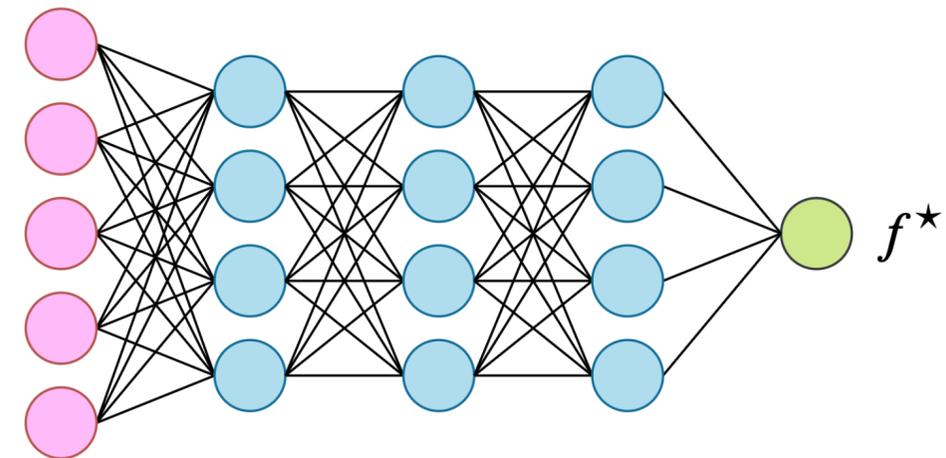
$$Z_i^\mu = \sigma \left( \frac{\langle \Theta_i, X_\mu \rangle}{\sqrt{D}} \right) \in \mathbb{R}^{N \times P}, \quad \Sigma = \frac{1}{N} Z^T Z \in \mathbb{R}^{P \times P}$$

**BAD CONDITIONING CAUSES PEAKS**

DNN MODEL

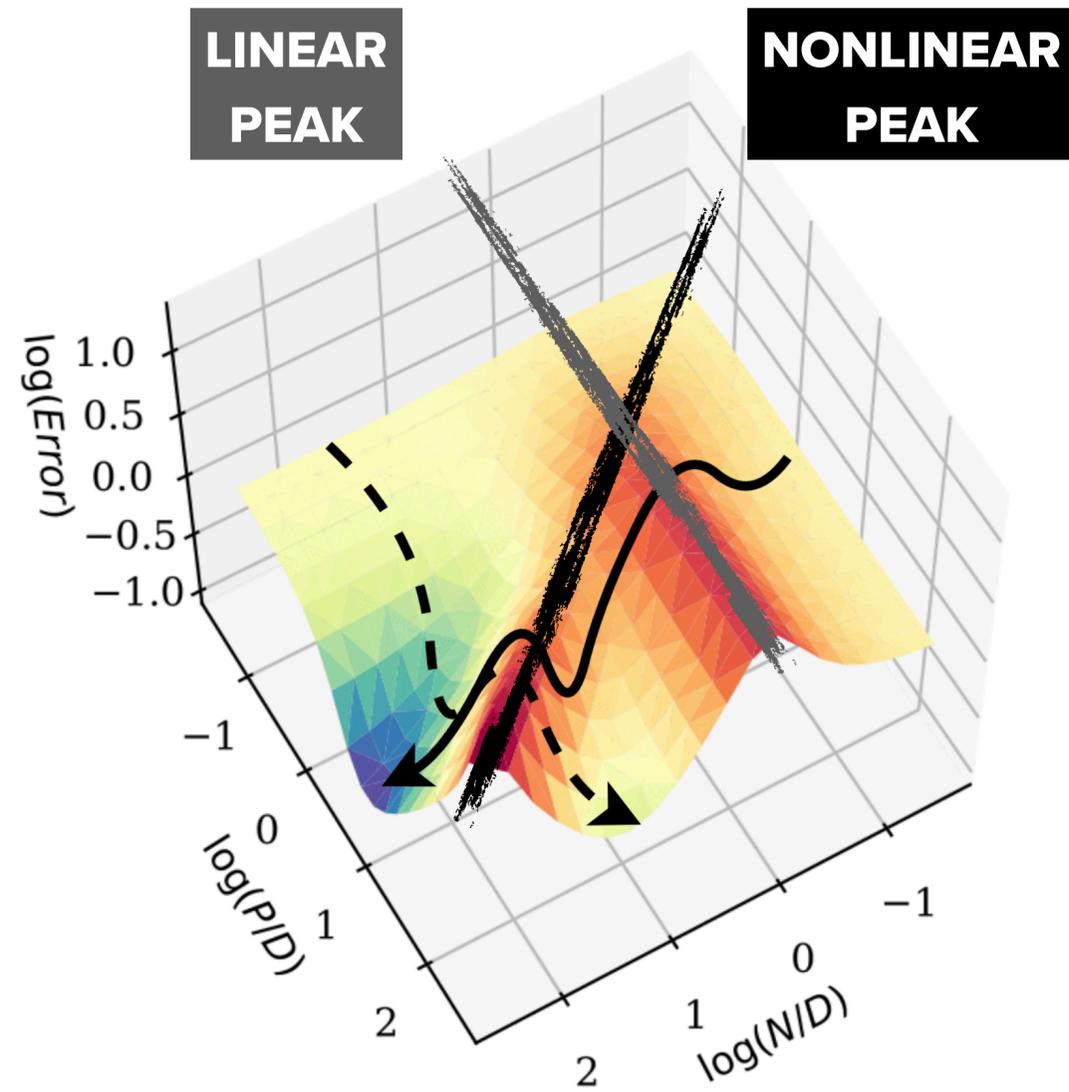


GD TRAINED STUDENT

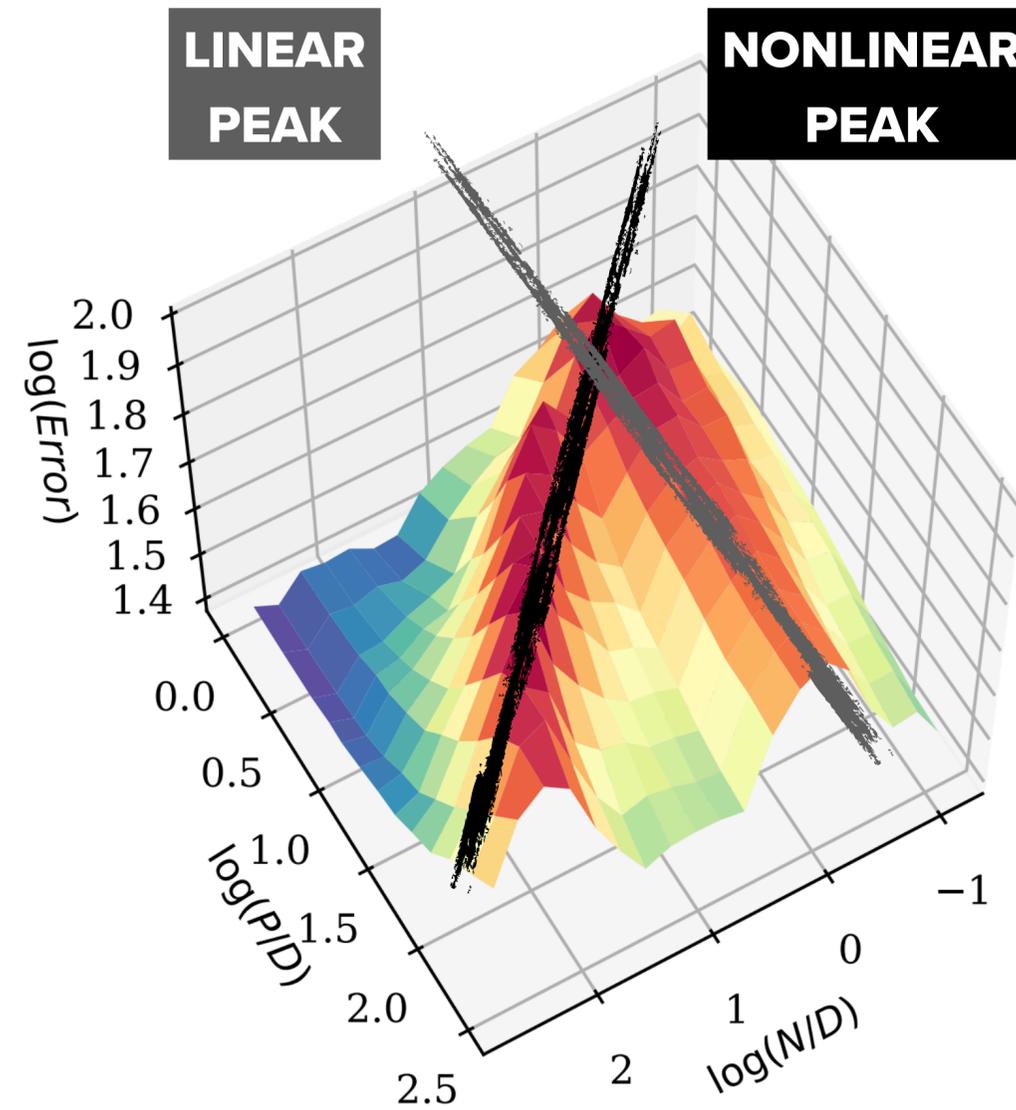


RANDOM TEACHER

# EVIDENCE OF TRIPLE DESCENT



**RF  
MODEL**



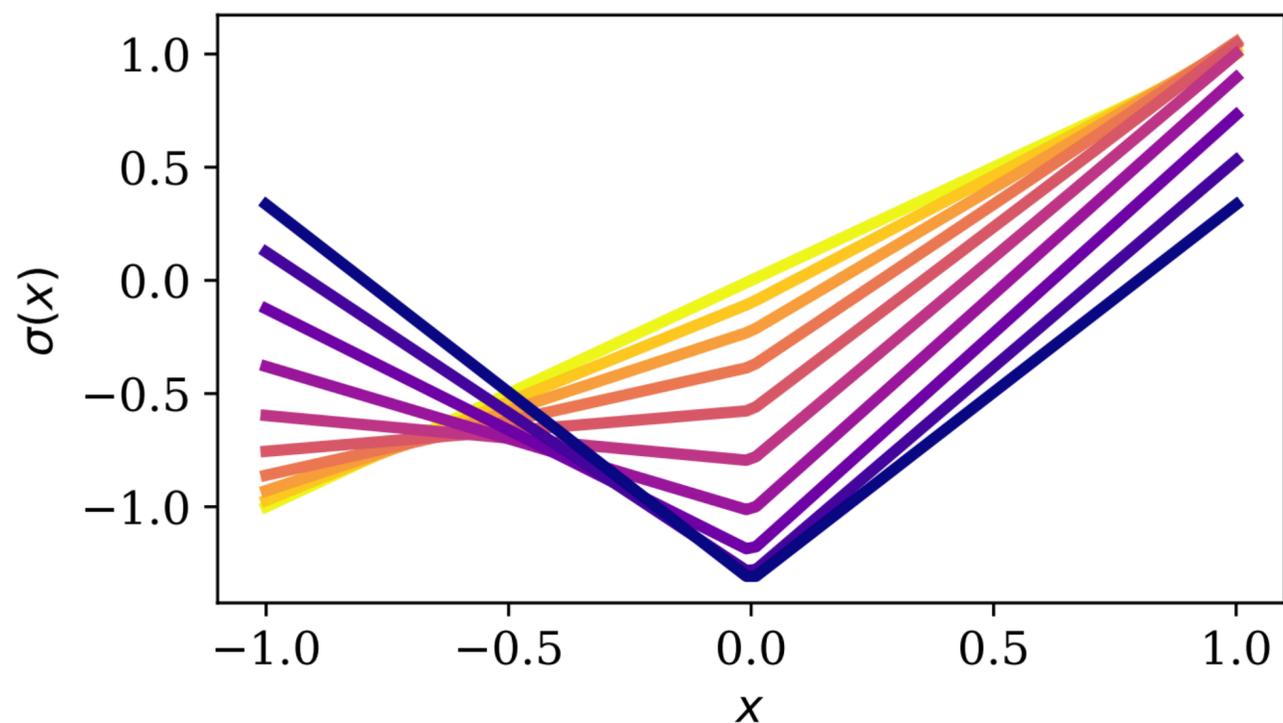
**DNN  
MODEL**

# ANALYTICAL DESCRIPTION

## HIGH-DIMENSIONAL LIMIT

$$N, D, P \rightarrow \infty, \quad \frac{D}{P} = \mathcal{O}(1), \quad \frac{D}{N} = \mathcal{O}(1)$$

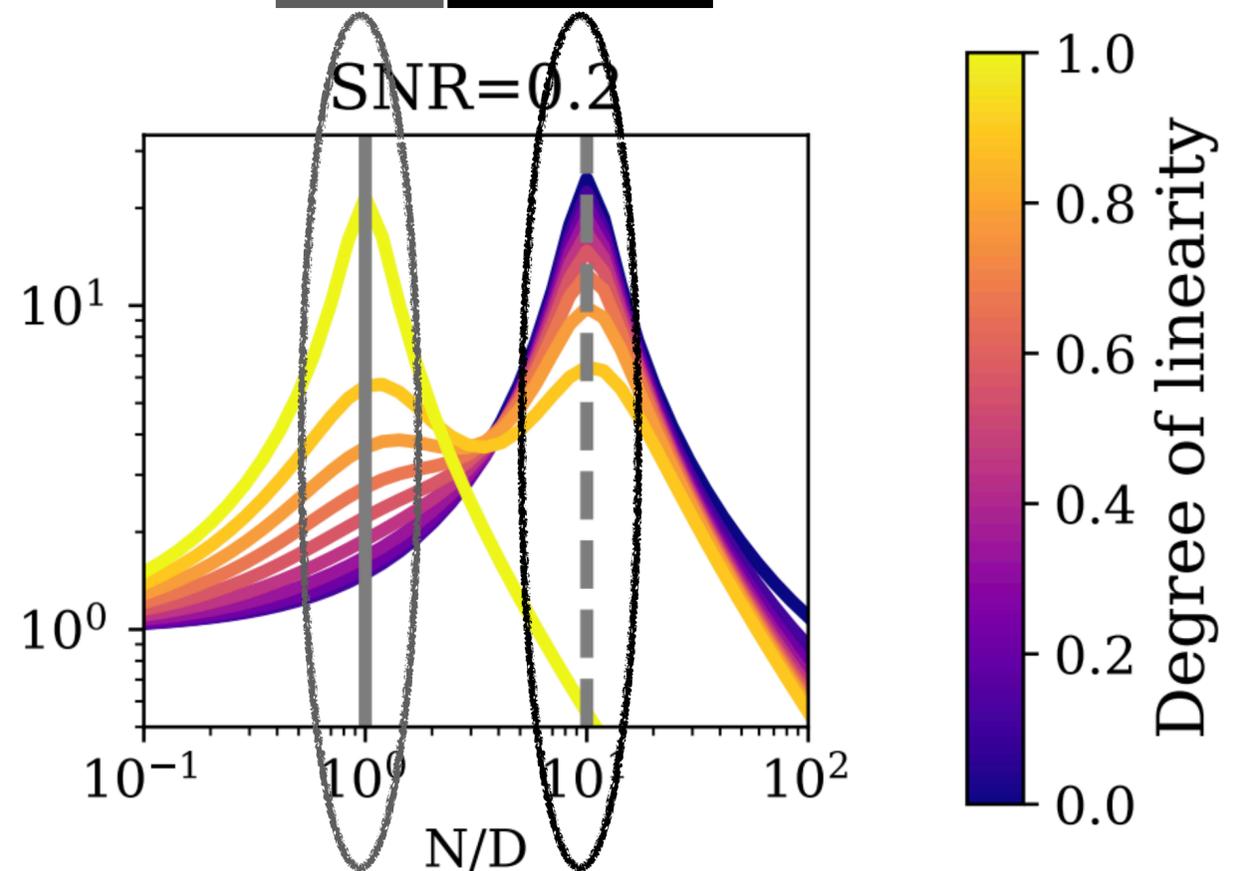
$$\eta = \int dz \frac{e^{-z^2/2}}{\sqrt{2\pi}} \sigma^2(z), \quad \zeta = \left[ \int dz \frac{e^{-z^2/2}}{\sqrt{2\pi}} \sigma'(z) \right]^2$$



## DEGREE OF LINEARITY

$$r = \frac{\zeta}{\eta}$$

**LINEAR PEAK**   **NONLINEAR PEAK**



# ANALYTICAL SPECTRUM

NONLINEAR = LINEAR + NOISE

LINEAR  
PART

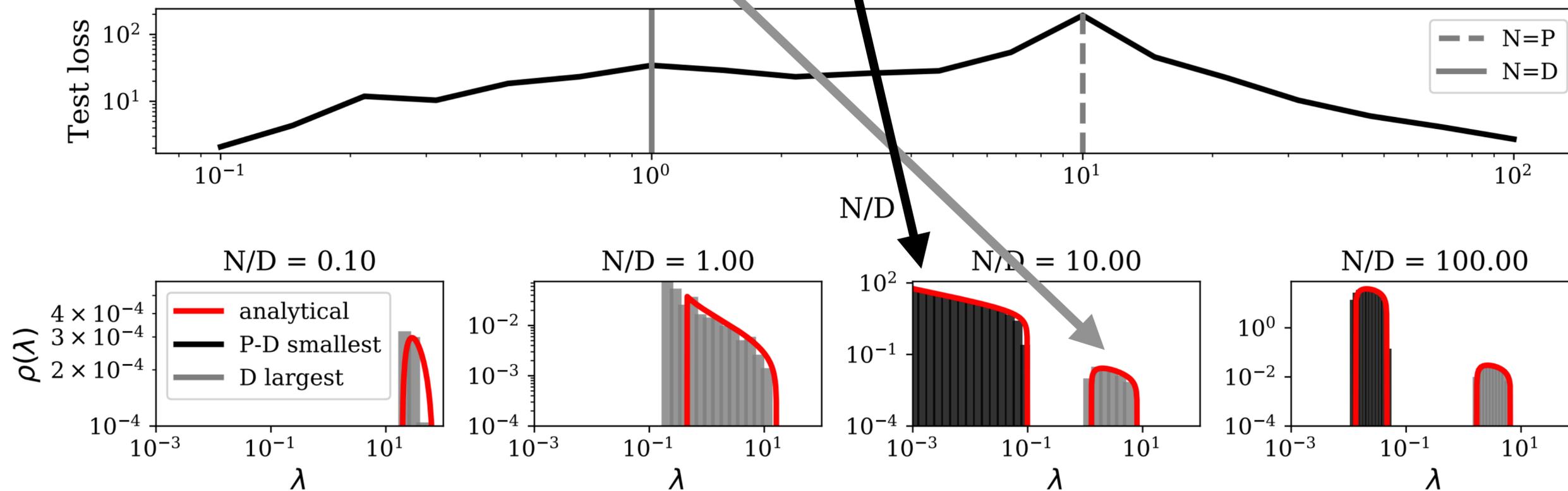
NONLINEAR  
PART

$$\rho(\lambda) = \frac{1}{\pi} \lim_{\epsilon \rightarrow 0^+} \text{Im} G(\lambda - i\epsilon), \quad G(z) = \frac{\psi}{z} A \left( \frac{1}{z\psi} \right) + \frac{1-\psi}{z}$$

$$A(t) = 1 + (\eta - \zeta)tA_\phi(t)A_\psi(t) + \frac{A_\phi(t)A_\psi(t)t\zeta}{1 - A_\phi(t)A_\psi(t)t\zeta}$$

$$Z = \sigma \left( \frac{X\Theta^\top}{\sqrt{D}} \right) \rightarrow \sqrt{\zeta} \frac{X\Theta^\top}{\sqrt{D}} + \sqrt{\eta - \zeta} W, \quad W \sim \mathcal{N}(0,1)$$

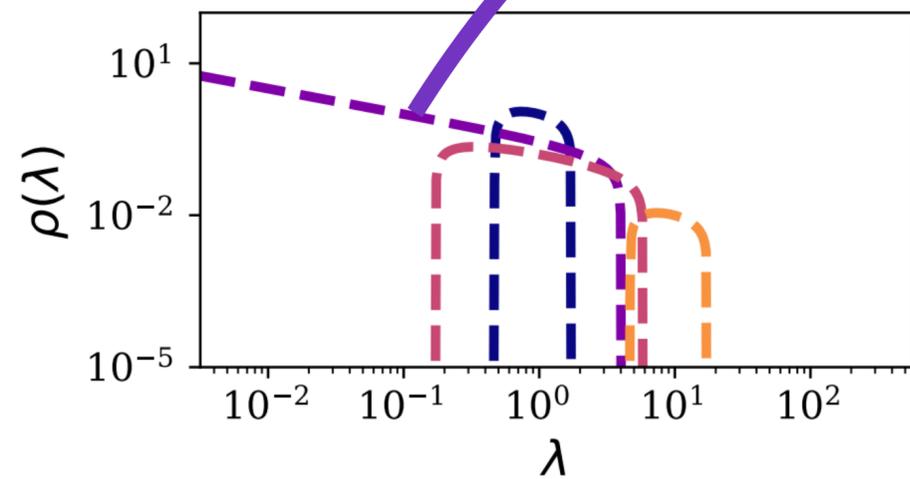
[Pennington & Worah 2017]



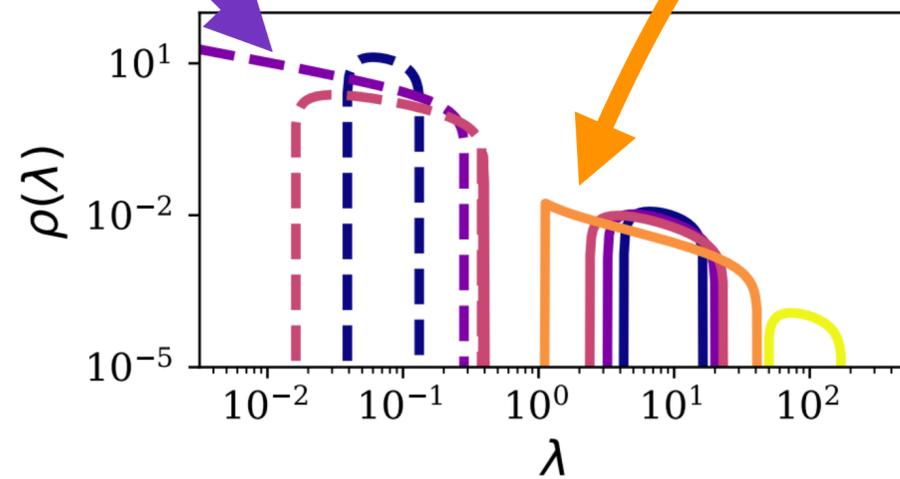
# ANALYTICAL SPECTRUM

**N=P GAP SURVIVES**

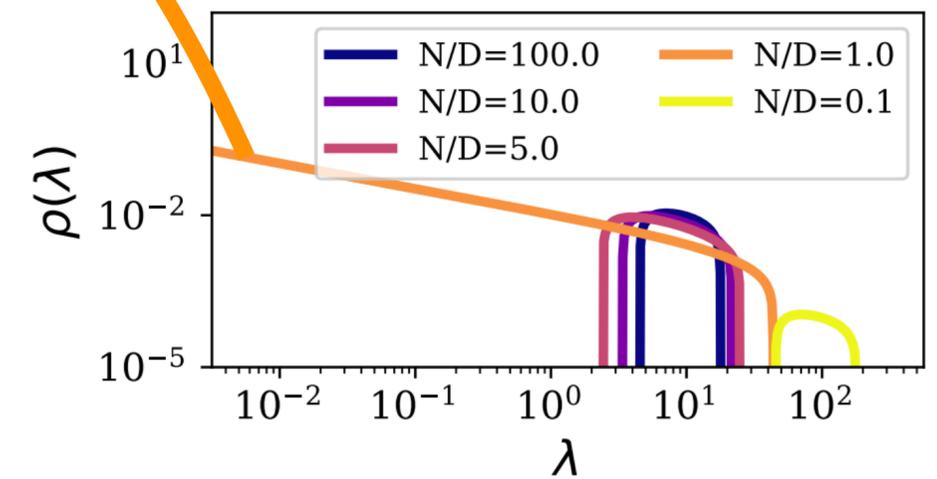
**N=D GAP IS REGULARISED**



(a) Absolute value ( $r=0$ )



(b) Tanh ( $r \simeq 0.92$ )



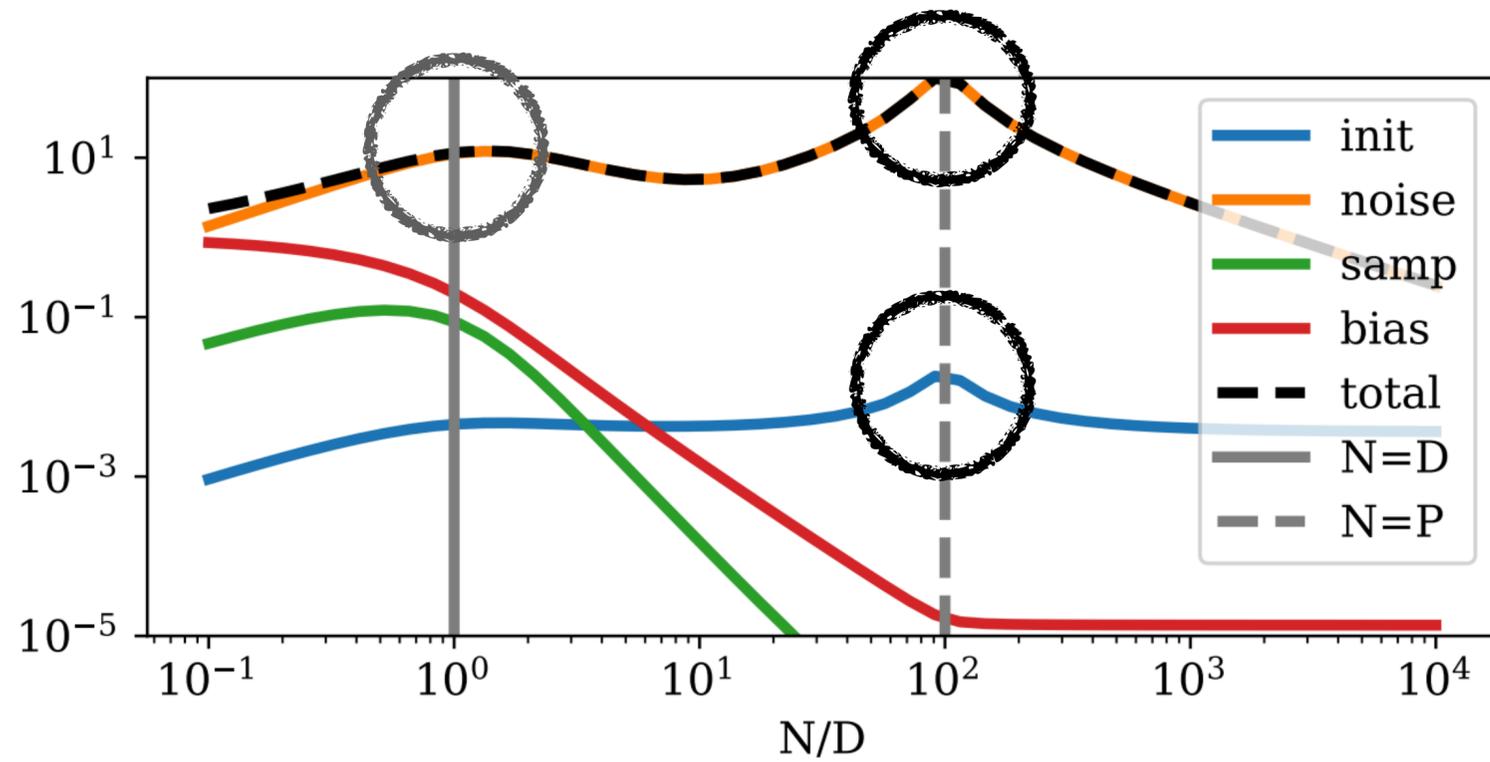
(c) Linear ( $r=1$ )

# BIAS AND VARIANCES

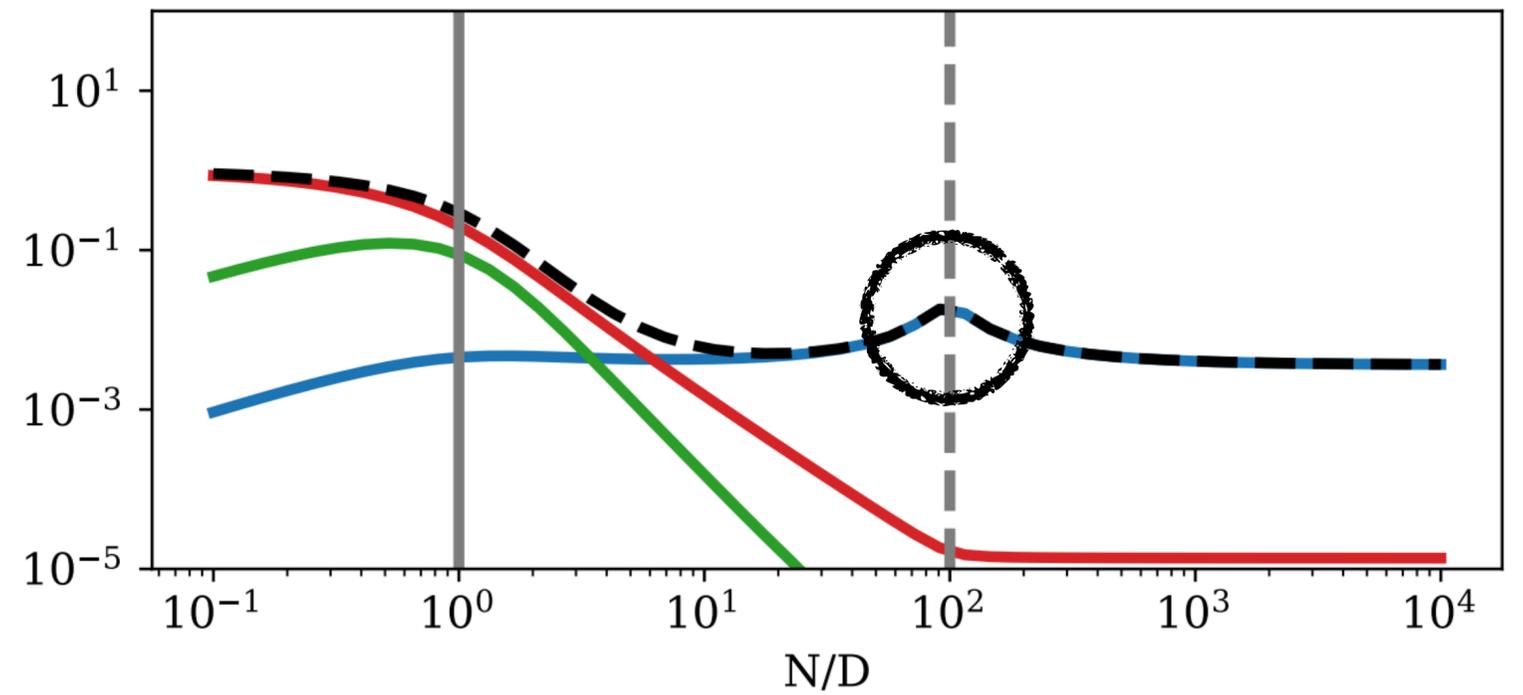
**LINEAR PEAK  
CAUSED BY NOISE**

**NONLINEAR PEAK  
CAUSED BY NOISE & INIT**

**NONLINEAR PEAK  
SURVIVES IN ABSENCE OF NOISE**



**NOISY**



**NOISELESS**

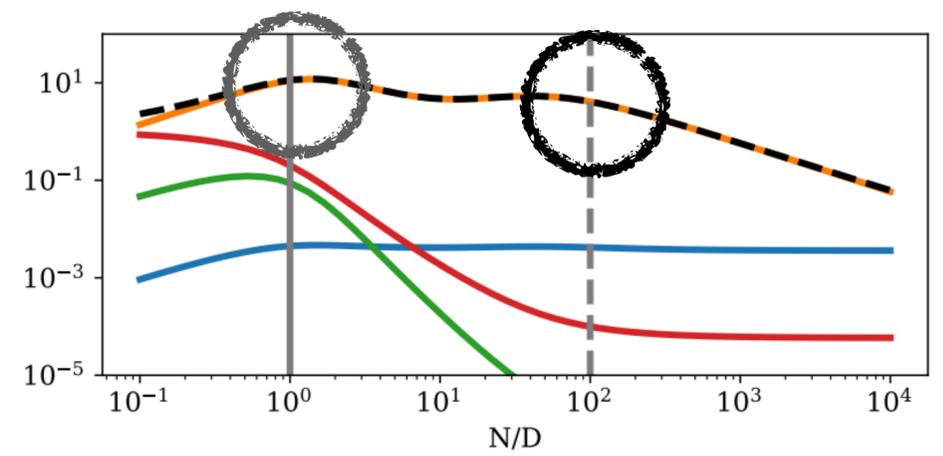
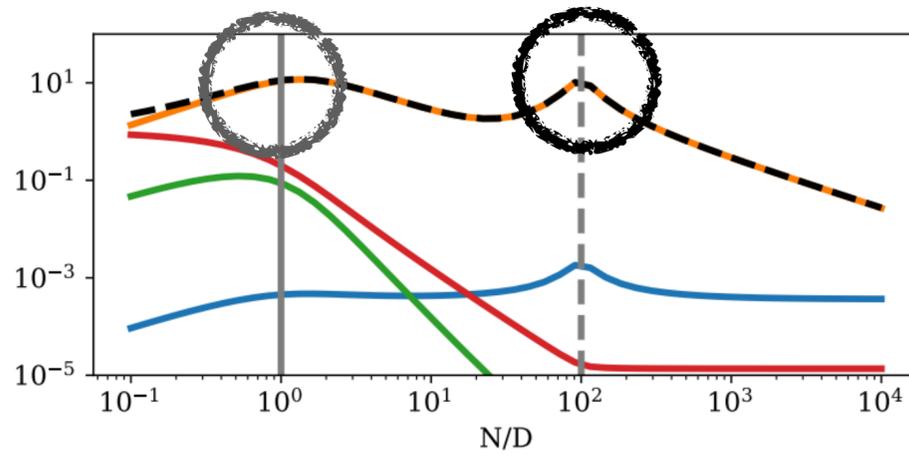
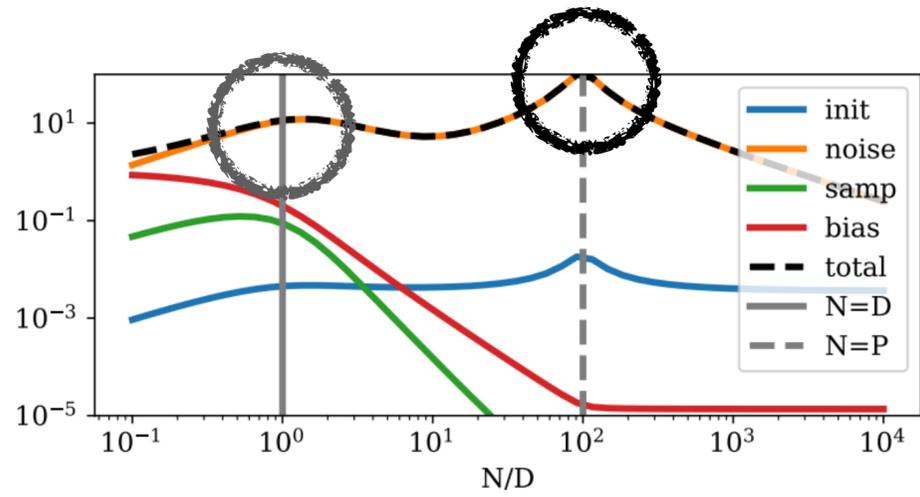
# EFFECT OF REGULARISATION

VANILLA

ENSEMBLING

REGULARIZING

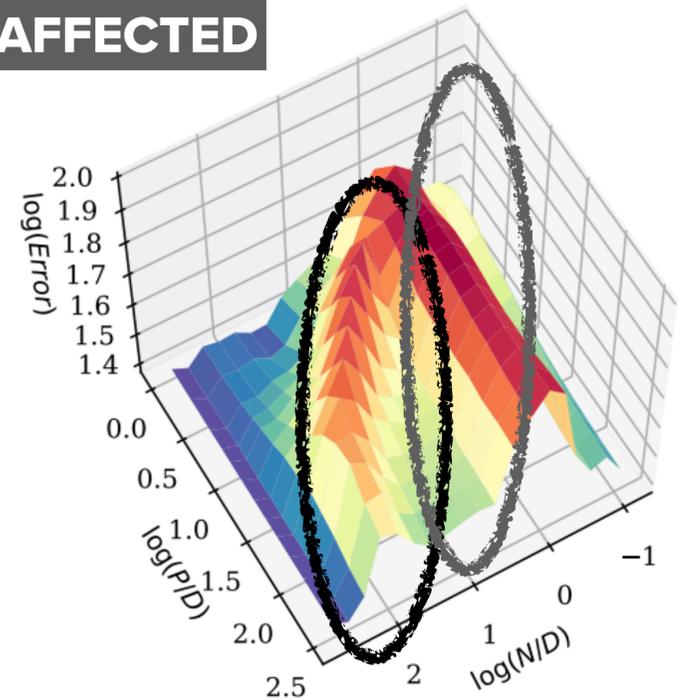
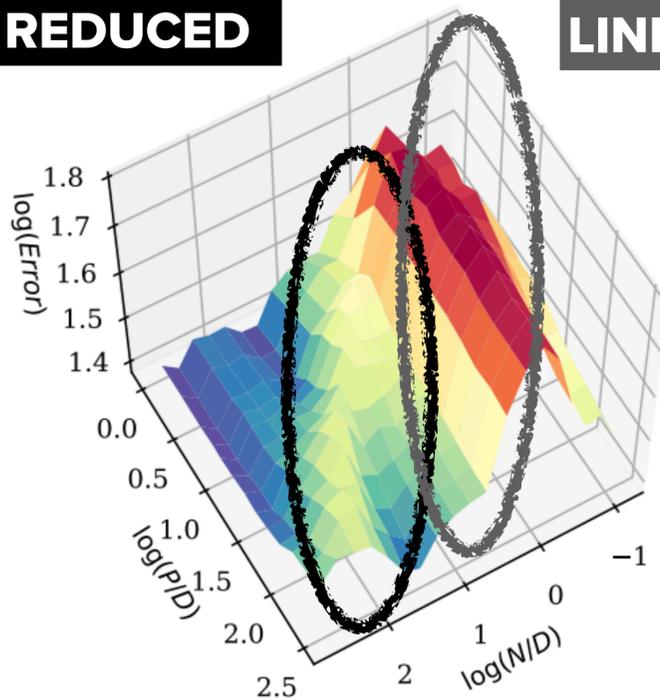
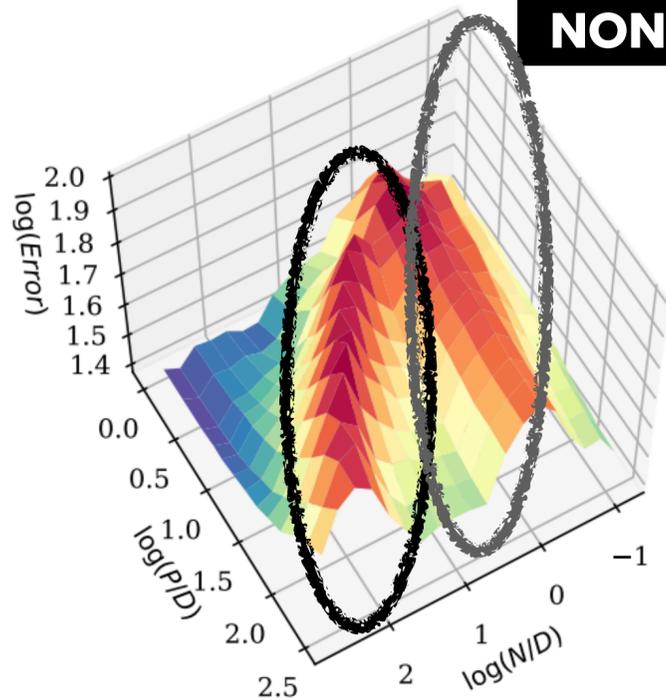
RF  
MODEL



DNN  
MODEL

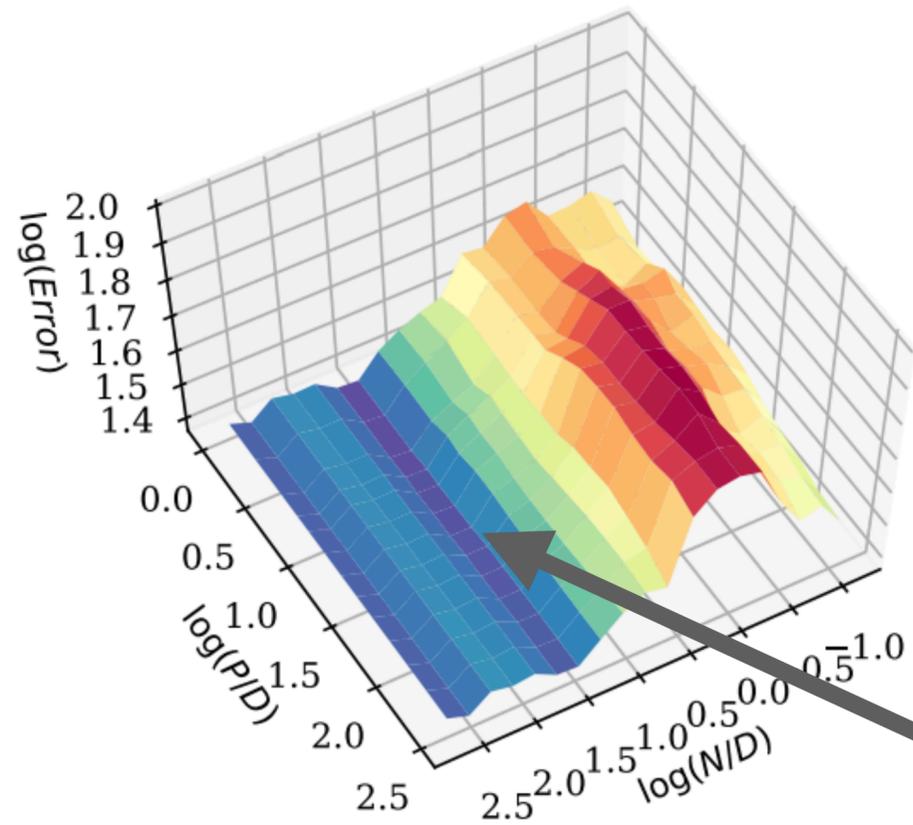
NONLINEAR PEAK REDUCED

LINEAR PEAK UNAFFECTED

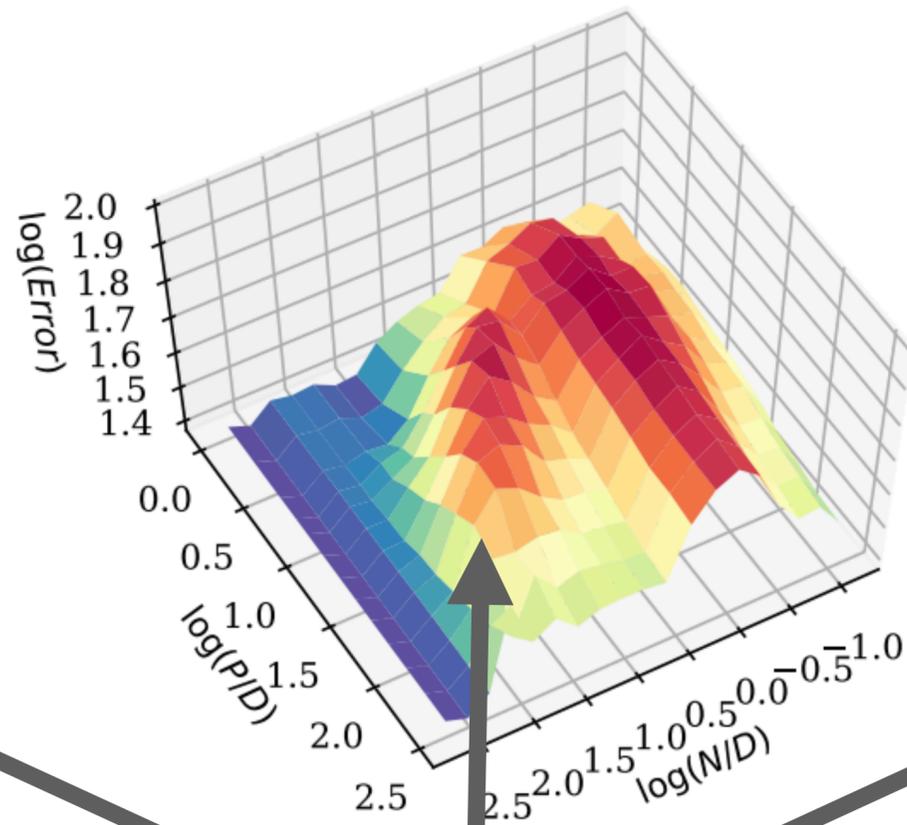


# TIME DEPENDENCE

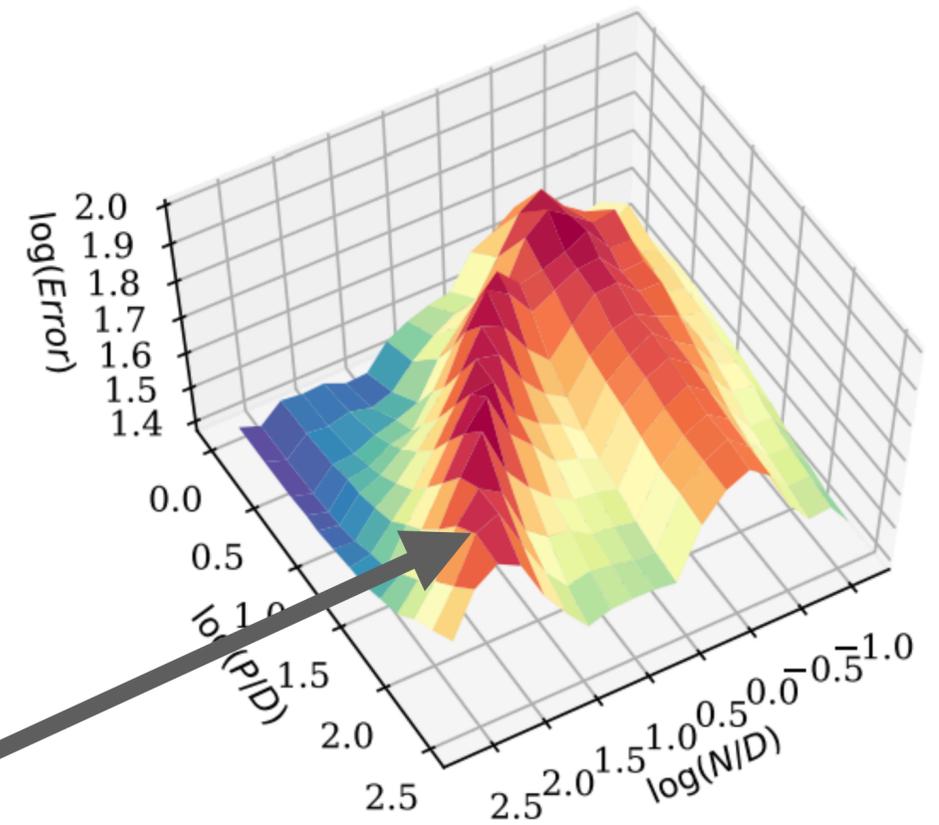
t=37 epochs



t=162 epochs



t=695 epochs



**THE NONLINEAR PEAK FORMS AT LATE TIMES**

# THANK YOU



**Stéphane d'Ascoli**



**Levent Sagun**



**Maria Refinetti**



**Giulio Biroli**



**Florent Krzakala**

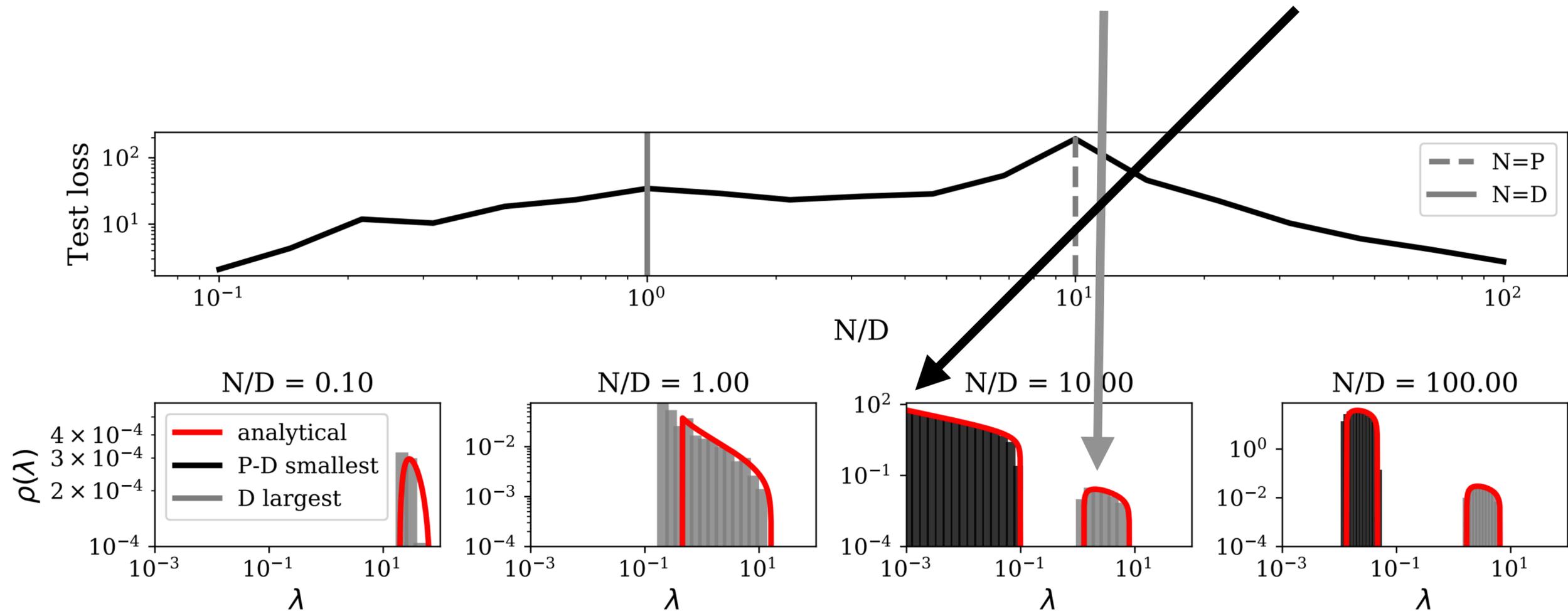
# ANALYTICAL SPECTRUM

**NONLINEAR = LINEAR + NOISE**

$$Z = \sigma \left( \frac{X\Theta^T}{\sqrt{D}} \right) \rightarrow \sqrt{\zeta} \frac{X\Theta^T}{\sqrt{D}} + \sqrt{\eta - \zeta} W, \quad W \sim \mathcal{N}(0,1)$$

**LINEAR  
PART**

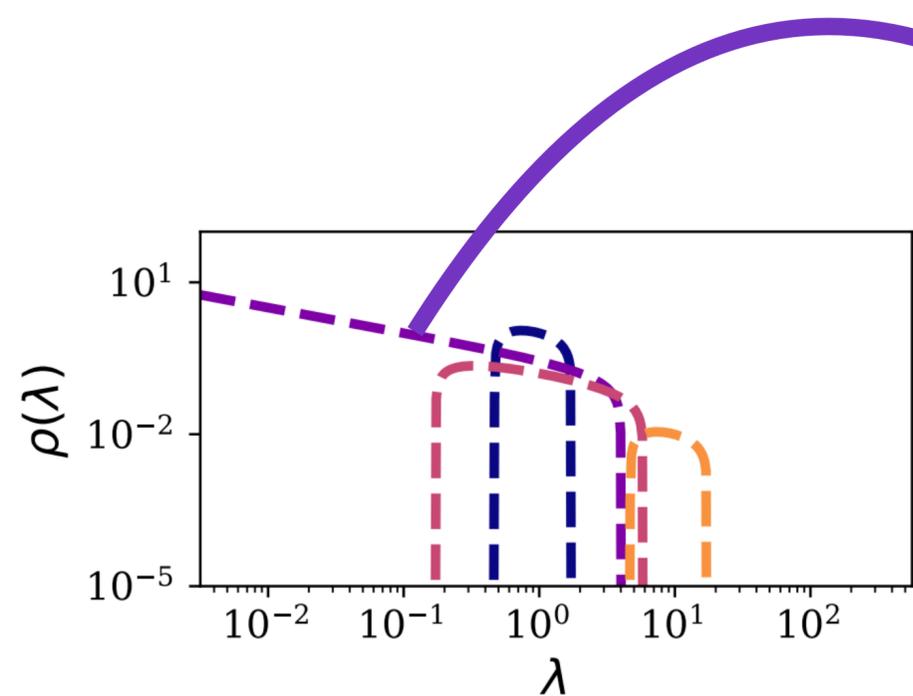
**NONLINEAR  
PART**



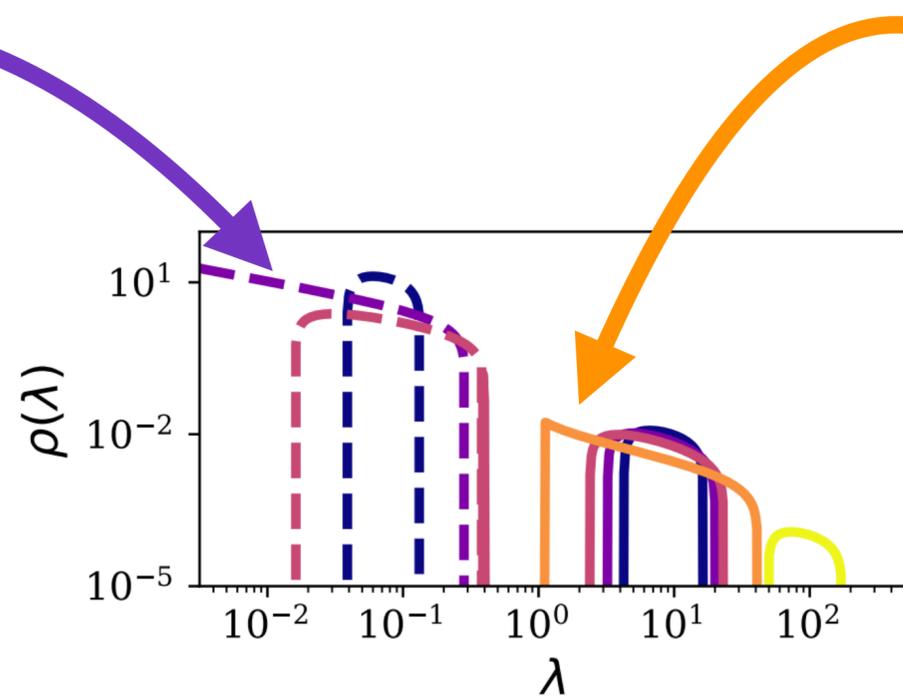
# ANALYTICAL SPECTRUM

**N=P GAP SURVIVES**

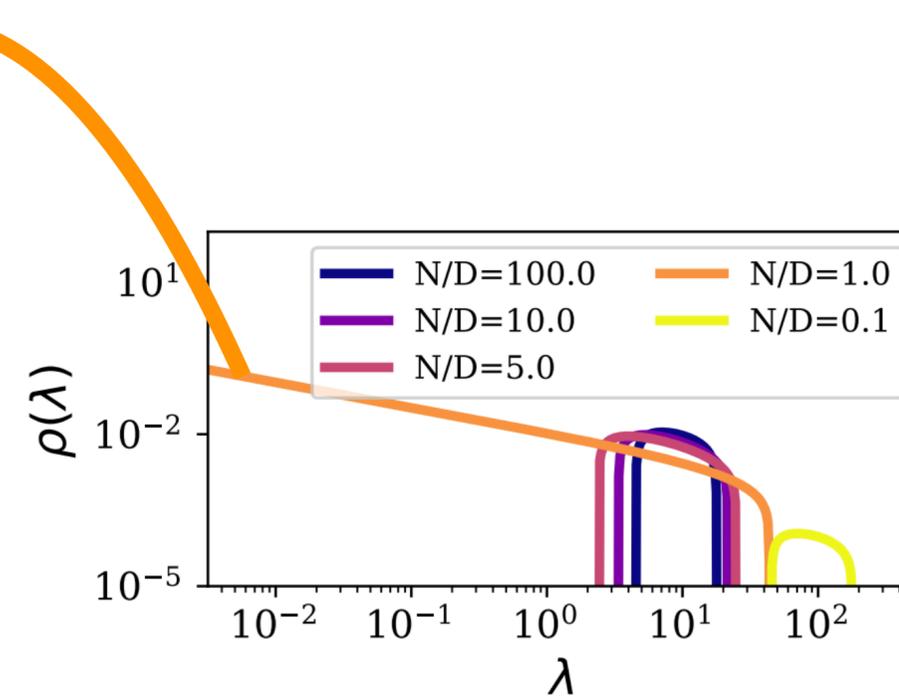
**N=D GAP IS REGULARISED**



(a) Absolute value ( $r=0$ )

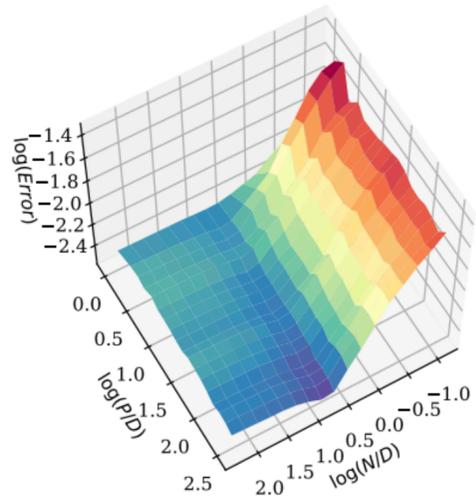


(b) Tanh ( $r \simeq 0.92$ )



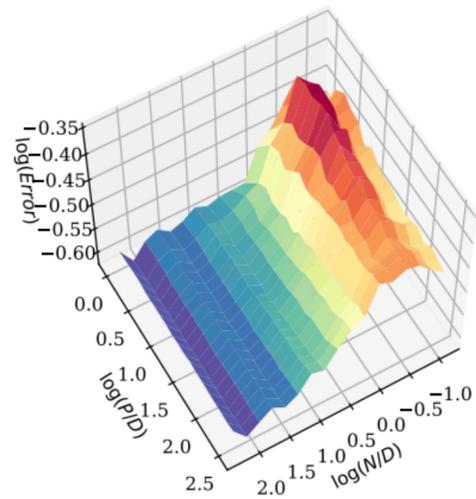
(c) Linear ( $r=1$ )

# STRUCTURED DATASETS



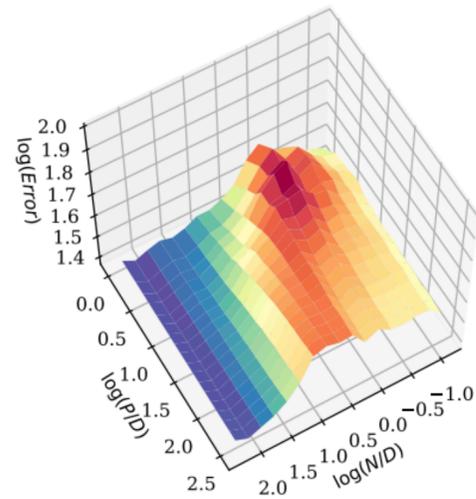
(a) MNIST,  $SNR = \infty$

t=37 epochs



(b) MNIST,  $SNR = 2$

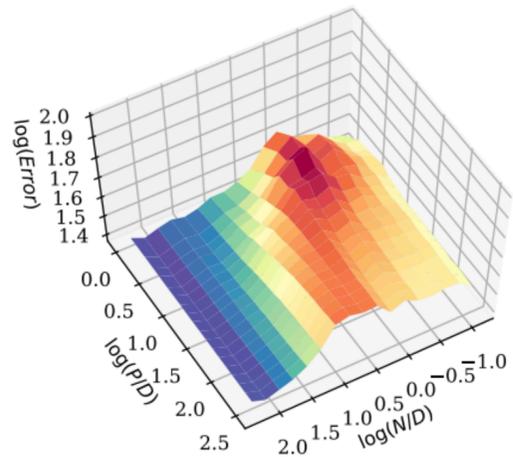
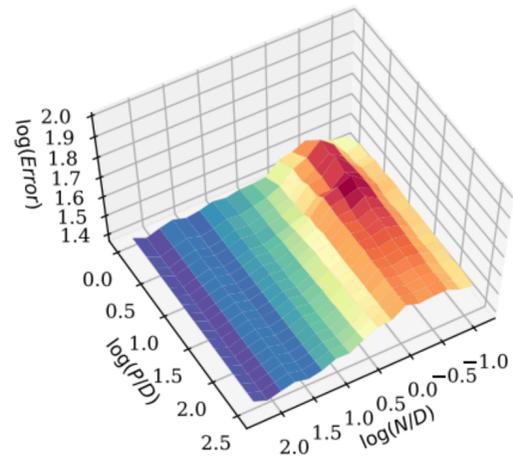
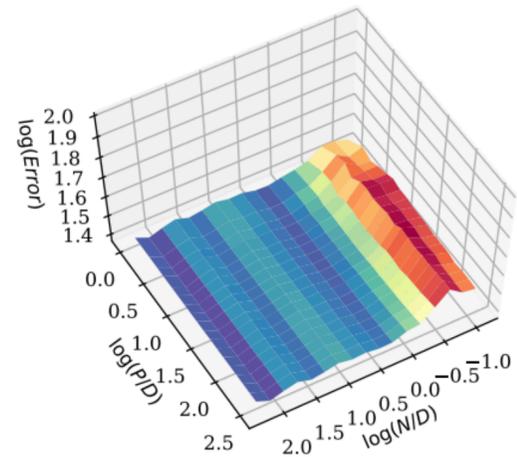
t=162 epochs



(c) MNIST,  $SNR = 0.2$

t=695 epochs

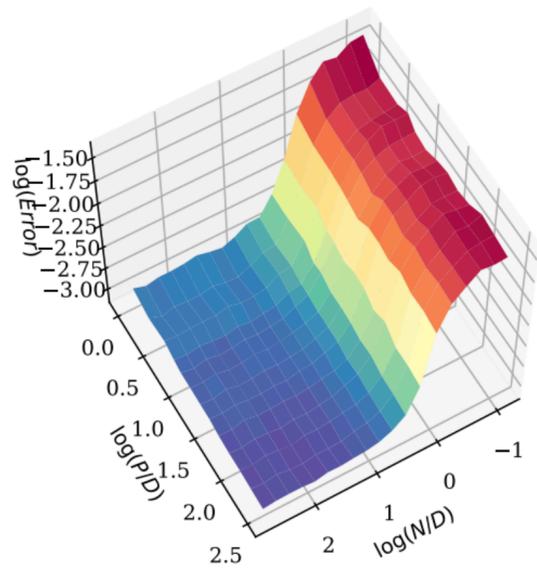
**LINEAR AND NONLINEAR PEAK  
ARE MERGED TOGETHER**



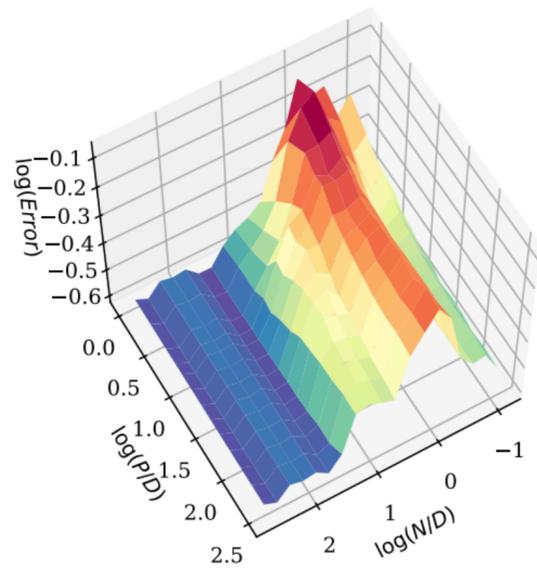
(d) Dynamics on MNIST at  $SNR = 0.2$

**SHIFT FROM LINEAR TO NONLINEAR  
DURING TRAINING**

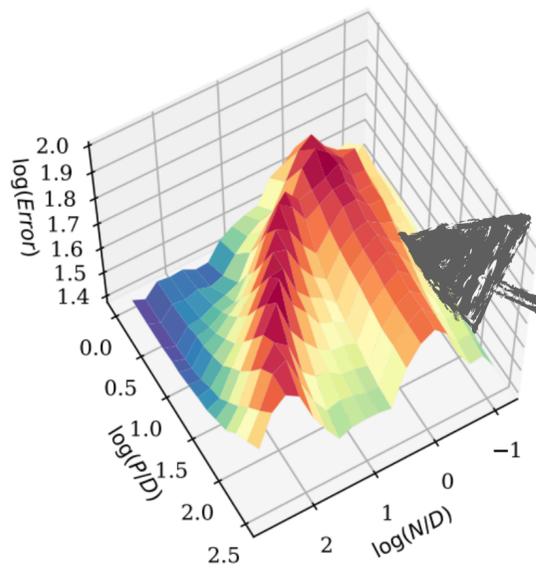
# EFFECT OF NOISE AND NONLINEARITY



(a) Tanh,  $SNR = \infty$

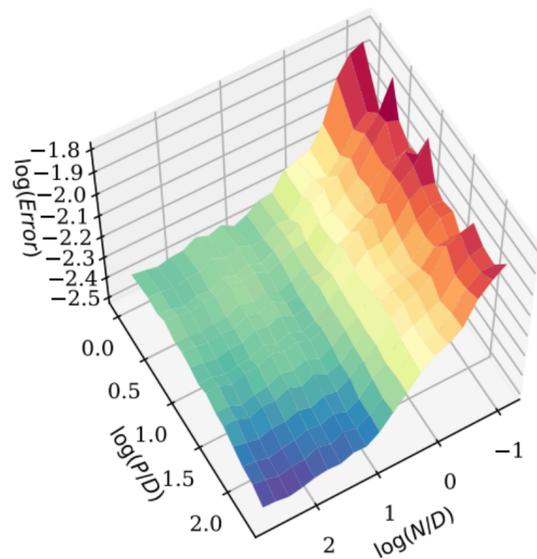


(b) Tanh,  $SNR = 2$

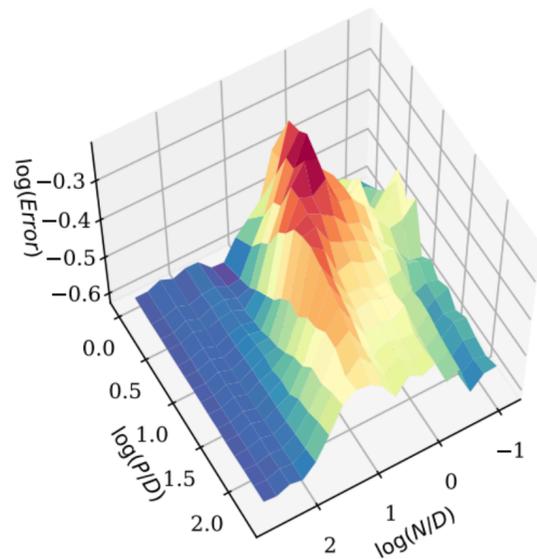


(c) Tanh,  $SNR = 0.2$

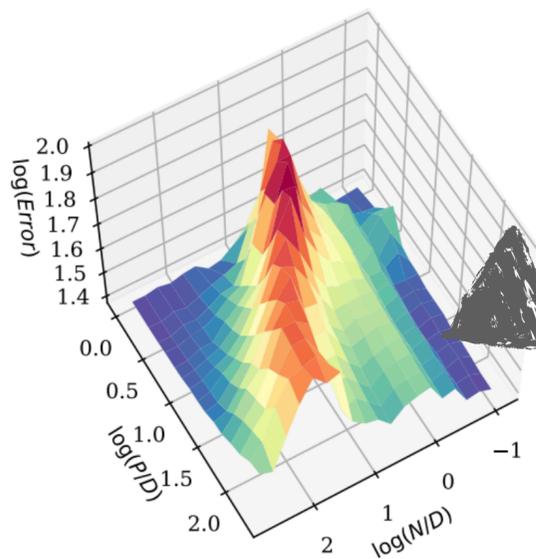
**LINEAR PEAK IS WEAKER FOR RELU**



(d) ReLU,  $SNR = \infty$



(e) ReLU,  $SNR = 2$



(f) ReLU,  $SNR = 0.2$