

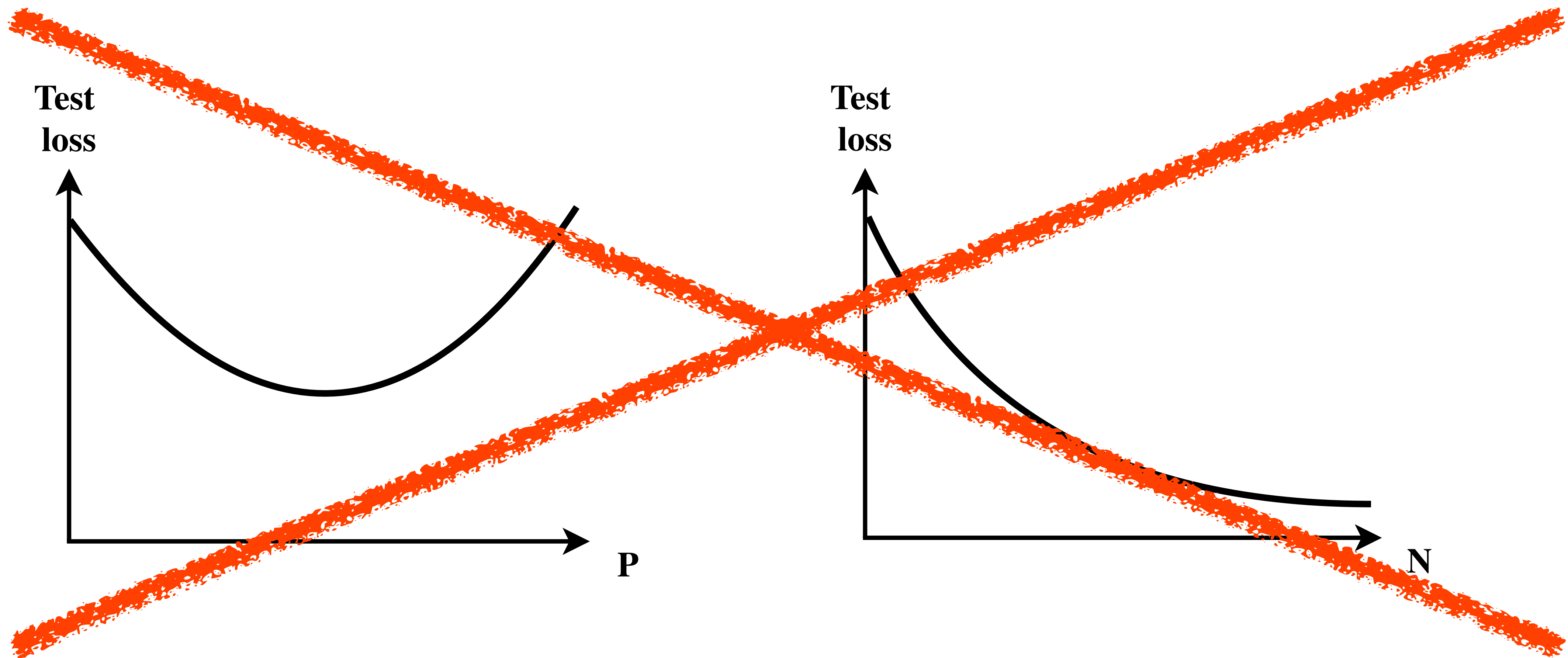


TRIPLE DESCENT : THE TWO KINDS OF OVERFITTING

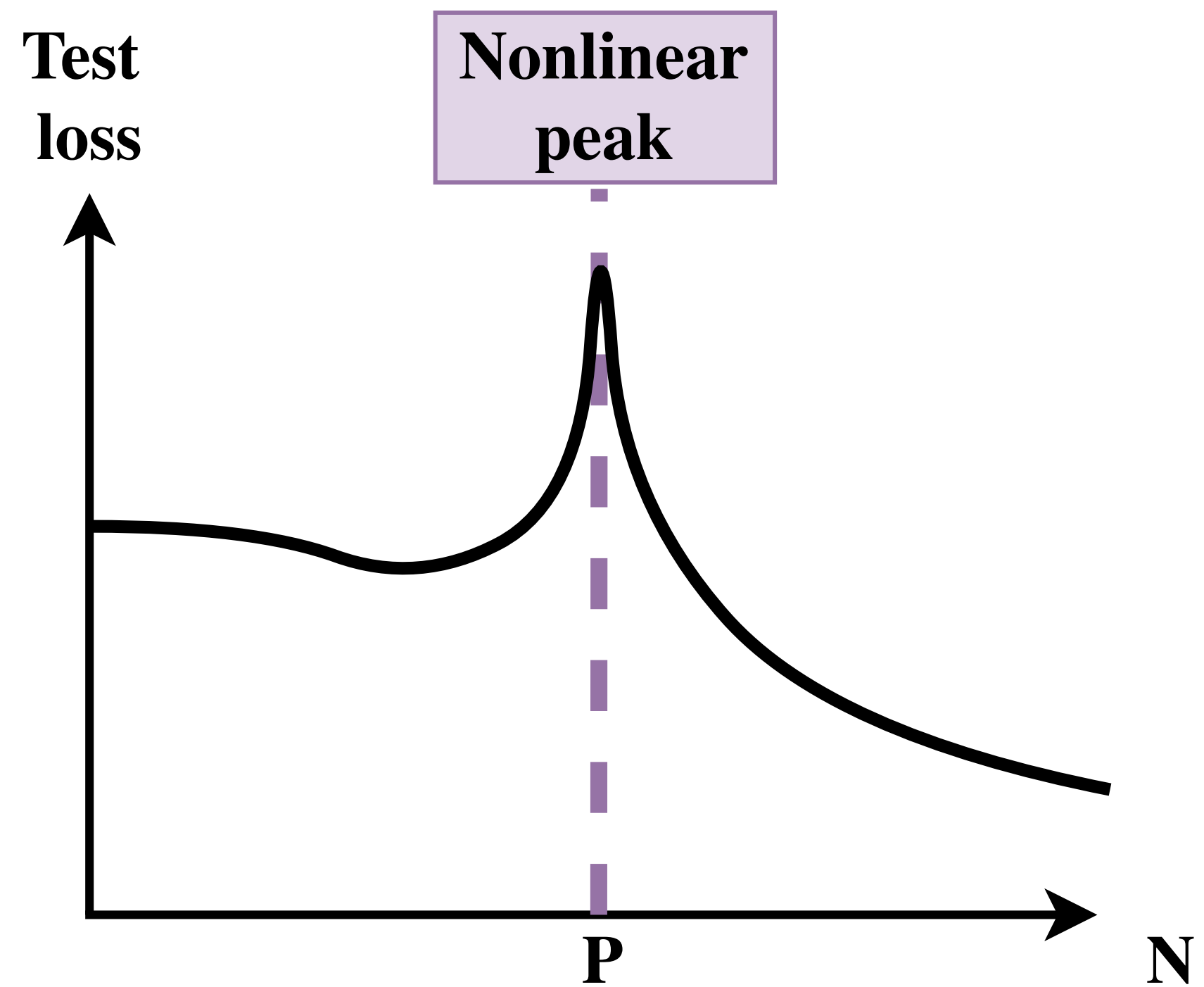
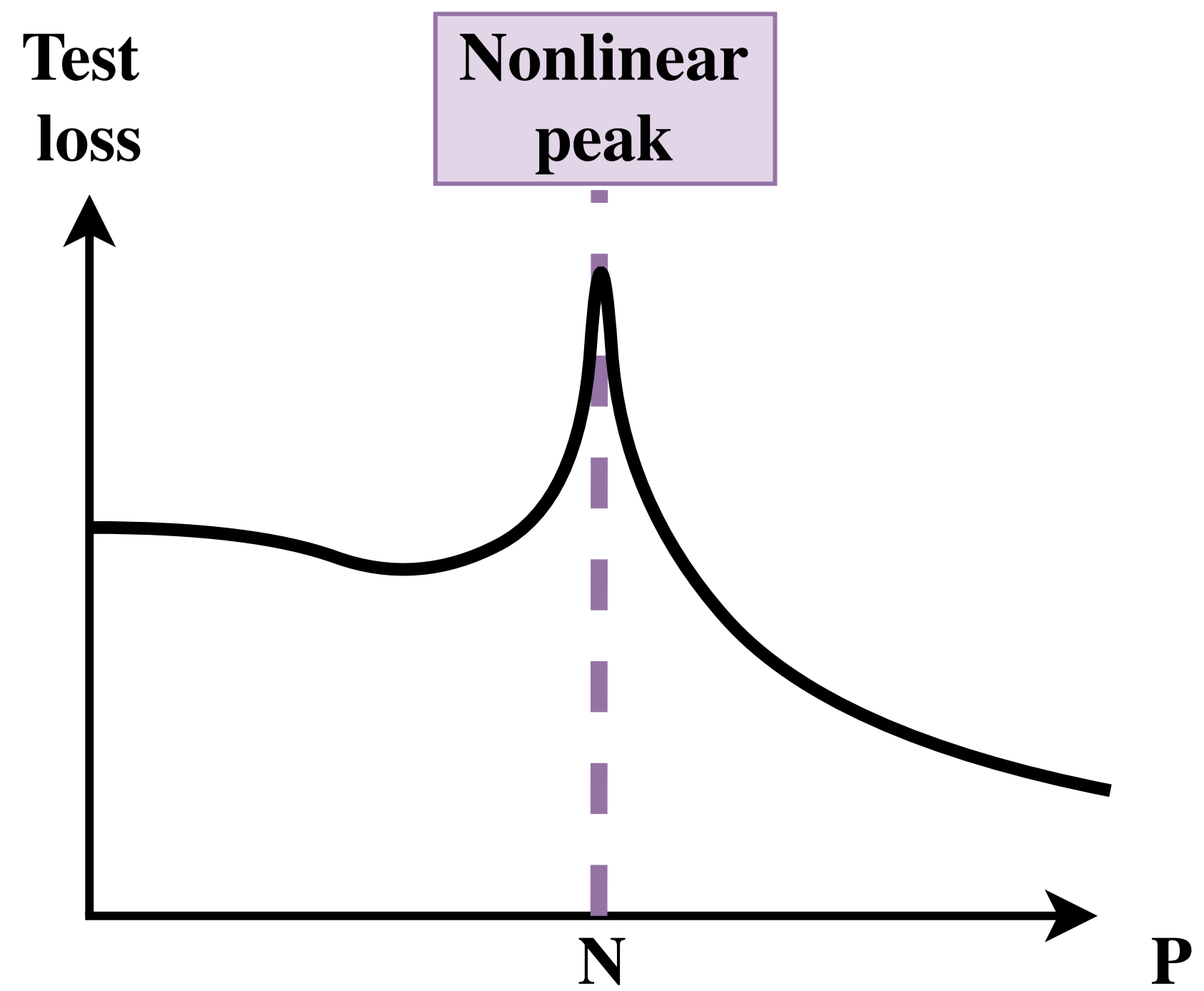
STÉPHANE D'ASCOLI, LEVENT SAGUN, GIULIO BIROLI

ÉCOLE NORMALE SUPÉRIEURE & FACEBOOK AI RESEARCH

PARAMETER-WISE AND SAMPLE-WISE



PARAMETER-WISE AND SAMPLE-WISE

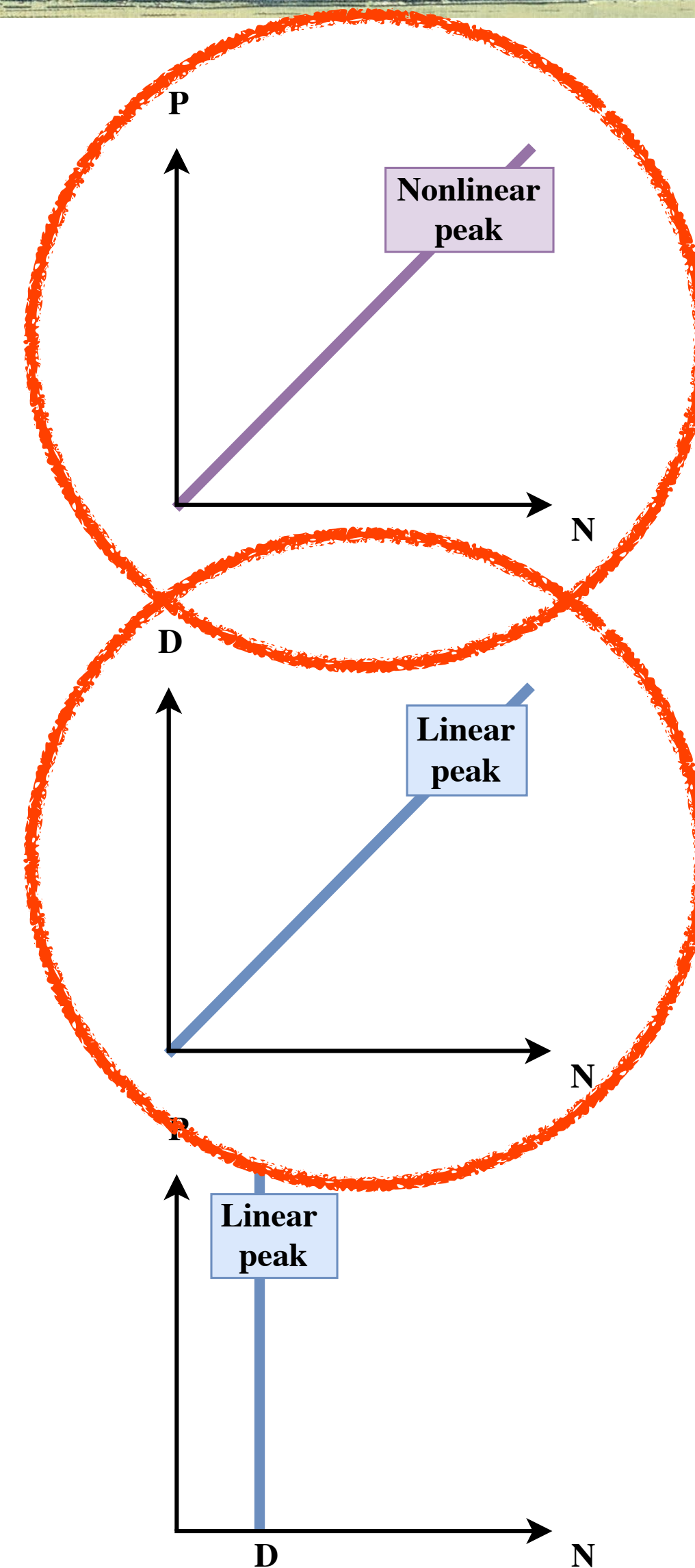
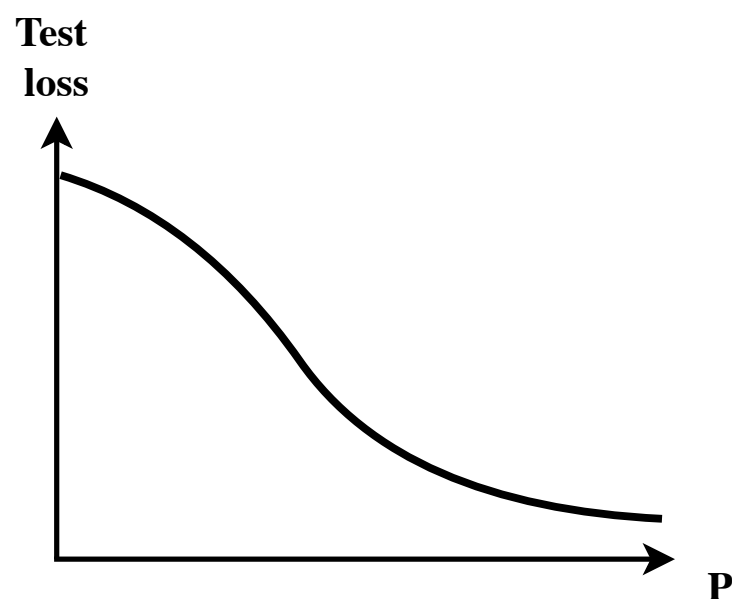
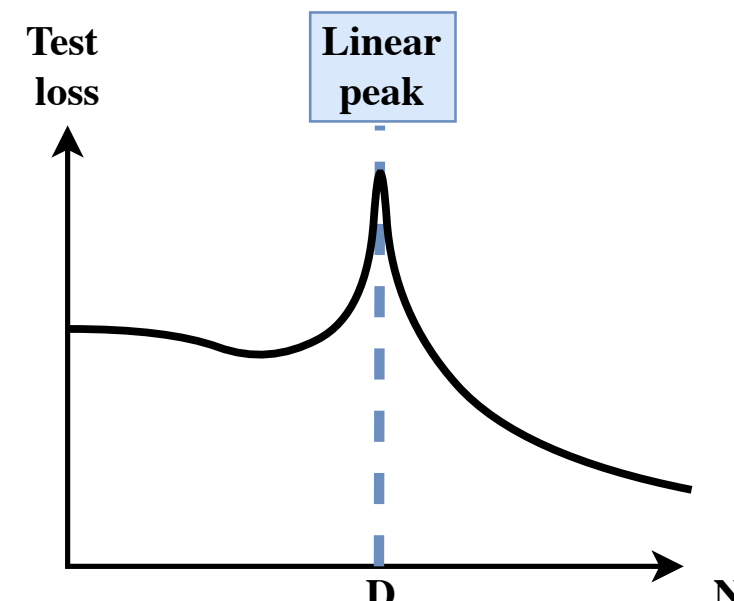
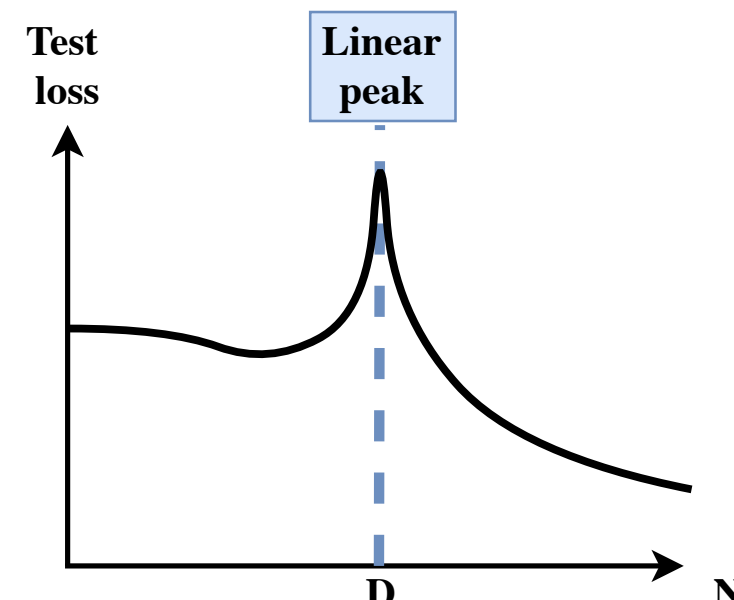
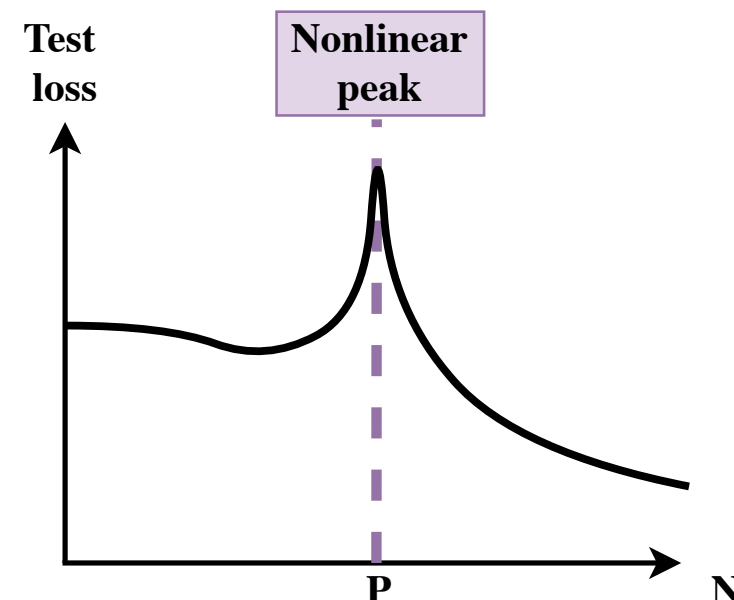
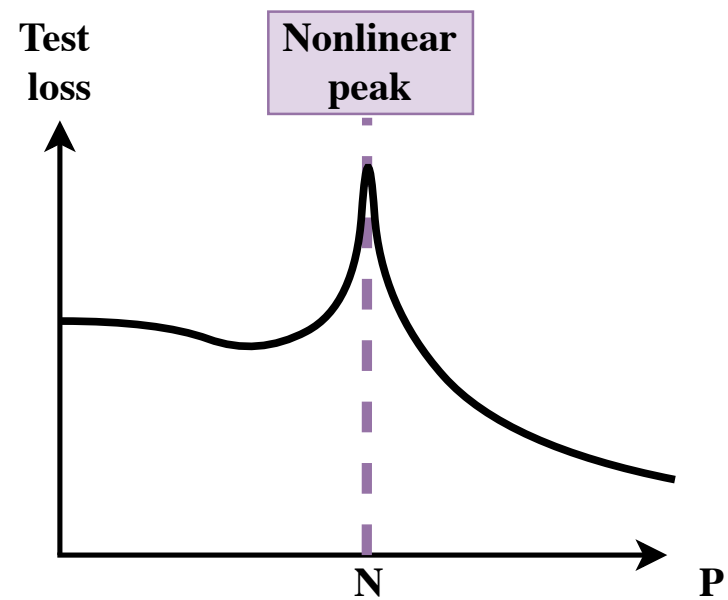


FOR LINEAR TO NONLINEAR

NONLINEAR NETWORKS

LINEAR MODELS

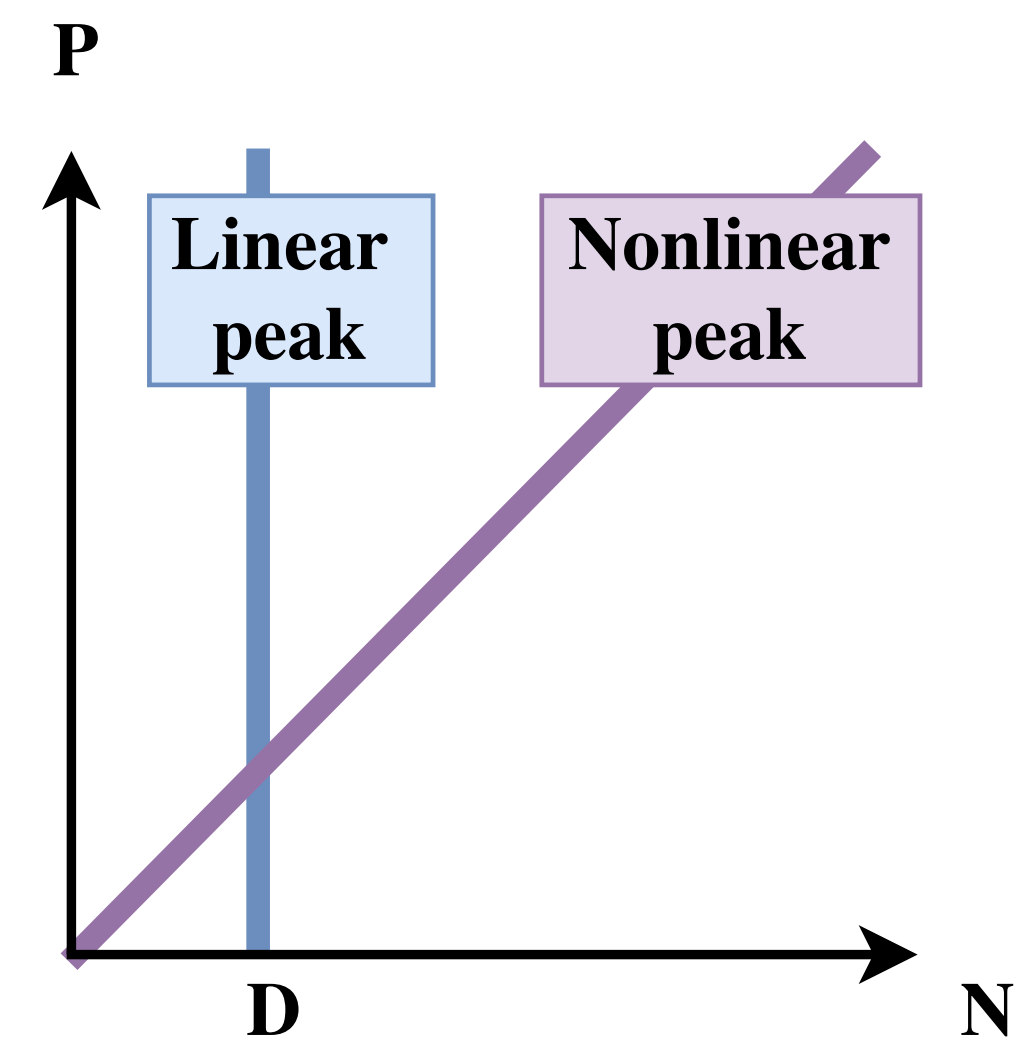
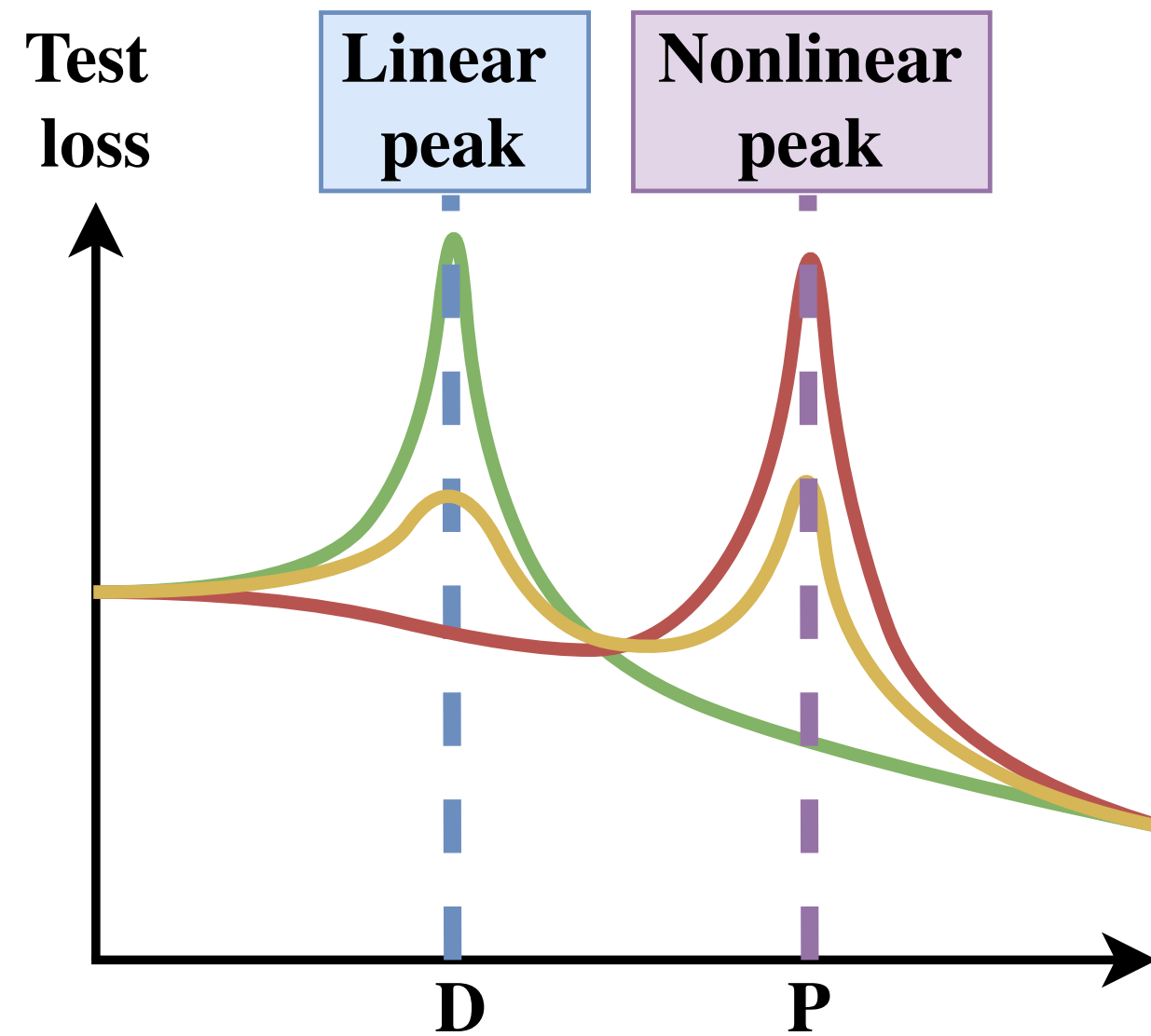
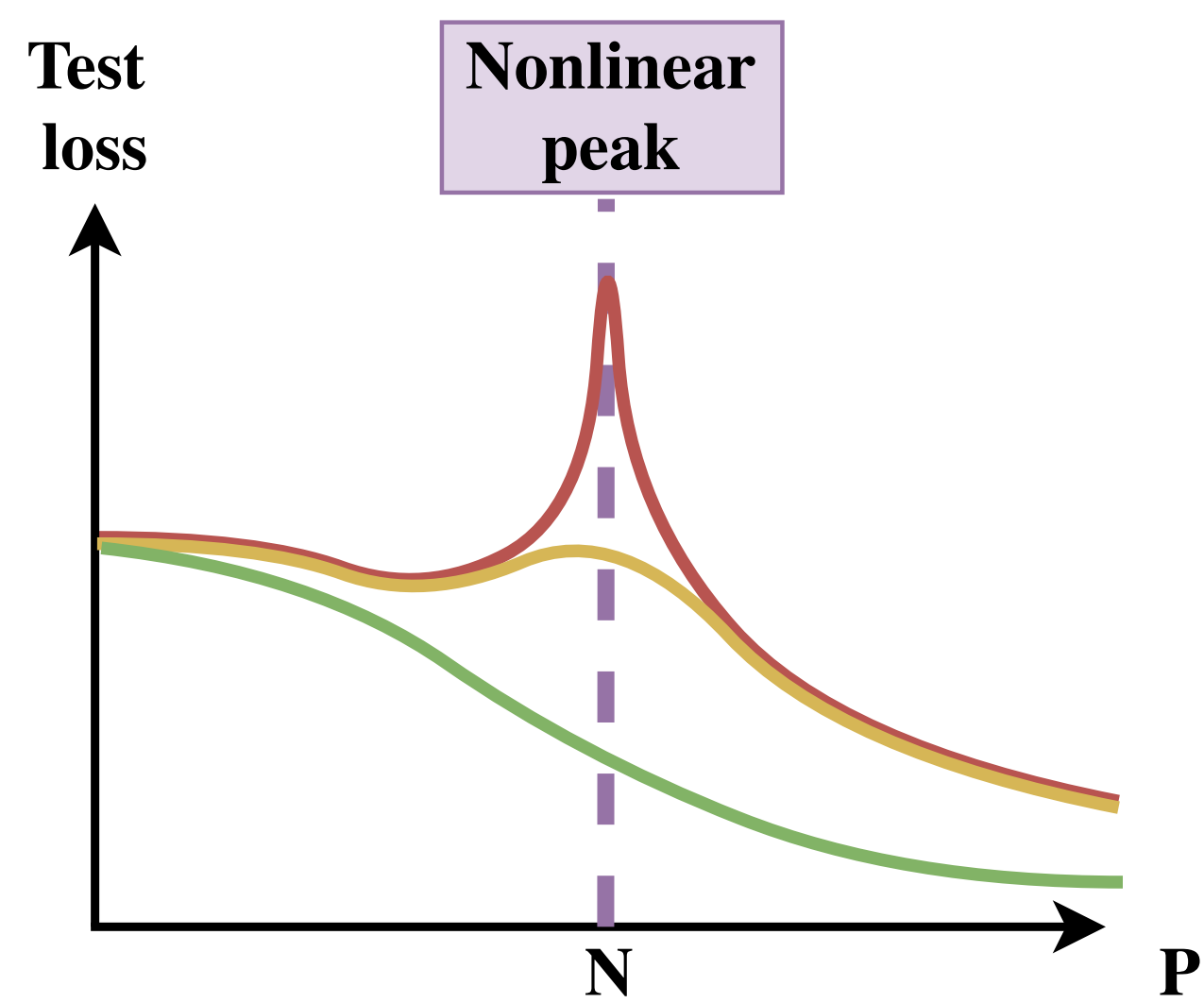
LINEAR NETWORKS



ARE THESE THE SAME ?

ANSWER : NO !

FROM LINEAR TO LINEAR



Activation function	
—	Strongly nonlinear
—	Weakly nonlinear
—	Linear

WHAT MECHANISMS UNDERLIE THESE PEAKS ?
HOW ARE THEY DIFFERENT ?

THE TWO MODELS

DATASET

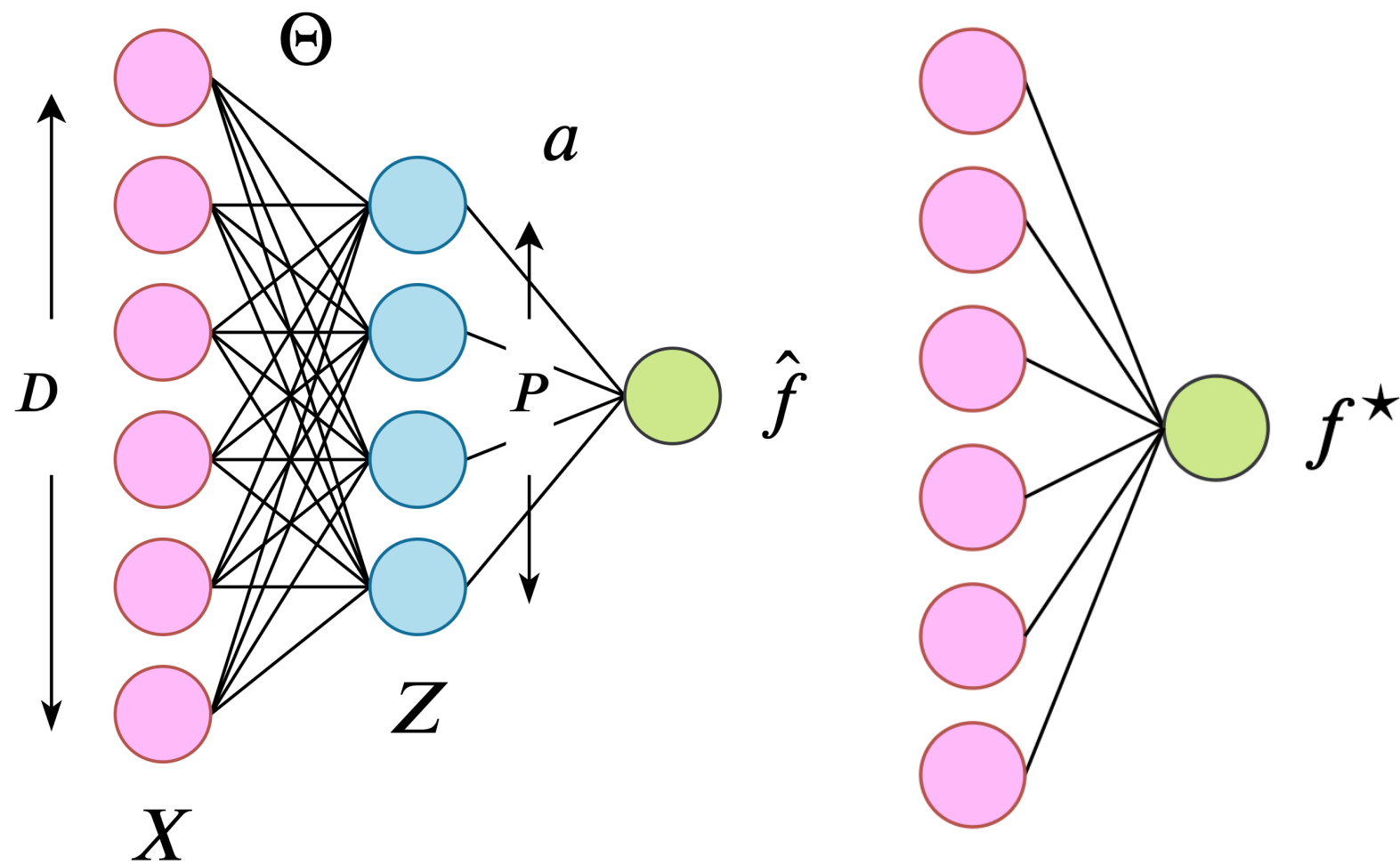
$$\mathbf{X} \sim \mathcal{N}(0,1) \in \mathbb{R}^{N \times D}$$

$$y = f^*(\mathbf{x}) + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, 1/\text{SNR})$$

$$\mathcal{L}_g = \mathbb{E}_x \left[\left(f(x) - \hat{f}(x) \right)^2 \right]$$

RF MODEL

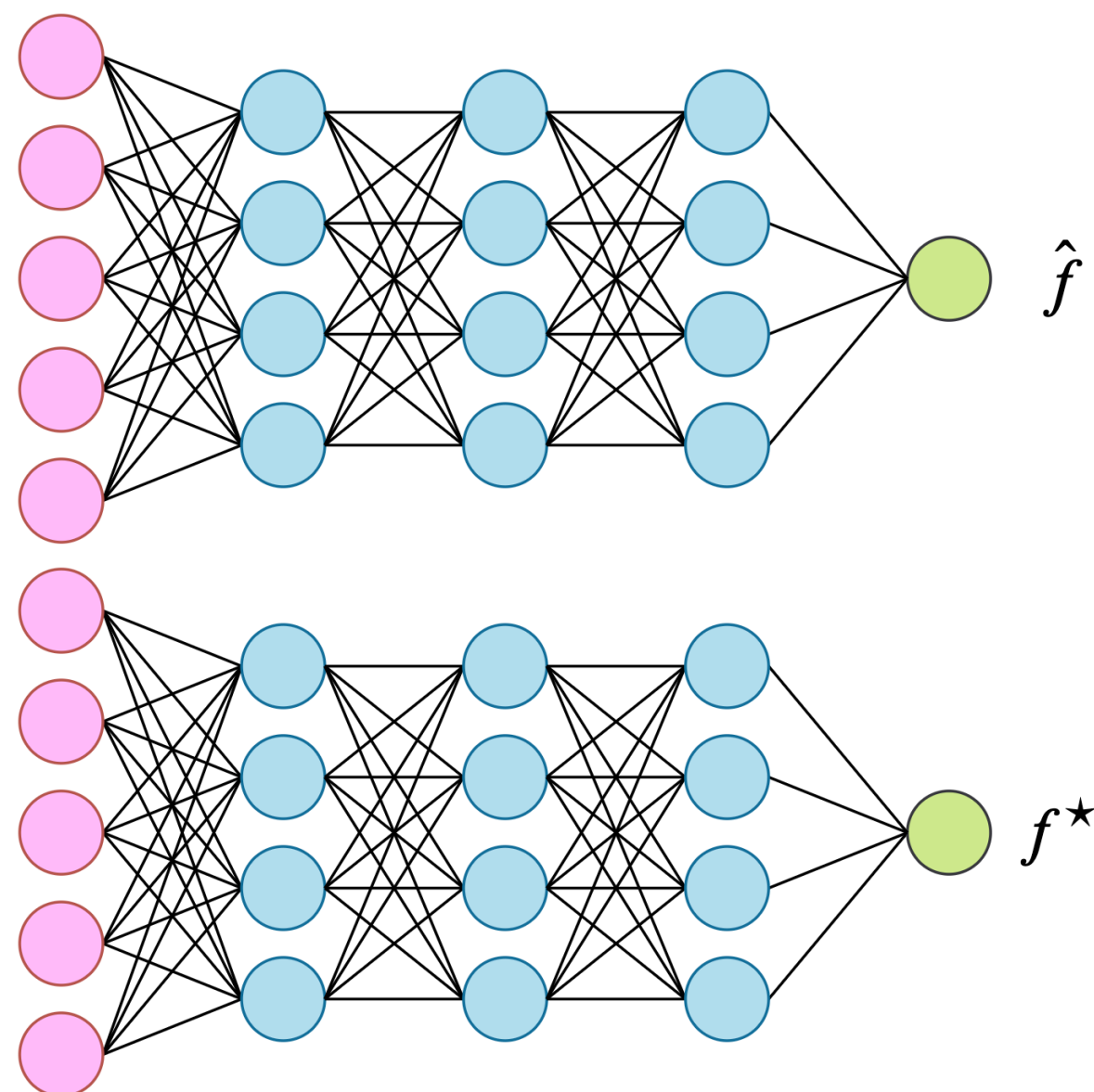


$$\hat{f}(x) = \sum_{i=1}^P a_i \sigma \left(\frac{\langle \Theta_i, x \rangle}{\sqrt{D}} \right)$$

$$\hat{a} = \arg \min_{a \in \mathbb{R}^P} \left[\frac{1}{N} (y - aZ^T)^2 + \frac{P\gamma}{D} \|a\|_2^2 \right]$$

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

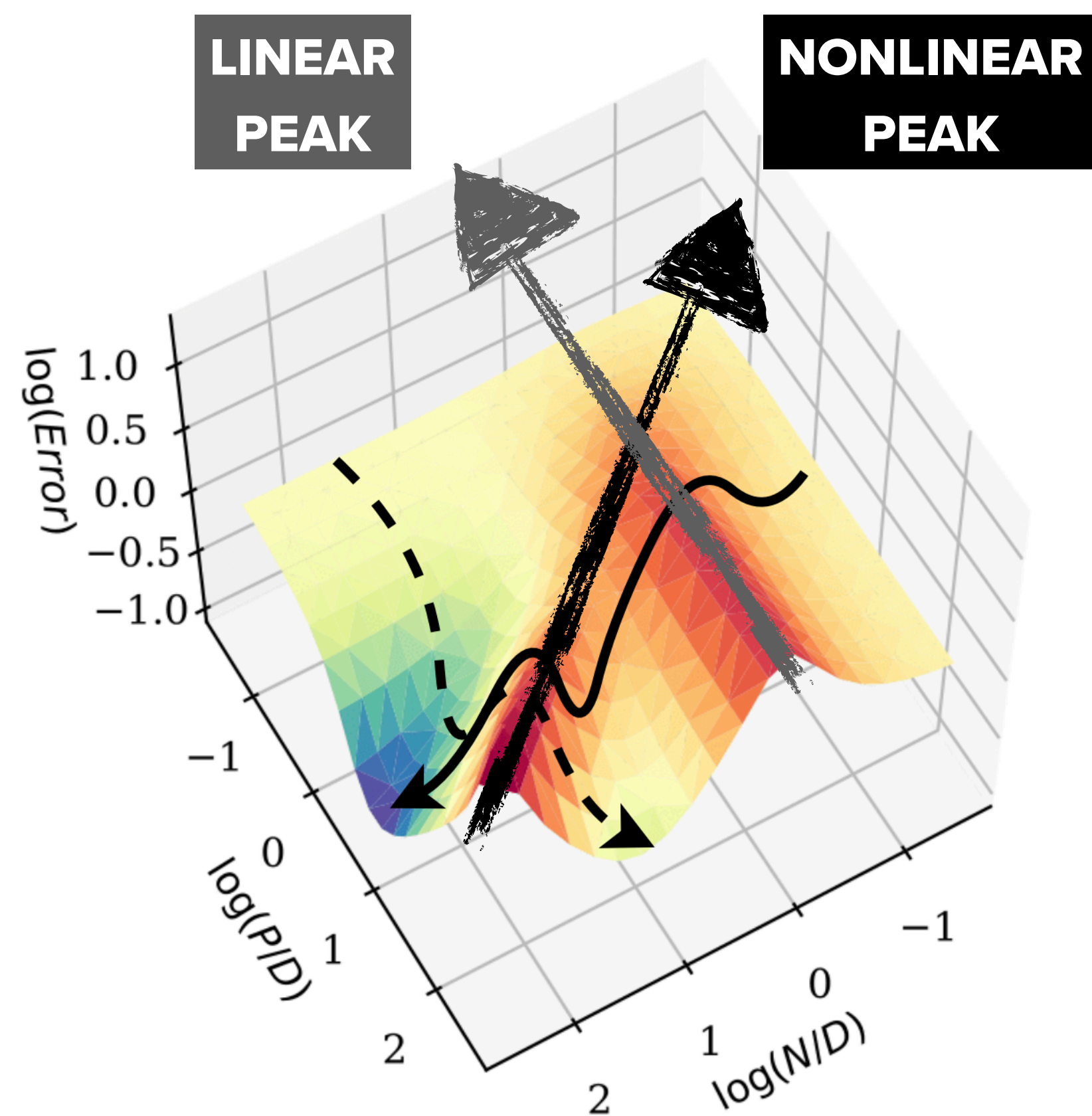
DNN MODEL



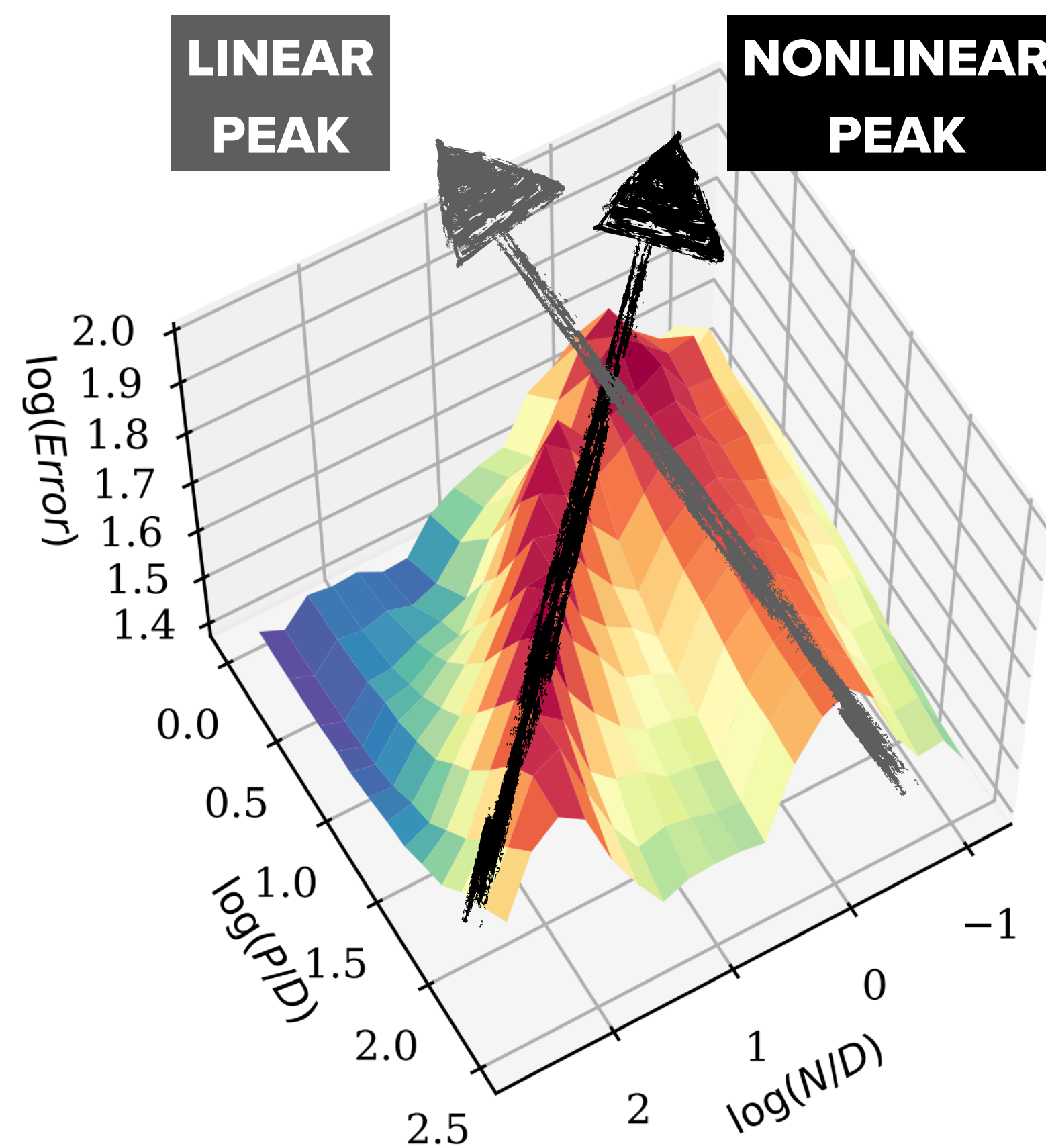
Random teacher

Student trained by GD

EVIDENCE OF TRIPLE DESCENT



**RF
MODEL**



**DNN
MODEL**

ANALYTICAL DESCRIPTION

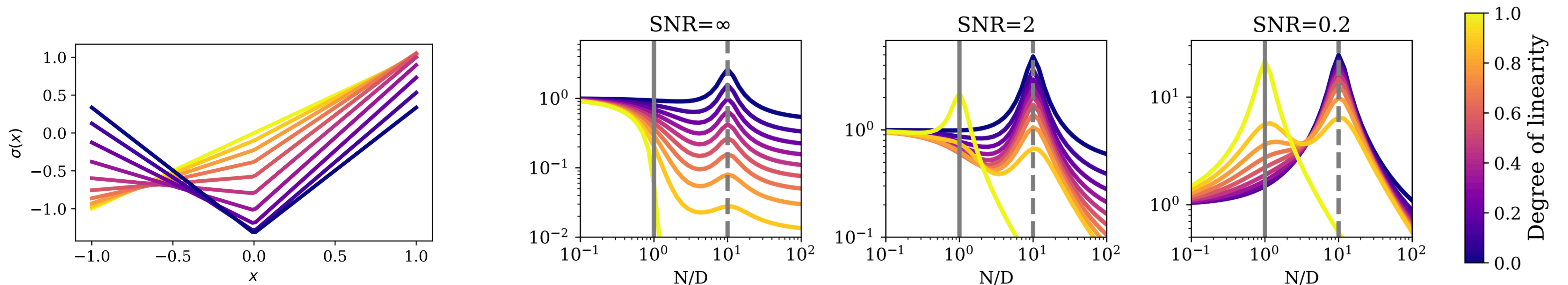
HIGH-DIMENSIONAL LIMIT

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

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

DEGREE OF LINEARITY

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



ANALYTICAL SPECTRUM

**LINEAR
PART**

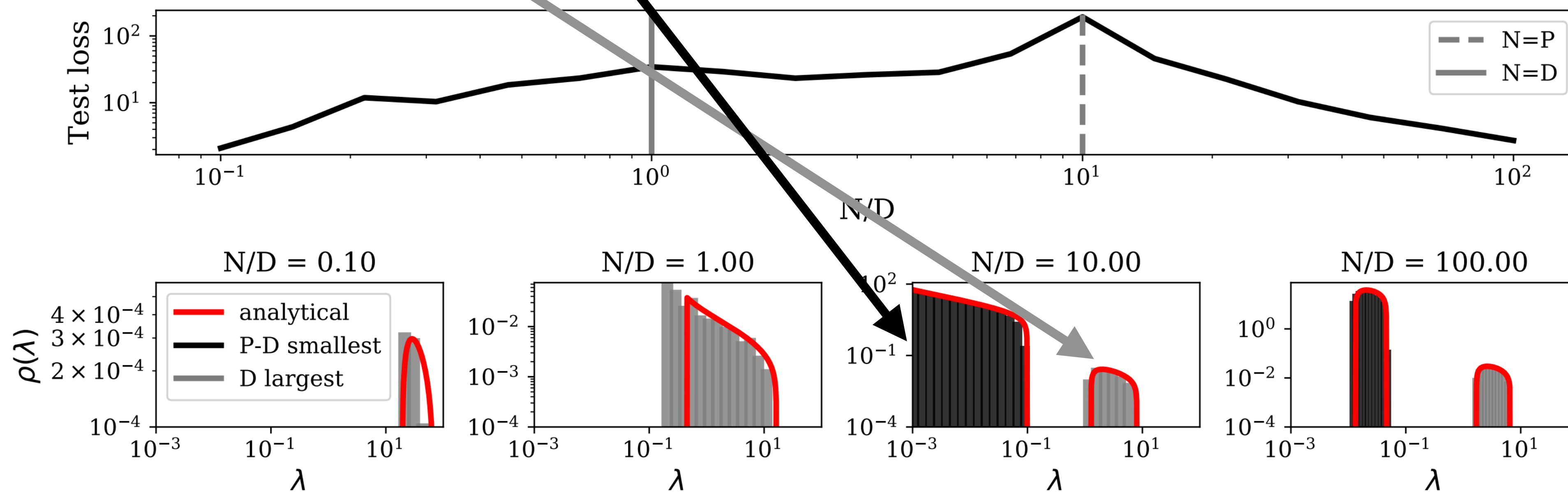
**NONLINEAR
PART**

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

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

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

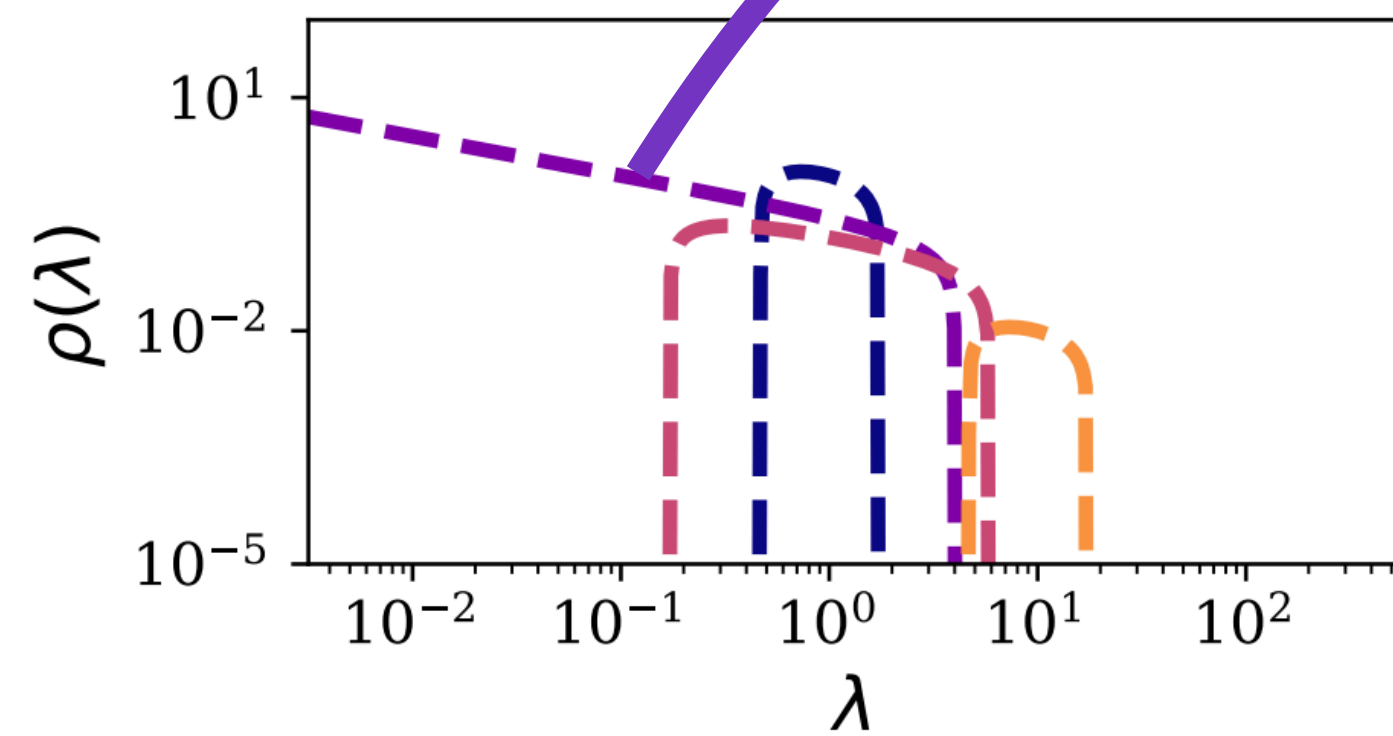
[Pennington & Worah 2017]



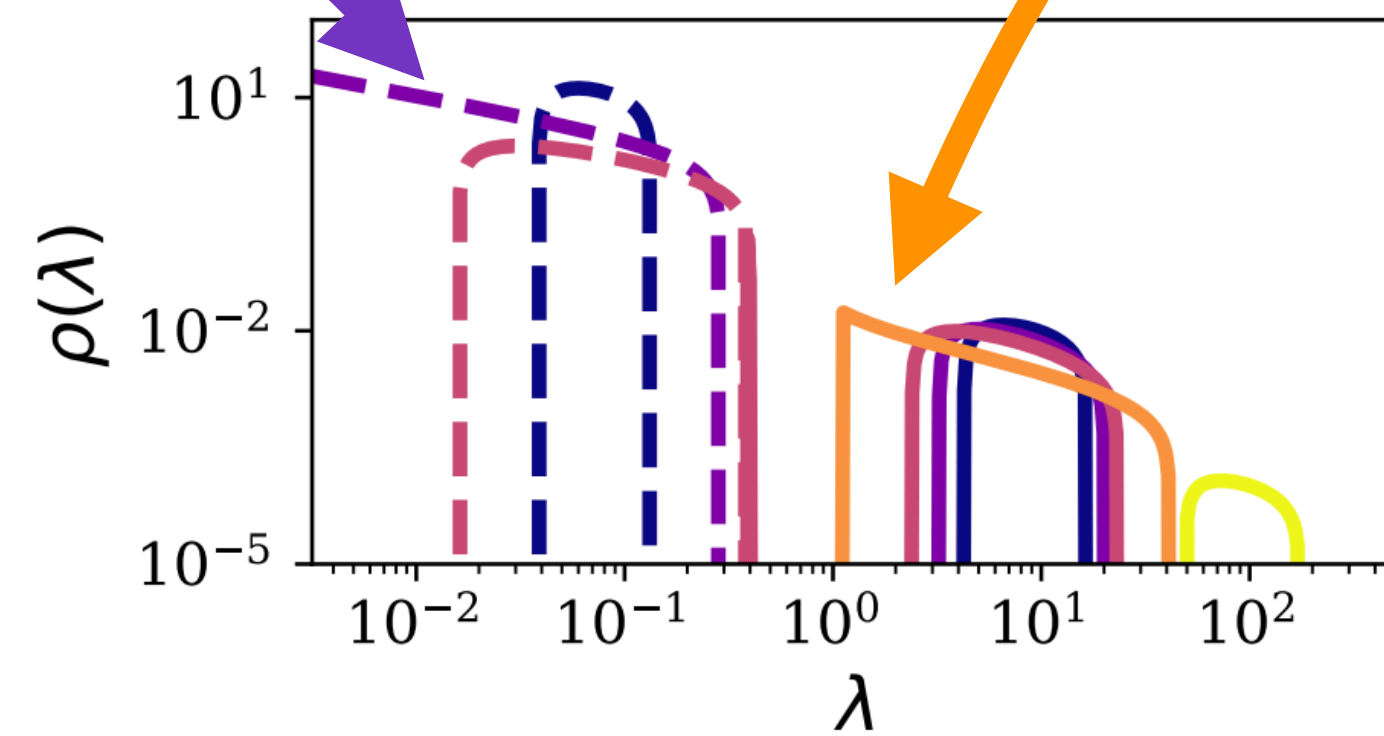
ANALYTICAL SPECTRUM

N=P GAP SURVIVES

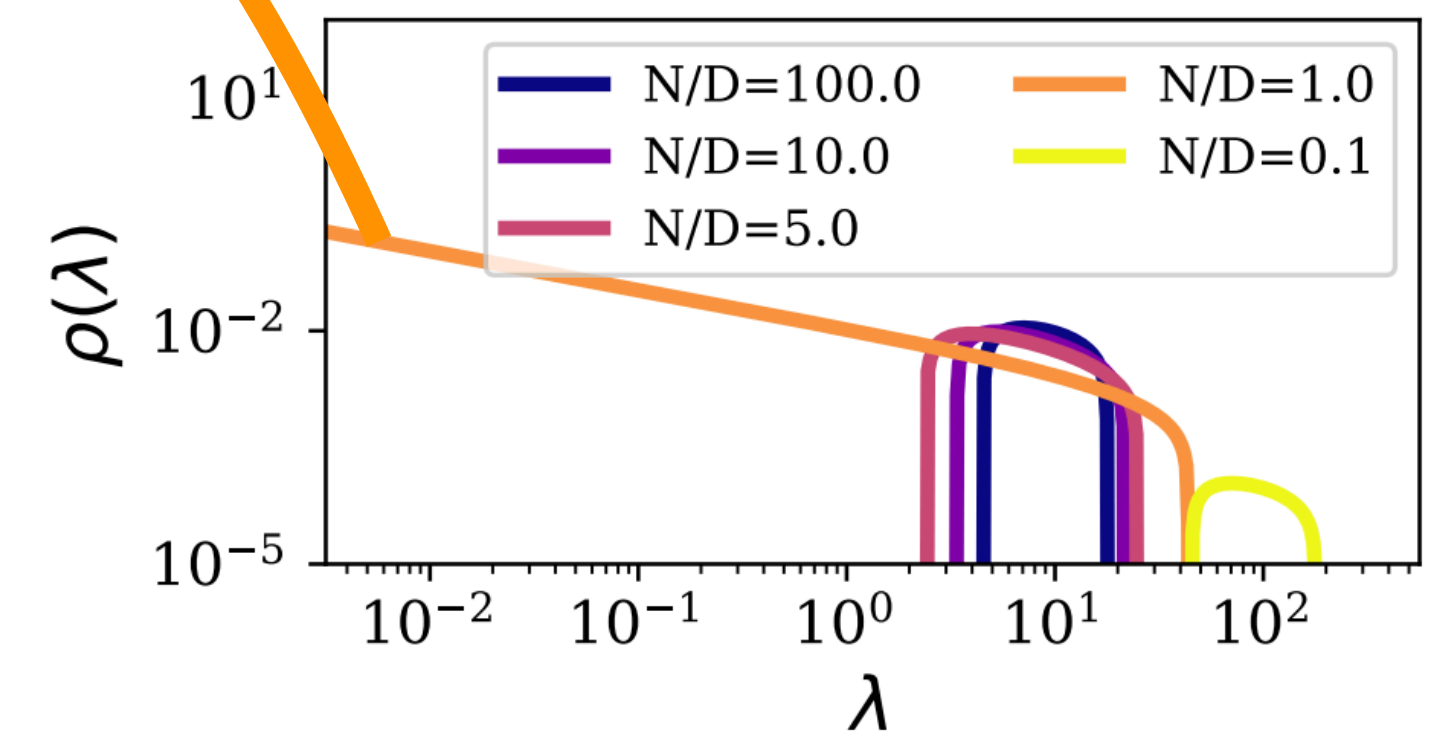
N=D GAP IS REGULARISED



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



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



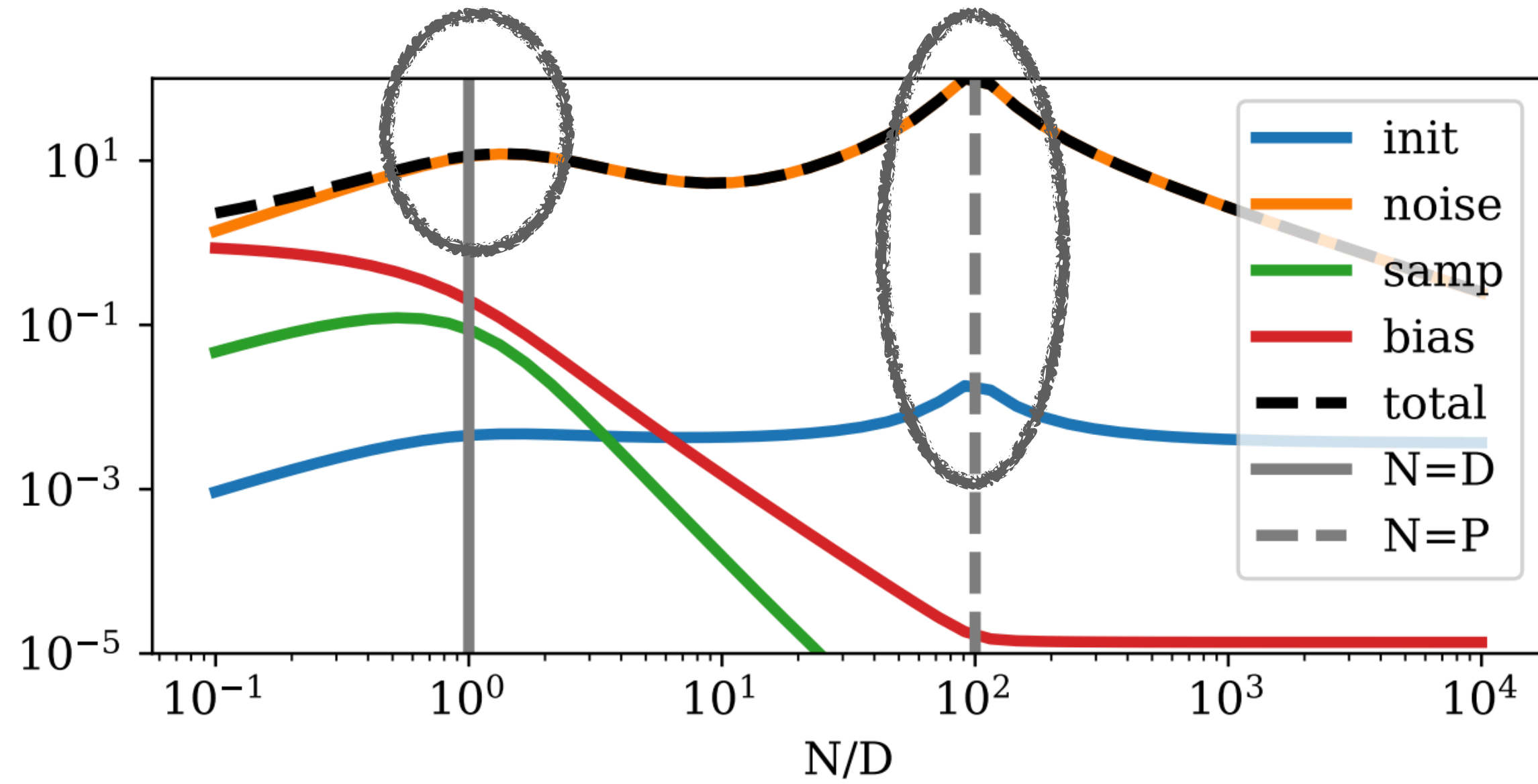
(c) Linear ($r=1$)

BIAS AND VARIANCES

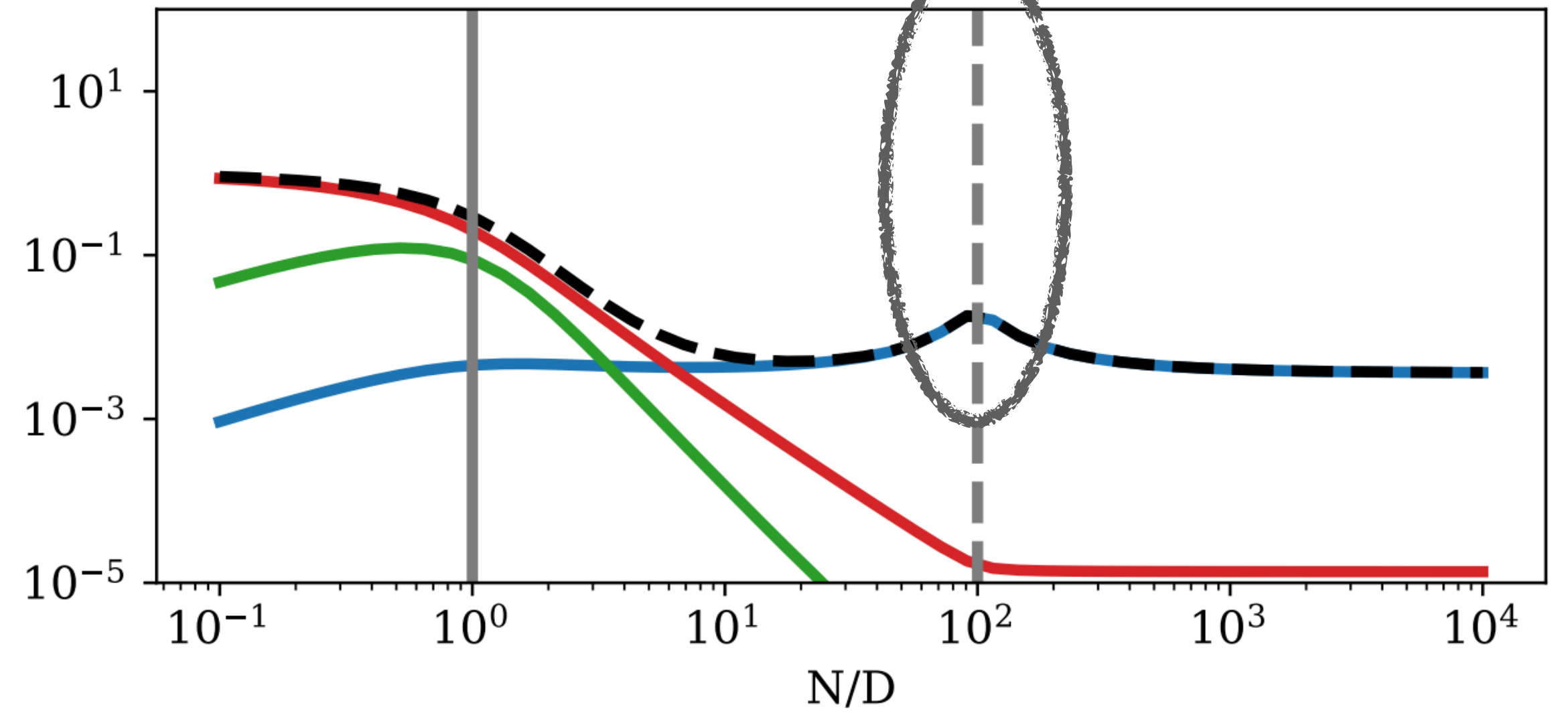
**LINEAR PEAK
CAUSED BY NOISE**

**NONLINEAR PEAK
CAUSED BY NOISE & INIT**

**NONLINEAR PEAK
SURVIVES IN ABSENCE OF NOISE**

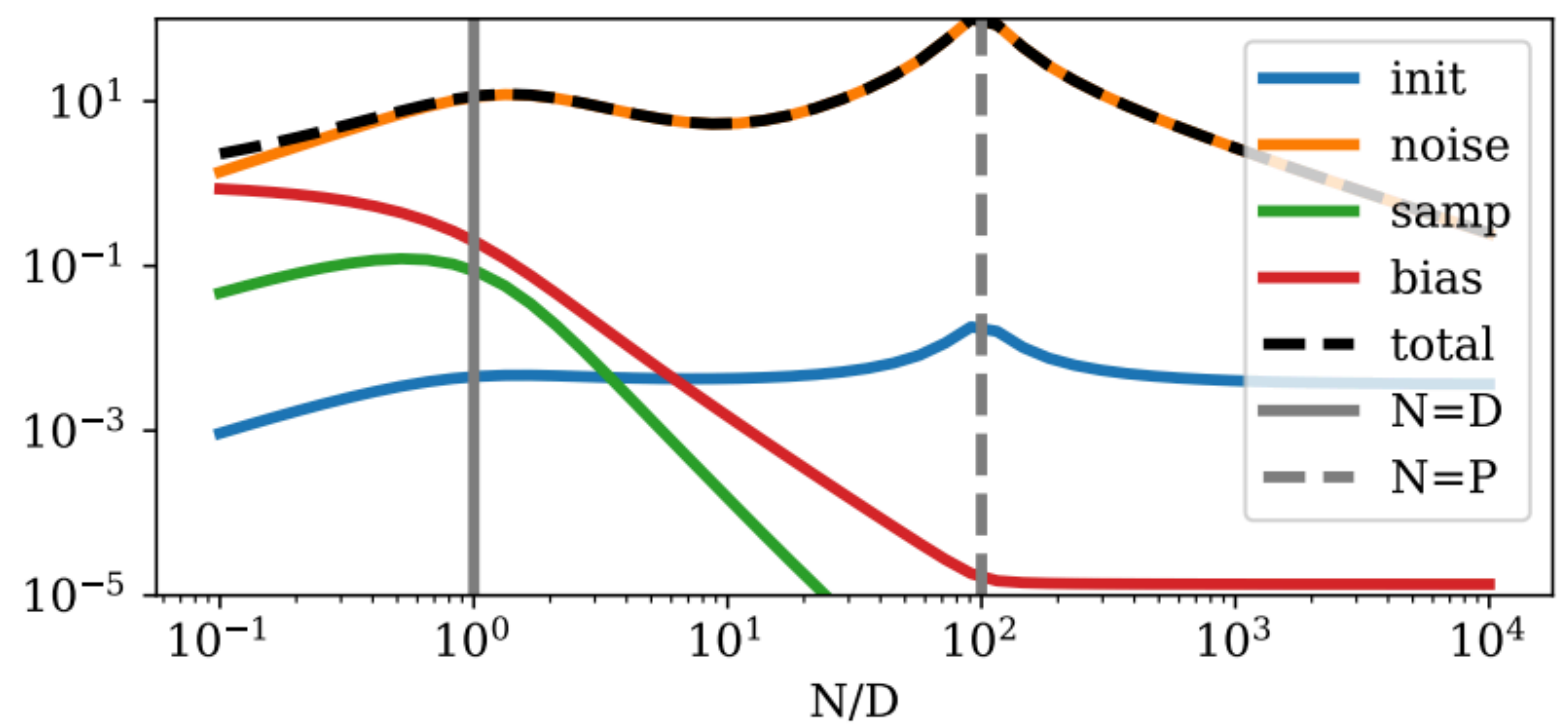


NOISY

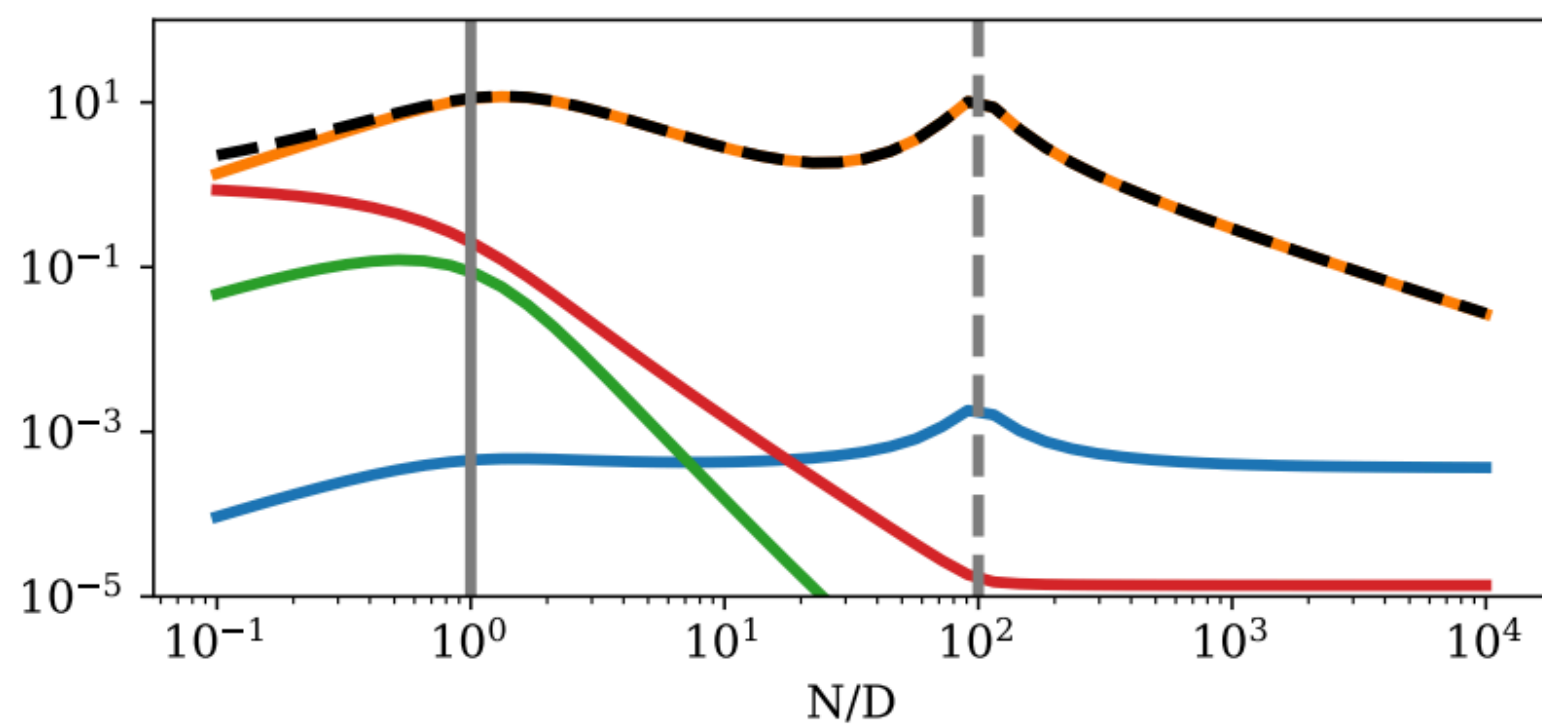
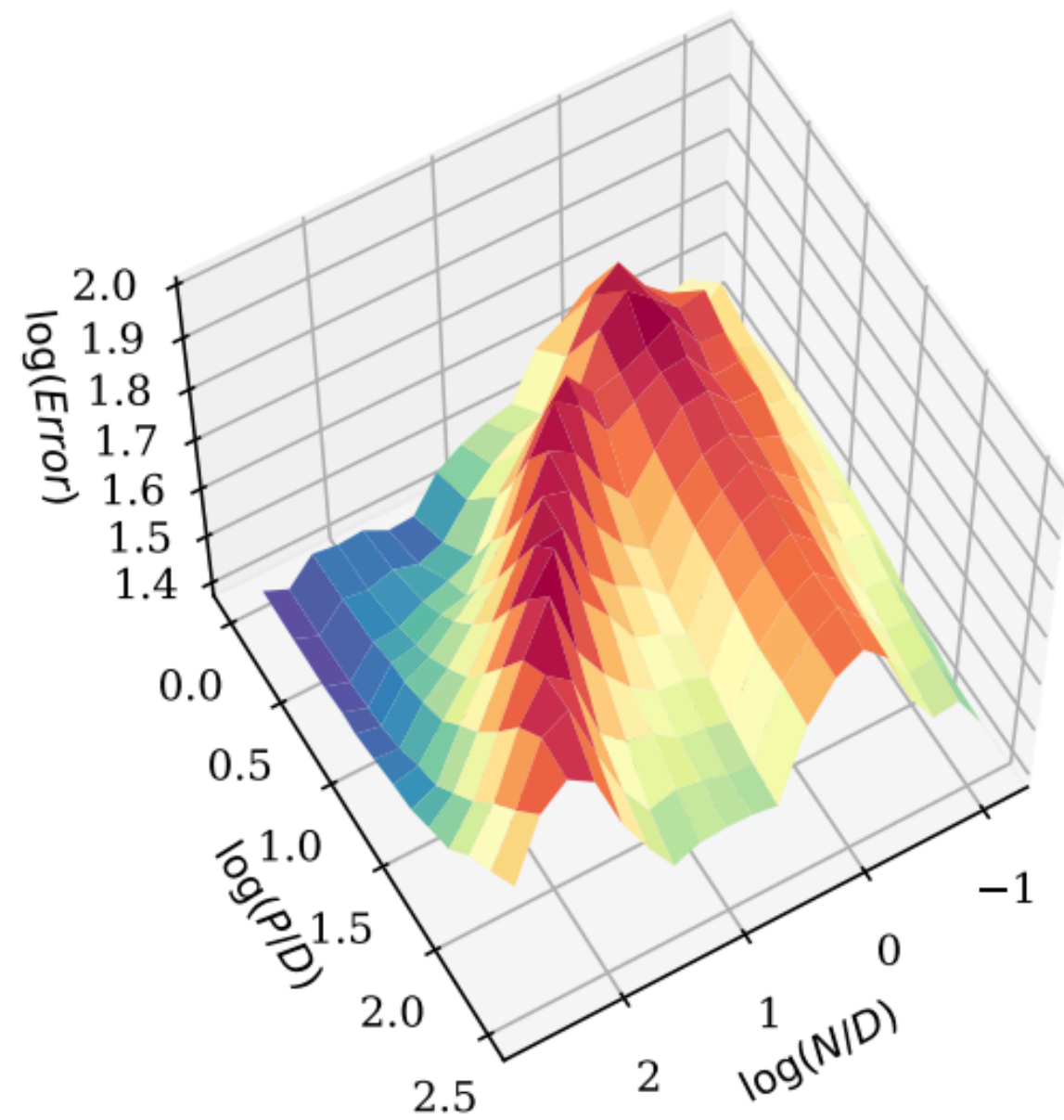


NOISELESS

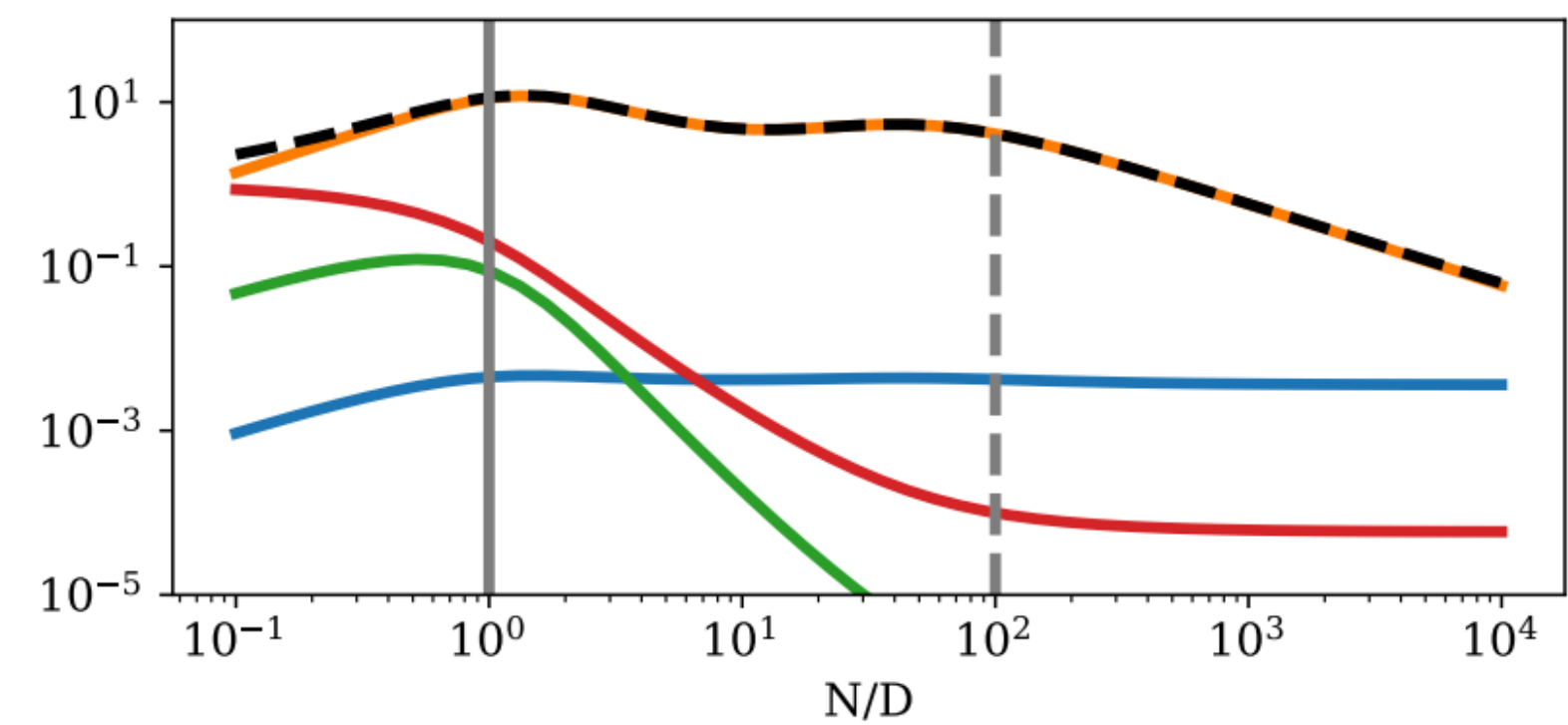
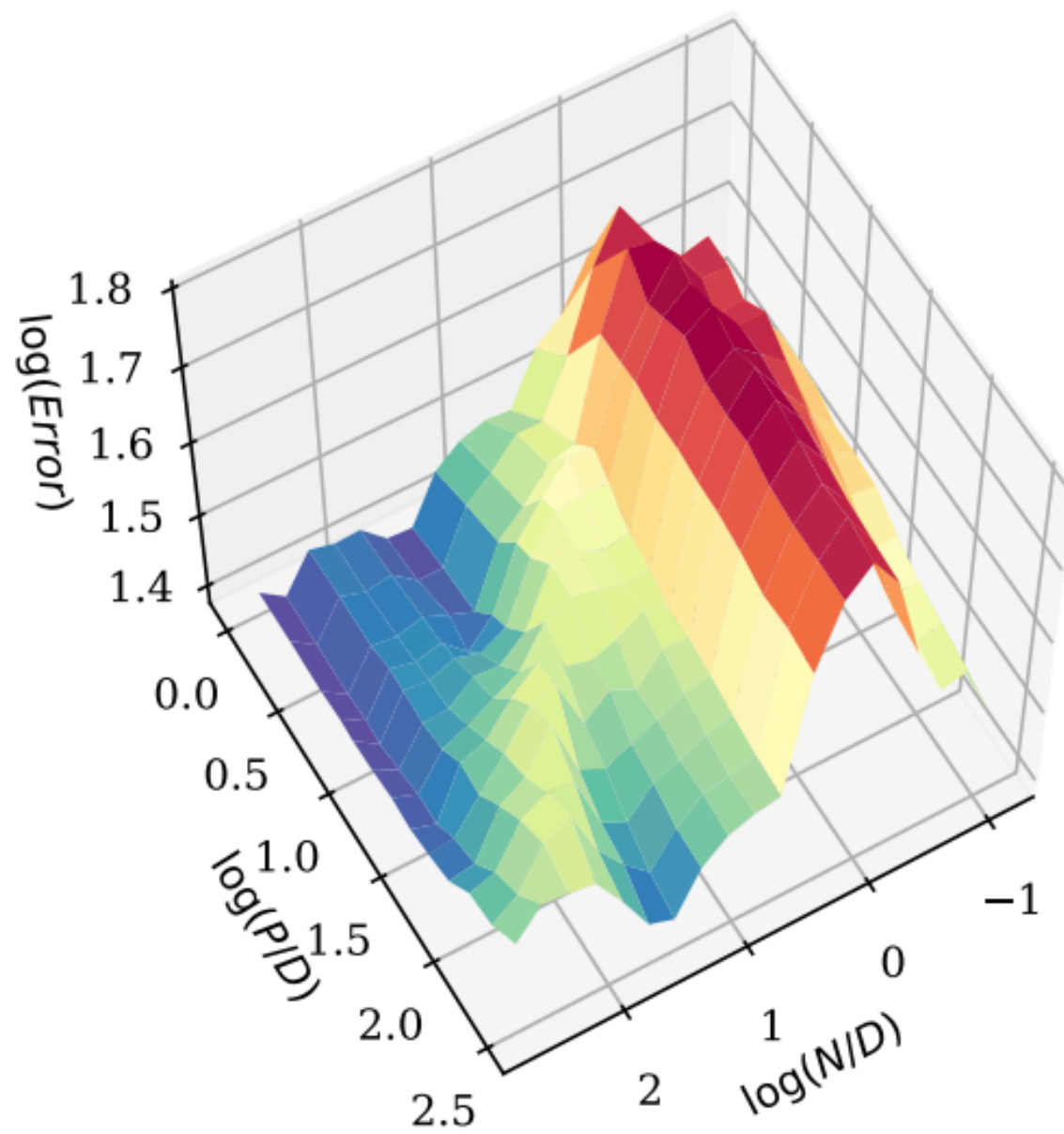
EFFECT OF REGULARISATION



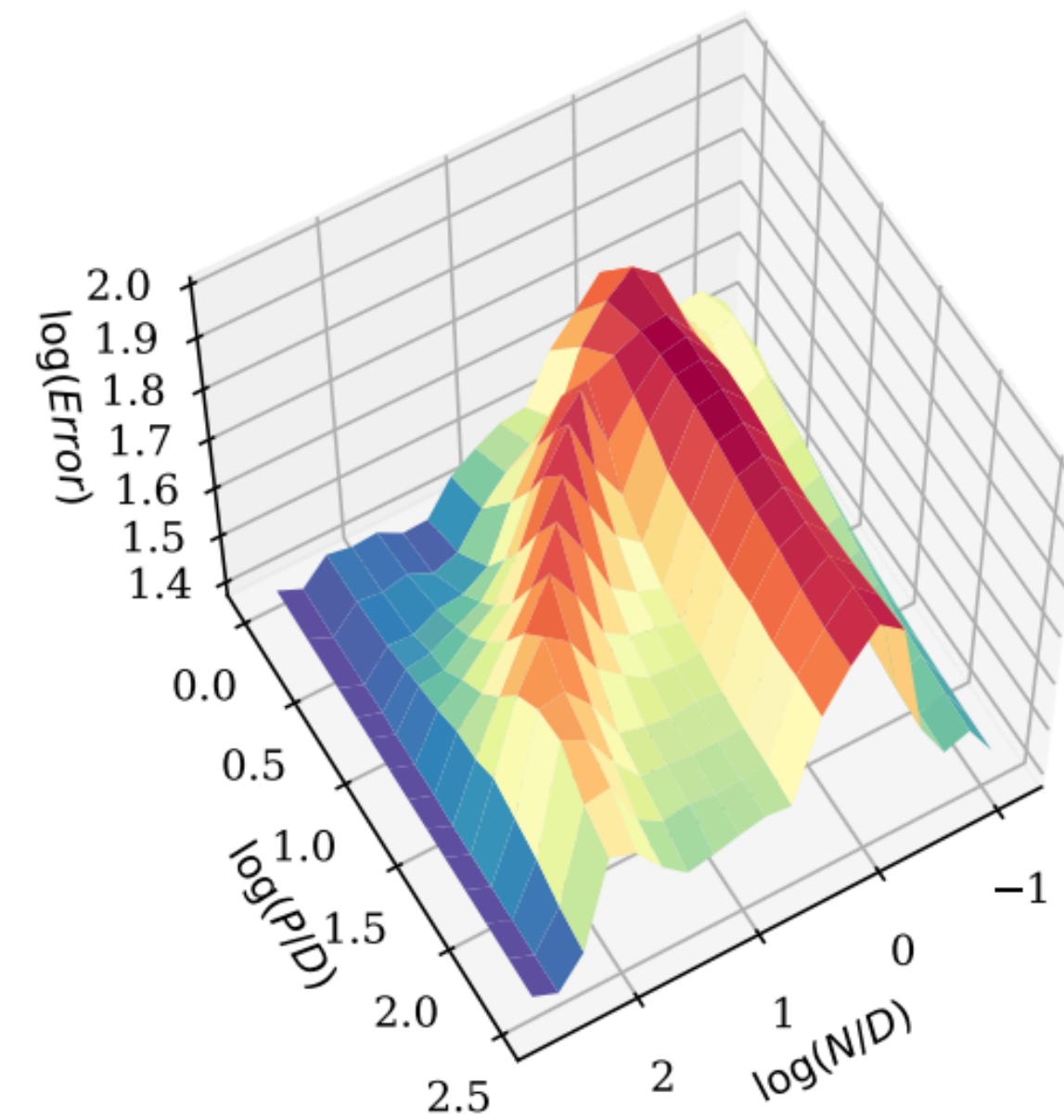
VANILLA



ENSEMBLING

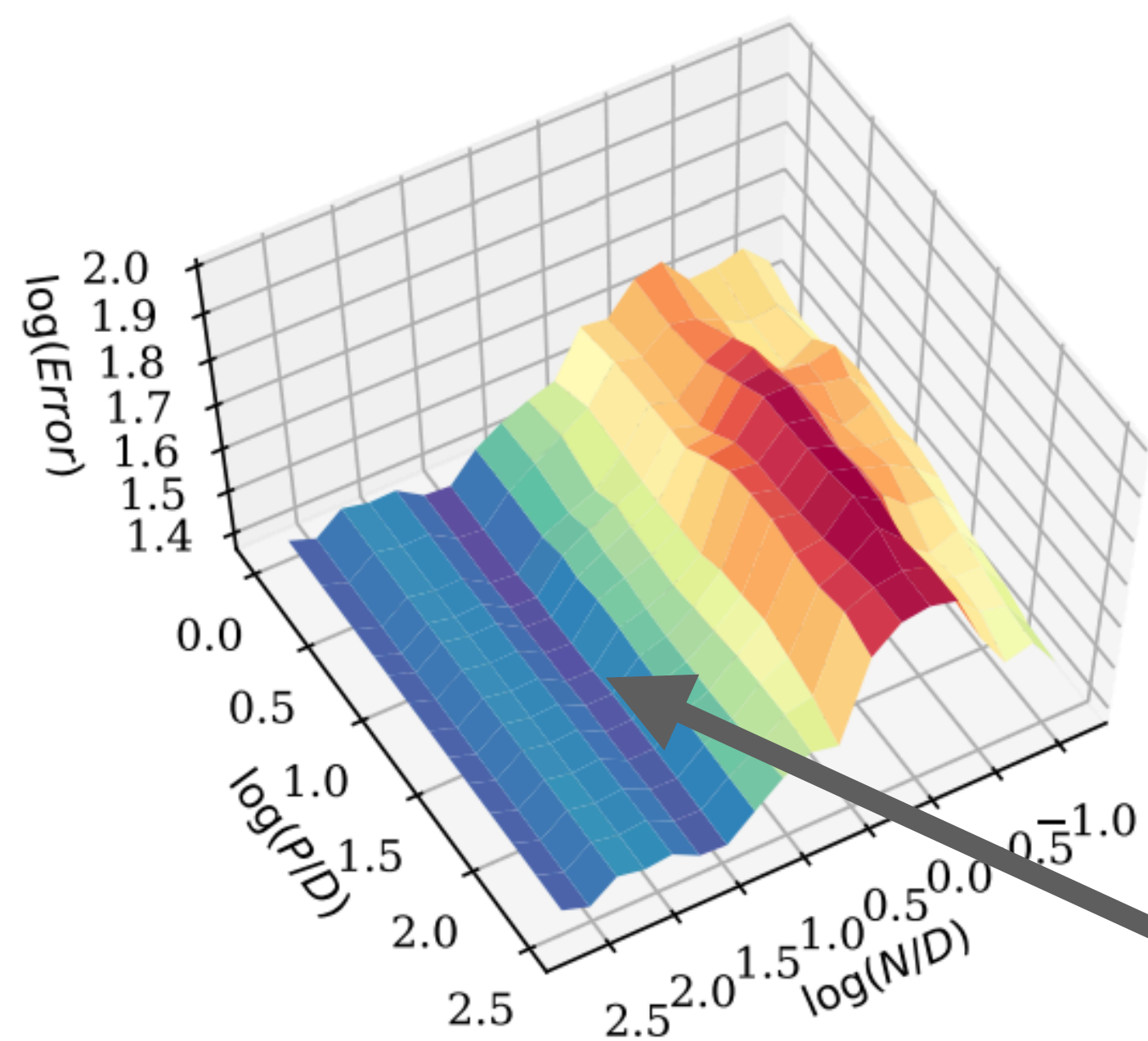


REGULARIZING

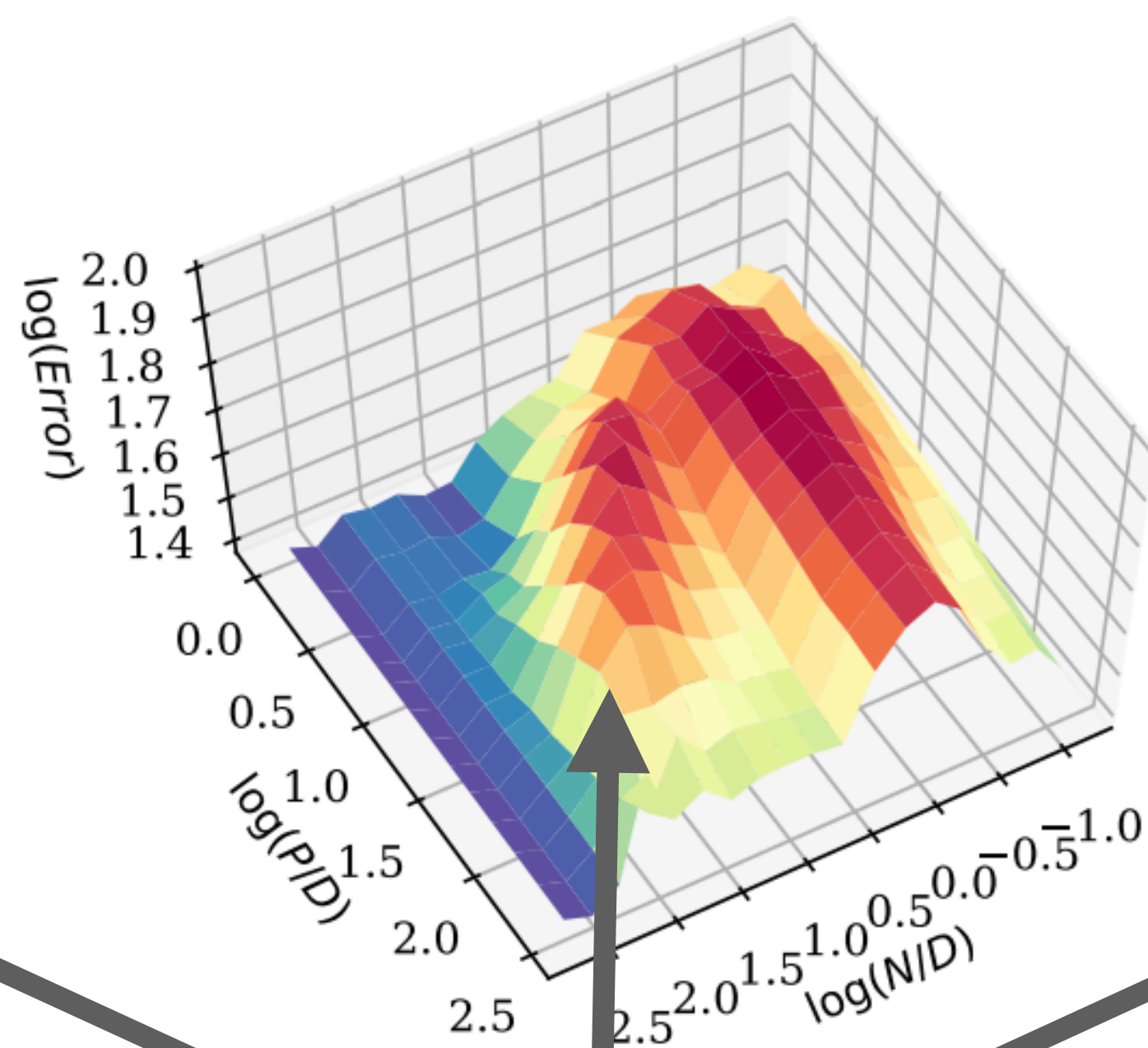


TIME DEPENDENCE

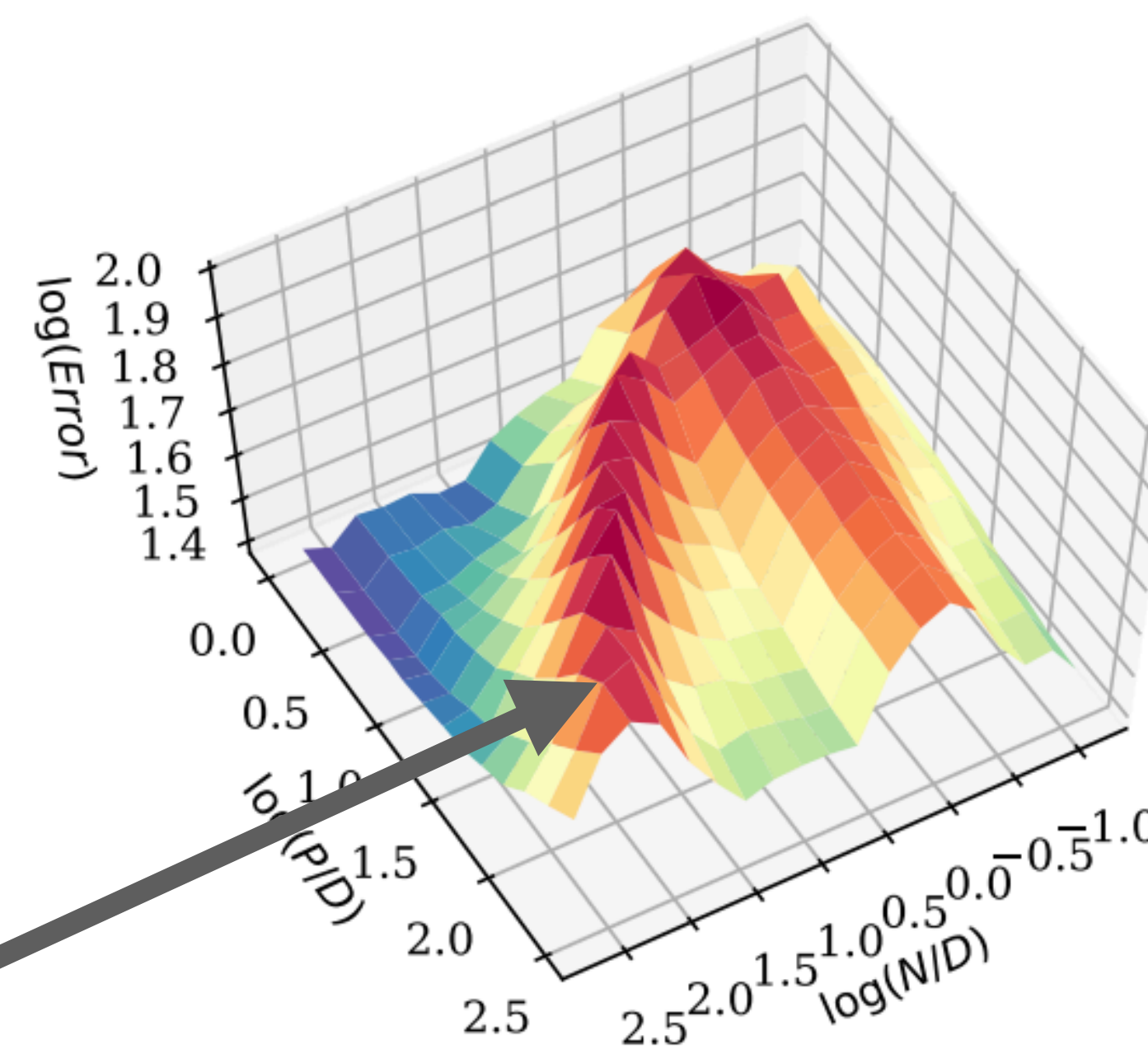
t=37 epochs



t=162 epochs

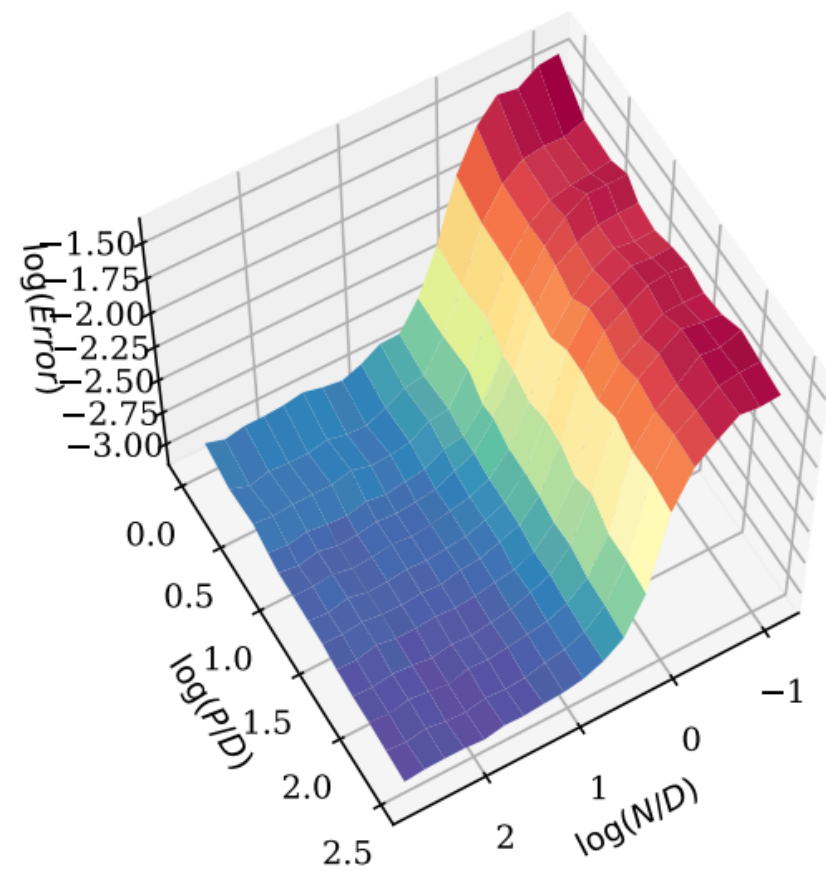


t=695 epochs

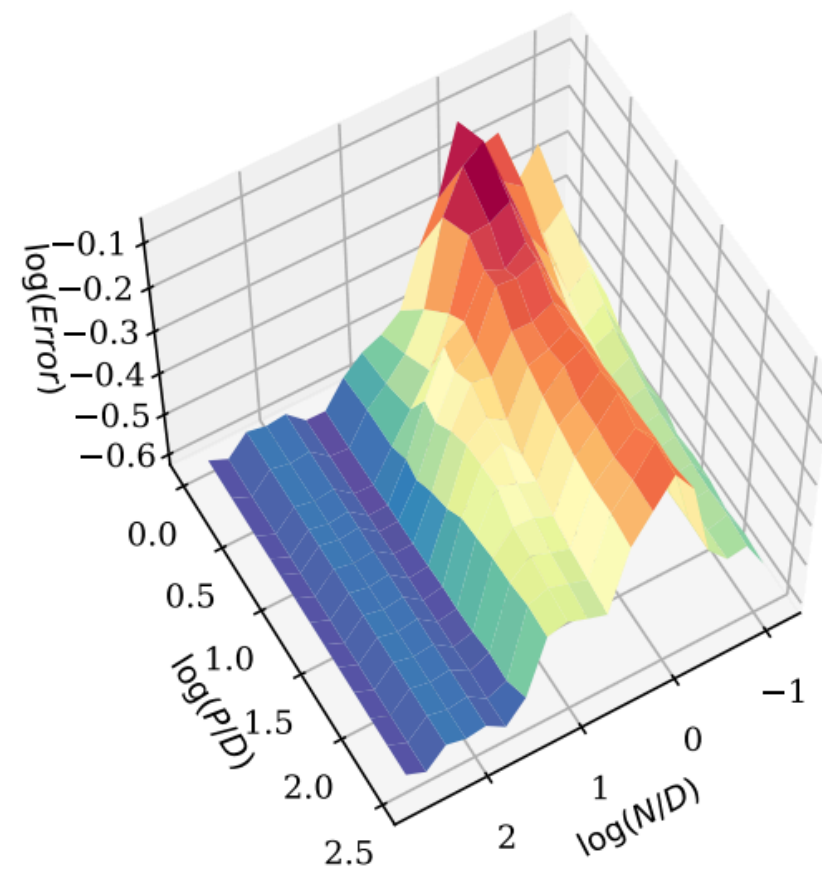


THE NONLINEAR PEAK FORMS AT LATE TIMES

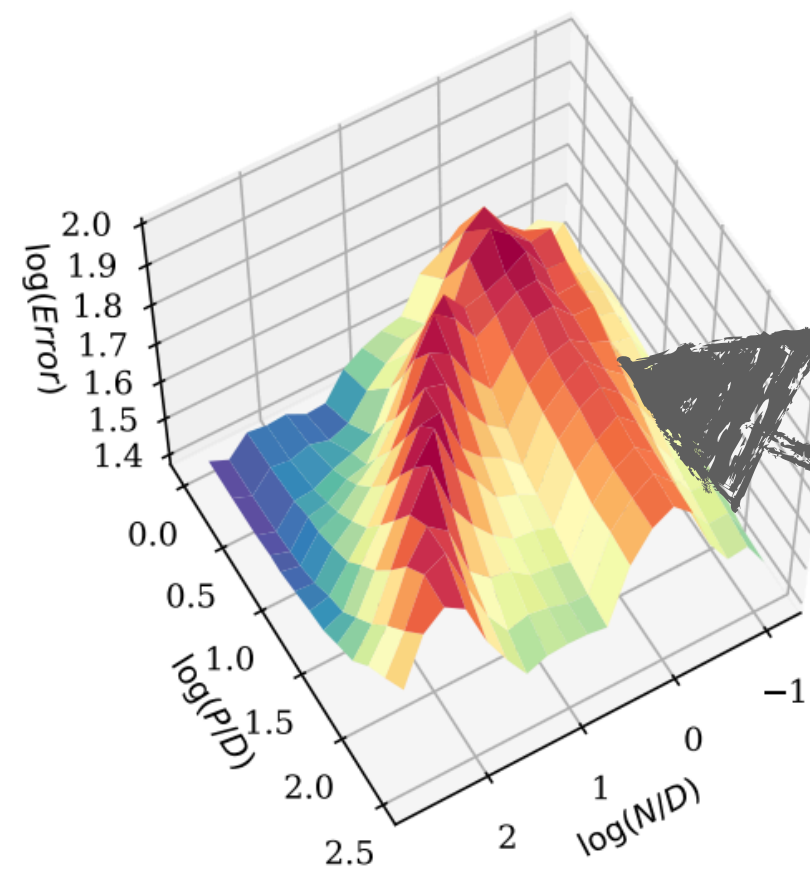
EFFECT OF NOISE AND NONLINEARITY



(a) Tanh, $SNR = \infty$

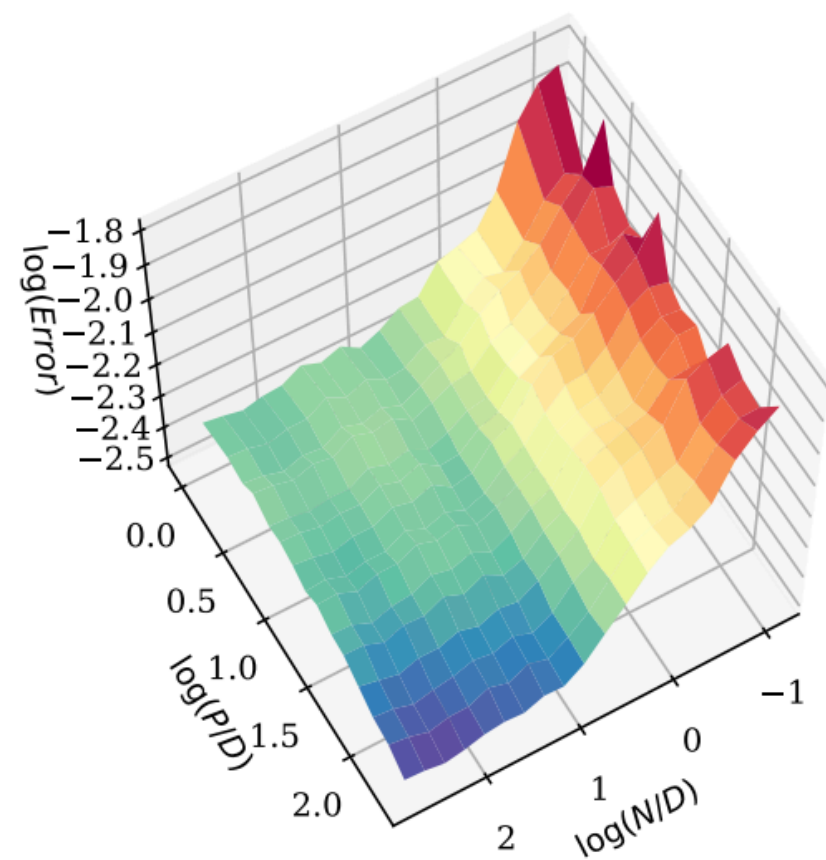


(b) Tanh, $SNR = 2$

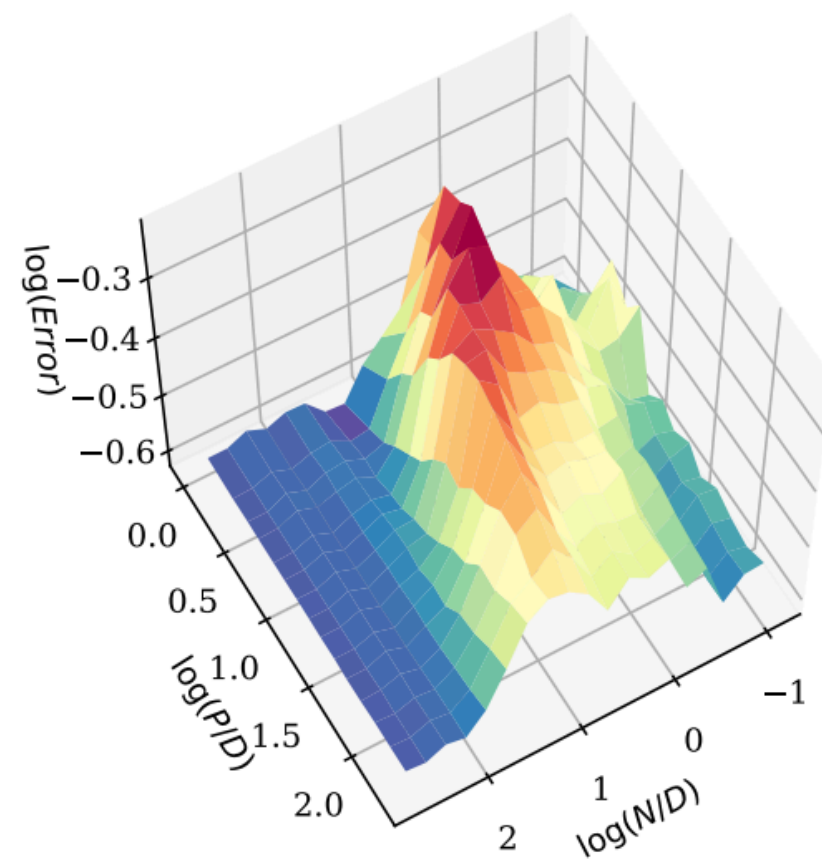


(c) Tanh, $SNR = 0.2$

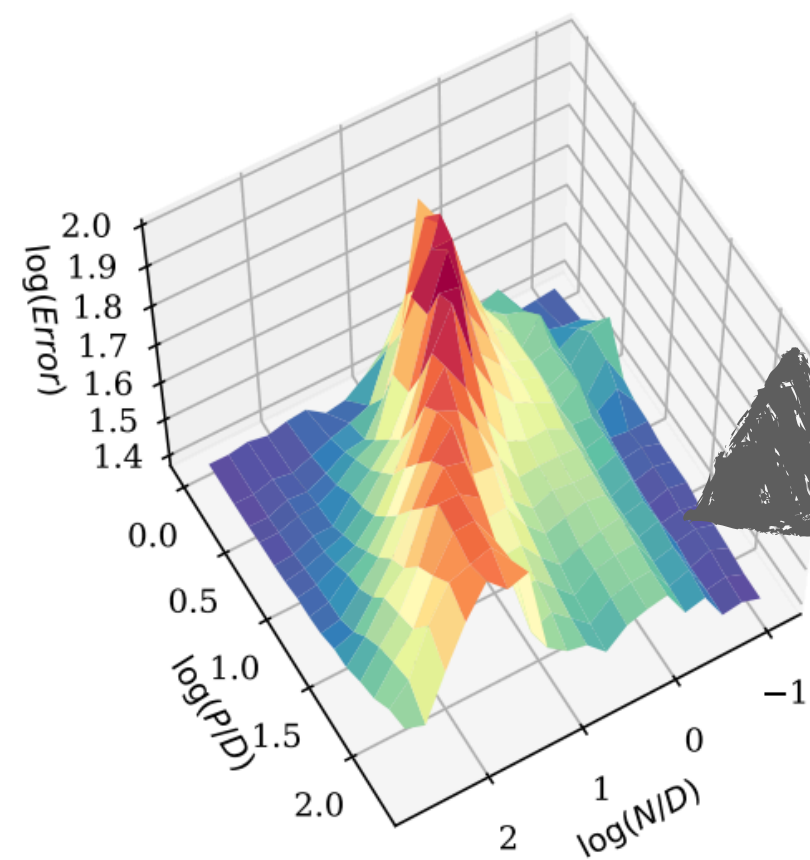
LINEAR PEAK IS WEAKER FOR RELU



(d) ReLU, $SNR = \infty$

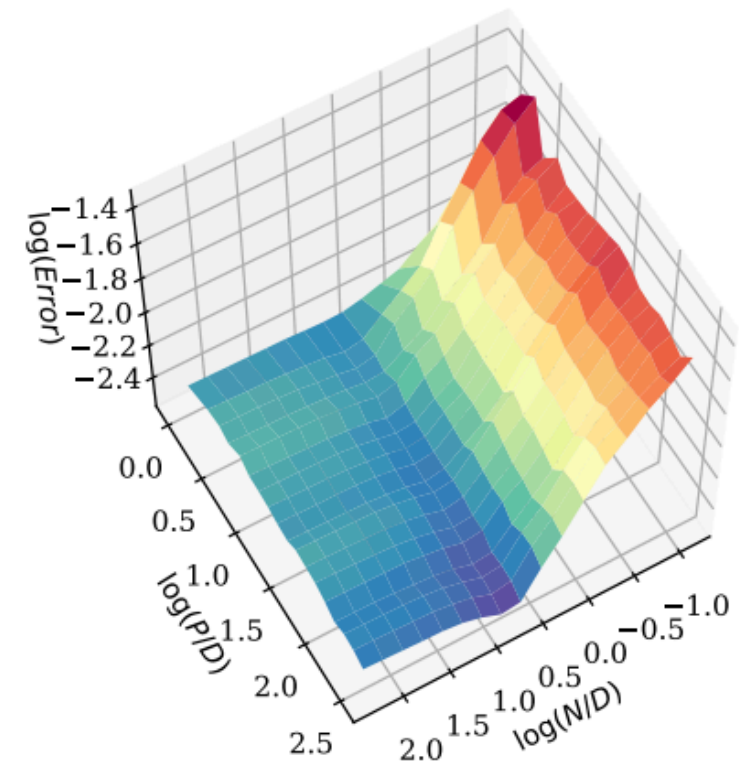


(e) ReLU, $SNR = 2$



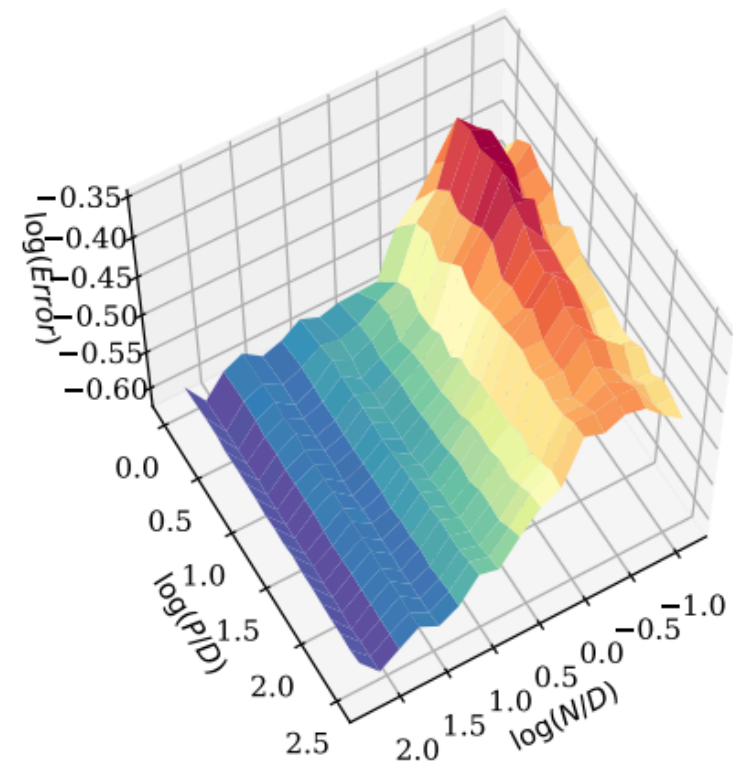
(f) ReLU, $SNR = 0.2$

STRUCTURED DATASETS



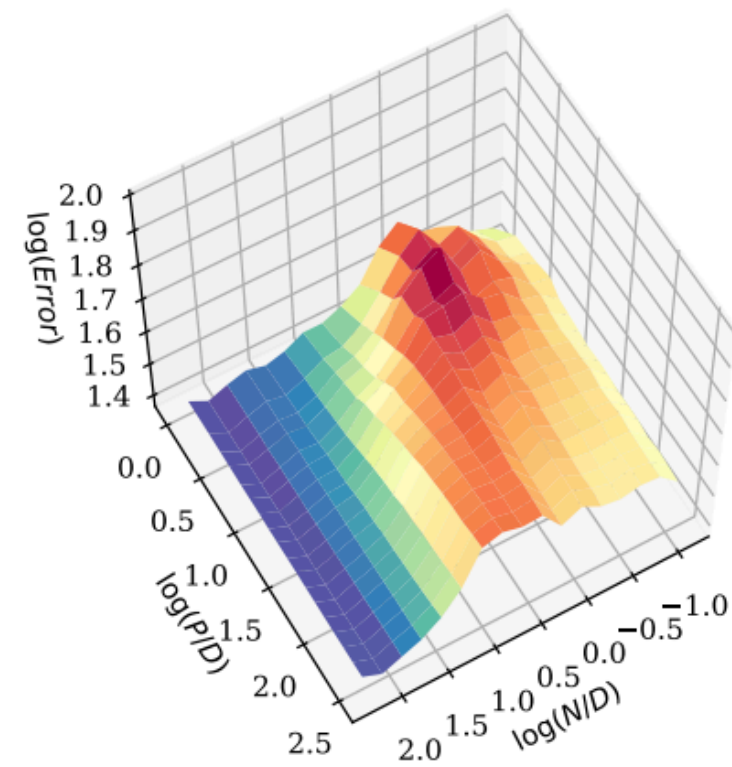
(a) MNIST, $SNR = \infty$

t=37 epochs



(b) MNIST, $SNR = 2$

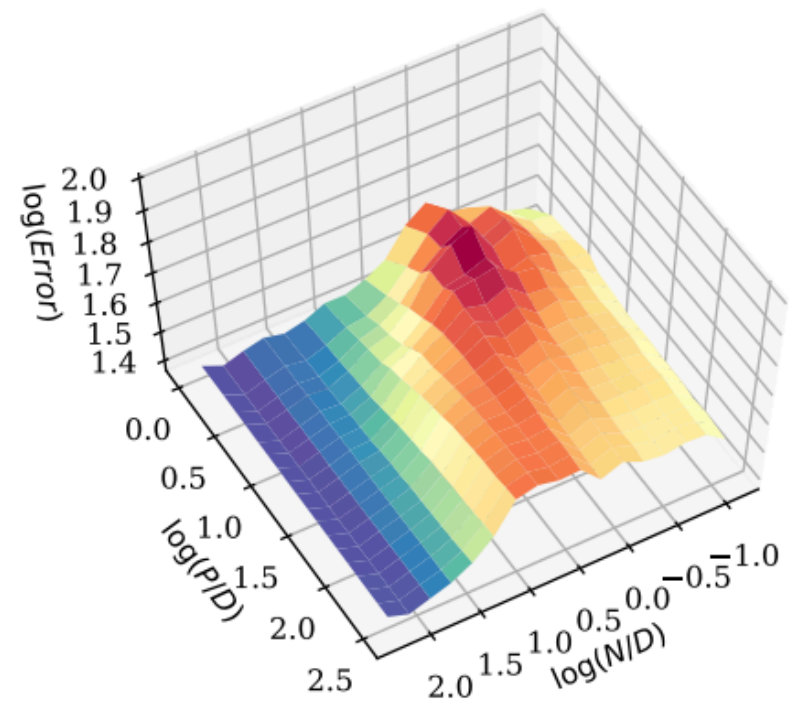
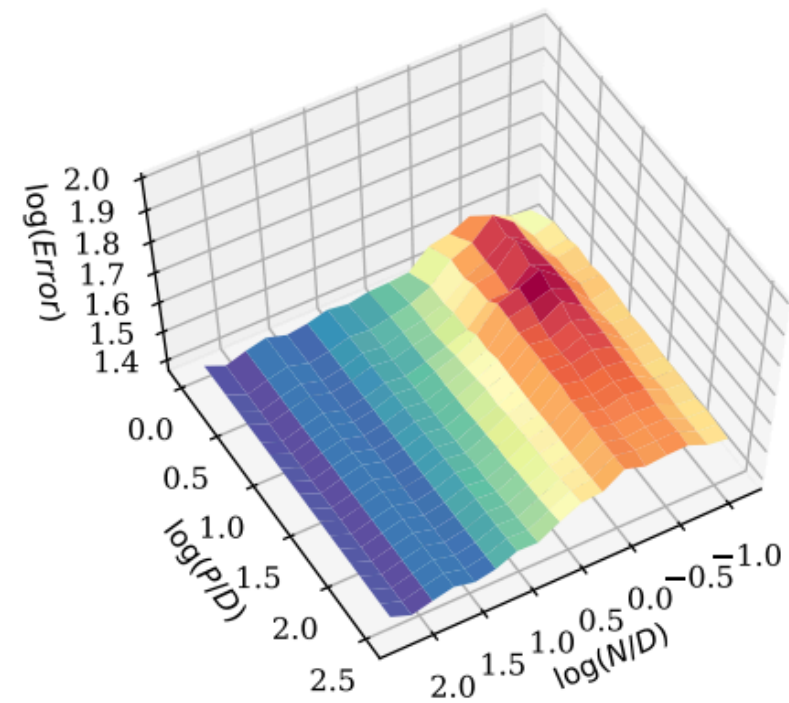
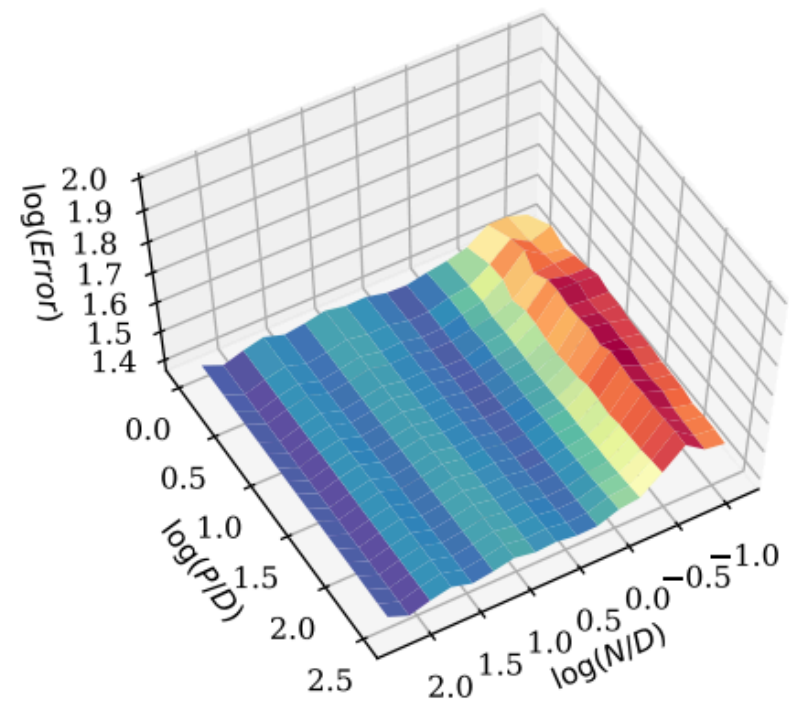
t=162 epochs



(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**