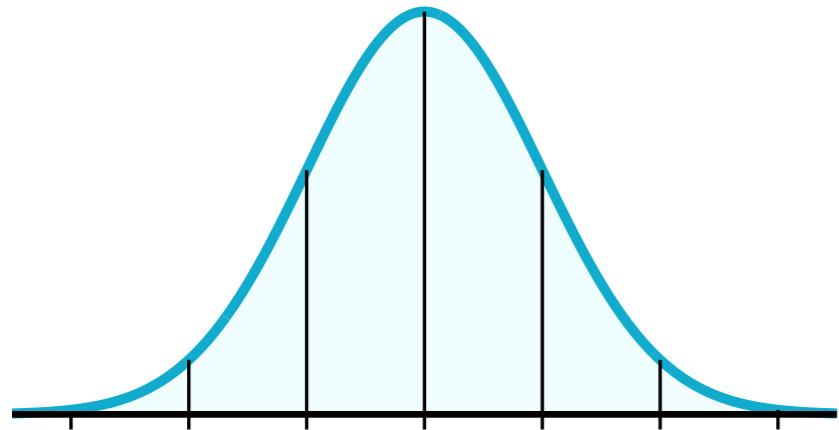


Attention-index models and how to solve them using tools for quadratic networks

Emanuele Troiani

Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian x_μ

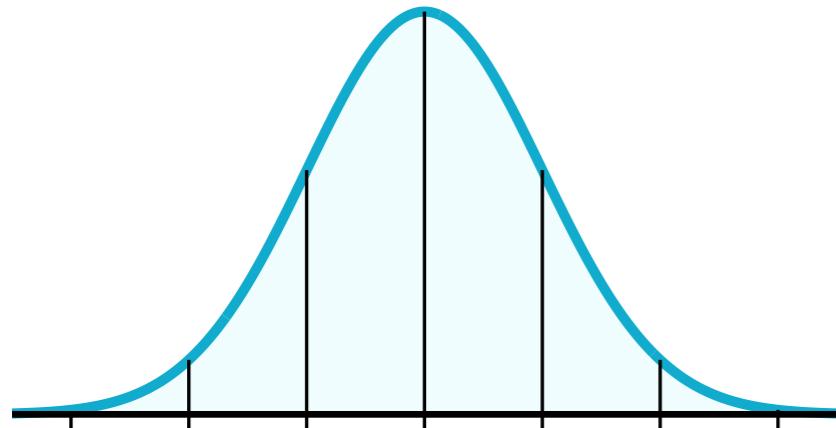


Projections
on $p = \mathcal{O}_d(1)$
random directions



Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian \mathbf{x}_μ



Projections
on $p = \mathcal{O}_d(1)$
random directions

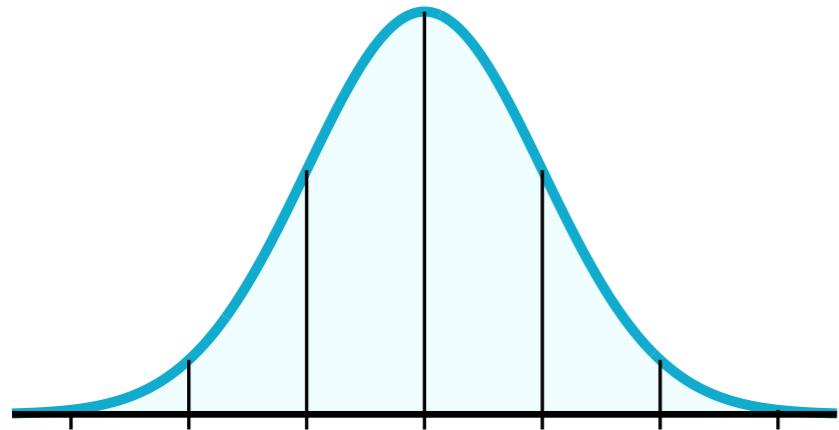


Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

1. Generate dataset $\mathcal{D} = \{\mathbf{x}_\mu, y_\mu\}_{\mu=1,\dots,n}$

$$y_\mu = g(w_1^\star \mathbf{x}_\mu^\top, \dots, w_p^\star \mathbf{x}_\mu^\top)$$

Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian \mathbf{x}_μ



Projections
on $p = \mathcal{O}_d(1)$
random directions



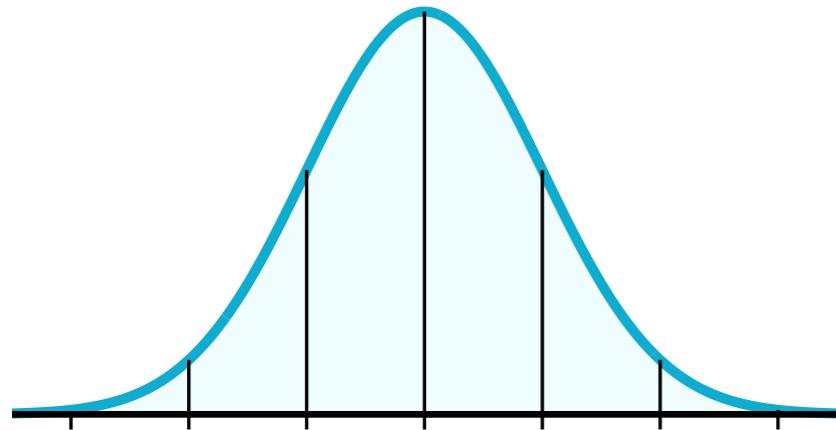
Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

1. Generate dataset $\mathcal{D} = \{\mathbf{x}_\mu, y_\mu\}_{\mu=1,\dots,n}$

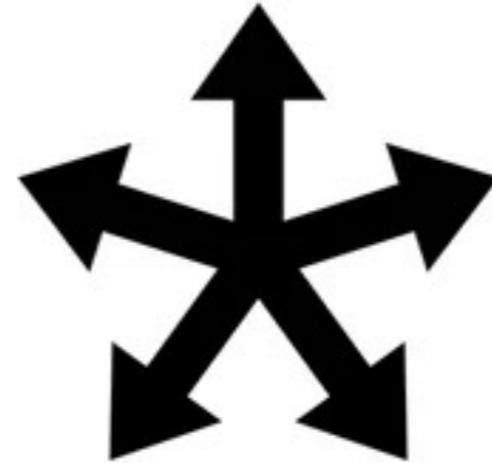
$$y_\mu = g(w_1^\star \mathbf{x}_\mu^\top, \dots, w_p^\star \mathbf{x}_\mu^\top)$$

2. Learn the function $x \mapsto y$

Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian \mathbf{x}_μ



Projections
on $p = \mathcal{O}_d(1)$
random directions



Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

1. Generate dataset $\mathcal{D} = \{\mathbf{x}_\mu, y_\mu\}_{\mu=1,\dots,n}$

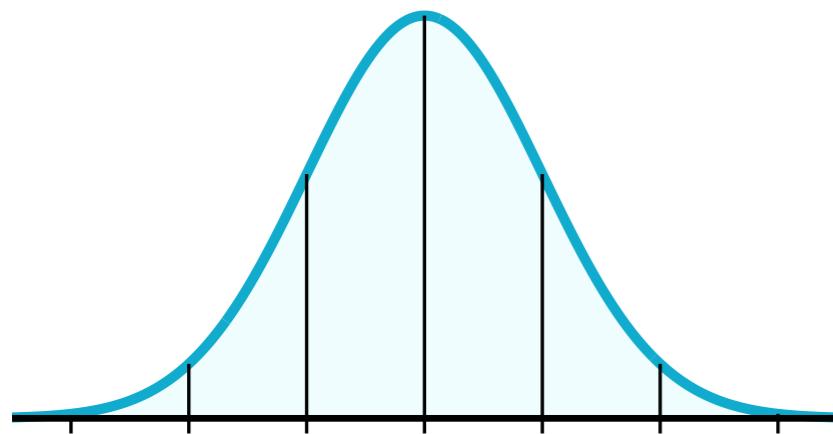
$$y_\mu = g(w_1^\star \mathbf{x}_\mu^\top, \dots, w_p^\star \mathbf{x}_\mu^\top)$$

2. Learn the function $x \mapsto y$

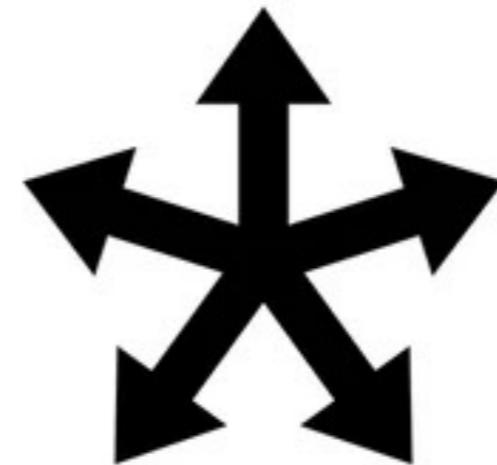


Bayes Optimal

Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian x_μ



Projections
on $p = \mathcal{O}_d(1)$
random directions

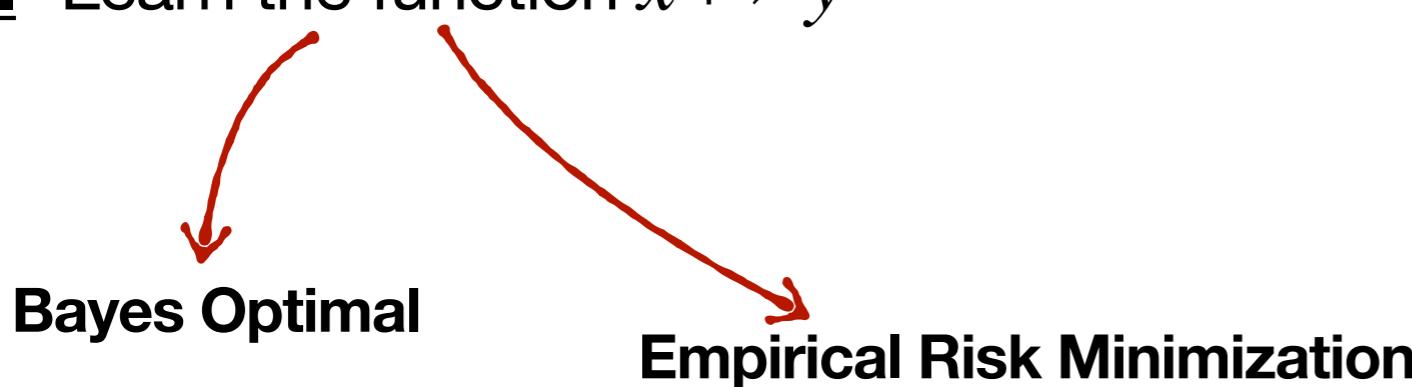


Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

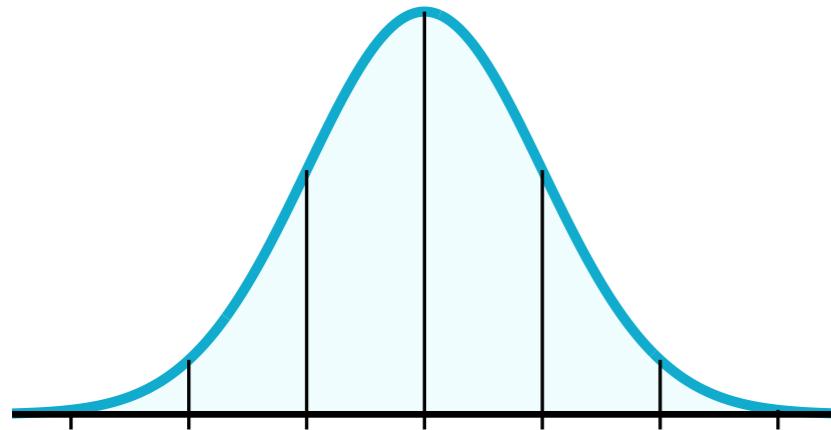
1. Generate dataset $\mathcal{D} = \{x_\mu, y_\mu\}_{\mu=1,\dots,n}$

$$y_\mu = g(w_1^\star x_\mu^\top, \dots, w_p^\star x_\mu^\top)$$

2. Learn the function $x \mapsto y$



Multi-index models...



Input data

$d \gg 1$ dimensional Gaussian x_μ

Projections
on $p = \mathcal{O}_d(1)$
random directions

Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

1. Generate dataset $\mathcal{D} = \{x_\mu, y_\mu\}_{\mu=1,\dots,n}$

$$y_\mu = g(w_1^\star x_\mu^\top, \dots, w_p^\star x_\mu^\top)$$

2. Learn the function $x \mapsto y$



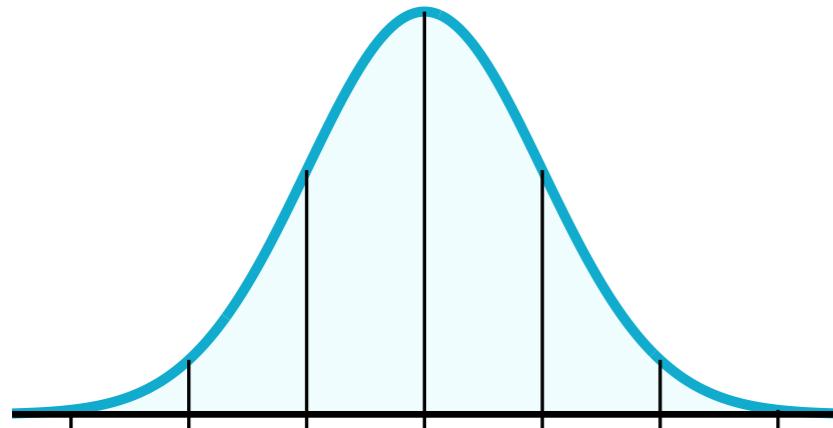
Bayes Optimal

Empirical Risk Minimization

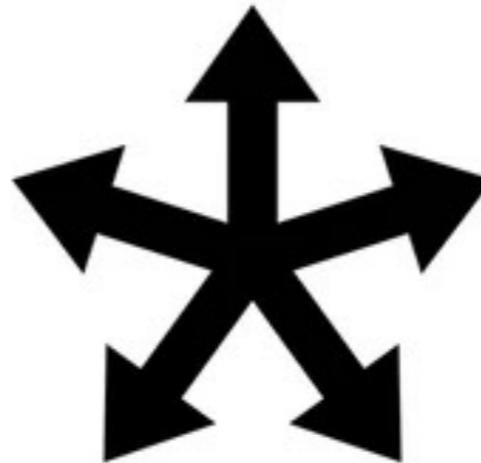


Many functions in this class

Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian x_μ



Projections
on $p = \mathcal{O}_d(1)$
random directions



Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

1. Generate dataset $\mathcal{D} = \{x_\mu, y_\mu\}_{\mu=1,\dots,n}$

$$y_\mu = g(w_1^\star x_\mu^\top, \dots, w_p^\star x_\mu^\top)$$

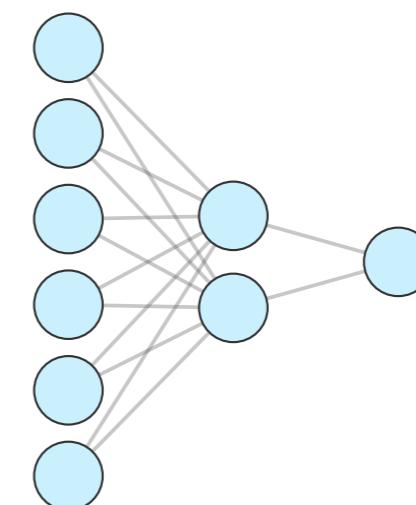
2. Learn the function $x \mapsto y$



Bayes Optimal

Empirical Risk Minimization

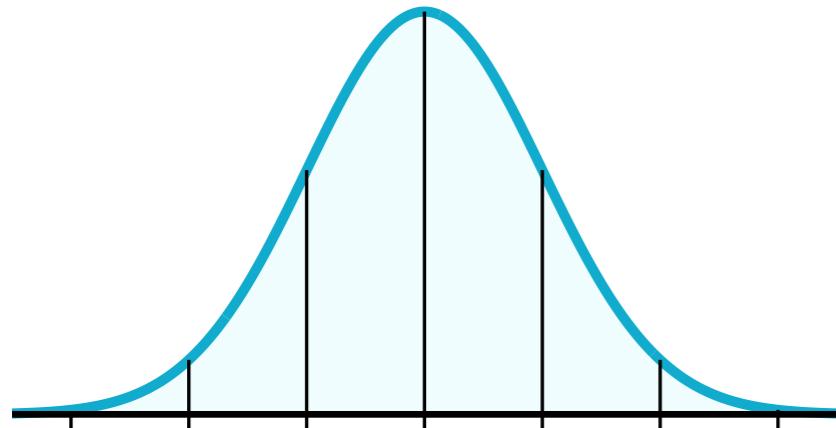
Many functions in this class



Narrow
two-layer networks...

$$y_\mu = \sum_{i=1}^p \sigma(w_i^\top x_\mu)$$

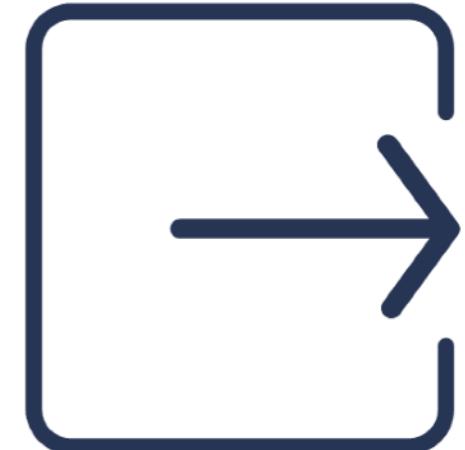
Multi-index models...



Input data
 $d \gg 1$ dimensional Gaussian x_μ



Projections
on $p = \mathcal{O}_d(1)$
random directions



Link function
using $g : \mathbb{R}^p \rightarrow \mathbb{R}$

1. Generate dataset $\mathcal{D} = \{x_\mu, y_\mu\}_{\mu=1,\dots,n}$

$$y_\mu = g(w_1^\star x_\mu^\top, \dots, w_p^\star x_\mu^\top)$$

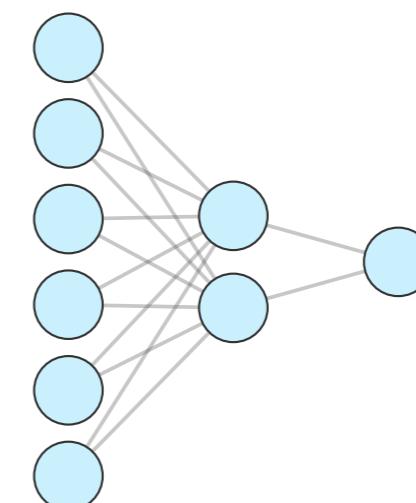
2. Learn the function $x \mapsto y$



Bayes Optimal

Empirical Risk Minimization

Many functions in this class

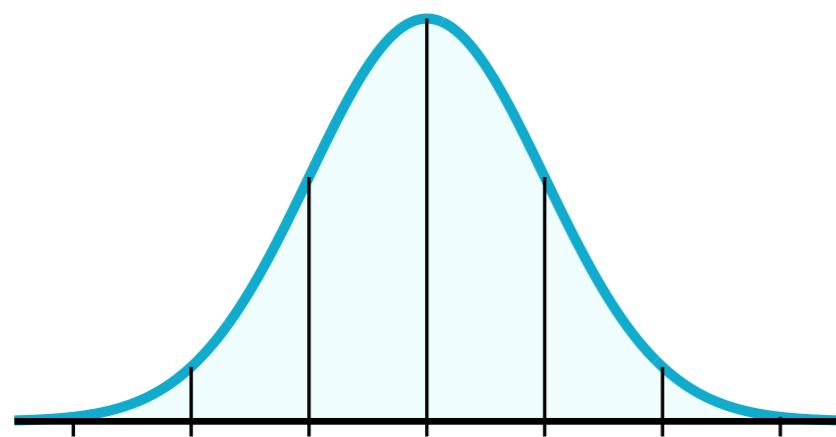


Narrow
two-layer networks...

$$y_\mu = \sum_{i=1}^p \sigma(w_i^\top x_\mu)$$

But also **Finite-rank**
attention networks!

Multi-index models... are ideal for attention



Input data

Collections of Gaussian tokens $X_\mu \in \mathbb{R}^{T \times d}$, $T = \mathcal{O}_d(1)$, $d \gg 1$

Projections

on $p = \mathcal{O}_d(1)$ random directions

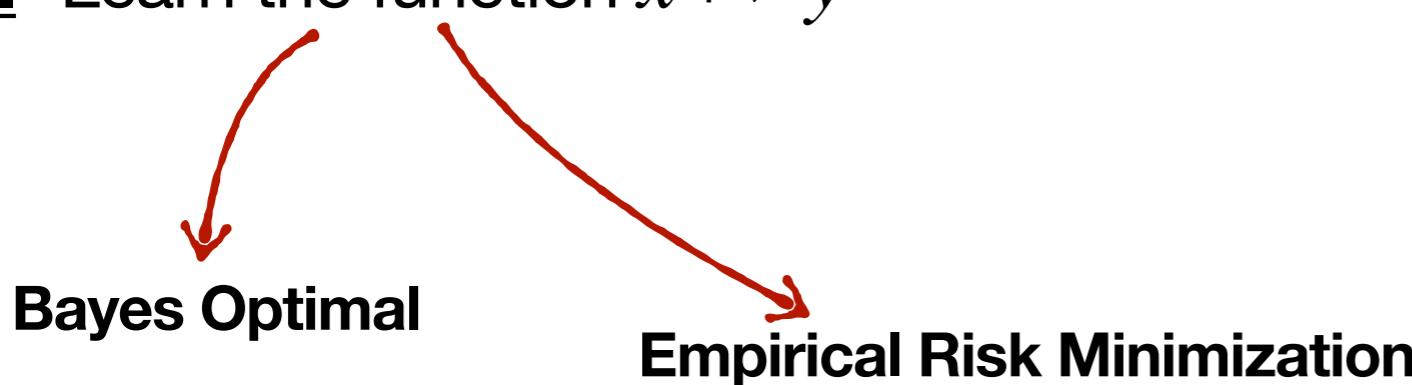
Link function

using $g : \mathbb{R}^p \rightarrow \mathbb{R}^{T \times T}$

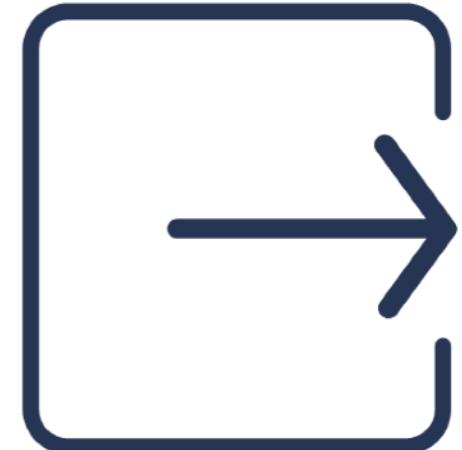
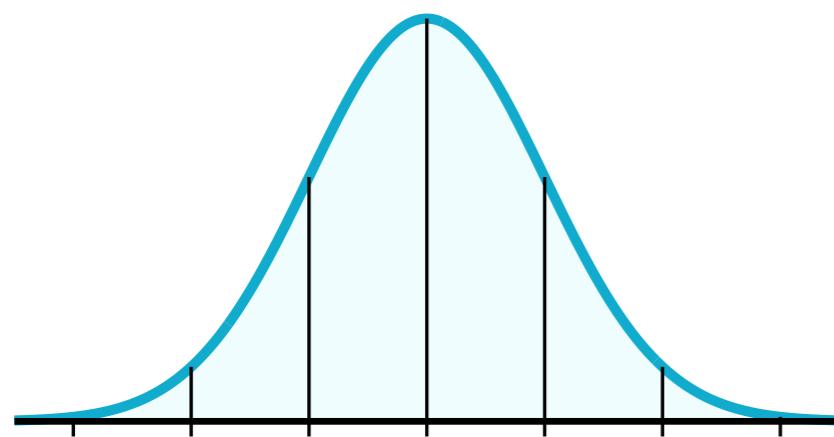
1. Generate dataset $\mathcal{D} = \{X_\mu, y_\mu\}_{\mu=1,\dots,n}$

$$y_\mu = g(w_1^\star X_\mu^\top, \dots, w_p^\star X_\mu^\top)$$

2. Learn the function $x \mapsto y$



Multi-index models... are ideal for attention



Input data

Collections of Gaussian
tokens $X_\mu \in \mathbb{R}^{T \times d}$, $T = \mathcal{O}_d(1)$, $d \gg 1$

Projections

on $p = \mathcal{O}_d(1)$
random directions

Link function

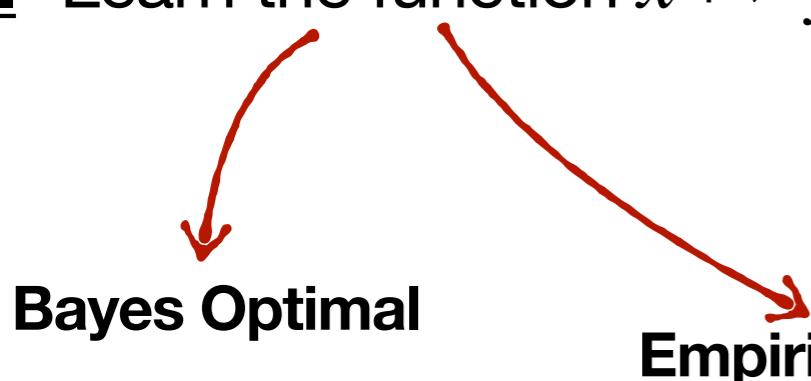
using $g : \mathbb{R}^p \rightarrow \mathbb{R}^{T \times T}$

1. Generate dataset $\mathcal{D} = \{X_\mu, y_\mu\}_{\mu=1,\dots,n}$
 $y_\mu = g(w_1^\star X_\mu^\top, \dots, w_p^\star X_\mu^\top)$

$$y_\mu = \sigma \left(X_\mu \sum_{i=1}^p w_i w_i^\top X_\mu^\top \right)$$

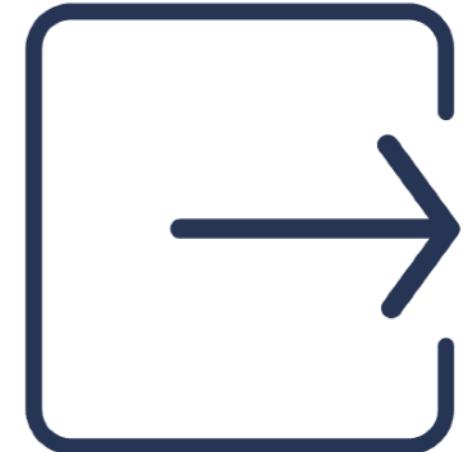
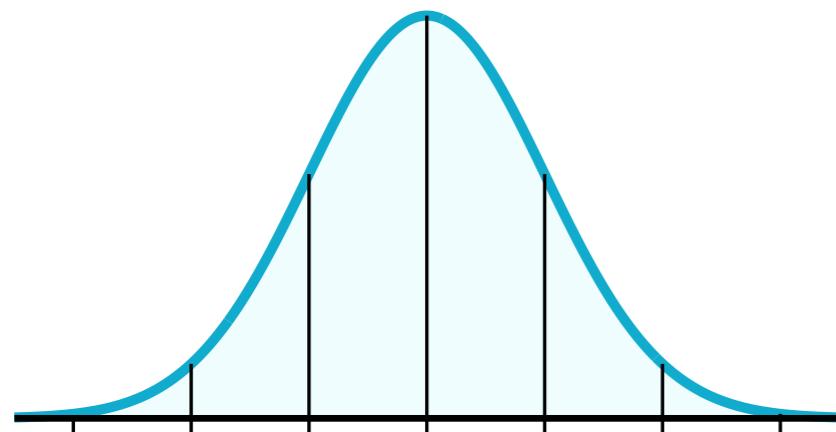
Tied keys-queries

2. Learn the function $x \mapsto y$



[H. Cui, F. Behrens, F. Krzakala, L. Zdeborová, '24]
[L. Arnaboldi, B. Loureiro, L. Stephan, F. Krzakala, L. Zdeborová, '25]
[H. Cui, '25]

Multi-index models... are ideal for attention



Input data

Collections of Gaussian
tokens $X_\mu \in \mathbb{R}^{T \times d}$, $T = \mathcal{O}_d(1)$, $d \gg 1$

Projections

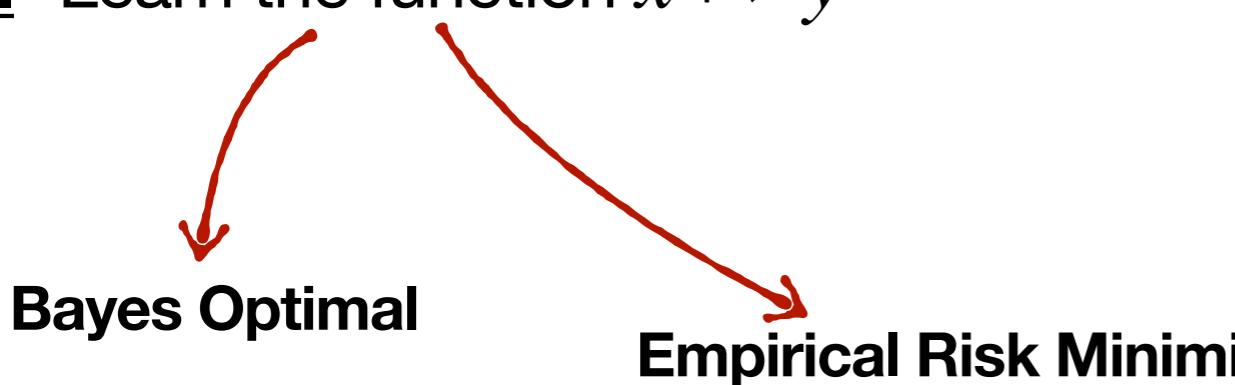
on $p = \mathcal{O}_d(1)$
random directions

Link function

using $g : \mathbb{R}^p \rightarrow \mathbb{R}^{T \times T}$

1. Generate dataset $\mathcal{D} = \{X_\mu, y_\mu\}_{\mu=1,\dots,n}$
 $y_\mu = g(w_1^\star X_\mu^\top, \dots, w_p^\star X_\mu^\top)$

2. Learn the function $x \mapsto y$



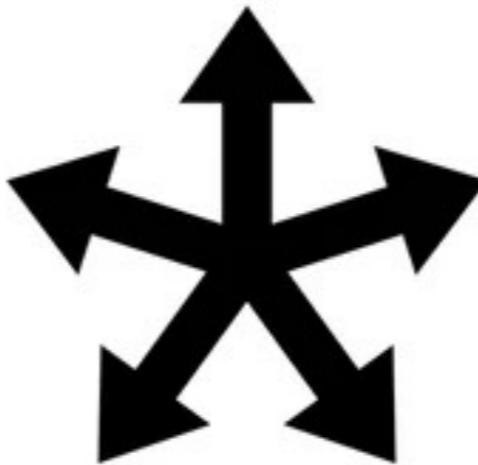
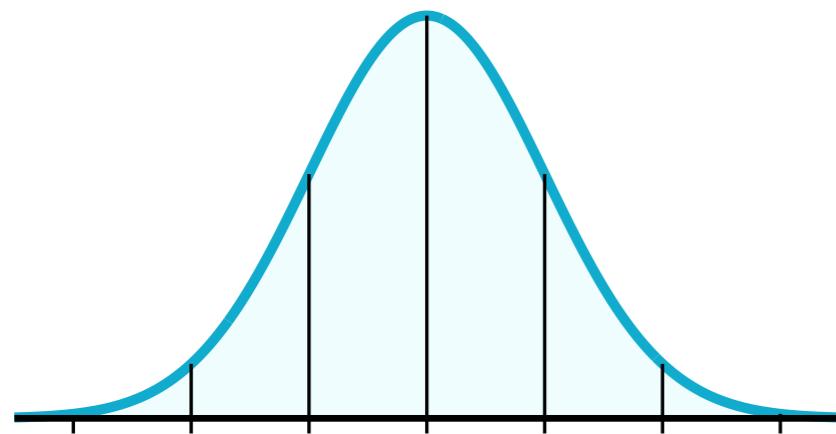
$$y_\mu = \sigma \left(X_\mu \sum_{i=1}^p w_i w_i^\top X_\mu^\top \right)$$

Tied keys-queries

[H. Cui, F. Behrens, F. Krzakala, L. Zdeborová, '24]
[L. Arnaboldi, B. Loureiro, L. Stephan, F. Krzakala, L. Zdeborová, '25]
[H. Cui, '25]

Also applies to
multi-layer attention!

Multi-index models... are ideal for attention



Input data

Collections of Gaussian tokens $X_\mu \in \mathbb{R}^{T \times d}$, $T = \mathcal{O}_d(1)$, $d \gg 1$

Projections

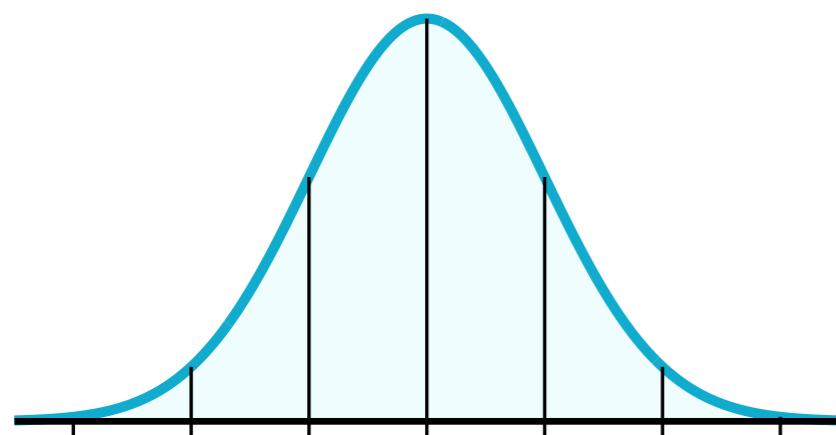
$P = p_1 + \dots + p_L = \mathcal{O}_d(1)$
random directions

Link function

using $g : \mathbb{R}^P \rightarrow \mathbb{R}^{T \times T}$

$$y_\mu = \sigma \left(X_\mu^L \sum_{i=1}^{p_L} \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_\mu^{L\top} \right) \quad X_\mu^{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\mu^\ell \sum_{i=1}^{p_\ell} \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\mu^{\ell\top} \right) \right] X_\mu^\ell$$

Multi-index models... are ideal for attention



Input data

Collections of Gaussian tokens $X_\mu \in \mathbb{R}^{T \times d}$, $T = \mathcal{O}_d(1)$, $d \gg 1$

Projections

$P = p_1 + \dots + p_L = \mathcal{O}_d(1)$ random directions

Link function

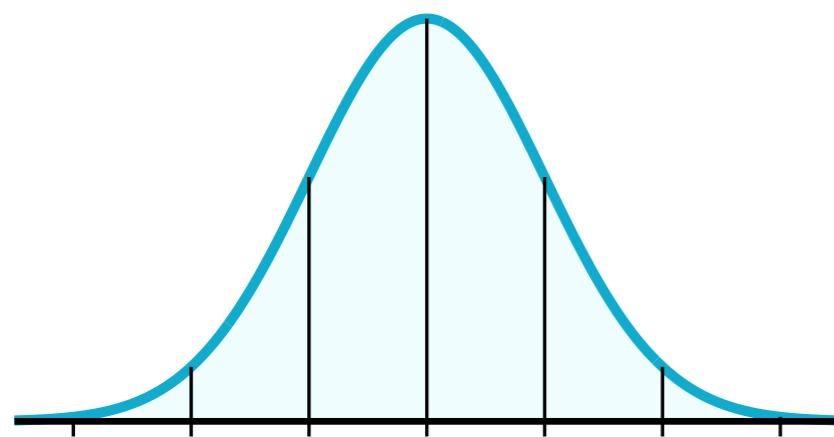
using $g : \mathbb{R}^P \rightarrow \mathbb{R}^{T \times T}$

$$y_\mu = \sigma \left(X_\mu^L \sum_{i=1}^{p_L} \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_\mu^{L\top} \right) \quad X_\mu^{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\mu^\ell \sum_{i=1}^{p_\ell} \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\mu^{\ell\top} \right) \right] X_\mu^\ell$$

$$y_\mu = g \left(X_\mu \sum_{i=1}^{p_1} \mathbf{w}_i^1 \mathbf{w}_i^{1\top} X_\mu^\top, \dots, X_\mu \sum_{i=1}^{p_L} \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_\mu^\top \right)$$

Exact g found by recursion

Multi-index models... are ideal for attention



Input data

Collections of Gaussian tokens $X_\mu \in \mathbb{R}^{T \times d}$, $T = \mathcal{O}_d(1)$, $d \gg 1$

Projections

$P = p_1 + \dots + p_L = \mathcal{O}_d(1)$
random directions

Link function

using $g : \mathbb{R}^P \rightarrow \mathbb{R}^{T \times T}$

$$y_\mu = \sigma \left(X_\mu^L \sum_{i=1}^{p_L} \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_\mu^{L\top} \right) \quad X_\mu^{\ell+1} = \left[c\mathbf{1}_T + \left(\dots \sum_{i=1}^{p_\ell} \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\mu^{\ell\top} \right) \right]_{T \times T}$$

$$y_\mu = g \left(X_\mu \sum_{i=1}^{p_1} \mathbf{w}_i^1 \mathbf{w}_i^{1\top} X_\mu^\top, \dots, X_\mu \sum_{i=1}^{p_L} \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_\mu^\top \right)$$

Exact g found by recursion

Also multiple heads,
untied keys-queries, ecc...

[E. Troiani, H. Cui, Y. Dandi, F. Krzakala, L. Zdeborová, '25]

yet these are **finite rank...**
attention networks!

Multi-index models... are limited

Can we remove the $P \ll d$ assumption?

Multi-index models... are limited

Can we remove the $P \ll d$ assumption?



IDEA: the attention is an index

Define a new class of models

Multi-index models... are limited

Can we remove the $P \ll d$ assumption?



IDEA: the attention is an index

Define a new class of models

$$y_\mu = g(\boldsymbol{w}^\top \boldsymbol{x}_\mu)$$

“Linear” index model

$$\boldsymbol{x}_\mu \in \mathbb{R}^d$$

$$d \rightarrow \infty$$

Multi-index models... are limited

Can we remove the $P \ll d$ assumption?



IDEA: the attention is an index

Define a new class of models

$$y_\mu = g(w^\top x_\mu)$$

“Linear” index model

$$\begin{aligned}x_\mu &\in \mathbb{R}^d \\ d &\rightarrow \infty\end{aligned}$$

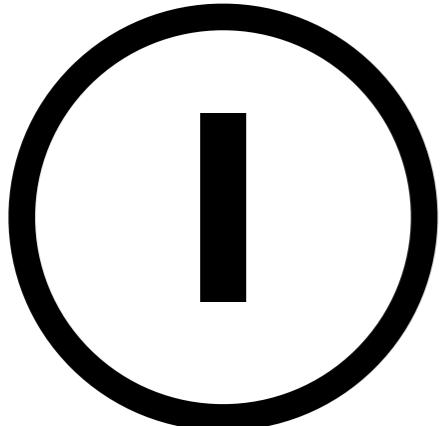


$$y_\mu = g(X_\mu^\top A X_\mu)$$

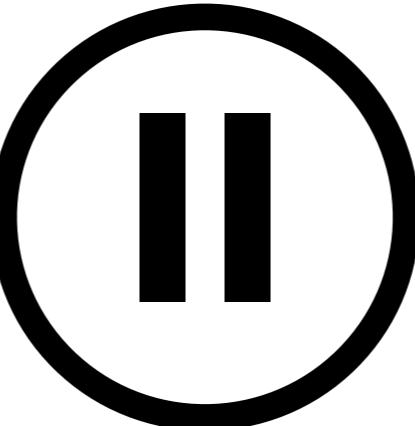
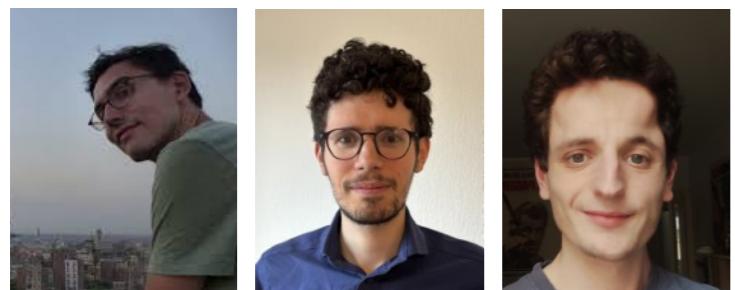
“Attention” index model

$$\begin{aligned}X_\mu &\in \mathbb{R}^{T \times d} & \text{rank}(A) &= \rho d \\ d &\rightarrow \infty\end{aligned}$$

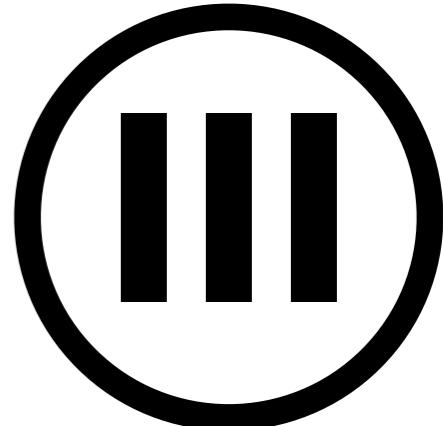
Structure of the talk



**Bayes-Optimal learning
of quadratic networks**



**Asymptotics of ERM in
overparametrised
quadratic networks**



**Bayes-Optimal of
attention networks**



Simplest possible model: quadratic networks

Data model:

Proportional-width
two-layer networks with
centered quadratic
activations

$$y_\mu = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}_\mu}{\sqrt{d}}\right)$$

$$\sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \frac{(\mathbf{w}_i^\top \mathbf{x})^2}{d} - \frac{\|\mathbf{w}_i\|^2}{d}$$

Simplest possible model: quadratic networks

Data model:

Proportional-width
two-layer networks with
centered quadratic
activations

$$y_\mu = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}_\mu}{\sqrt{d}}\right) = \text{Tr} \left[\frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - \mathbf{1}}{\sqrt{d}} \sum_{i=1}^p \frac{\mathbf{w}_i \mathbf{w}_i^\top}{\sqrt{pd}} \right]$$

$$\sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \frac{(\mathbf{w}_i^\top \mathbf{x})^2}{d} - \frac{\|\mathbf{w}_i\|^2}{d}$$

Simplest possible model: quadratic networks

Data model:

Proportional-width
two-layer networks with
centered quadratic
activations

$$y_\mu = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}_\mu}{\sqrt{d}}\right) = \text{Tr}\left[\frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - \mathbf{1}}{\sqrt{d}} \sum_{i=1}^p \frac{\mathbf{w}_i \mathbf{w}_i^\top}{\sqrt{pd}}\right] = \text{Tr}[G_\mu A]$$

$$\sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \frac{(\mathbf{w}_i^\top \mathbf{x})^2}{d} - \frac{\|\mathbf{w}_i\|^2}{d}$$

Simplest possible model: quadratic networks

Data model:

Proportional-width
two-layer networks with
centered quadratic
activations

$$y_\mu = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}_\mu}{\sqrt{d}}\right) = \text{Tr}\left[\frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - \mathbf{1}}{\sqrt{d}} \sum_{i=1}^p \frac{\mathbf{w}_i \mathbf{w}_i^\top}{\sqrt{pd}}\right] = \text{Tr}[G_\mu A]$$

$$\sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \frac{(\mathbf{w}_i^\top \mathbf{x})^2}{d} - \frac{\|\mathbf{w}_i\|^2}{d}$$

Goal: Bayes optimal study of simplest index model

$$y_\mu = \mathbf{x}_\mu^\top A^\star \mathbf{x}_\mu \quad A^\star \in \text{Sym}(d), \quad \text{rank}(A^\star) = \rho d, \quad \mathbf{x} \in \mathbb{R}^d, \quad d \rightarrow \infty$$

Simplest possible model: quadratic networks

Data model:

Proportional-width
two-layer networks with
centered quadratic
activations

$$y_\mu = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}_\mu}{\sqrt{d}}\right) = \text{Tr}\left[\frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - \mathbf{1}}{\sqrt{d}} \sum_{i=1}^p \frac{\mathbf{w}_i \mathbf{w}_i^\top}{\sqrt{pd}}\right] = \text{Tr}[G_\mu A]$$

$$\sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \frac{(\mathbf{w}_i^\top \mathbf{x})^2}{d} - \frac{\|\mathbf{w}_i\|^2}{d}$$

Goal: Bayes optimal study of simplest index model

$$y_\mu = \mathbf{x}_\mu^\top A^\star \mathbf{x}_\mu \quad A^\star \in \text{Sym}(d), \quad \text{rank}(A^\star) = \rho d, \quad \mathbf{x} \in \mathbb{R}^d, \quad d \rightarrow \infty$$

Results:

Fundamental limits :

You have $n = \alpha d^2$ samples $\{\mathbf{x}_\mu, y_\mu\}$

What is the minimal $\mathcal{E} = \|\hat{A} - A^\star\|_F$ over choices of \hat{A} ?

Simplest possible model: quadratic networks

Data model:

Proportional-width
two-layer networks with
centered quadratic
activations

$$y_\mu = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}_\mu}{\sqrt{d}}\right) = \text{Tr}\left[\frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - \mathbf{1}}{\sqrt{d}} \sum_{i=1}^p \frac{\mathbf{w}_i \mathbf{w}_i^\top}{\sqrt{pd}}\right] = \text{Tr}[G_\mu A]$$

$$\sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \frac{(\mathbf{w}_i^\top \mathbf{x})^2}{d} - \frac{\|\mathbf{w}_i\|^2}{d}$$

Goal: Bayes optimal study of simplest index model

$$y_\mu = \mathbf{x}_\mu^\top A^\star \mathbf{x}_\mu \quad A^\star \in \text{Sym}(d), \quad \text{rank}(A^\star) = \rho d, \quad \mathbf{x} \in \mathbb{R}^d, \quad d \rightarrow \infty$$

Results:

Fundamental limits :

You have $n = \alpha d^2$ samples $\{\mathbf{x}_\mu, y_\mu\}$

What is the minimal $\mathcal{E} = \|\hat{A} - A^\star\|_F$ over choices of \hat{A} ?

Optimal algorithms :

Write an algorithm that achieves \mathcal{E}

Sketch of the analysis though AMP

1. Gaussian equivalence on input

$$G_\mu = \frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - 1}{\sqrt{d}} \sim \text{GOE}(d)$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Matrix compressed sensing
[Y. Xu, A. Maillard, L. Zdeborová, F. Krzakala]

Sketch of the analysis though AMP

1. Gaussian equivalence on input

$$G_\mu = \frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - 1}{\sqrt{d}} \sim \text{GOE}(d)$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Matrix compressed sensing
[Y. Xu, A. Maillard, L. Zdeborová, F. Krzakala]

2. Write an AMP iteration

$$A^{t+1} = \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

$$h_\mu^{t+1} = y_\mu - \text{Tr}[G_\mu A^{t+1}] + \nabla \cdot \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

Sketch of the analysis though AMP

1. Gaussian equivalence on input

$$G_\mu = \frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - 1}{\sqrt{d}} \sim \text{GOE}(d)$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Matrix compressed sensing
[Y. Xu, A. Maillard, L. Zdeborová, F. Krzakala]

2. Write an AMP iteration

$$A^{t+1} = \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

Denoising
subproblem

$$h_\mu^{t+1} = y_\mu - \text{Tr}[G_\mu A^{t+1}] + \nabla \cdot \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

Sketch of the analysis though AMP

1. Gaussian equivalence on input

$$G_\mu = \frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - 1}{\sqrt{d}} \sim \text{GOE}(d)$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Matrix compressed sensing
[Y. Xu, A. Maillard, L. Zdeborová, F. Krzakala]

2. Write an AMP iteration

$$\boxed{A^{t+1} = \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)}$$

Denoising subproblem

$$h_\mu^{t+1} = y_\mu - \text{Tr}[G_\mu A^{t+1}] + \nabla \cdot \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

3. Study BO denoising

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant

Sketch of the analysis though AMP

1. Gaussian equivalence on input

$$G_\mu = \frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - 1}{\sqrt{d}} \sim \text{GOE}(d)$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Matrix compressed sensing
[Y. Xu, A. Maillard, L. Zdeborová, F. Krzakala]

2. Write an AMP iteration

$$\boxed{A^{t+1} = \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)}$$

Denoising subproblem

$$h_\mu^{t+1} = y_\mu - \text{Tr}[G_\mu A^{t+1}] + \nabla \cdot \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

3. Study BO denoising

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant



Estimate \mathcal{S} !

Sketch of the analysis though AMP

1. Gaussian equivalence on input

$$G_\mu = \frac{\mathbf{x}_\mu \mathbf{x}_\mu^\top - 1}{\sqrt{d}} \sim \text{GOE}(d)$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Matrix compressed sensing
[Y. Xu, A. Maillard, L. Zdeborová, F. Krzakala]

2. Write an AMP iteration

$$A^{t+1} = \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

Denoising
subproblem

$$h_\mu^{t+1} = y_\mu - \text{Tr}[G_\mu A^{t+1}] + \nabla \cdot \eta \left(A^t + \sum_{\mu=1}^N G_\mu h_\mu^t \right)$$

3. Study BO denoising

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant



Estimate \mathcal{S} !

Not quite a trivial problem...

[A. Maillard, F. Krzakala, M. Mézard, L. Zdeborová '21]
[E. Troiani, V. Erba, F. Krzakala, A. Maillard, L. Zdeborová '22]
[F. Pourkamali, N. Macris '23]
[G. Semerjian '24]

BO denoising and HClZ integrals

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant

Minimise MSE on \mathcal{S} !

BO denoising and HClZ integrals

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant

Minimise MSE on \mathcal{S} !



Rotationally Invariant Estimator (RIE)

[O. Ledoit, S. Péché '11]

[J. Bun, J.P. Bouchaud, M. Potters '16]

BO denoising and HClZ integrals

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant

Minimise MSE on \mathcal{S} !



Rotationally Invariant Estimator (RIE)

[O. Ledoit, S. Péché '11]
[J. Bun, J.P. Bouchaud, M. Potters '16]

$$\mathcal{Y} = U\Lambda_{\mathcal{Y}}U^{\top}$$



$$\hat{\mathcal{S}} = UF(\Lambda_{\mathcal{Y}})U^{\top}$$

$$F[\Lambda_{\mathcal{Y}}]_i = f_{\mathcal{Y}}[(\Lambda_{\mathcal{Y}})_i]$$

BO denoising and HClZ integrals

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant

Minimise MSE on \mathcal{S} !



Rotationally Invariant Estimator (RIE)

[O. Ledoit, S. Péché '11]
[J. Bun, J.P. Bouchaud, M. Potters '16]

$$\mathcal{Y} = U\Lambda_{\mathcal{Y}}U^{\top}$$



$$\hat{\mathcal{S}} = UF(\Lambda_{\mathcal{Y}})U^{\top}$$

$$F[\Lambda_{\mathcal{Y}}]_i = f_{\mathcal{Y}}[(\Lambda_{\mathcal{Y}})_i]$$

Integrate posterior over rotations

$$P(\mathcal{S} | \mathcal{Y}) \propto e^{-\frac{\|\mathcal{Y} - \mathcal{S}\|^2}{2\delta^2}} P(\mathcal{S}) = e^{-\frac{\|\mathcal{Y}\|^2 + \|\mathcal{S}\|^2}{2\delta^2}} P(\mathcal{S}) e^{-\frac{\text{Tr}[\mathcal{Y}\mathcal{S}]}{\delta^2}}$$

BO denoising and HCIZ integrals

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant
Minimise MSE on \mathcal{S} !



Rotationally Invariant Estimator (RIE)

[O. Ledoit, S. Péché '11]
[J. Bun, J.P. Bouchaud, M. Potters '16]

$$\mathcal{Y} = U\Lambda_{\mathcal{Y}}U^{\top}$$



$$\hat{\mathcal{S}} = UF(\Lambda_{\mathcal{Y}})U^{\top}$$

$$F[\Lambda_{\mathcal{Y}}]_i = f_{\mathcal{Y}}[(\Lambda_{\mathcal{Y}})_i]$$

Integrate posterior over rotations

$$P(\mathcal{S} | \mathcal{Y}) \propto e^{-\frac{\|\mathcal{Y} - \mathcal{S}\|^2}{2\delta^2}} P(\mathcal{S}) = e^{-\frac{\|\mathcal{Y}\|^2 + \|\mathcal{S}\|^2}{2\delta^2}} P(\mathcal{S}) e^{-\frac{\text{Tr}[\mathcal{Y}\mathcal{S}]}{\delta^2}}$$

A red oval highlights the term $e^{-\frac{\text{Tr}[\mathcal{Y}\mathcal{S}]}{\delta^2}}$, and a red arrow points from it to a yellow starburst icon labeled "HCIZ!".

$$\int e^{-\frac{\text{Tr}[U\Lambda_{\mathcal{Y}}U^{\top}\Lambda_{\mathcal{S}}]}{\delta^2}} dU$$

Recall A. Maillard's talk...

BO denoising and HCIZ integrals

Observe $\mathcal{Y} = \mathcal{S} + \delta Z$ with $Z \sim \text{GOE}(d)$ and $P_{\mathcal{S}}$ rotationally invariant

Minimise MSE on \mathcal{S} !



Rotationally Invariant Estimator (RIE)

[O. Ledoit, S. Péché '11]
[J. Bun, J.P. Bouchaud, M. Potters '16]

$$\mathcal{Y} = U\Lambda_{\mathcal{Y}}U^{\top}$$

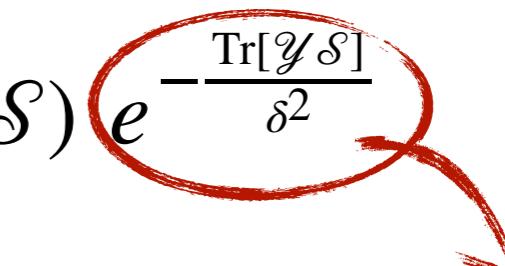


$$\hat{\mathcal{S}} = UF(\Lambda_{\mathcal{Y}})U^{\top}$$

$$F[\Lambda_{\mathcal{Y}}]_i = f_{\mathcal{Y}}[(\Lambda_{\mathcal{Y}})_i]$$

Integrate posterior over rotations

$$P(\mathcal{S} | \mathcal{Y}) \propto e^{-\frac{\|\mathcal{Y} - \mathcal{S}\|^2}{2\delta^2}} P(\mathcal{S}) = e^{-\frac{\|\mathcal{Y}\|^2 + \|\mathcal{S}\|^2}{2\delta^2}} P(\mathcal{S})$$



HCIZ!

Asymptotic Minimum MSE

$$\|\hat{\mathcal{S}} - \mathcal{S}\|_F \xrightarrow{d \rightarrow \infty} \delta^2 - \frac{4\pi^2\delta^4}{3} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$

$$\int e^{-\frac{\text{Tr}[U\Lambda_{\mathcal{Y}}U^{\top}\Lambda_{\mathcal{S}}]}{\delta^2}} dU$$

Recall A. Maillard's talk...

Learning curves for quadratic networks

Minimal MMSE

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} \frac{2\alpha\rho}{\hat{q}}$$

$$1 - 2\alpha = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$

$\mathcal{Y} = A + \hat{q}^{-1/2}G$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

$$\begin{aligned} G &\sim \text{GOE}(d) & p &= \rho d \\ A &\sim \mathcal{W}_{p,d} & n &= \alpha d^2 \end{aligned}$$

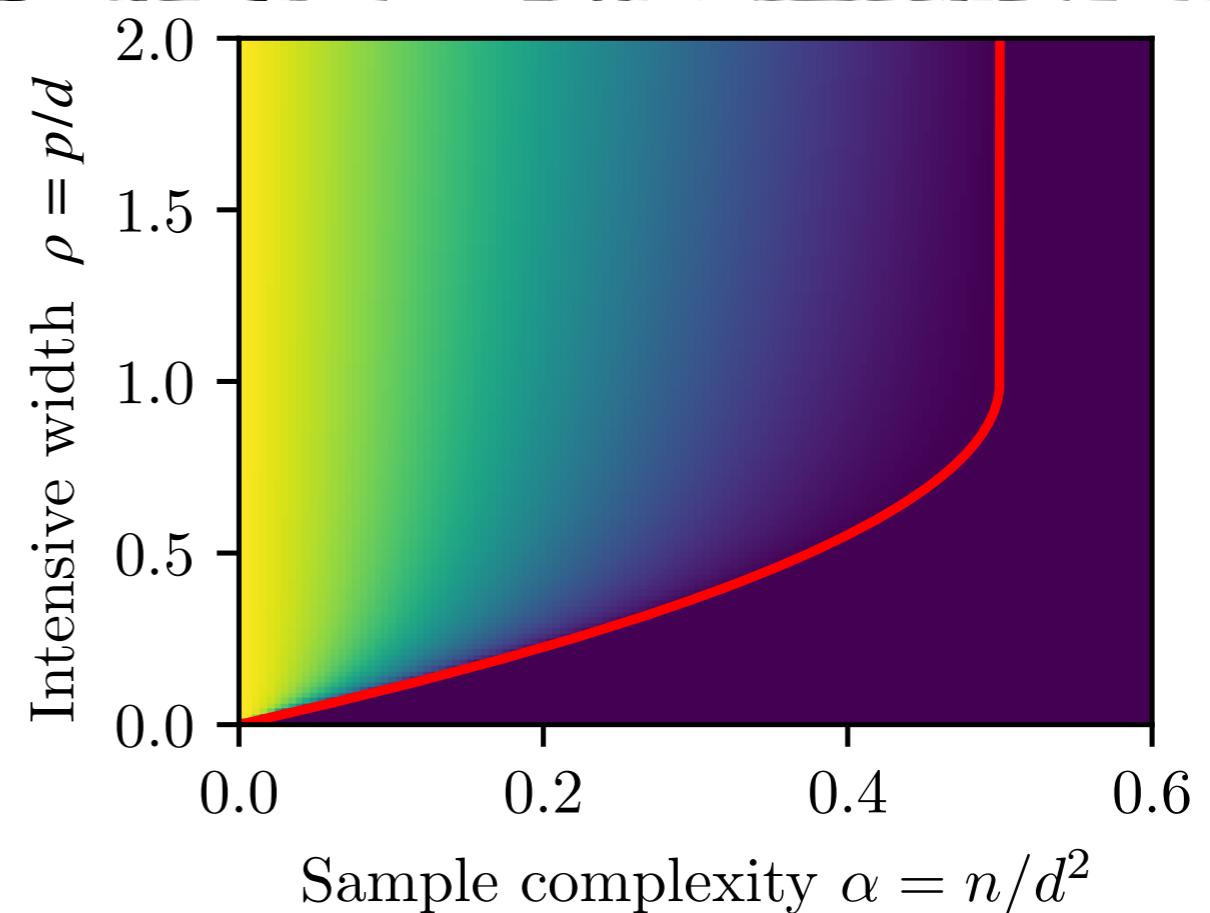
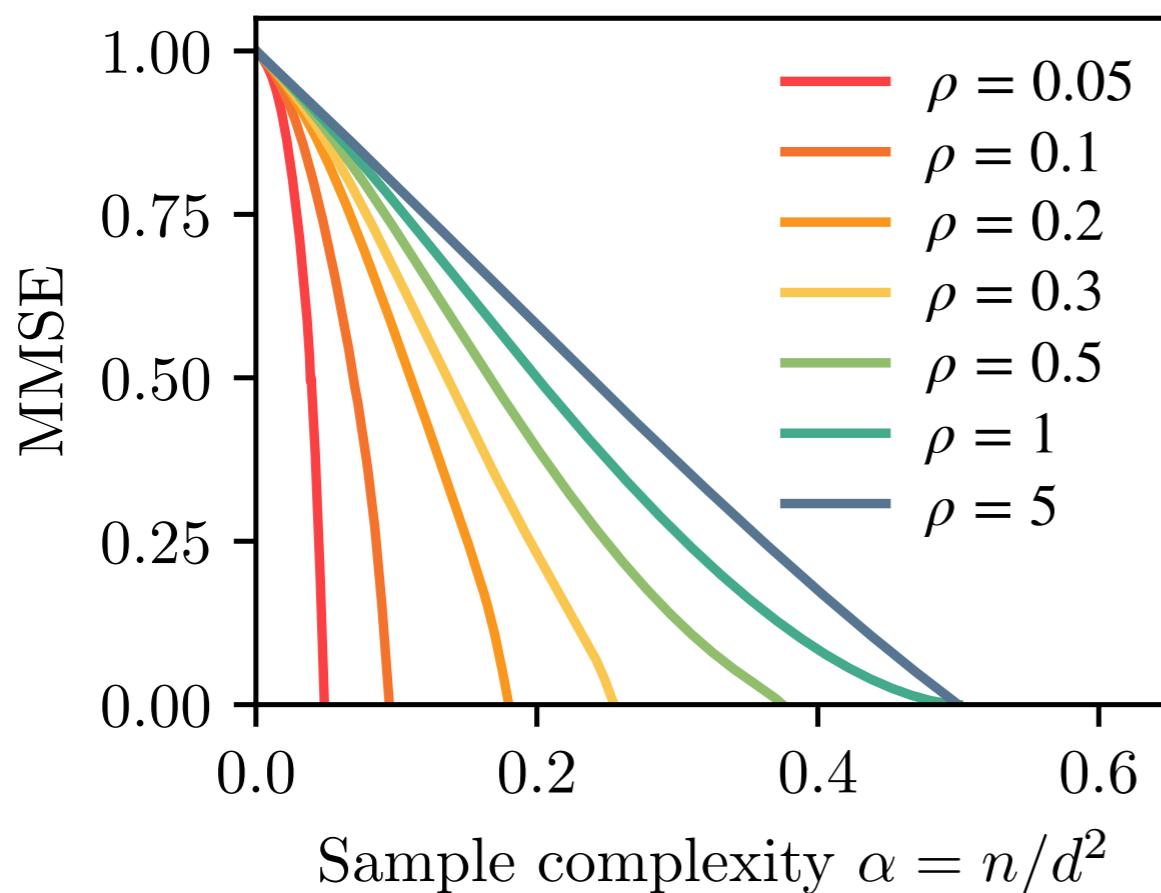
Learning curves for quadratic networks

Minimal MMSE

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} \frac{2\alpha\rho}{\hat{q}}$$

$$1 - 2\alpha = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$

$\mathcal{Y} = A + \hat{q}^{-1/2}G$



$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

$$G \sim \text{GOE}(d)$$

$$A \sim \mathcal{W}_{p,d}$$

$$p = \rho d$$

$$n = \alpha d^2$$

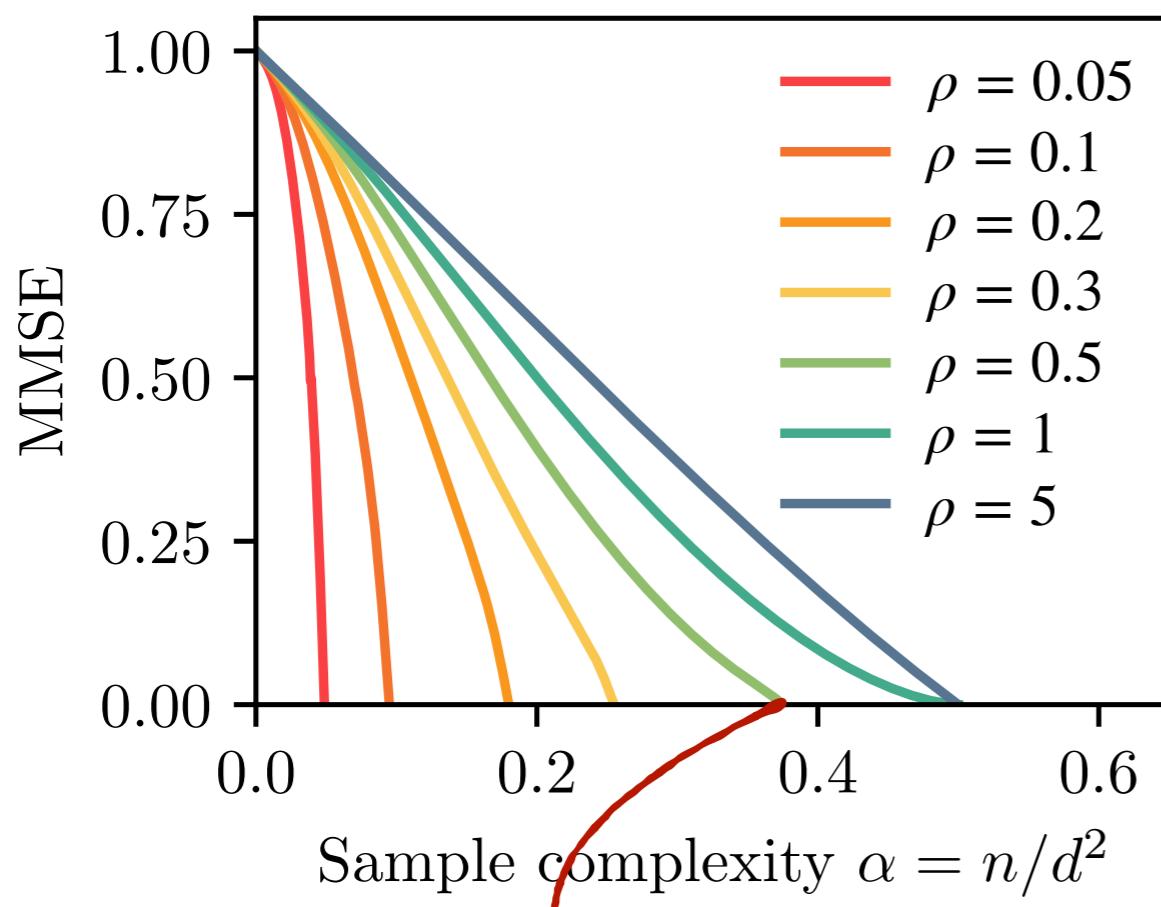
Learning curves for quadratic networks

Minimal MMSE

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} \frac{2\alpha\rho}{\hat{q}}$$

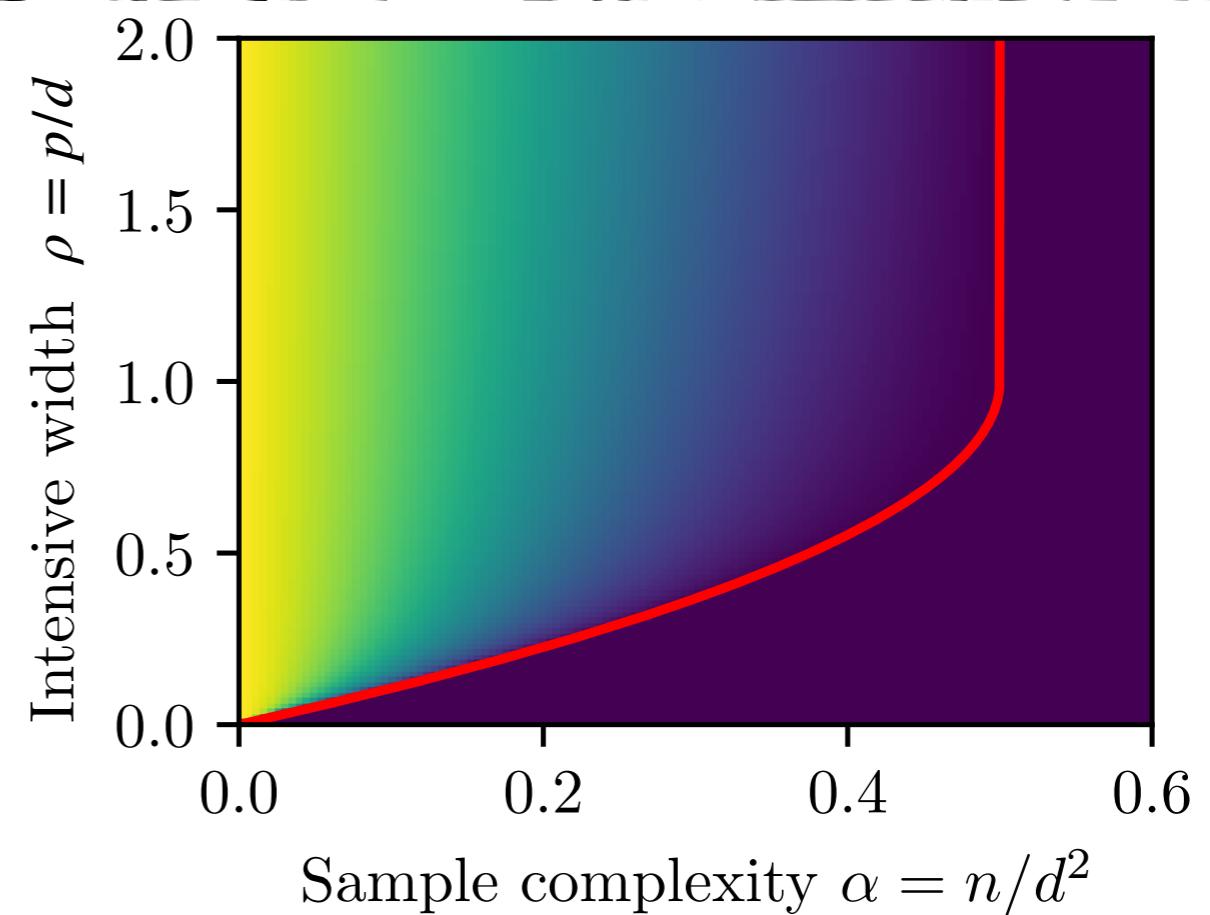
$$1 - 2\alpha = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$

$\mathcal{Y} = A + \hat{q}^{-1/2}G$



$$\alpha_c = \rho - \rho^2/2$$

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$



Sample complexity $\alpha = n/d^2$

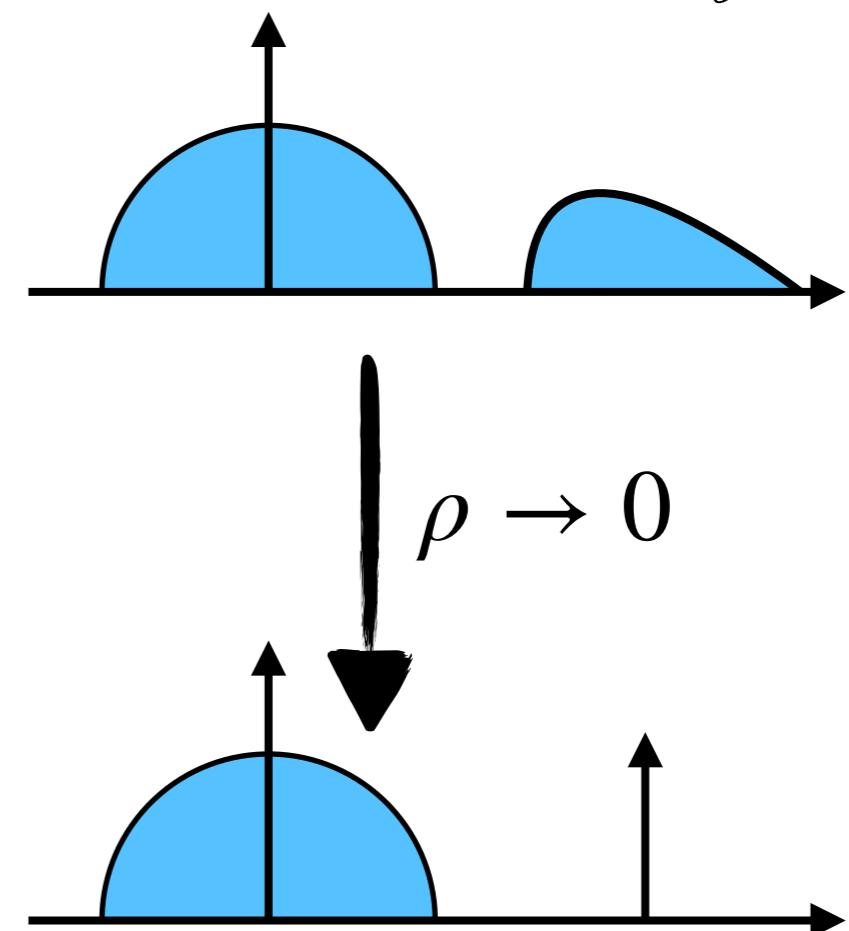
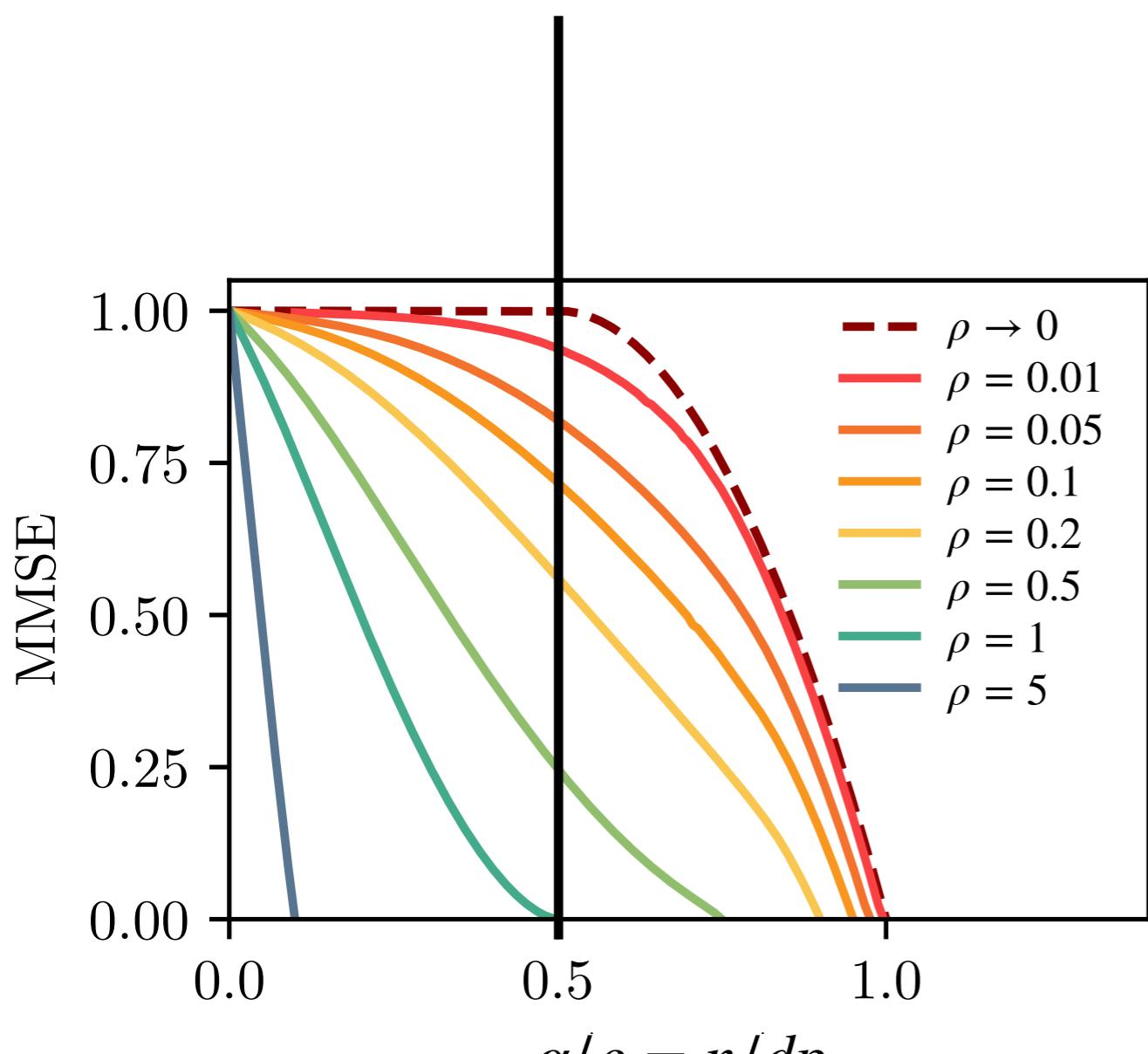
$$\begin{aligned} G &\sim \text{GOE}(d) & p &= \rho d \\ A &\sim \mathcal{W}_{p,d} & n &= \alpha d^2 \end{aligned}$$

Narrow limit: a BBP transition for weak recovery

Narrow width limit $\rho \rightarrow 0$

$$1 - 2\alpha = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$

$\mathcal{Y} = A + \hat{q}^{-1/2}G$



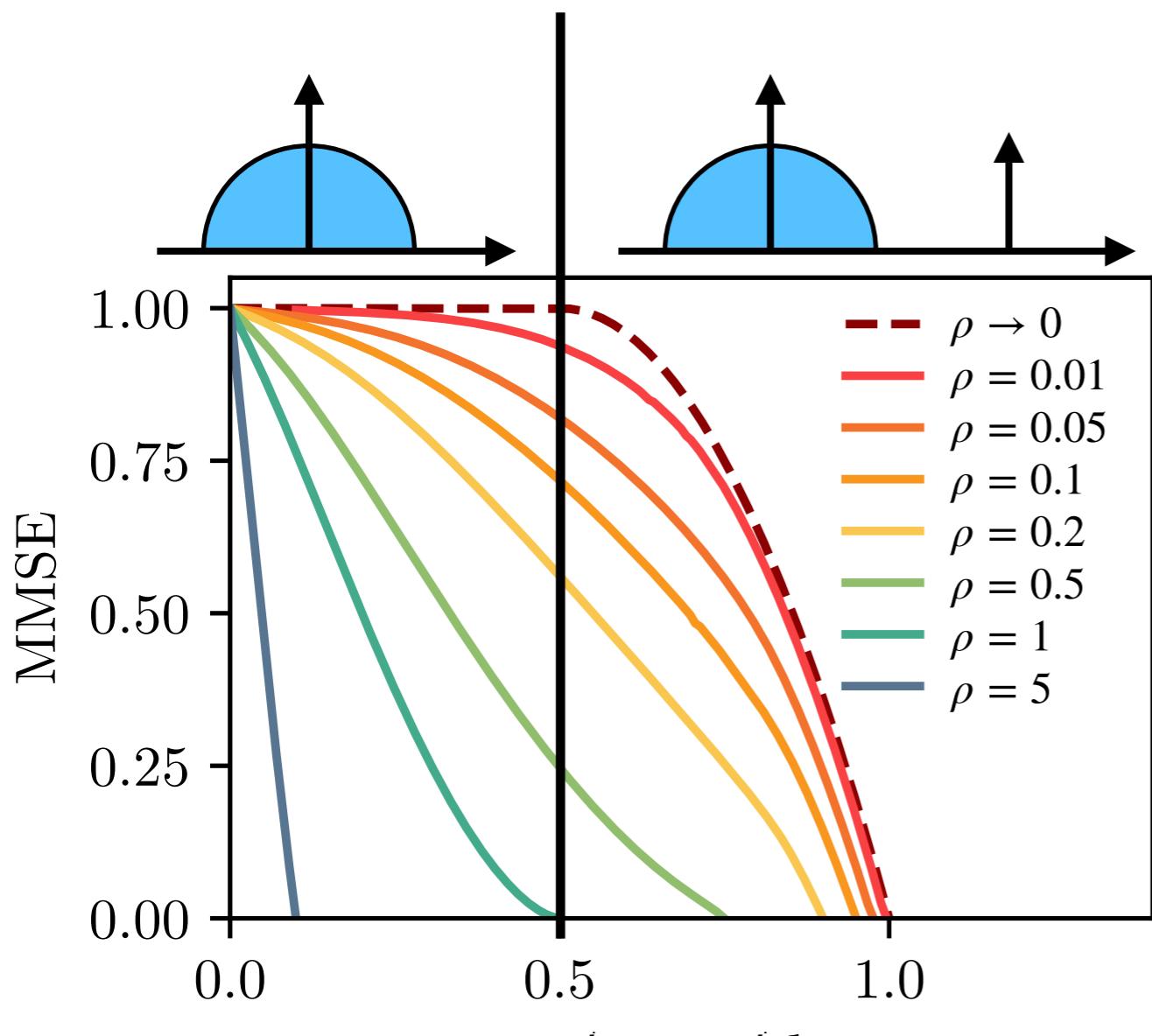
$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A &\sim \mathcal{W}_{p,d} \end{aligned}$$

$$\begin{aligned} p &= \rho d \\ n &= \alpha d^2 \end{aligned}$$

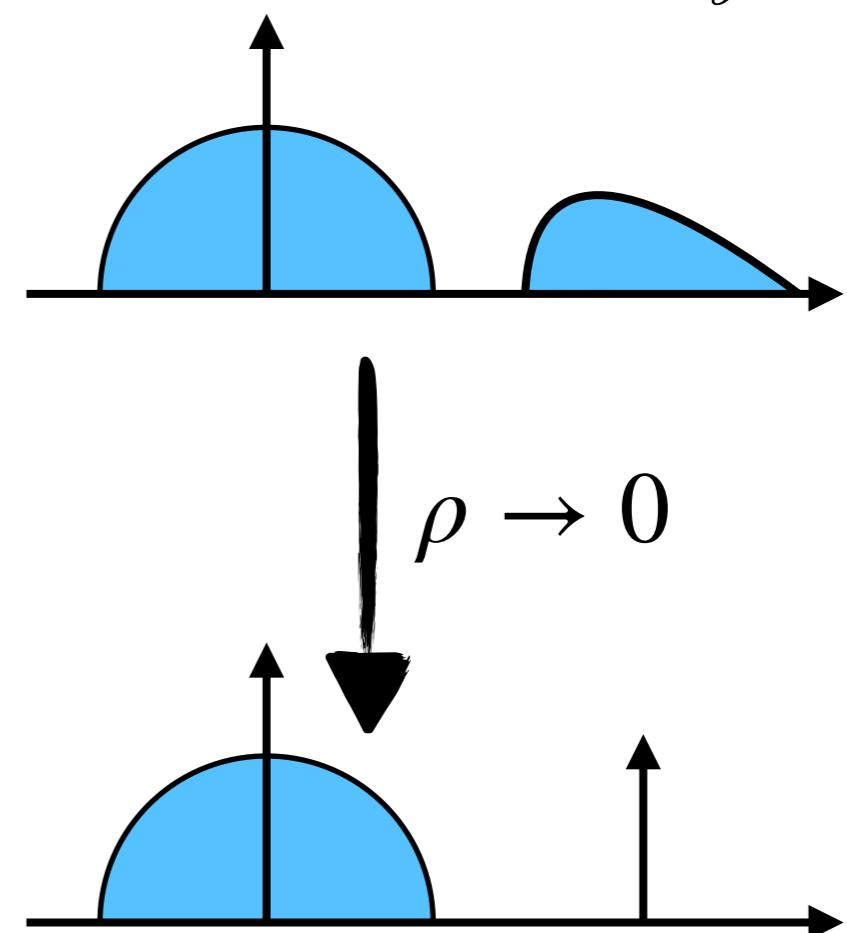
Narrow limit: a BBP transition for weak recovery

Narrow width limit $\rho \rightarrow 0$



$$1 - 2\alpha = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$

$\mathcal{Y} = A + \hat{q}^{-1/2}G$



$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA]$$

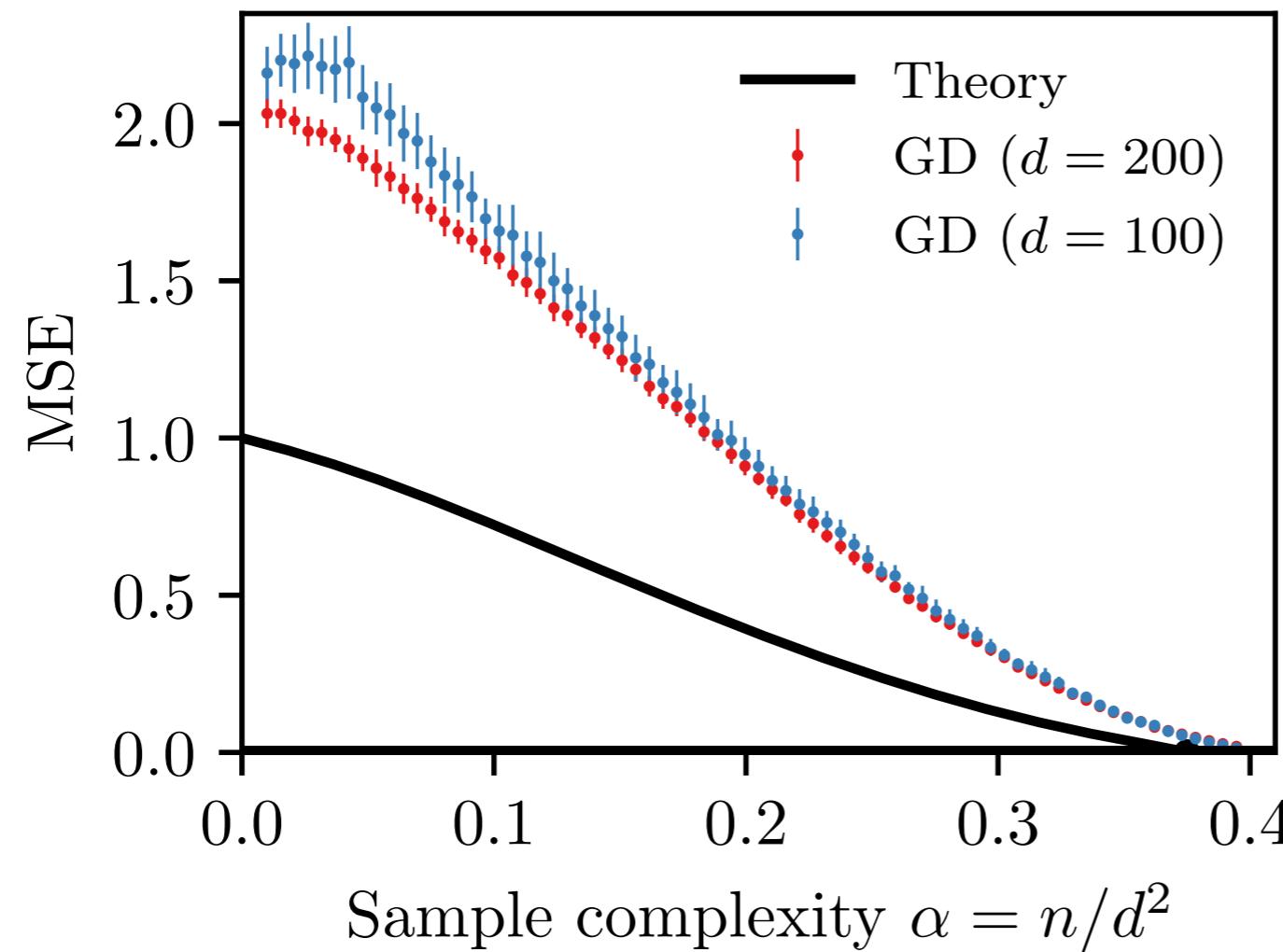
$$\begin{aligned} G &\sim \text{GOE}(d) \\ A &\sim \mathcal{W}_{p,d} \end{aligned}$$

$$\begin{aligned} p &= \rho d \\ n &= \alpha d^2 \end{aligned}$$

(Some form of) GD can have BO performance

Everyone: I don't care about AMP, what about GD?

Minimise $\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2$



$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA^\star]$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{\mathbf{w}}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$G \sim \text{GOE}(d)$$

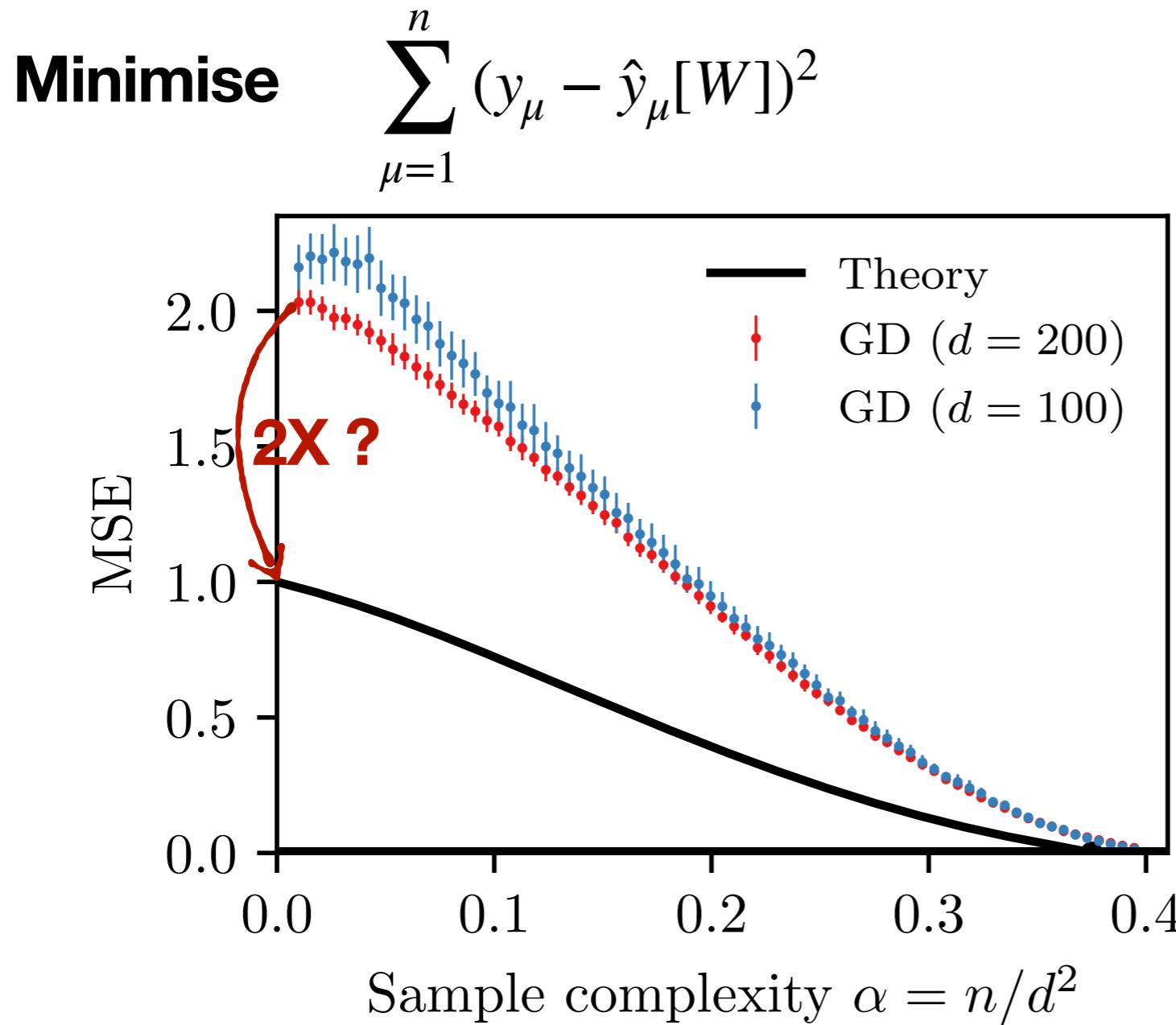
$$A^\star \sim \mathcal{W}_{p,d}$$

$$p = \rho d$$

$$n = \alpha d^2$$

(Some form of) GD can have BO performance

Everyone: I don't care about AMP, what about GD?



$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA^\star]$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{\mathbf{w}}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$G \sim \text{GOE}(d)$$

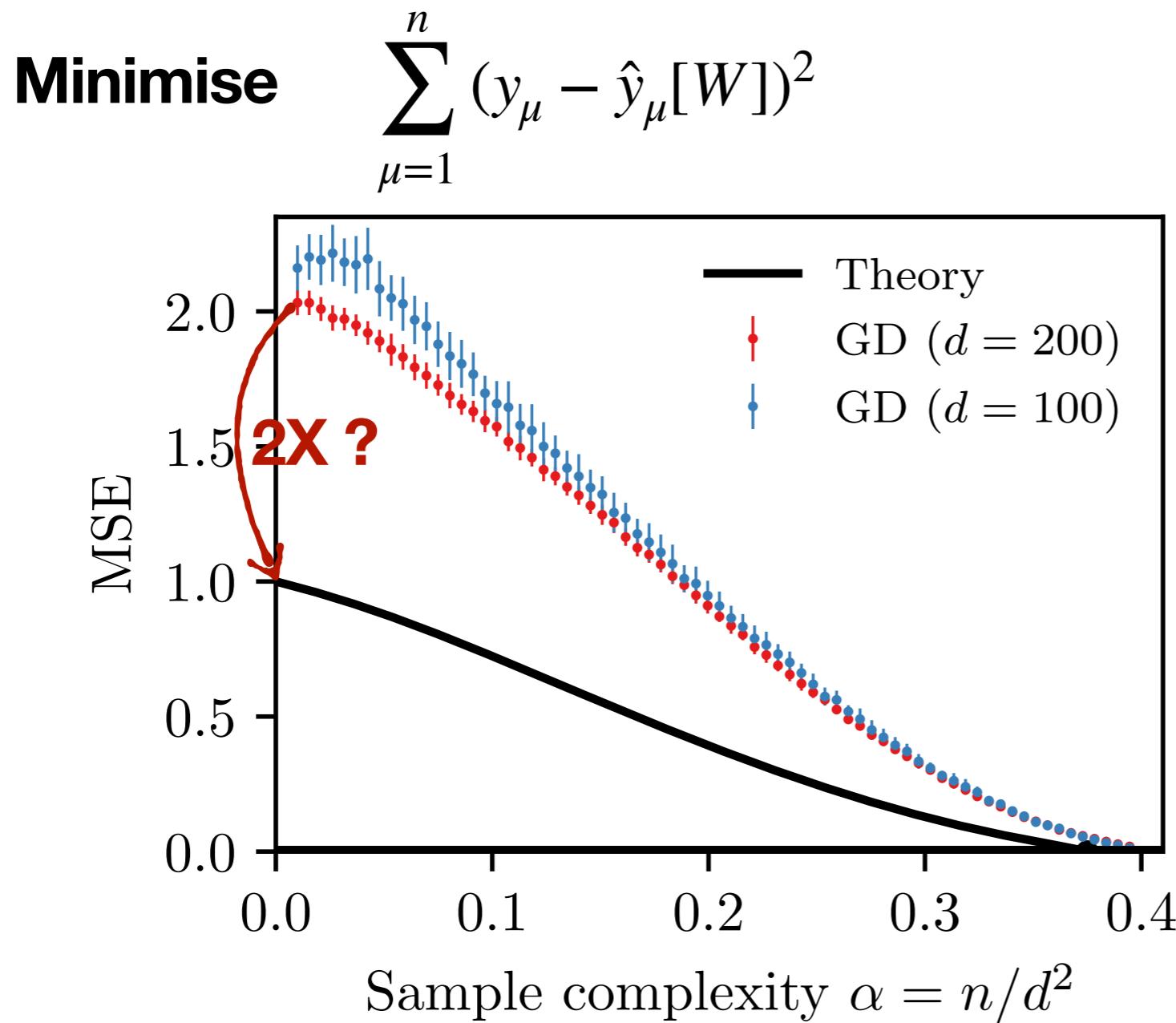
$$A^\star \sim \mathcal{W}_{p,d}$$

$$p = \rho d$$

$$n = \alpha d^2$$

(Some form of) GD can have BO performance

Everyone: I don't care about AMP, what about GD?



Averaged GD:

1. Run multiple times initializing in prior;
2. Average the labels at convergence

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA^\star]$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{\mathbf{w}}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$G \sim \text{GOE}(d)$$

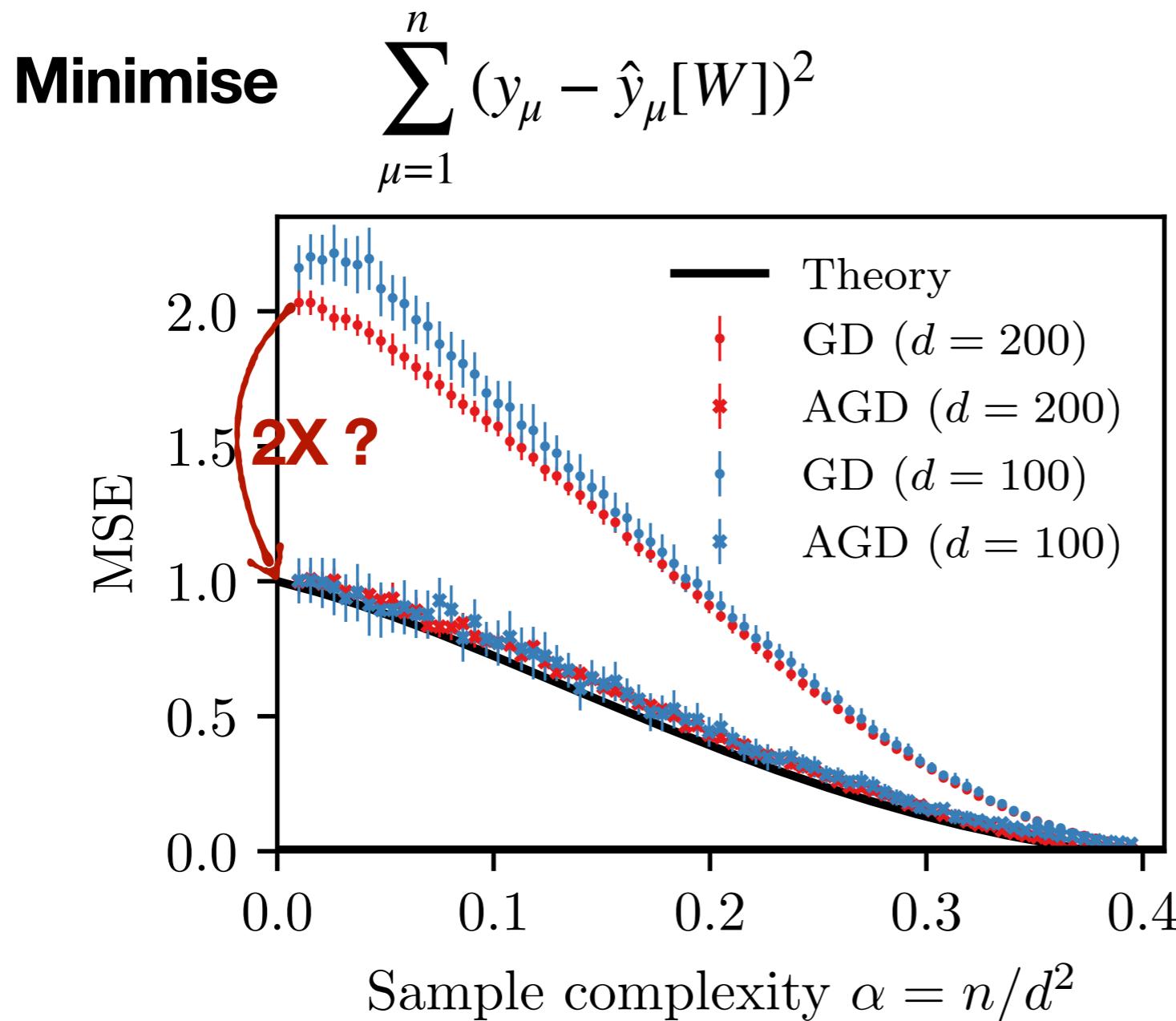
$$A^\star \sim \mathcal{W}_{p,d}$$

$$p = \rho d$$

$$n = \alpha d^2$$

(Some form of) GD can have BO performance

Everyone: I don't care about AMP, what about GD?



Averaged GD:

1. Run multiple times initializing in prior;
2. Average the labels at convergence

BO performance!

(hard) open problem

$$y = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[GA^\star]$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{\mathbf{w}}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$G \sim \text{GOE}(d)$$

$$A^\star \sim \mathcal{W}_{p,d}$$

$$p = \rho d$$

$$n = \alpha d^2$$

L2 regularized quadratic networks

Let's add L2 regularization and learn with overparametrized
 $p \gg p^*$ network

Minimise $\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 + \lambda \|W\|_F^2$

$$y = \text{Tr}[GA^\star] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^{\textcolor{red}{p}} \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^\star &\sim \mathcal{W}_{p^\star, d} \end{aligned}$$

$$\begin{aligned} \textcolor{red}{p} &\geq d \\ \textcolor{red}{p^\star} &= \rho^\star d \\ n &= \alpha d^2 \end{aligned}$$

L2 regularized quadratic networks

Let's add L2 regularization and learn with overparametrized
 $p \gg p^*$ network

Minimise $\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 + \lambda \|W\|_F^2$

(i) Training: Optimisation with **GD is easy**, all minima are global minima

[Venturi, Bandeira, Bruna '19]

$$y = \text{Tr}[GA^*] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^{\textcolor{red}{p}} \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^* &\sim \mathcal{W}_{p^*, d} \end{aligned}$$

$$\begin{aligned} \textcolor{red}{p} &\geq d \\ \textcolor{red}{p^*} &= \rho^* d \\ n &= \alpha d^2 \end{aligned}$$

L2 regularized quadratic networks

Let's add L2 regularization and learn with overparametrized $p \gg p^*$ network

Minimise
$$\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 + \lambda \|W\|_F^2$$

(i) Training: Optimisation with **GD is easy**, all minima are global minima

[Venturi, Bandeira, Bruna '19]

(ii) Equivalent to a (convex) matrix compressed sensing with a **nuclear norm regularisation**

[Fazel, Candes, Recht, Parrilo '08; Donoho, Gavish, Montanari '13; Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro '17]

$$\sum_{\mu=1}^n \|(A^* - A)G_\mu + \sqrt{\Delta}\zeta\|_F^2 + \lambda \|A\|_*$$

$$y = \text{Tr}[GA^*] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^{\textcolor{red}{p}} \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^* &\sim \mathcal{W}_{p^*, d} \end{aligned}$$

$$\begin{aligned} \textcolor{red}{p} &\geq d \\ \textcolor{red}{p^*} &= \rho^* d \\ n &= \alpha d^2 \end{aligned}$$

L2 regularized quadratic networks

Let's add L2 regularization and learn with overparametrized $p \gg p^*$ network

Minimise
$$\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 + \lambda \|W\|_F^2$$

(i) Training: Optimisation with **GD is easy**, all minima are global minima

[Venturi, Bandeira, Bruna '19]

(ii) Equivalent to a (convex) matrix compressed sensing with a **nuclear norm regularisation**

[Fazel, Candes, Recht, Parrilo '08; Donoho, Gavish, Montanari '13; Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro '17]

$$\sum_{\mu=1}^n \|(A^* - A)G_\mu + \sqrt{\Delta}\zeta\|_F^2 + \lambda \|A\|_*$$

(iii) NTK/Lazy is instead equivalent to matrix compressed sensing with a Frobenius **norm regularisation**

$$\sum_{\mu=1}^n \|(A^* - A)G_\mu + \sqrt{\Delta}\zeta\|_F^2 + \lambda \|A\|_F^2$$

$$y = \text{Tr}[GA^*] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^{\textcolor{red}{p}} \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^* &\sim \mathcal{W}_{p^*, d} \end{aligned}$$

$$\begin{aligned} \textcolor{red}{p} &\geq d \\ \textcolor{red}{p^*} &= \rho^* d \\ n &= ad^2 \end{aligned}$$

Asymptotics of overparametrized quadratic nets

Generalization error

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} 2\alpha\delta^2 - \frac{\Delta}{2}$$

Training loss

$$\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 \xrightarrow{d \rightarrow \infty} \frac{\delta^2}{4\epsilon^2} - \frac{\lambda}{2}\partial_2 J(\delta, \lambda\epsilon)$$

$$y = \text{Tr}[GA^\star] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{\mathbf{w}}_i^\top \mathbf{x}}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^\star &\sim \mathcal{W}_{p^\star, d} \end{aligned}$$

$$\begin{aligned} p &\geq d \\ p^\star &= \rho^\star d \\ n &= \alpha d^2 \end{aligned}$$

Asymptotics of overparametrized quadratic nets

Generalization error

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} 2\alpha\delta^2 - \frac{\Delta}{2}$$

Training loss

$$\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 \xrightarrow{d \rightarrow \infty} \frac{\delta^2}{4\epsilon^2} - \frac{\lambda}{2}\partial_2 J(\delta, \lambda\epsilon)$$

$$4\alpha\delta - \frac{\delta}{\epsilon} = \partial_1 J(\delta, \lambda\epsilon)$$

$$Q^\star + \frac{\Delta}{2} + 2\alpha\delta^2 - \frac{\delta^2}{\epsilon} = (1 - \epsilon\lambda\partial_2)J(\delta, \lambda\epsilon)$$

$$J(a, b) = \int_b^{+\infty} dx \mu_y(x) (x - b)^2$$

\downarrow

$$y = A^\star + a^{-1/2}G$$

$$y = \text{Tr}[GA^\star] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^\star &\sim \mathcal{W}_{p^\star, d} \end{aligned}$$

$$\begin{aligned} p &\geq d \\ p^\star &= \rho^\star d \\ n &= \alpha d^2 \end{aligned}$$

Asymptotics of overparametrized quadratic nets

Generalization error

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} 2\alpha\delta^2 - \frac{\Delta}{2}$$

Training loss

$$\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 \xrightarrow{d \rightarrow \infty} \frac{\delta^2}{4\epsilon^2} - \frac{\lambda}{2}\partial_2 J(\delta, \lambda\epsilon)$$

$$4\alpha\delta - \frac{\delta}{\epsilon} = \partial_1 J(\delta, \lambda\epsilon)$$

$$Q^\star + \frac{\Delta}{2} + 2\alpha\delta^2 - \frac{\delta^2}{\epsilon} = (1 - \epsilon\lambda\partial_2)J(\delta, \lambda\epsilon)$$

$$J(a, b) = \int_b^{+\infty} dx \mu_{\mathcal{Y}}(x) (x - b)^2$$

$\mathcal{Y} = A^\star + a^{-1/2}G$

Spectrum of \hat{A} at convergence

$$\mu_{\hat{W}}(x) = C\delta(x) + I(x > 0)[2x\mu_{\mathcal{Y}}(x^2 + \lambda\epsilon)]$$

$$y = \text{Tr}[GA^\star] + \sqrt{\Delta}\zeta$$

$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^\star &\sim \mathcal{W}_{p^\star, d} \end{aligned}$$

$$\begin{aligned} p &\geq d \\ p^\star &= \rho^\star d \\ n &= \alpha d^2 \end{aligned}$$

Asymptotics of overparametrized quadratic nets

Generalization error

$$\sum_{\mu=1}^n (\hat{y}[x_\mu^{\text{new}}] - y_\mu^{\text{new}})^2 \xrightarrow{d \rightarrow \infty} 2\alpha\delta^2 - \frac{\Delta}{2}$$

Training loss

$$\sum_{\mu=1}^n (y_\mu - \hat{y}_\mu[W])^2 \xrightarrow{d \rightarrow \infty} \frac{\delta^2}{4\epsilon^2} - \frac{\lambda}{2}\partial_2 J(\delta, \lambda\epsilon)$$

$$4\alpha\delta - \frac{\delta}{\epsilon} = \partial_1 J(\delta, \lambda\epsilon)$$

$$Q^\star + \frac{\Delta}{2} + 2\alpha\delta^2 - \frac{\delta^2}{\epsilon} = (1 - \epsilon\lambda\partial_2)J(\delta, \lambda\epsilon)$$

$$J(a, b) = \int_b^{+\infty} dx \mu_{\mathcal{Y}}(x) (x - b)^2$$

$\mathcal{Y} = A^\star + a^{-1/2}G$

Spectrum of \hat{A} at convergence $\mu_{\hat{W}}(x) = C\delta(x) + I(x > 0)[2x\mu_{\mathcal{Y}}(x^2 + \lambda\epsilon)]$

Analysis similar to BO:

Different spectral denoiser: shifted ReLU on eigenvalues

$$y = \text{Tr}[GA^\star] + \sqrt{\Delta}\zeta$$

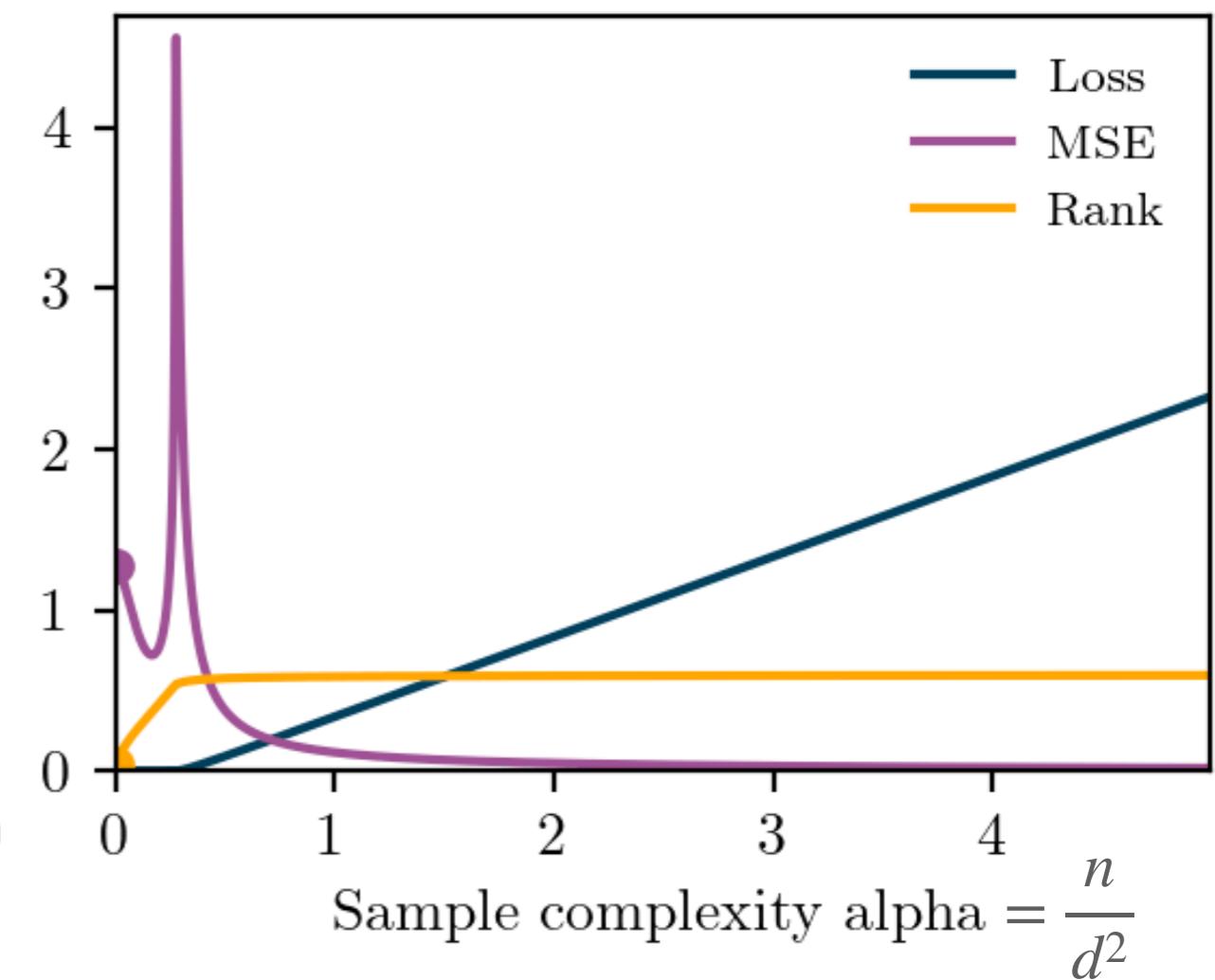
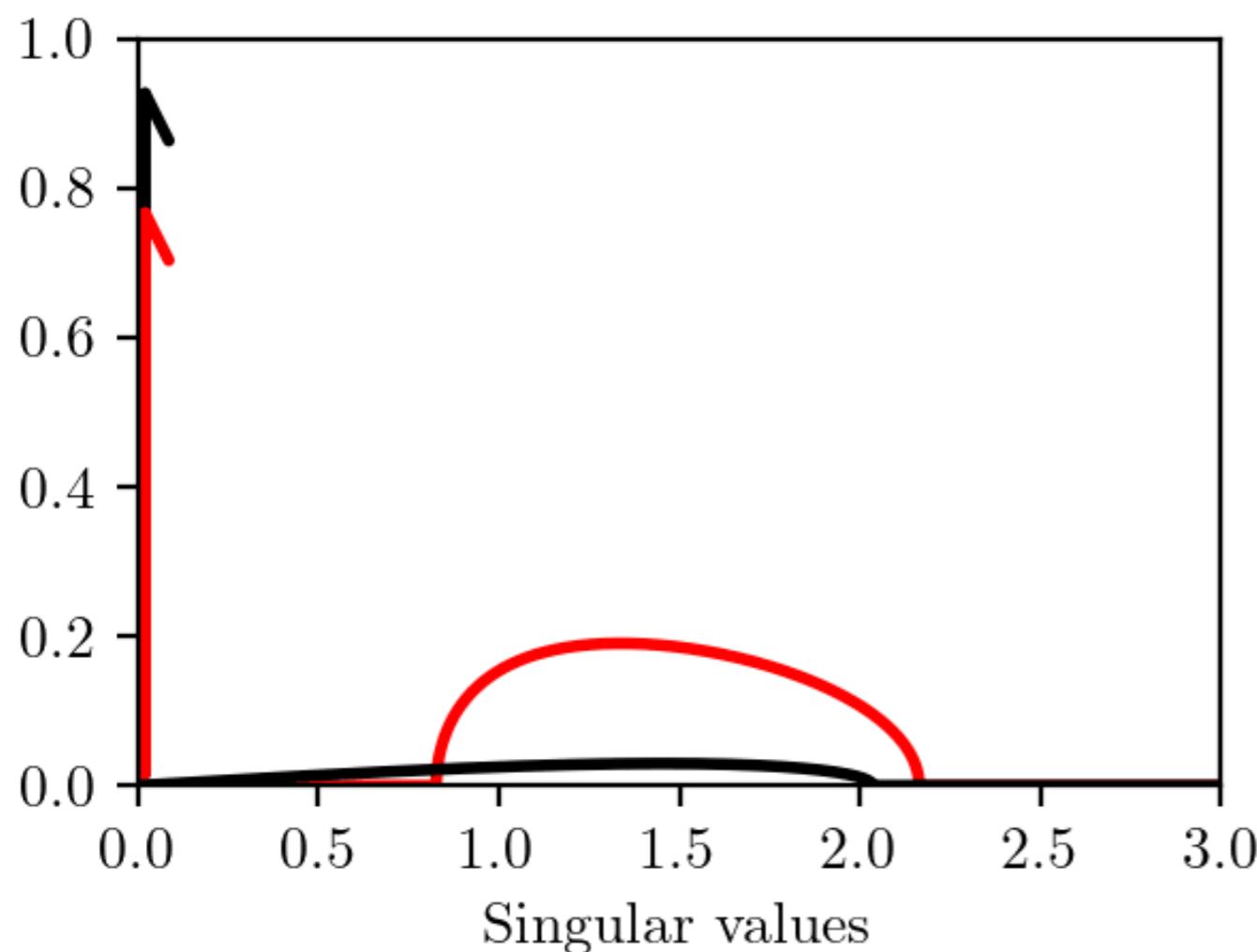
$$\hat{y} = \frac{1}{\sqrt{p}} \sum_{i=1}^p \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[G\hat{A}]$$

$$\begin{aligned} G &\sim \text{GOE}(d) \\ A^\star &\sim \mathcal{W}_{p^\star, d} \end{aligned}$$

$$\begin{aligned} p &\geq d \\ p^\star &= \rho^\star d \\ n &= \alpha d^2 \end{aligned}$$

Vignette I : spectrum of matrix \hat{W} , lower noise

Noisy data, low regularisation



Memorization
(But min- ℓ_1 norm interpolator)

$$\text{Double descent} = \frac{1}{\sqrt{P}} \sum_{i=1}^p \frac{\langle \hat{w}_i^\top, r_i \rangle}{\sqrt{d}}$$

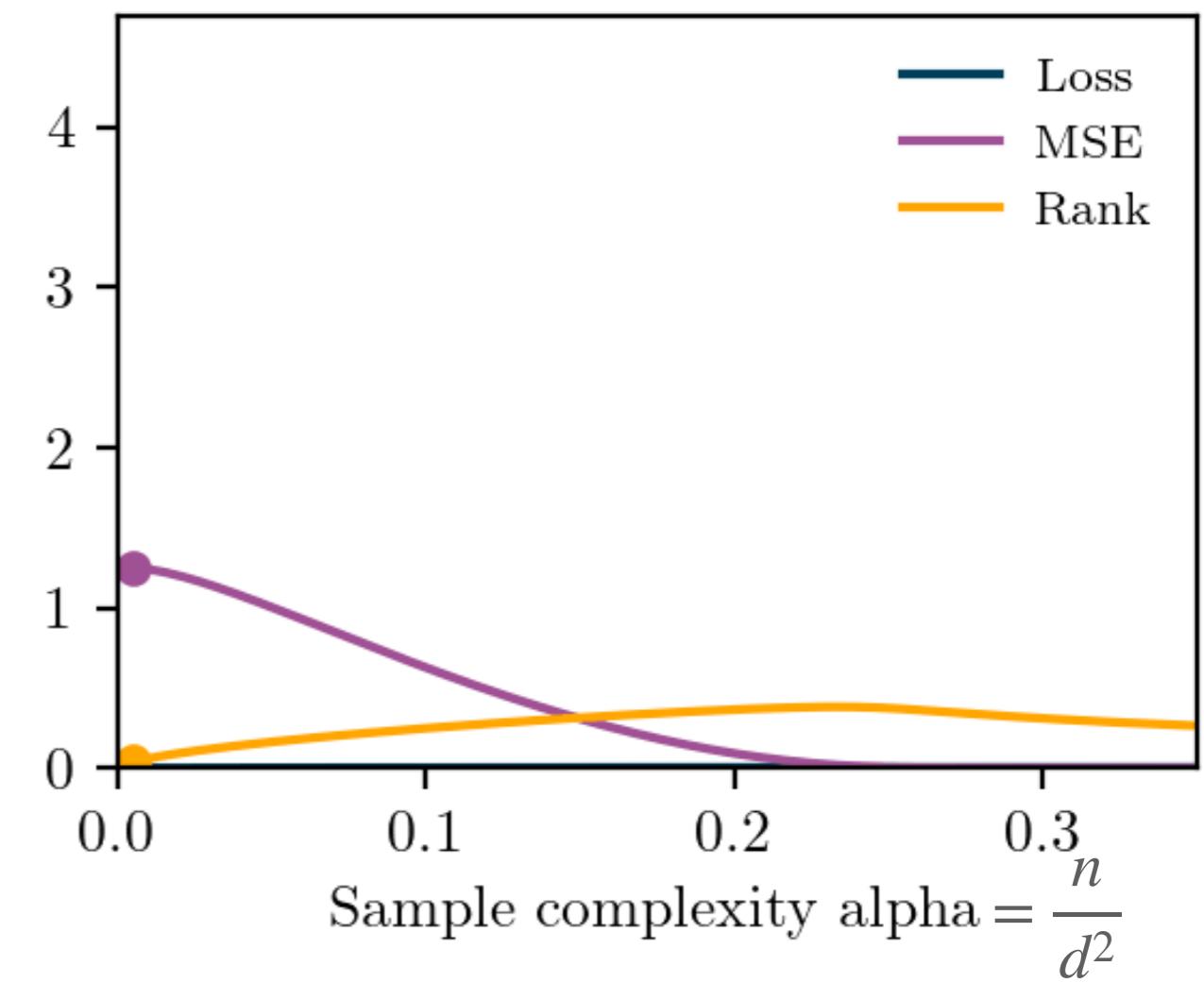
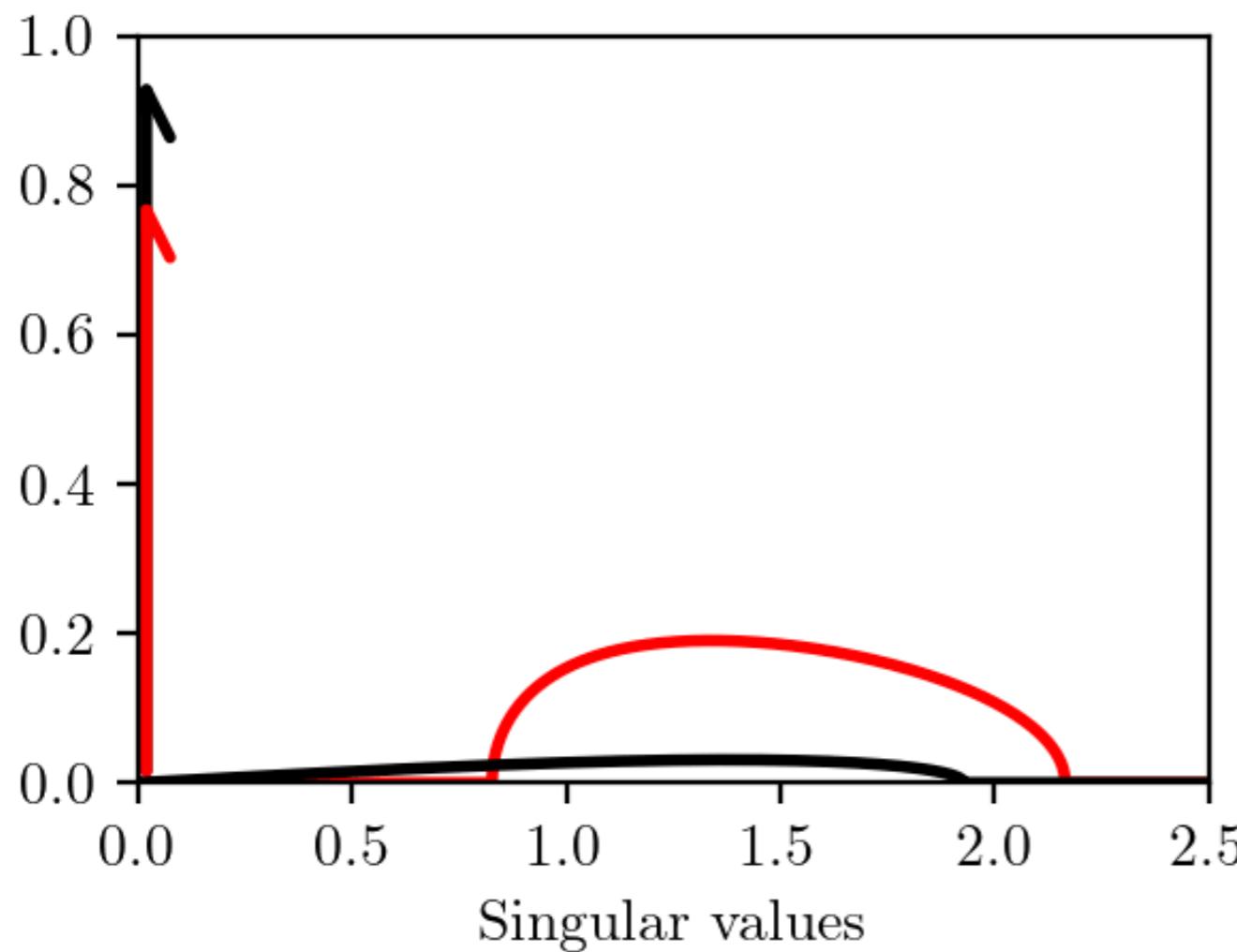
Interpolation peak $\text{Tr}[G\hat{A}]$

Generalization

$$p \geq d$$
$$p^* = \rho^* d$$
$$n = \alpha d^2$$

Vignette I : spectrum of matrix \hat{W} , lower noise

Less noisy data, low regularisation



Memorization
 $y = \min_{\|\hat{w}\|_1} \|y - \hat{w}^\top x\|$ (But min- ℓ_1 norm interpolator)

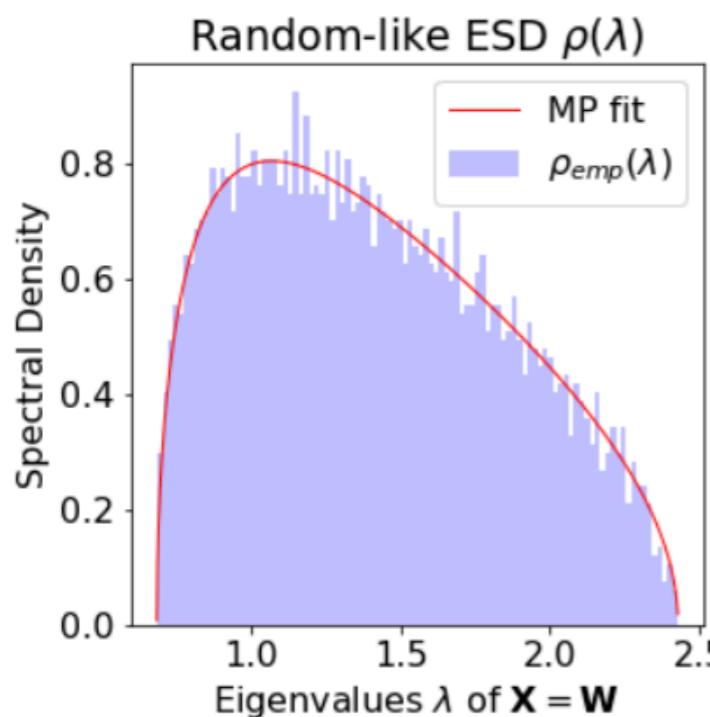
$$\sum_{i=1}^p \sigma\left(\frac{\hat{w}_i^\top x}{\sqrt{d}}\right) = \text{Tr}[GA]$$

Generalization

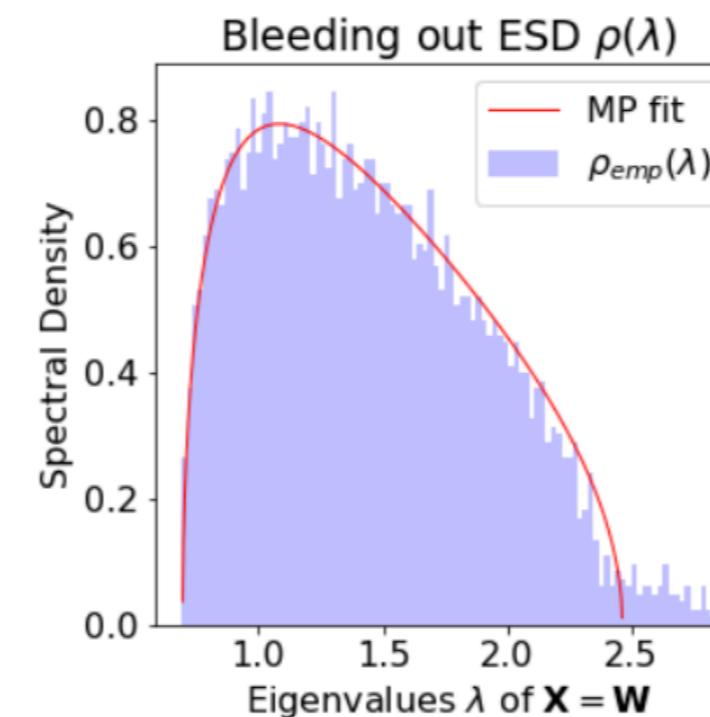
$$\begin{aligned} p &\geq d \\ p^* &= \rho^* d \\ n &= \alpha d^2 \end{aligned}$$

Vignette I : spectrum reproduces deep learning ones!

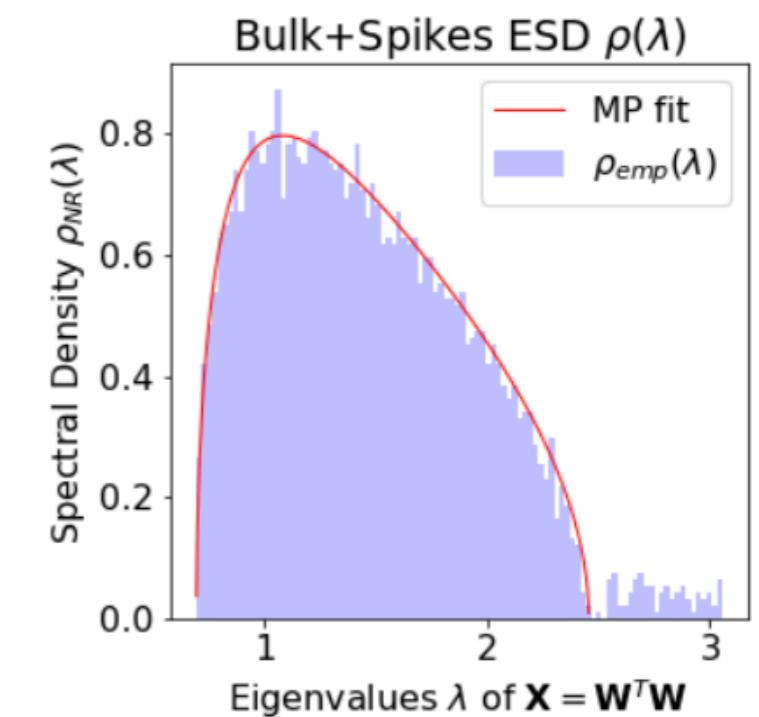
M. Michael W. Mahoney '2020



(a) RANDOM-LIKE.

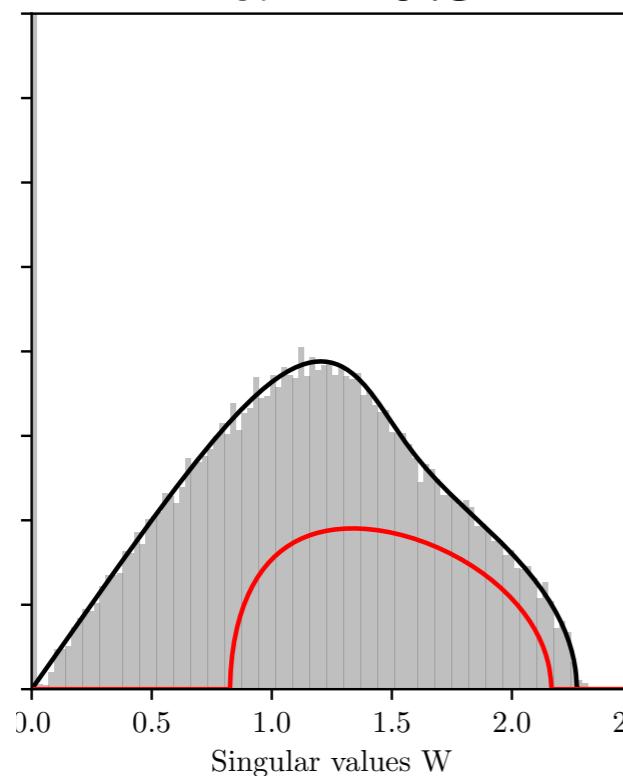


(b) BLEEDING-OUT.

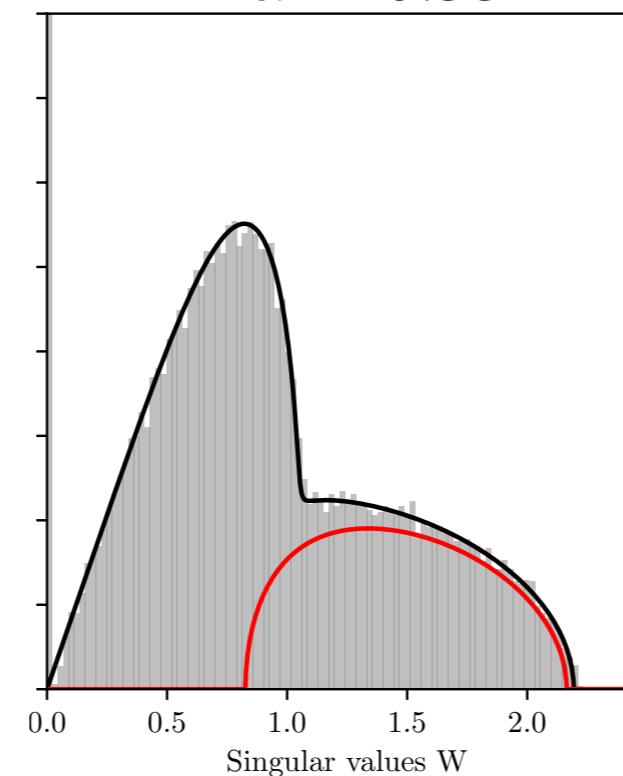


(c) BULK+SPIKES.

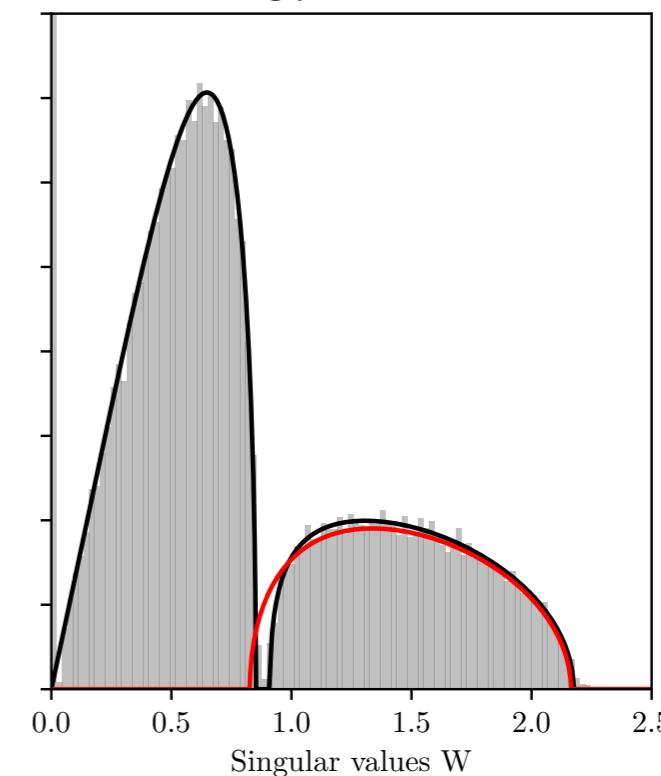
$$\alpha = 0.3$$



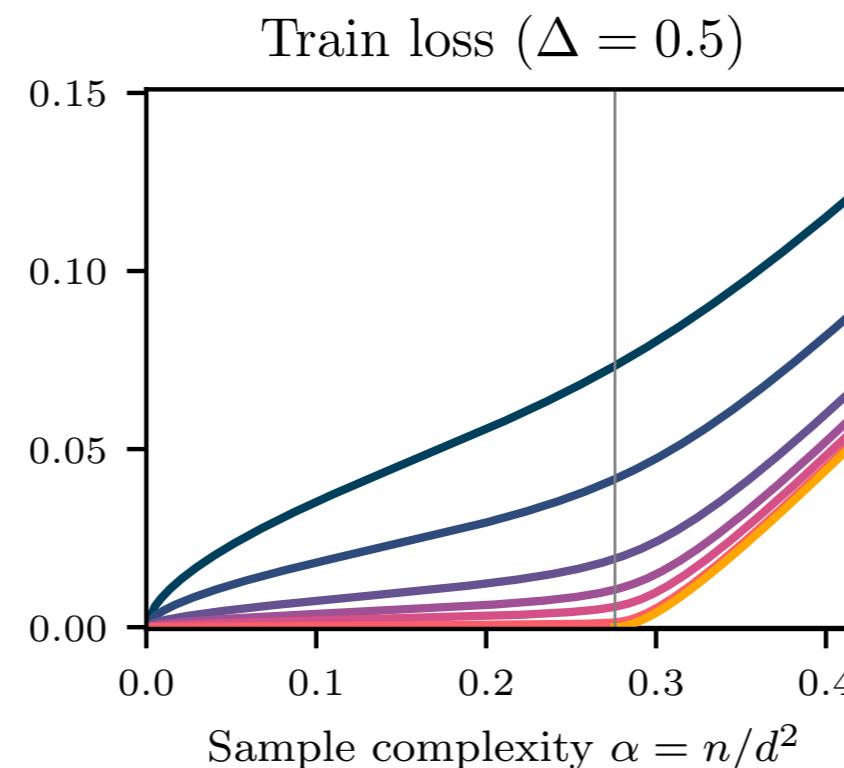
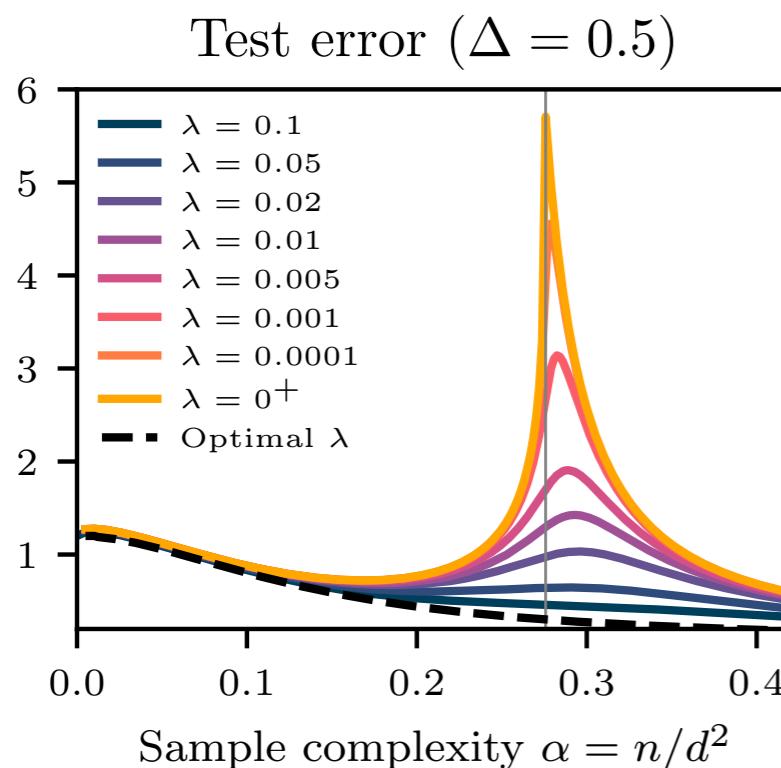
$$\alpha = 0.55$$



$$\alpha = 1$$



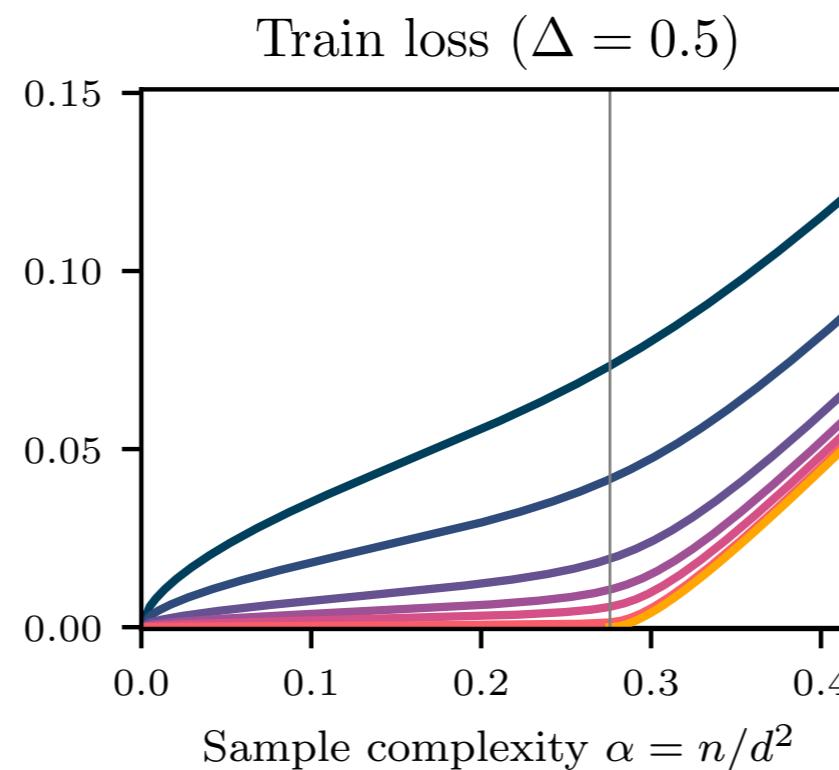
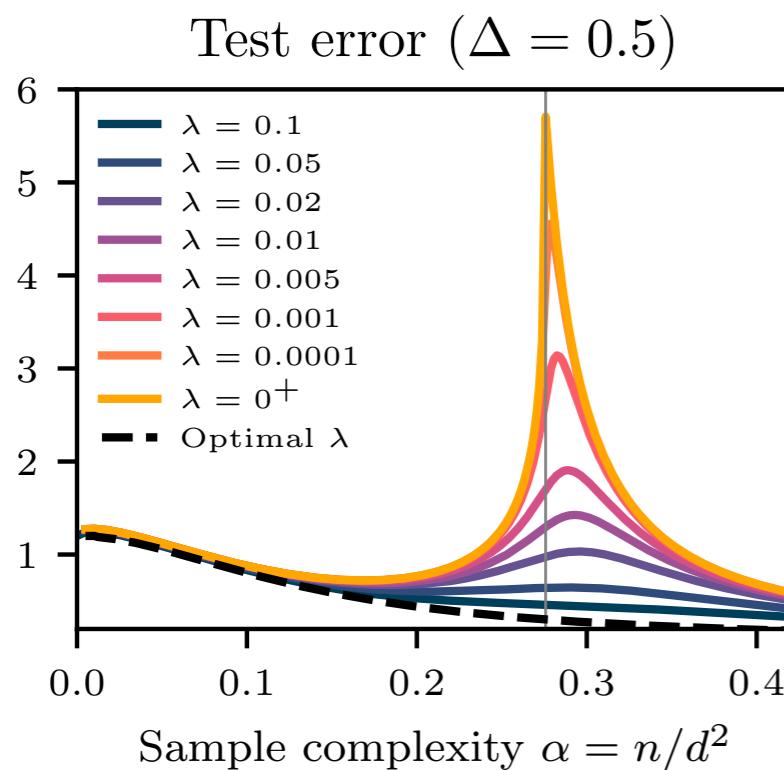
Vignette II : Interpolation threshold & double descent



Interpolation threshold

Largest value of $\alpha = n/d^2$ such that a perfect fit with “zero training loss” (aka, an interpolator) solution exists

Vignette II : Interpolation threshold & double descent



Interpolation threshold

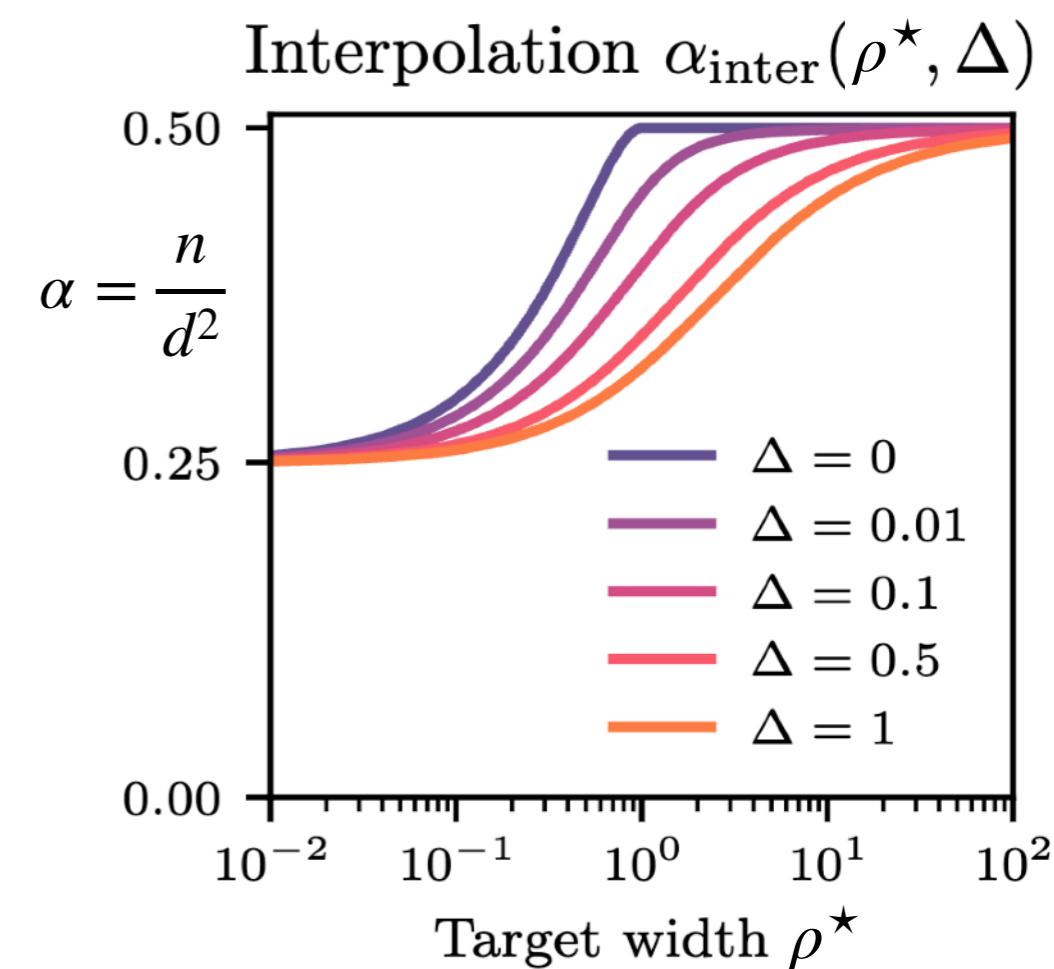
Largest value of $\alpha = n/d^2$ such that a perfect fit with “zero training loss” (aka, an interpolator) solution exists

Location of the interpolation threshold

Here with the objective a Marcenko-Pastur with parameter ρ^*

$$\lim_{\Delta \rightarrow 0^+} \alpha_{\text{inter}}(\rho^*, \Delta) = \frac{1}{4} \begin{cases} 1 + 2\rho^* - \rho^{*2}, & \text{if } 0 < \rho^* < 1, \\ 2, & \text{if } \rho^* > 1. \end{cases}$$

$$\lim_{\Delta \rightarrow \infty} \alpha_{\text{inter}}(\rho^*, \Delta) = \frac{1}{4} \quad \text{as in [Maillard and Bandeira '24]}$$

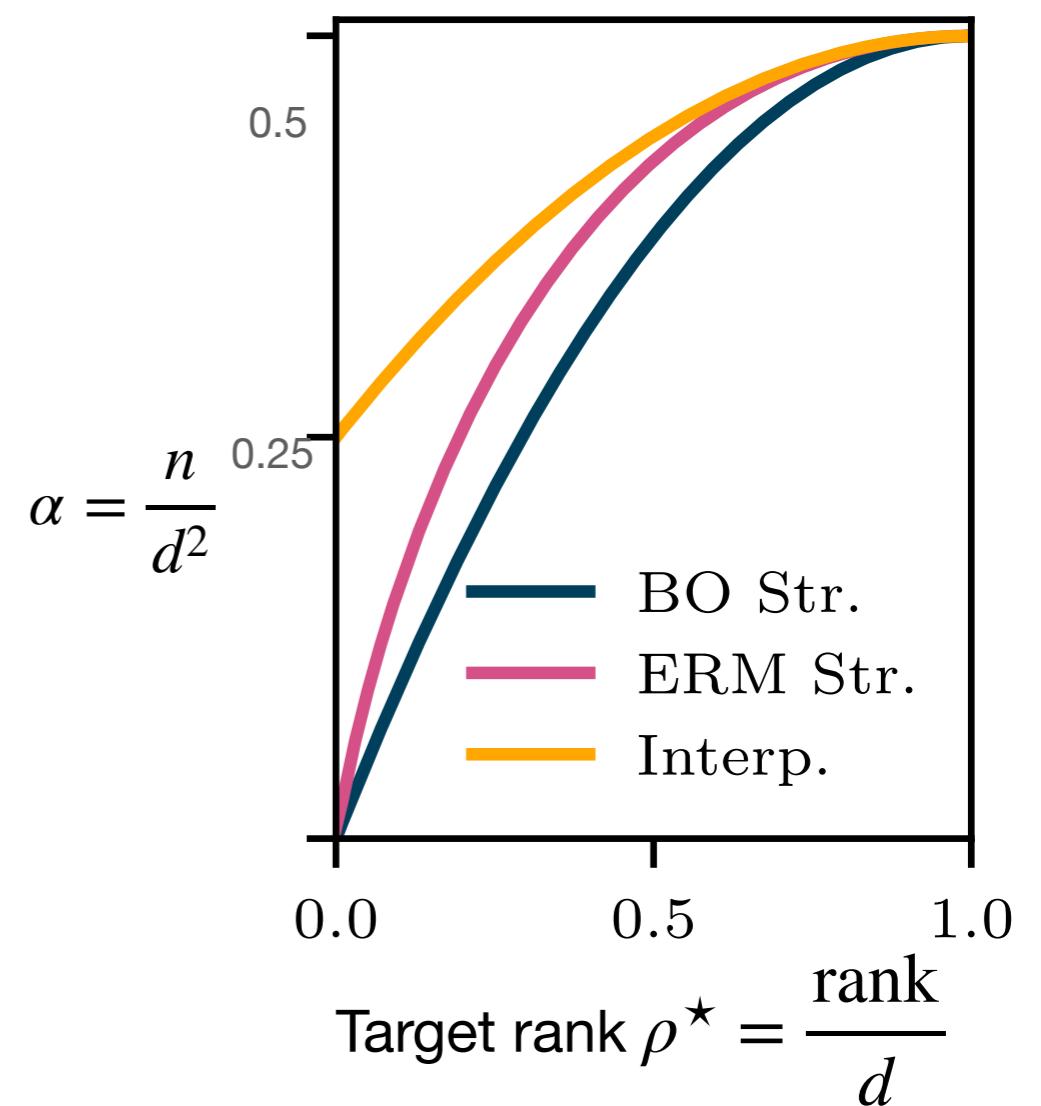


Vignette III : Perfect recovery ℓ_1 versus ℓ_2

Strong recovery:

Number of measurement to reach perfect recovery with noiseless measurements

Thresholds at $\Delta = 0$

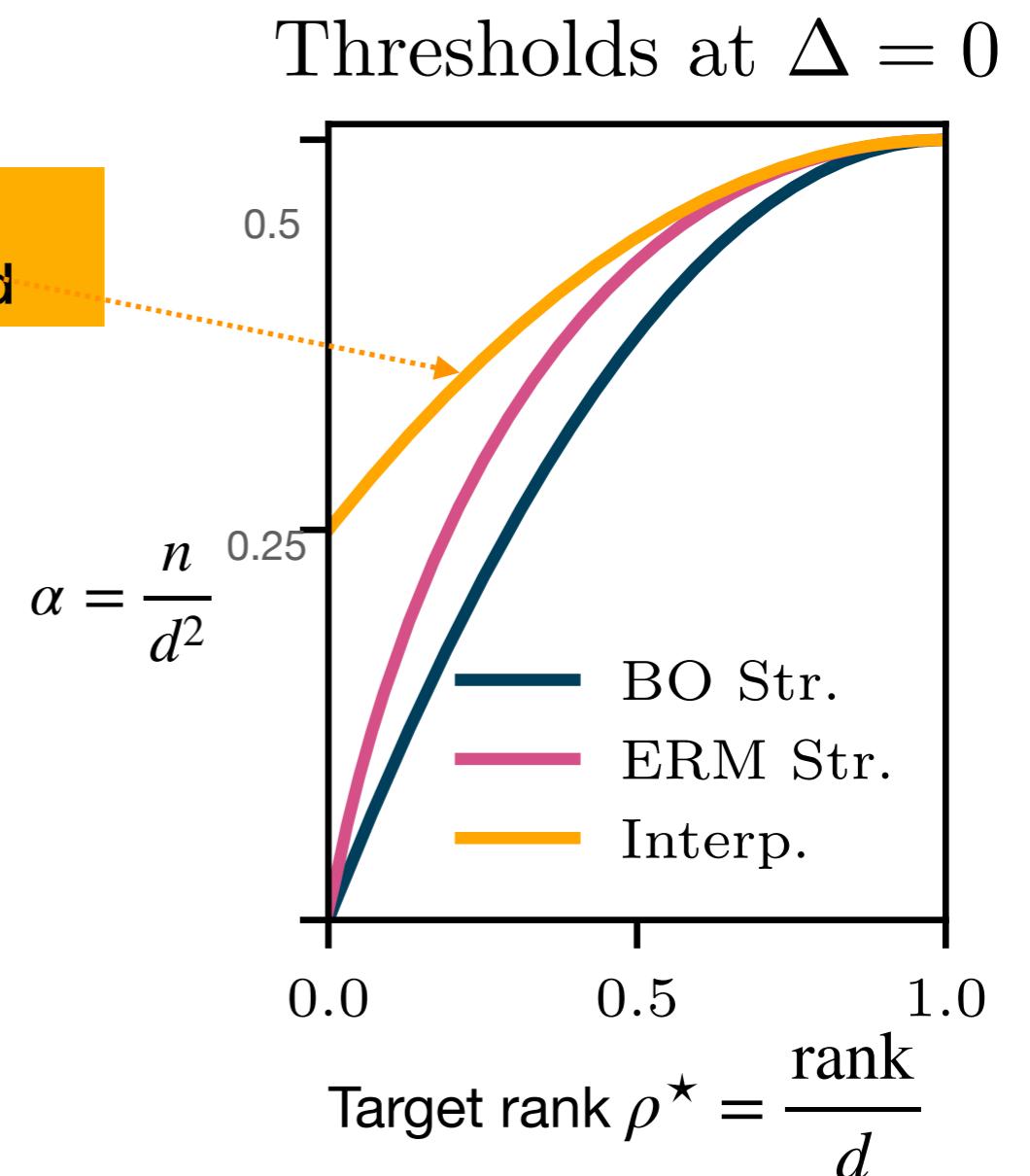


Vignette III : Perfect recovery ℓ_1 versus ℓ_2

Strong recovery:

Number of measurement to reach perfect recovery with noiseless measurements

NTK, KernelLazy...
Interpolation Threshold



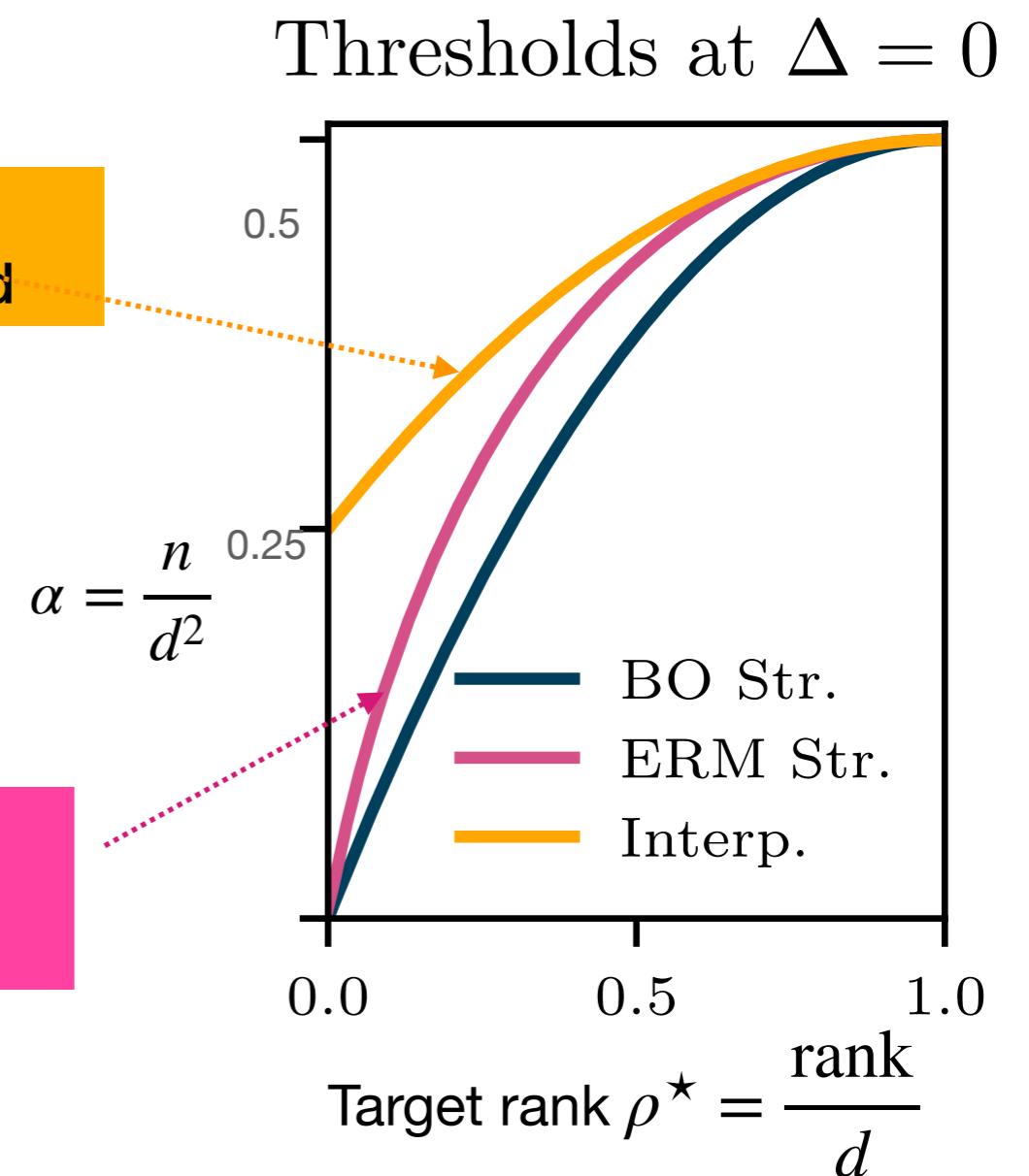
Vignette III : Perfect recovery ℓ_1 versus ℓ_2

Strong recovery:

Number of measurement to reach perfect recovery with noiseless measurements

NTK, KernelLazy...
Interpolation Threshold

Neural net:
Min ℓ_1 interpolator



Vignette III : Perfect recovery ℓ_1 versus ℓ_2

Strong recovery:

Number of measurement to reach perfect recovery with noiseless measurements

NTK, Kernel
(aka ℓ_2)

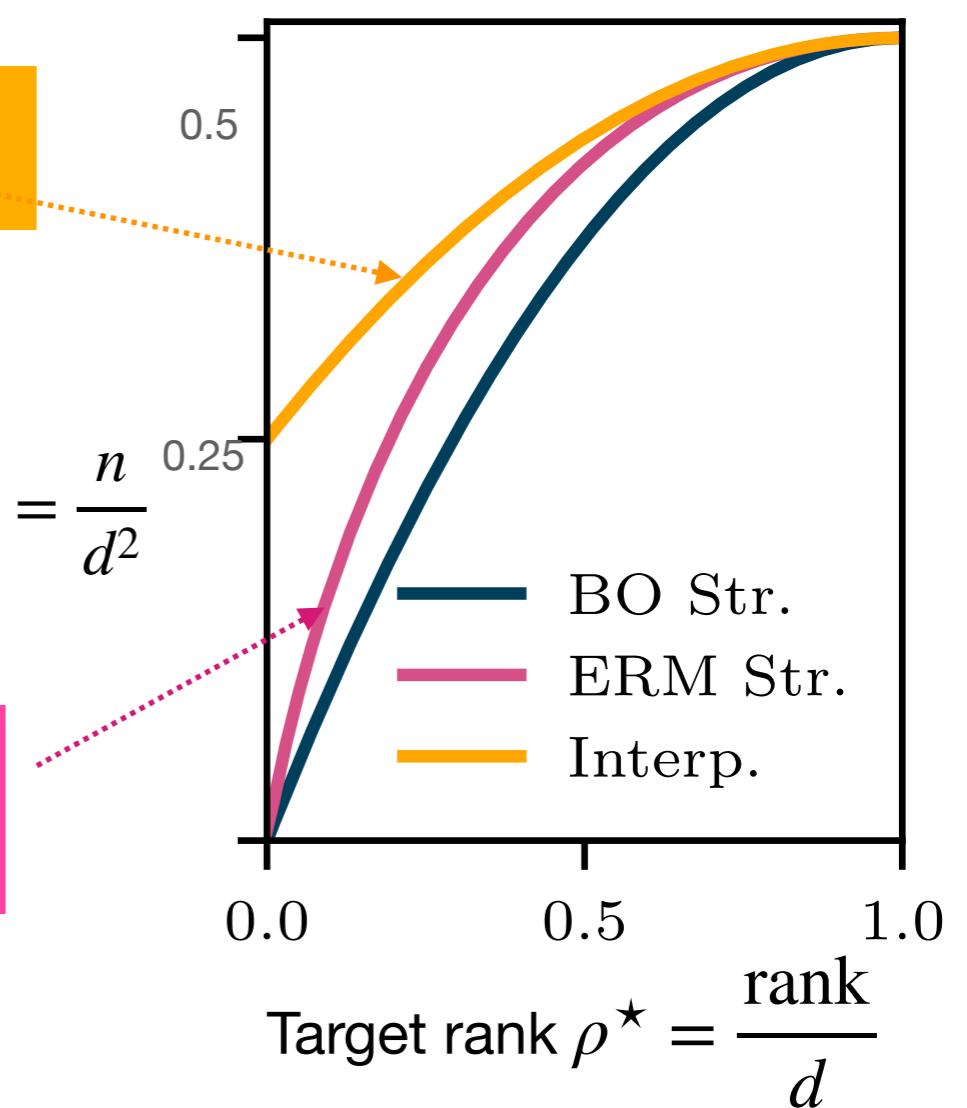
$$\alpha_C \rightarrow \frac{1}{4}$$

Neural net:
Min ℓ_1 interpolator

NTK, KernelLazy...
Interpolation Threshold

$$\alpha = \frac{n}{d^2}$$

Thresholds at $\Delta = 0$



Vignette III : Perfect recovery ℓ_1 versus ℓ_2

Strong recovery:

Number of measurement to reach perfect recovery with noiseless measurements

NTK, Kernel
(aka ℓ_2)

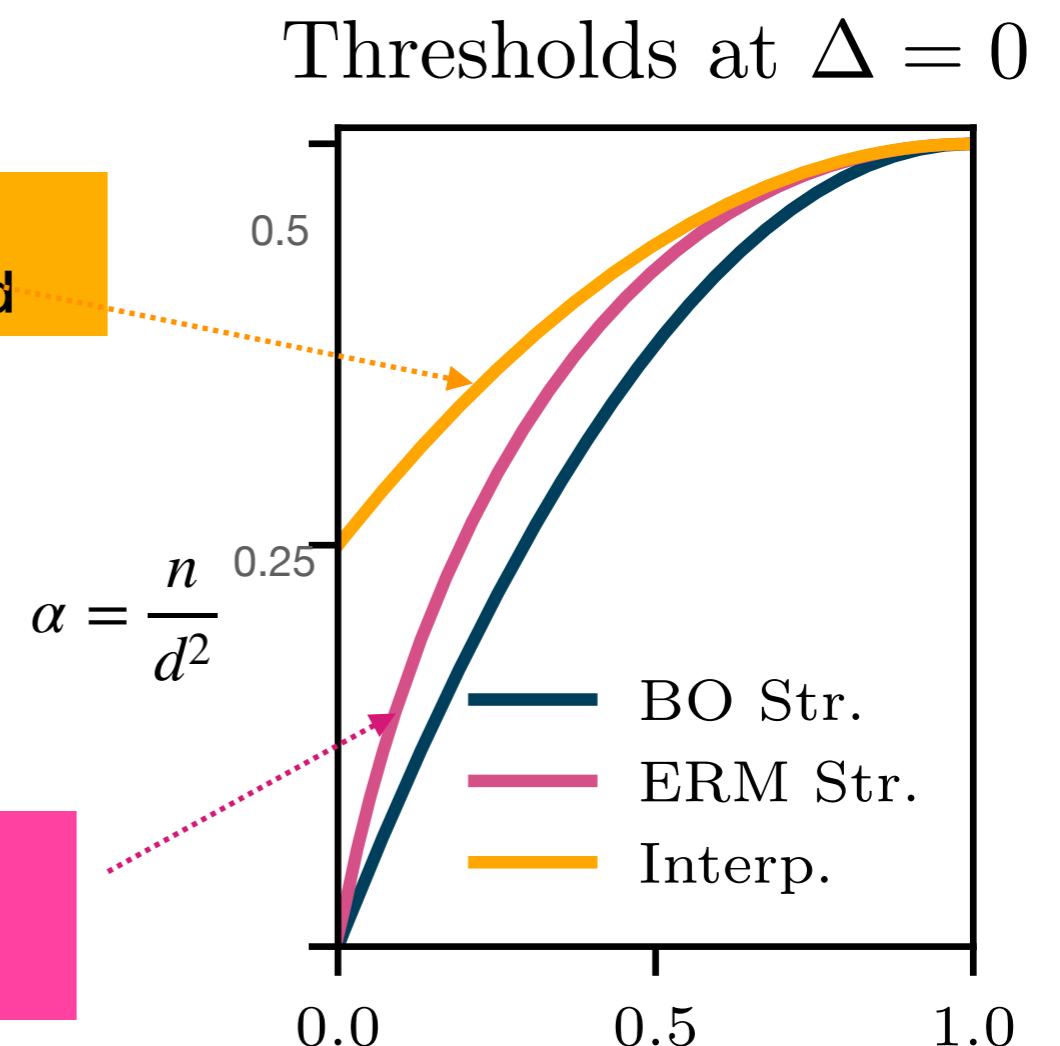
$$\alpha_C \rightarrow \frac{1}{4}$$

Neural net:
Min ℓ_1 interpolator

Neural Networks
(aka ℓ_1)

$$\alpha_C \rightarrow 3\rho^* \quad \rho_{\text{eff}} \rightarrow \rho^{\star^{3/5}} \quad \text{Target rank } \rho^* = \frac{\text{rank}}{d}$$

NTK, KernelLazy...
Interpolation Threshold



Methodologically BO attention is analogous to quadratic network

Back to attention $X \in \mathbb{R}^{T \times d}, X_{a\mu} \sim \mathcal{N}(0, 1_d)$

$$y_\mu = \text{hardmax}(X_\mu A X_\mu^\top)$$



$$y = \begin{pmatrix} E & L & V & F \\ E & 0 & 1 & 0 & 0 \\ L & 0 & 0 & 0 & 1 \\ V & 1 & 0 & 0 & 0 \\ F & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$\begin{array}{ll} X \in \mathbb{R}^{T \times d} & p = \rho d \\ \text{rank}(A) = p & d \rightarrow \infty \end{array}$$

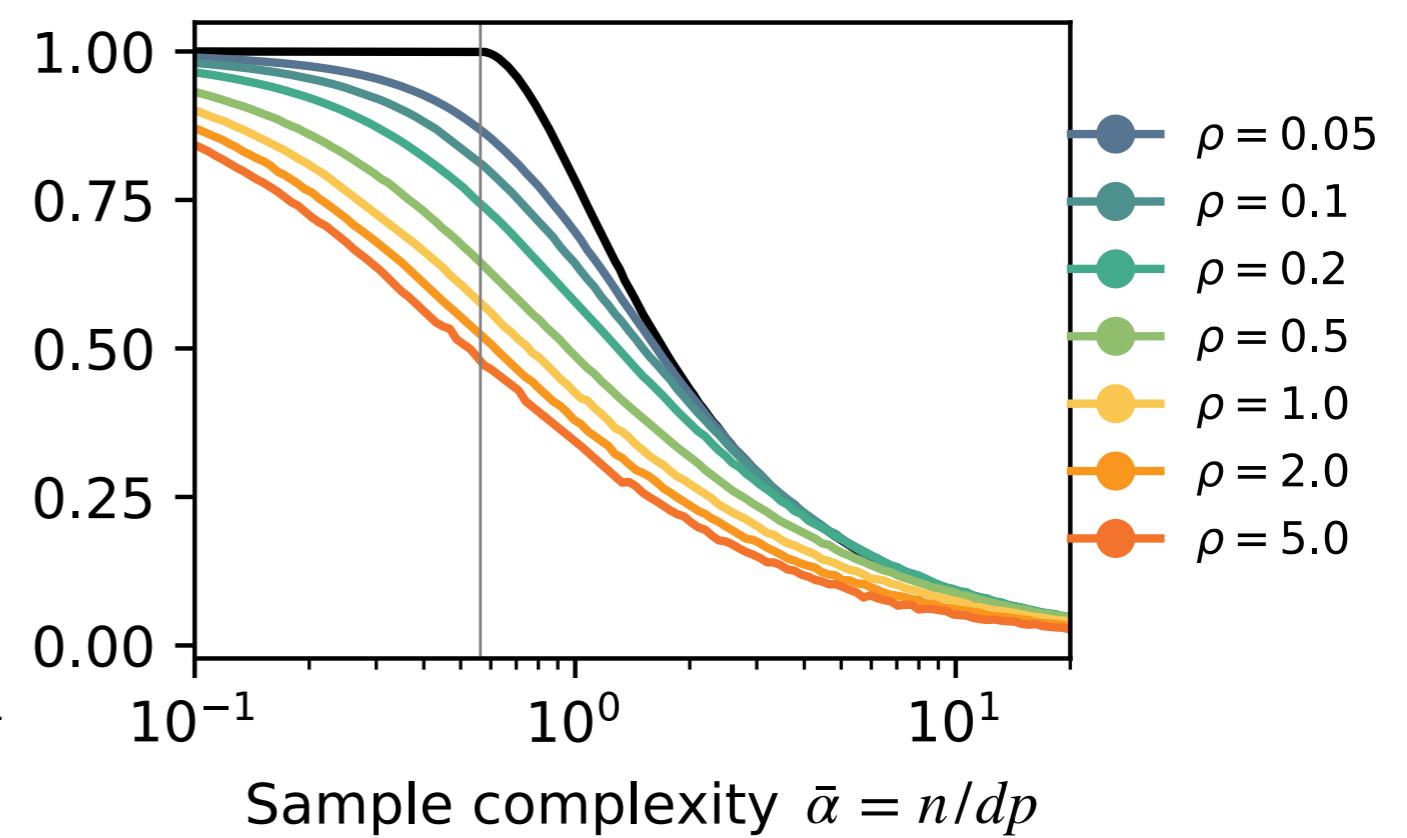
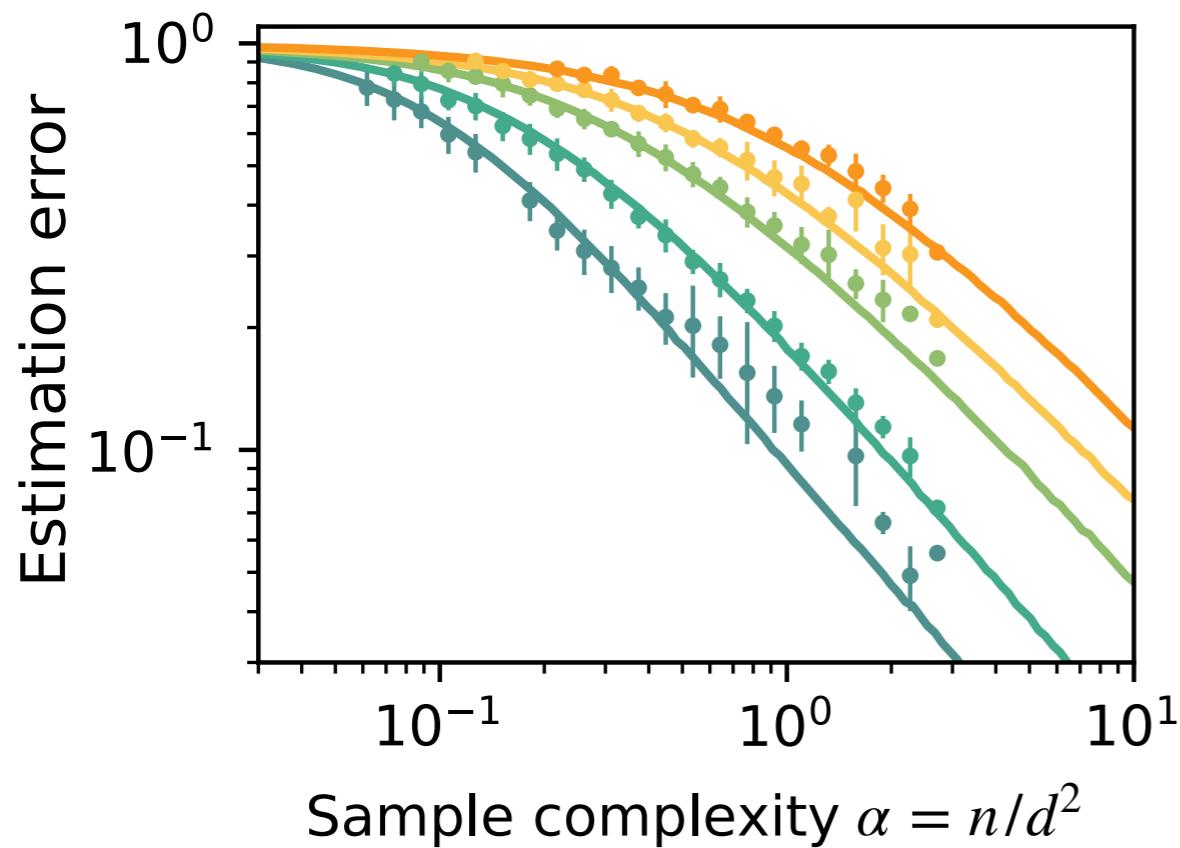
Methodologically BO attention is analogous to quadratic network

Back to attention $X \in \mathbb{R}^{T \times d}, X_{a\mu} \sim \mathcal{N}(0, 1_d)$

$$y_\mu = \text{hardmax}(X_\mu A X_\mu^\top)$$



$$y = \begin{pmatrix} E & L & V & F \\ E & 0 & 1 & 0 & 0 \\ L & 0 & 0 & 0 & 1 \\ V & 1 & 0 & 0 & 0 \\ F & 0 & 0 & 1 & 0 \end{pmatrix}$$



$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$

$$\text{rank}(A) = p \quad d \rightarrow \infty$$

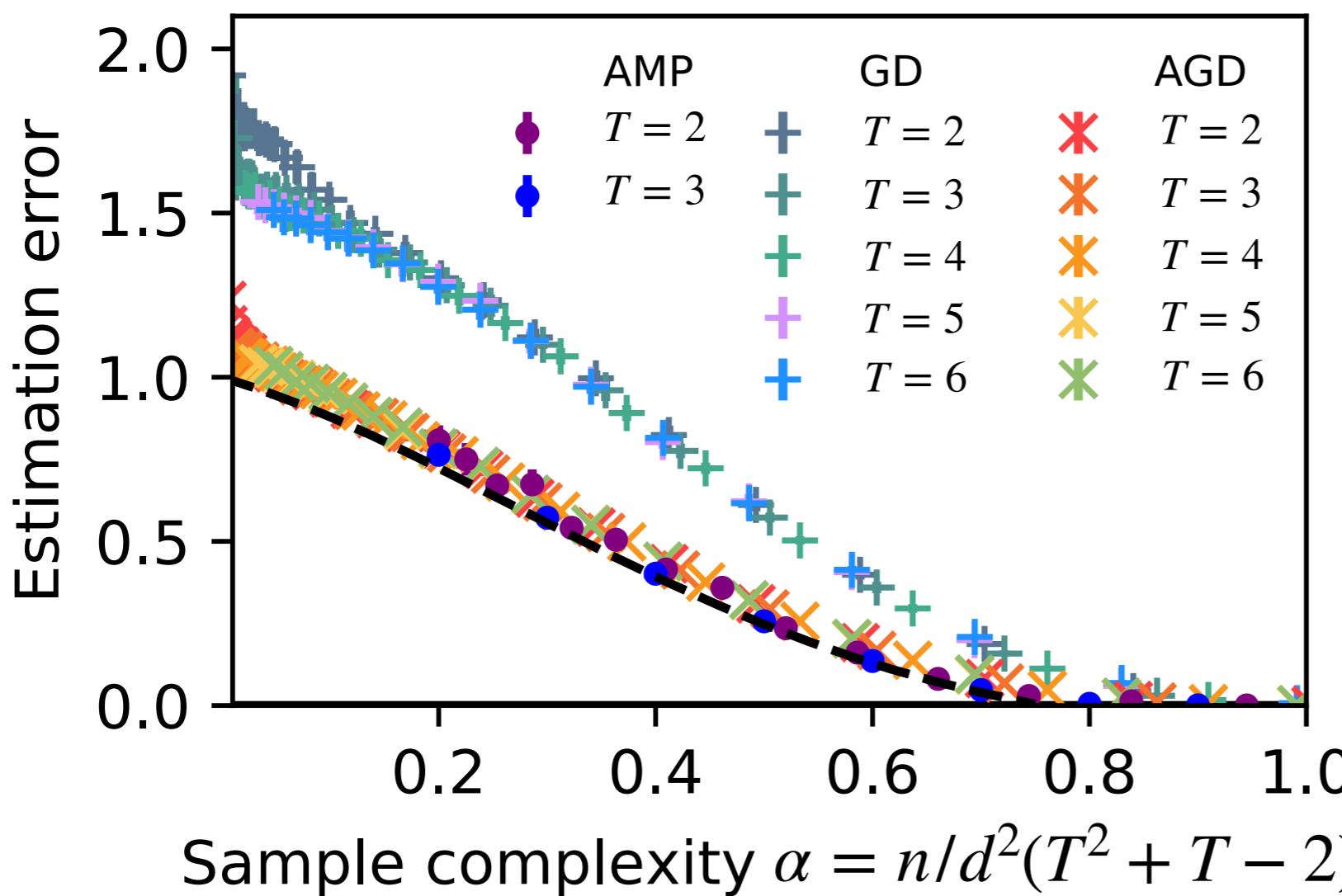
Methodologically BO attention is analogous to quadratic network

Back to attention

$$y_\mu = \text{softmax}(X_\mu A X_\mu^\top)$$

$$\text{Estimation error} \xrightarrow{d \rightarrow \infty} \frac{2\alpha\rho}{\hat{q}}$$

$$1 - \alpha(T^2 + T - 2) = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$



$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$
$$\text{rank}(A) = p \quad d \rightarrow \infty$$

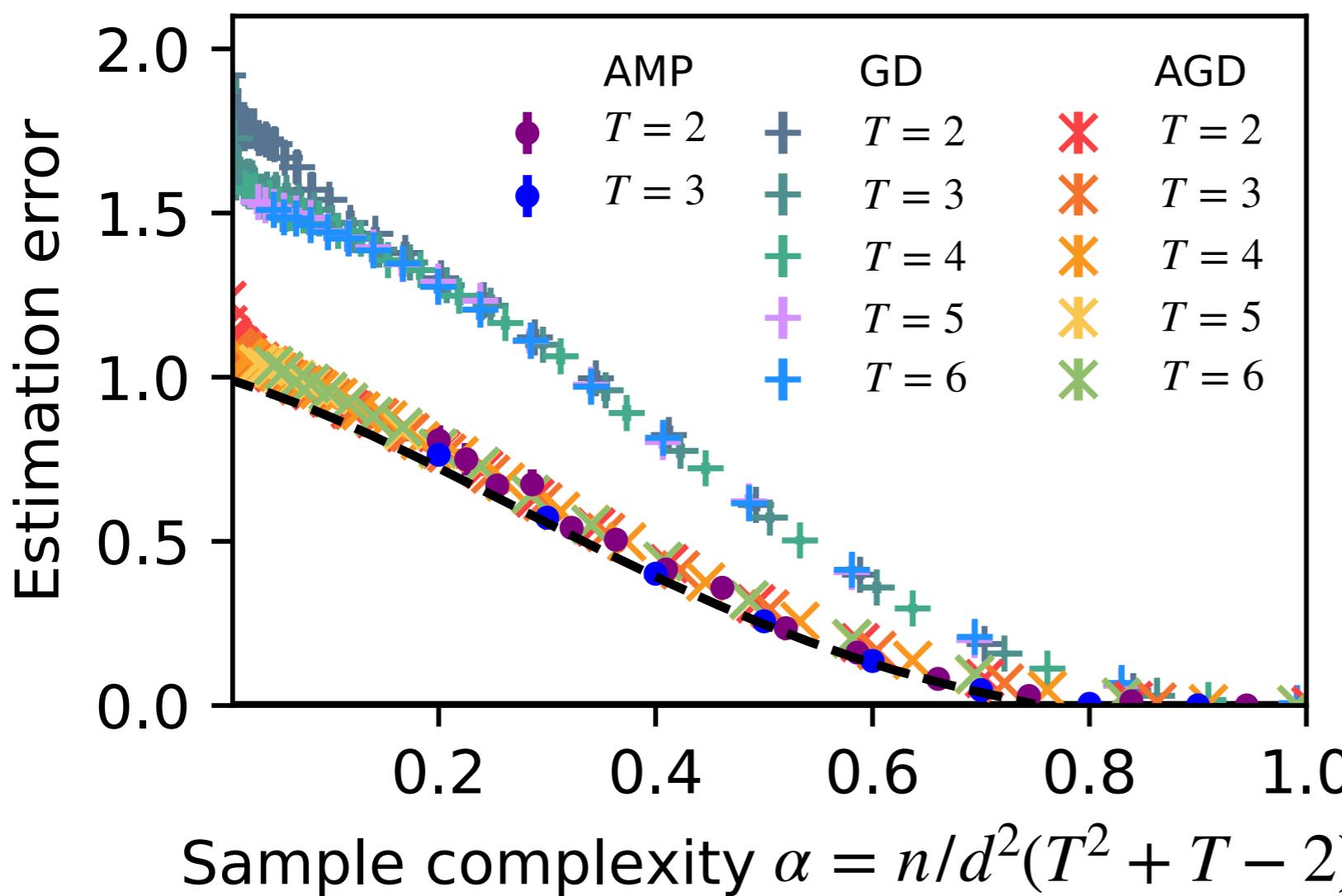
Methodologically BO attention is analogous to quadratic network

Back to attention

$$y_\mu = \text{softmax}(X_\mu A X_\mu^\top)$$

$$\text{Estimation error} \xrightarrow{d \rightarrow \infty} \frac{2\alpha\rho}{\hat{q}}$$

$$1 - \alpha(T^2 + T - 2) = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$



Simple scaling for **tokens**
 $n \rightarrow n/(T^2 + T - 2)$

$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$
$$\text{rank}(A) = p \quad d \rightarrow \infty$$

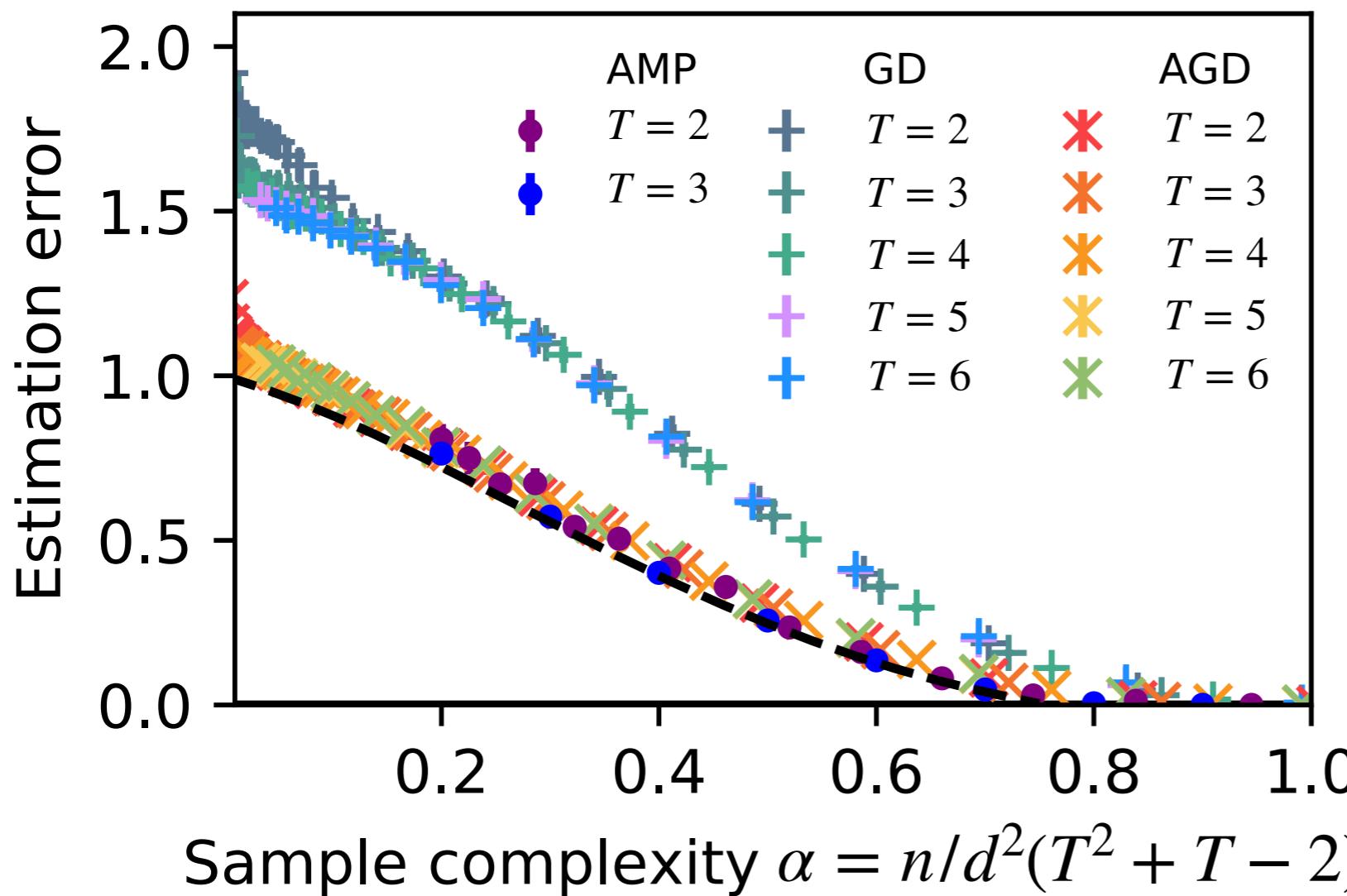
Methodologically BO attention is analogous to quadratic network

Back to attention

$$y_\mu = \text{softmax}(X_\mu A X_\mu^\top)$$

$$\text{Estimation error} \xrightarrow{d \rightarrow \infty} \frac{2\alpha\rho}{\hat{q}}$$

$$1 - \alpha(T^2 + T - 2) = \frac{4\pi^2}{3\hat{q}} \int \mu_{\mathcal{Y}}(\lambda)^3 d\lambda$$



Simple scaling for **tokens**
 $n \rightarrow n/(T^2 + T - 2)$

No effect of $\beta < \infty$



It's a BO artifact

$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$
$$\text{rank}(A) = p \quad d \rightarrow \infty$$

What about multiple layers ?

Single head, no MLP layer

$$y = \sigma \left(X_L \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_L^\top \right)$$

$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell \sum_{i=1}^p \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\ell^\top \right) \right] X_\ell$$

$$\begin{aligned} X &\in \mathbb{R}^{T \times d} & p &= \rho d \\ \text{rank}(A) &= p & d &\rightarrow \infty \end{aligned}$$

What about multiple layers ?

Single head, no MLP layer

$$y = \sigma \left(X_L \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_L^\top \right)$$

$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell \sum_{i=1}^p \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\ell^\top \right) \right] X_\ell$$

$$y = g \left(X \sum_{i=1}^p \mathbf{w}_i^1 \mathbf{w}_i^{1\top} X^\top, \dots, X \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X^\top \right)$$



$$\begin{aligned} X &\in \mathbb{R}^{T \times d} & p &= \rho d \\ \text{rank}(A) &= p & d &\rightarrow \infty \end{aligned}$$

What about multiple layers ?

Single head, no MLP layer

$$y = \sigma \left(X_L \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_L^\top \right)$$
$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell \sum_{i=1}^p \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\ell^\top \right) \right] X_\ell$$
$$y = g \left(X \sum_{i=1}^p \mathbf{w}_i^1 \mathbf{w}_i^{1\top} X^\top, \dots, X \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X^\top \right) = g \left(X A^1 X^\top, \dots, X A^L X^\top \right)$$

$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$
$$\text{rank}(A) = p \quad d \rightarrow \infty$$

What about multiple layers ?

Single head, no MLP layer

$$y = \sigma \left(X_L \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_L^\top \right)$$
$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell \sum_{i=1}^p \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\ell^\top \right) \right] X_\ell$$
$$y = g \left(X \sum_{i=1}^p \mathbf{w}_i^1 \mathbf{w}_i^{1\top} X^\top, \dots, X \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X^\top \right) = g \left(X A^1 X^\top, \dots, X A^L X^\top \right)$$

Shown by induction: assume $X_\ell = g_\ell \left(X A^1 X^\top, \dots, X A^\ell X^\top \right) X$

$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell A^\ell X_\ell^\top \right) \right] X_\ell = \left[c\mathbf{1}_T + \sigma \left(g_\ell X A^\ell X^\top g_\ell^\top \right) \right] g_\ell X$$

$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$
$$\text{rank}(A) = p \quad d \rightarrow \infty$$

What about multiple layers ?

Single head, no MLP layer

$$y = \sigma \left(X_L \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X_L^\top \right)$$
$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell \sum_{i=1}^p \mathbf{w}_i^\ell \mathbf{w}_i^{\ell\top} X_\ell^\top \right) \right] X_\ell$$
$$y = g \left(X \sum_{i=1}^p \mathbf{w}_i^1 \mathbf{w}_i^{1\top} X^\top, \dots, X \sum_{i=1}^p \mathbf{w}_i^L \mathbf{w}_i^{L\top} X^\top \right) = g \left(X A^1 X^\top, \dots, X A^L X^\top \right)$$

Shown by induction: assume $X_\ell = g_\ell \left(X A^1 X^\top, \dots, X A^\ell X^\top \right) X$

$$X_{\ell+1} = \left[c\mathbf{1}_T + \sigma \left(X_\ell A^\ell X_\ell^\top \right) \right] X_\ell = \left[c\mathbf{1}_T + \sigma \left(g_\ell X A^\ell X^\top g_\ell^\top \right) \right] g_\ell X$$



Hinges on new problem:
Multiplexed denoising

$$\mathcal{Y}_\ell = \mathcal{S}_\ell + \sum_{k=1}^L \sqrt{\delta}_{\ell k} G_k$$

$$X \in \mathbb{R}^{T \times d} \quad p = \rho d$$
$$\text{rank}(A) = p \quad d \rightarrow \infty$$

Thanks for your attention!

Bayes-optimal learning of an extensive-width neural network from quadratically many samples

A. Maillard, E. Troiani, S. Martin, F. Krzakala, L. Zdeborová
arXiv:2408.03733

The Nuclear Route: Sharp Asymptotics of ERM in Overparameterized Quadratic Networks

V. Erba, E. Troiani, L. Zdeborová, F. Krzakala
arXiv:2505.17958

Bayes optimal learning of attention-indexed models

F. Boncoraglio, E. Troiani, V. Erba, L. Zdeborová
arXiv:2506.01582

Thanks for your attention!



Bayes-optimal learning of an extensive-width neural network from quadratically many samples

A. Maillard, E. Troiani, S. Martin, F. Krzakala, L. Zdeborová



The Nuclear Route: Sharp Asymptotics of ERM in Overparameterized Quadratic Networks

V. Erba, E. Troiani, L. Zdeborová, F. Krzakala



Bayes optimal learning of attention-indexed models

F. Boncoraglio, E. Troiani, V. Erba, L. Zdeborová