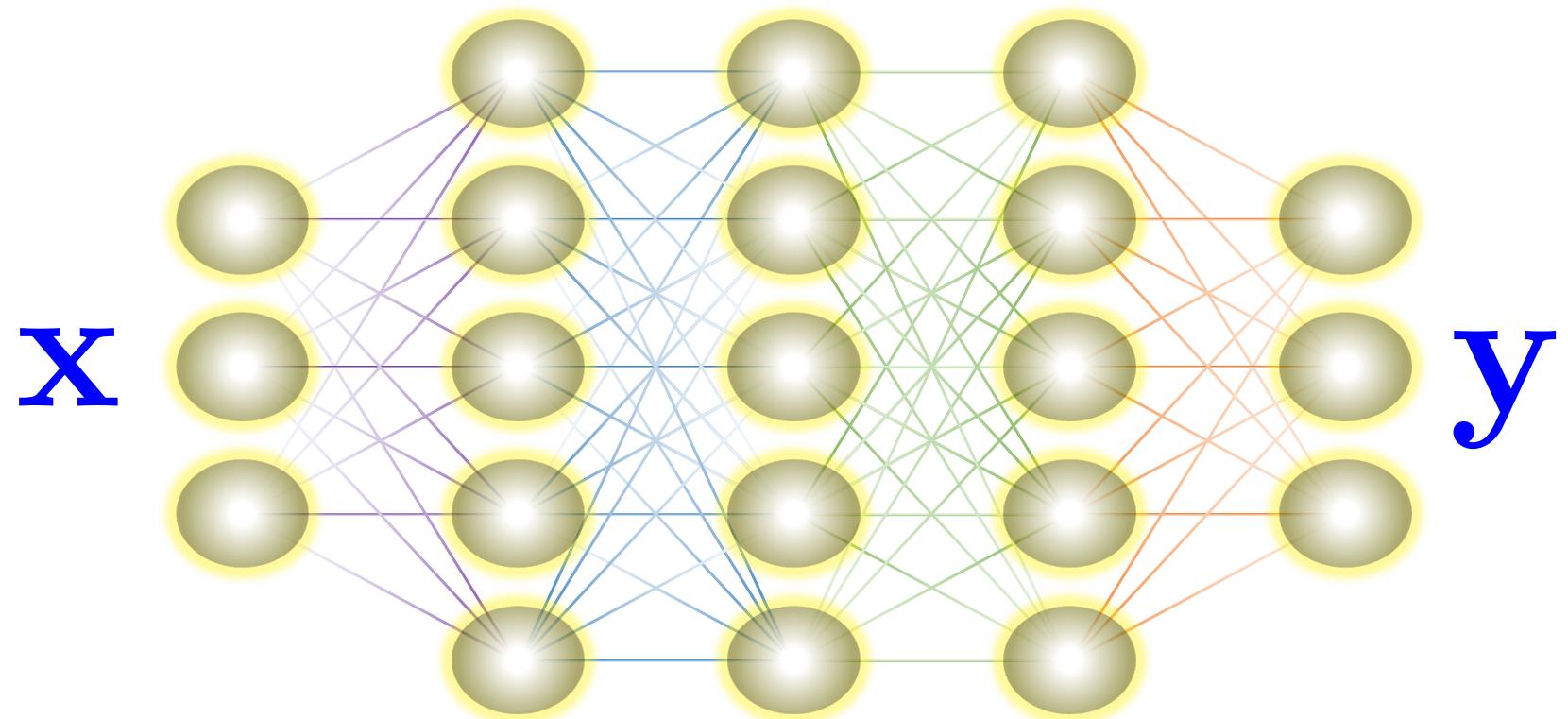


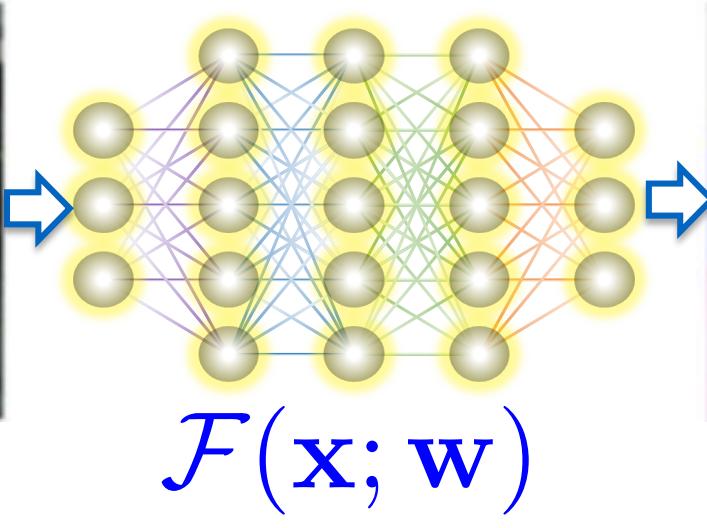
# **Optimization for Learning**

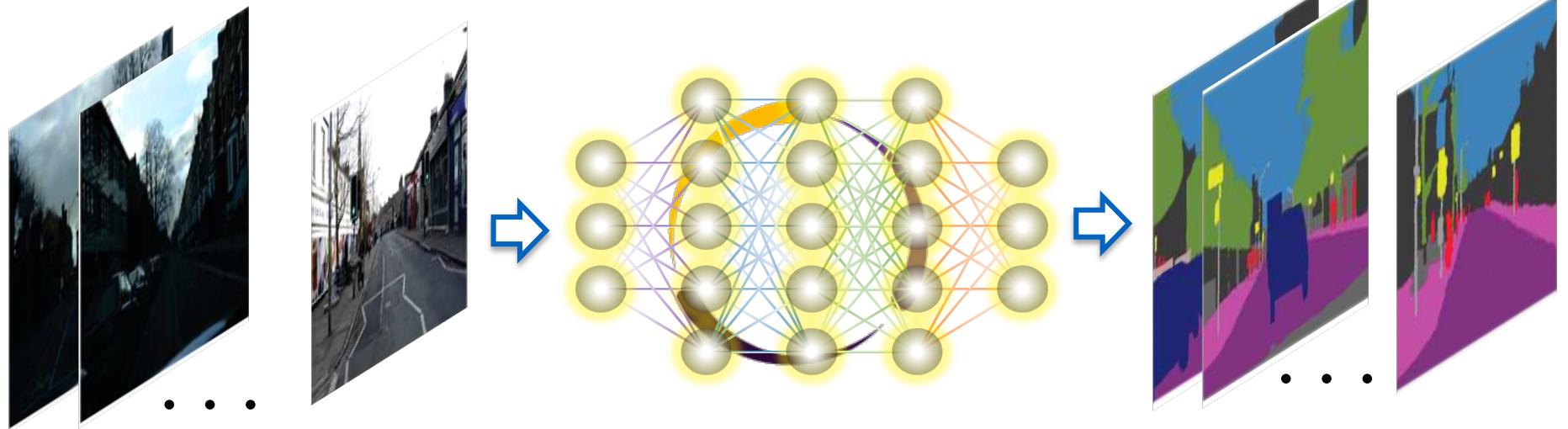
Instructor - Simon Lucey  
**Maths of AI - 2024**



**AUSTRALIAN  
INSTITUTE FOR  
MACHINE LEARNING**







$$N \propto 10^6 - 10^9$$

# Optimization and Learning

---

$$\min_{\mathbf{w}} \sum_{n=0}^{N-1} \mathcal{L}[y_n, \mathcal{F}(\mathbf{x}_n; \mathbf{w})]$$

**What is the simplest network and loss?**

$$\mathcal{L}[y, x] = ||y - x||_2^2 \quad \leftarrow \text{“least-squares loss”}$$

$$\mathcal{F}(\mathbf{x}; \mathbf{w}) = \mathbf{x}^T \mathbf{w} \quad \leftarrow \text{“linear network”}$$

# Today

---

- **Regularized Least-Squares**
- Gradient Descent
- Optimizing Deep Networks

# Least-Squares Linear Model

---

$$\min_{\mathbf{w}} \frac{1}{2} \left| \left| \mathbf{y} - \mathbf{X} \mathbf{w} \right| \right|_2^2$$

$$\mathbf{y} = [y_0, \dots, y_{N-1}]^T$$

$$\mathbf{X} = [\mathbf{x}_0, \dots, \mathbf{x}_{N-1}]^T$$

## Least-Squares Linear Model

---

$$\min_{\mathbf{w}} \frac{1}{2} \left| \left| \mathbf{y} - \mathbf{X}\mathbf{w} \right| \right|_2^2 + \Omega(\mathbf{w})$$

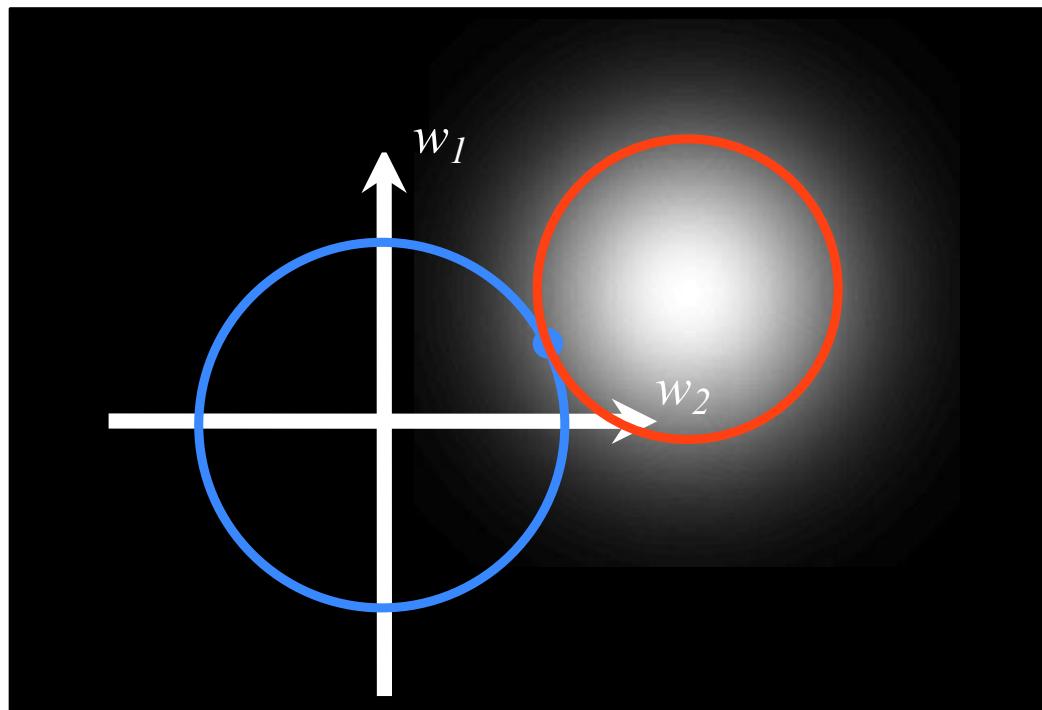
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$
 “closed form”

**What if  $(\mathbf{X}^T \mathbf{X})$ is non-singular?**

# Ridge Regression

---

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \text{ s.t. } \|\mathbf{w}\|_2^2 \leq \epsilon$$



## Ridge Regression – Closed Form

---

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$
 “closed form solution”

$$\mathbf{X} = \mathbf{U} \text{diag}(\mathbf{s}) \mathbf{V}^T$$
 “SVD”

$$\mathbf{U} \mathbf{U}^T = \mathbf{I} \quad \mathbf{V} \mathbf{V}^T = \mathbf{I} \quad \mathbf{s} \geq 0$$

## Ridge Shrinkage

---

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

$$\mathbf{w} = \mathbf{V} \text{diag}(\hat{\mathbf{s}}^{-1}) \mathbf{U}^T \mathbf{y}$$
 “closed form solution”

$$\hat{\mathbf{s}}_{\text{rdg}}^{-1} = \frac{s}{s^2 + \lambda}$$
 “ridge shrinkage”

**Drawback: attenuates all principal components!!**

## PCA Shrinkage

---

$$\mathbf{X} = \mathbf{U} \text{diag}(\mathbf{s}) \mathbf{V}^T \quad \text{"SVD"}$$

$$\mathbf{w} = \mathbf{V} \text{diag}(\hat{\mathbf{s}}^{-1}) \mathbf{U}^T \mathbf{y} \quad \text{"closed form solution"}$$

$$\hat{s}_{\text{pca}}^{-1} = m_{\text{pca}}(s, \kappa) \cdot s^{-1} \quad \text{"masked shrinkage"}$$

$$m_{\text{pca}}(s; \kappa) = \begin{cases} 1, & \text{if } s \geq \kappa \\ 0, & \text{otherwise.} \end{cases} \quad \text{"PCA mask"}$$

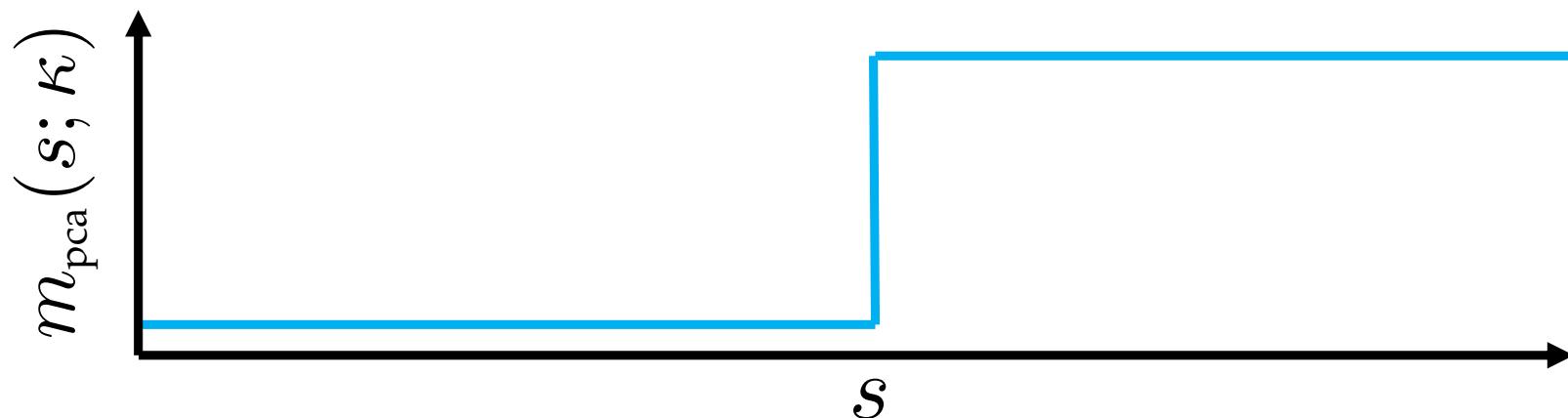
# PCA Shrinkage

---

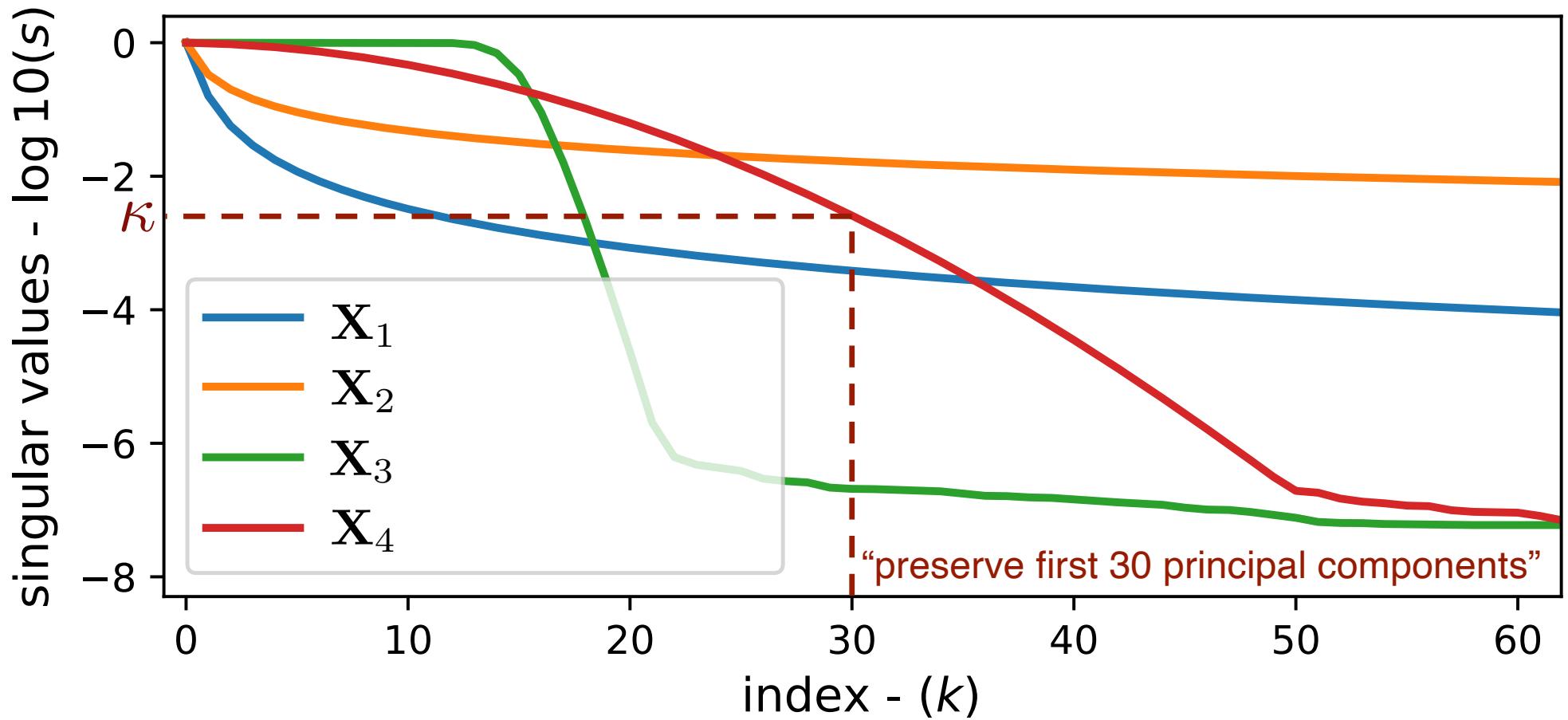
$$\mathbf{X} = \mathbf{U} \text{diag}(\mathbf{s}) \mathbf{V}^T \quad \text{"SVD"}$$

$$\mathbf{w} = \mathbf{V} \text{diag}(\hat{\mathbf{s}}^{-1}) \mathbf{U}^T \mathbf{y} \quad \text{"closed form solution"}$$

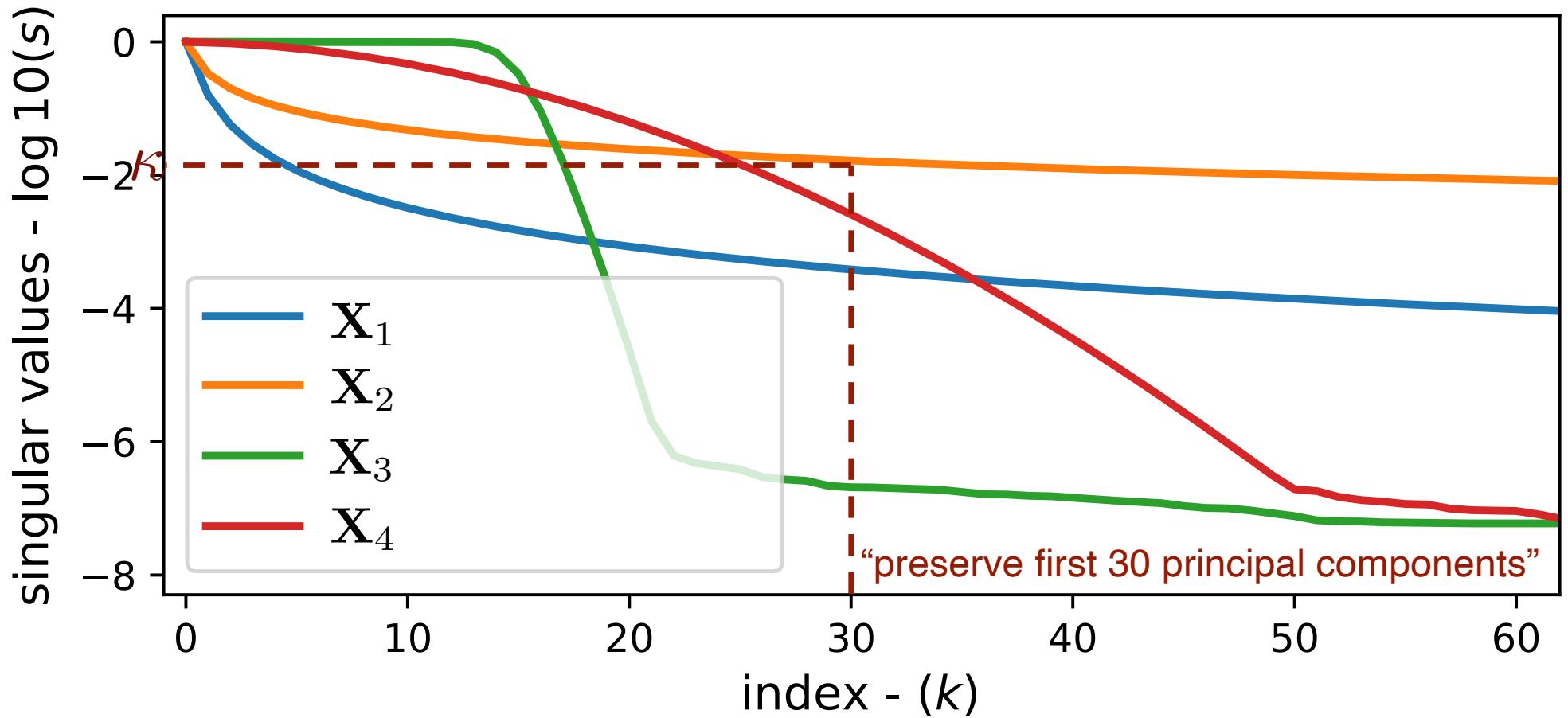
$$\hat{s}_{\text{pca}}^{-1} = m_{\text{pca}}(s, \kappa) \cdot s^{-1} \quad \text{"masked shrinkage"}$$



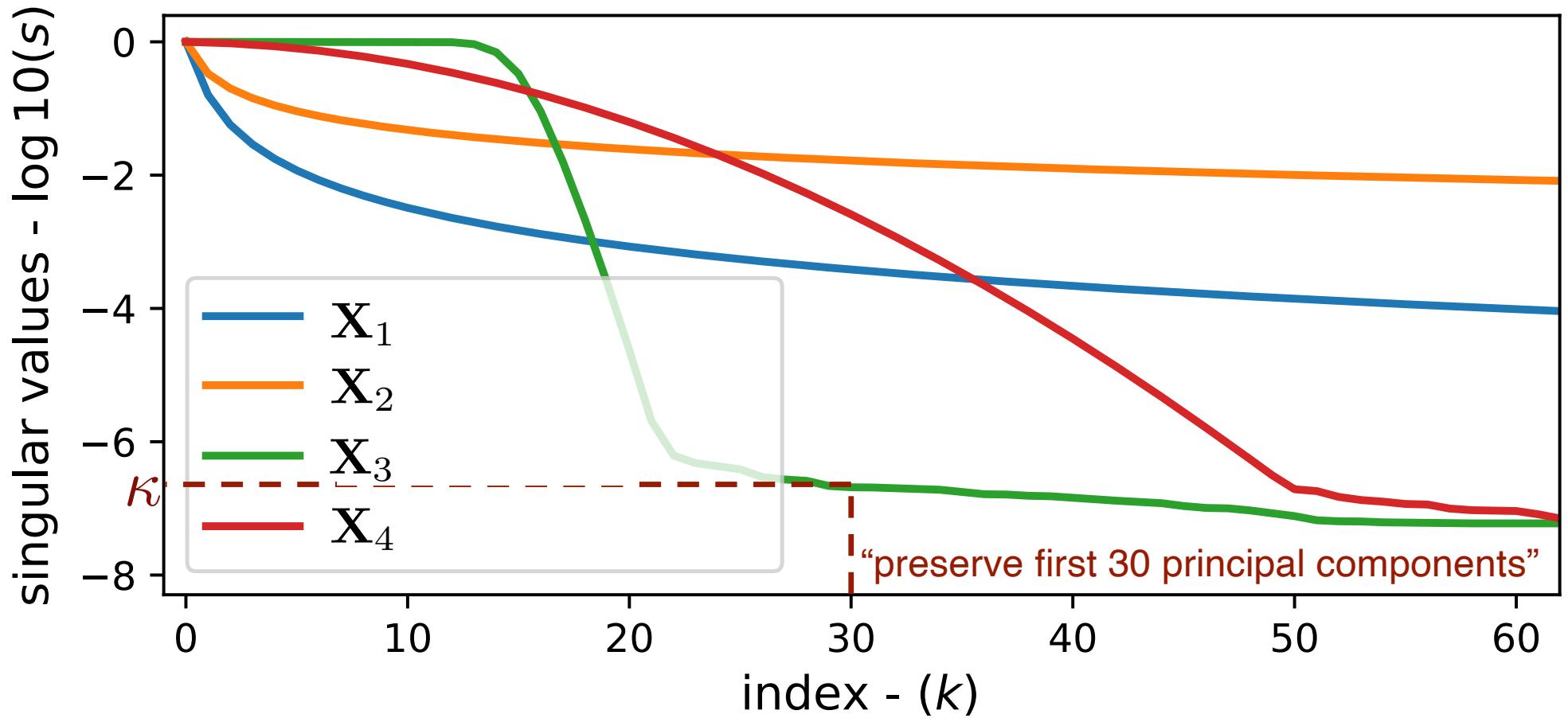
# PCA Shrinkage



# PCA Shrinkage



# PCA Shrinkage



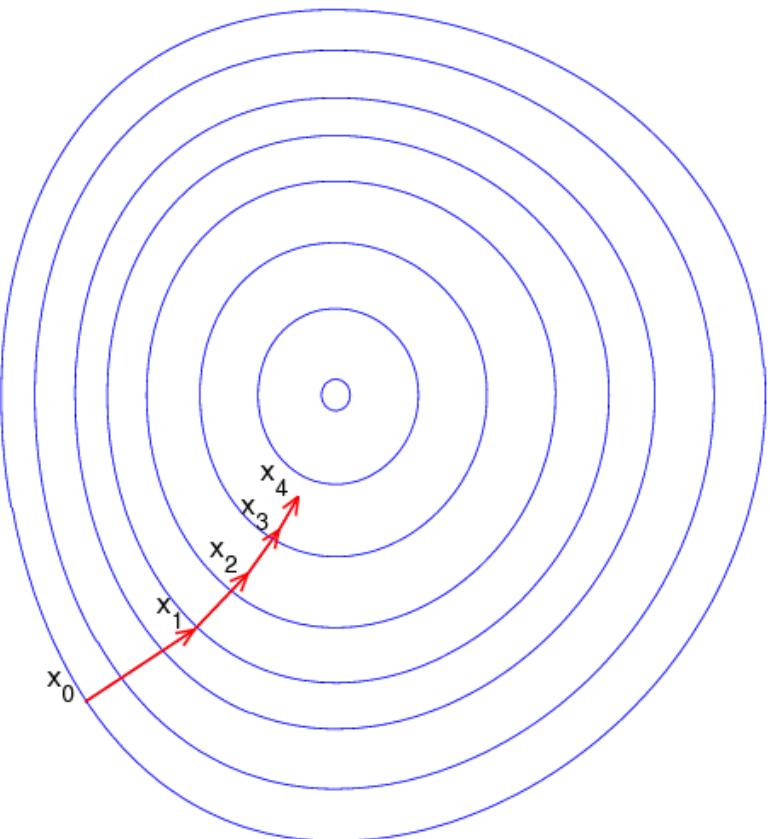
# Today

---

- Regularized Least-Squares
- **Gradient Descent**
- Optimizing Deep Networks

# Gradient Descent

---



**Idea:** just go downhill

**Pros:**

- Simple
- Always improves on current guess

**Cons:**

- Slow, (does not use higher order information)
- Requires objective to be differentiable.
- Heuristics on learning rate and number of iterations

# Gradient Descent

---

$$\mathbf{w} \rightarrow \mathbf{w} - \alpha \cdot \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}^T}$$

"learning rate"

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \left| \left| \mathbf{y} - \mathbf{Xw} \right| \right|_2^2 \quad \text{"loss function"}$$

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}^T} = -\mathbf{X}^T [\mathbf{y} - \mathbf{Xw}] \quad \text{"partial derivative"}$$

Does gradient descent  
implicitly regularize?

## Closed form - Gradient Descent

---

If we start with  $\mathbf{w}_0 = \mathbf{0}$

$$\mathbf{w}_1 = \alpha \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}_2 = (\mathbf{I} - \alpha \mathbf{X}^T \mathbf{X}) \alpha \mathbf{X}^T \mathbf{y} + \alpha \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}_3 = (\mathbf{I} - \alpha \mathbf{X}^T \mathbf{X})^2 \alpha \mathbf{X}^T \mathbf{y} + (\mathbf{I} - \alpha \mathbf{X}^T \mathbf{X}) \alpha \mathbf{X}^T \mathbf{y} + \alpha \mathbf{X}^T \mathbf{y}$$

at the  $q$ -th iteration generalizes to,

$$\mathbf{w}_q = \sum_{k=0}^{q-1} (\mathbf{I} - \alpha \mathbf{X}^T \mathbf{X})^{k-1} \alpha \mathbf{X}^T \mathbf{y}$$

## Closed form - Gradient Descent

---

Using the Neumann series identity,

$$\sum_{k=0}^{q-1} z^k = (1 - z^q)(1 - z)^{-1}$$

Closed form solution,

$$\mathbf{w}_q = [\mathbf{I} - (\mathbf{I} - \alpha \mathbf{X}^T \mathbf{X})^q](\alpha \mathbf{X}^T \mathbf{X})^{-1} \alpha \mathbf{X}^T \mathbf{y}$$

<https://www.cs.cornell.edu/courses/cs6241/2021sp/lec/2021-02-18.pdf>

# Gradient Descent Shrinkage

---

$$\mathbf{X} = \mathbf{U} \text{diag}(\mathbf{s}) \mathbf{V}^T \quad \text{"SVD"}$$

$$\mathbf{w} = \mathbf{V} \text{diag}(\hat{\mathbf{s}}^{-1}) \mathbf{U}^T \mathbf{y} \quad \text{"closed form solution"}$$

$$\hat{s}_{\text{gd}}^{-1} = m_{\text{gd}}(s; \alpha, q) \cdot s^{-1} \quad \text{"masked shrinkage"}$$

$$m_{\text{gd}}(s; \alpha, q) = [1 - (1 - \alpha \cdot s^2)^q] \quad \text{"GD mask"}$$

## Choosing Learning Rate and Iterations

---

Inspecting the masking function,

$$m_{\text{gd}}(s; \alpha, q) = [1 - (1 - \alpha \cdot s^2)^q]$$

we can see that,

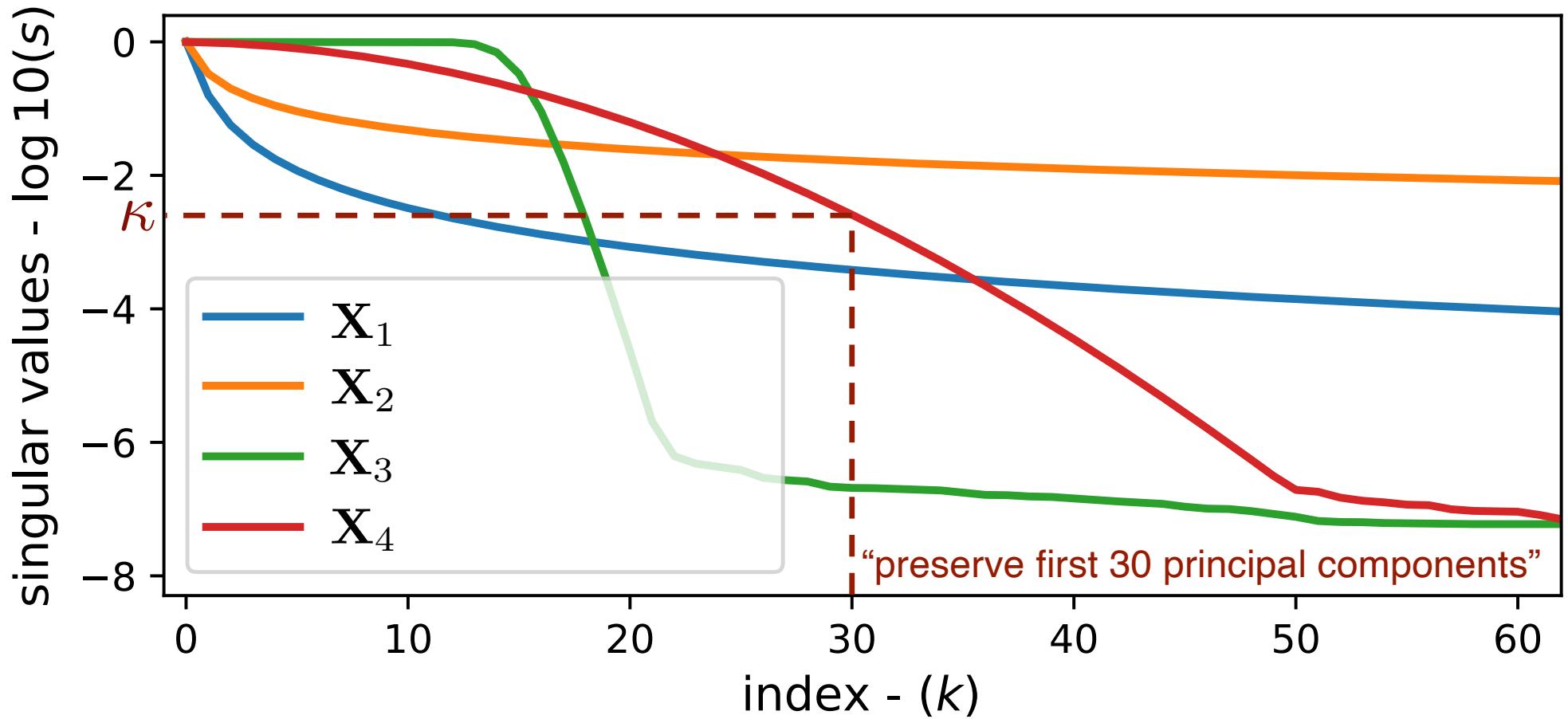
$$\alpha \leq 2 \cdot s_{\max}^{-2}$$

for stable efficient performance,

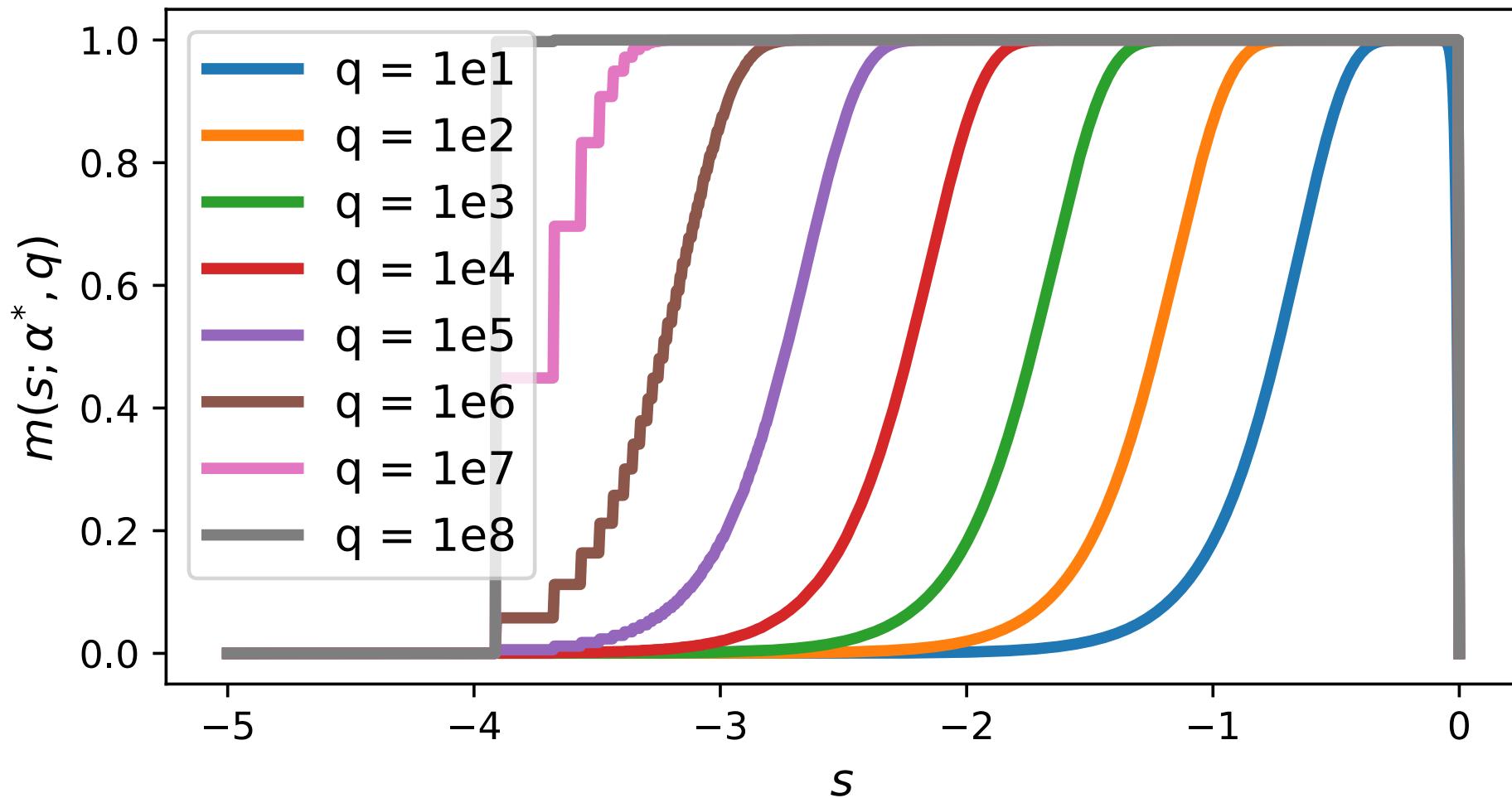
$$\alpha^* = s_{\max}^{-2} \quad \text{and} \quad q^* \approx \frac{s_{\max}^2}{\kappa^2}$$

“remember from PCA shrinkage”

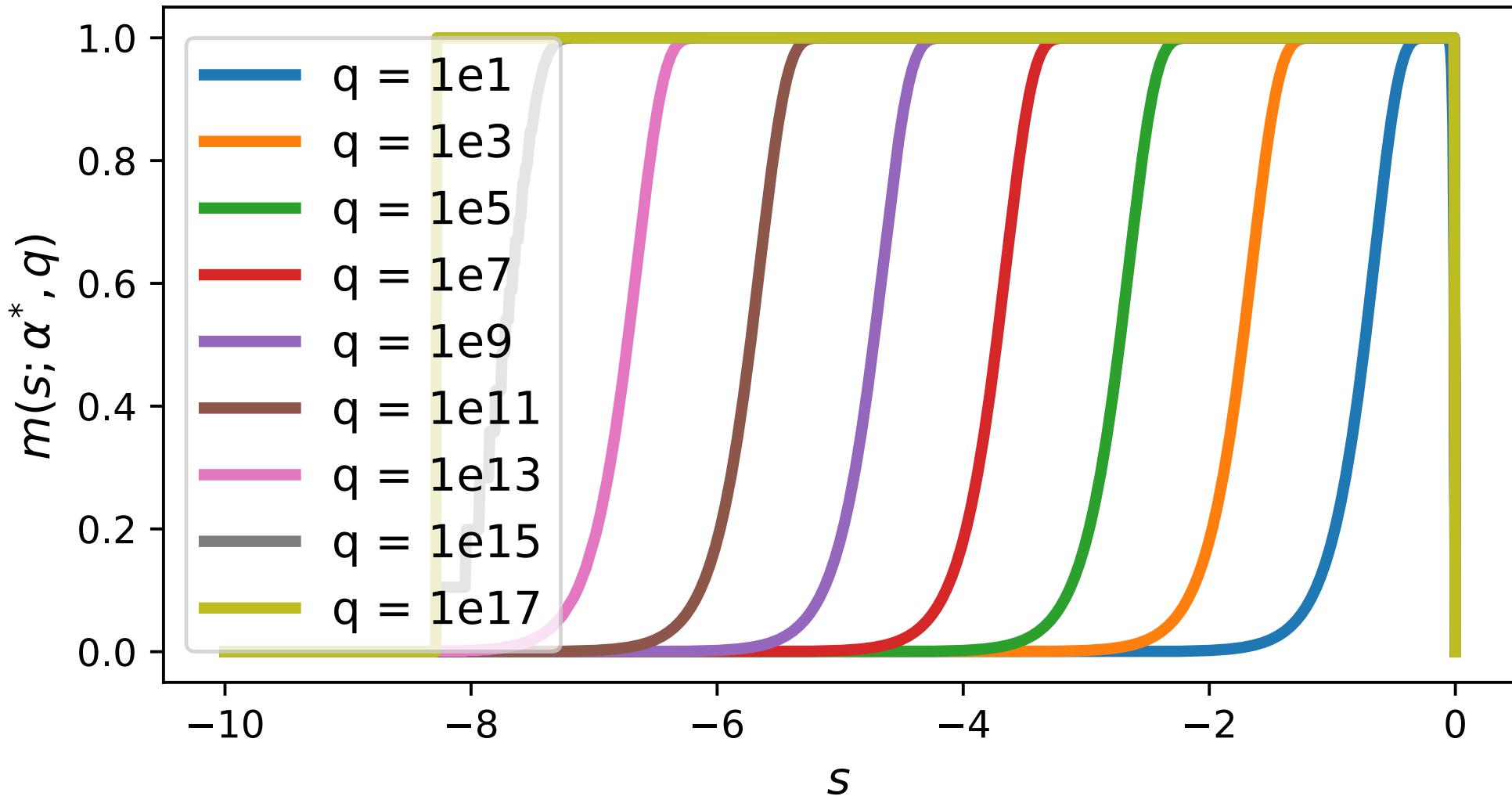
## Remember - PCA Shrinkage



# Gradient Descent – Shrinkage (Float 32 bit)



# Gradient Descent – Shrinkage (Float 64 bit)





# Gradient Descent Shrinkage - Review

---

- Gradient descent on linear models can be reinterpreted as a type of implicit shrinkage regularizer.
- Unlike explicit PCA shrinkage, it does not require any costly SVD step making it ideal for large scale problems.
- Number of iterations and learning rate approximately dictate the number of principal components to preserve,
- Floating point precision provides an upper bound on the number of principal components that can be preserved.

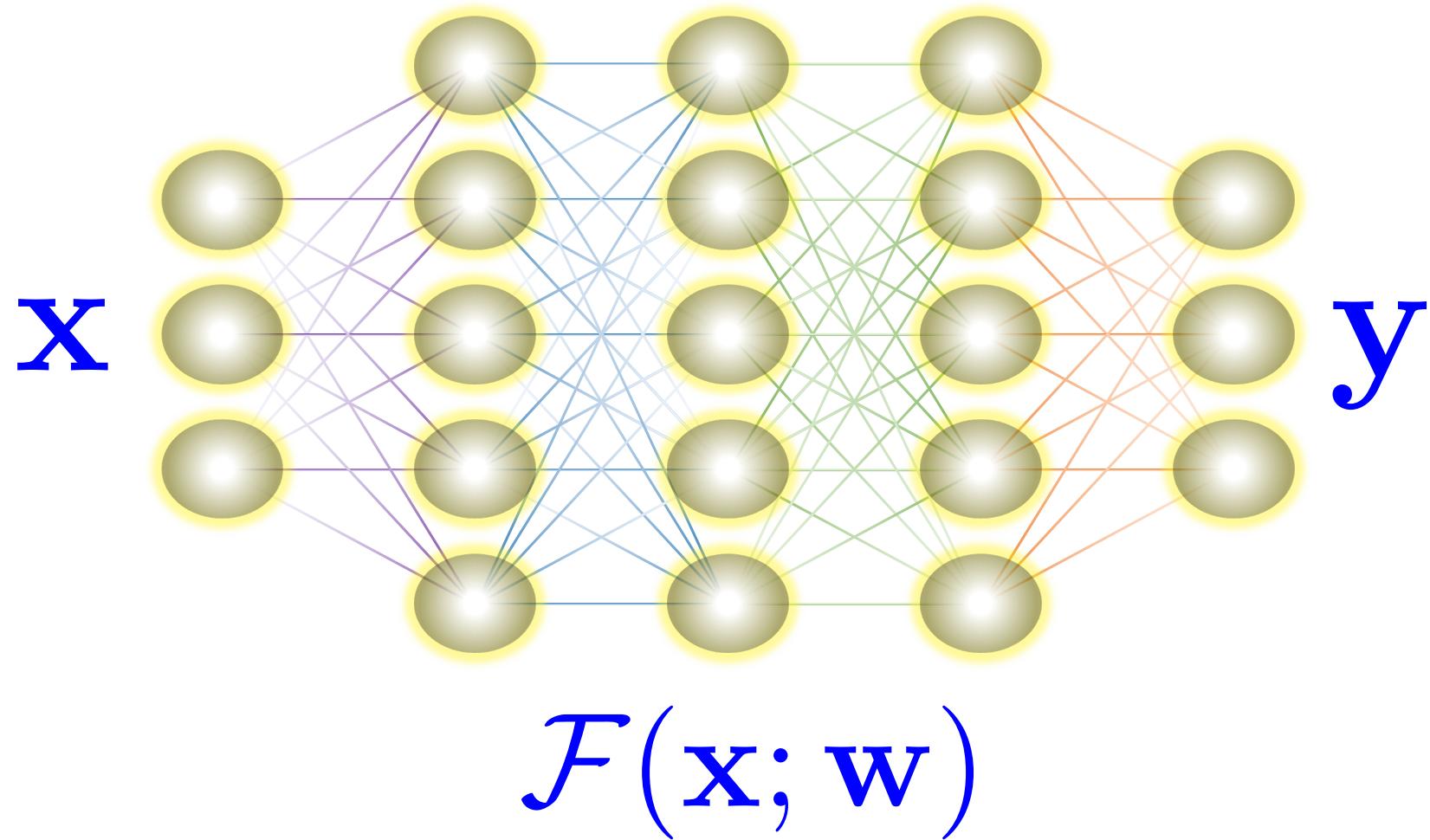
# Today

---

- Regularized Least-Squares
- Gradient Descent
- **Optimizing Deep Networks**

# Optimizing Deep Networks

---



# Linearizing Deep Networks

---

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathcal{F}(\mathbf{X}; \mathbf{w})\|_F^2$$

$$\mathcal{F}(\mathbf{X}; \mathbf{w}) \approx \mathcal{F}(\mathbf{X}; \mathbf{w}_0) + \underbrace{\frac{\partial \mathcal{F}(\mathbf{X}; \mathbf{w}_0)}{\partial \mathbf{w}^T} [\mathbf{w} = \mathbf{w}_0]}_{\text{"network Jacobian"}}$$

Could we apply the same trick?

## Reminder - Gradient Descent Shrinkage

---

$$\nabla \mathcal{F}(\mathbf{X}; \mathbf{w}) = \mathbf{U} \text{diag}(\mathbf{s}) \mathbf{V}^T \quad \text{"SVD"}$$

$$\mathbf{w} = \mathbf{V} \text{diag}(\hat{\mathbf{s}}^{-1}) \mathbf{U}^T \mathbf{y} \quad \text{"closed form solution"}$$

$$\hat{s}_{\text{gd}}^{-1} = m_{\text{gd}}(s; \alpha, q) \cdot s^{-1} \quad \text{"masked shrinkage"}$$

$$m_{\text{gd}}(s; \alpha, q) = [1 - (1 - \alpha \cdot s^2)^q] \quad \text{"GD mask"}$$

## Linear Network - Fixed Jacobian

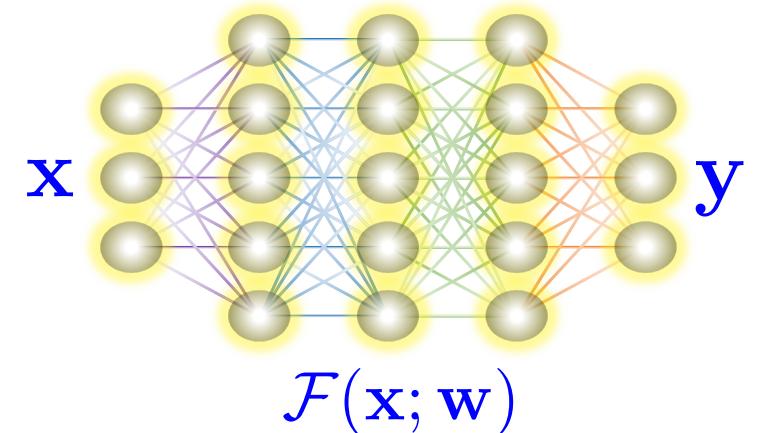
---

$$\frac{\partial \mathcal{F}(\mathbf{X}; \mathbf{w})}{\partial \mathbf{w}^T} = \text{constant}$$

# Deep Network - Dynamic Jacobian

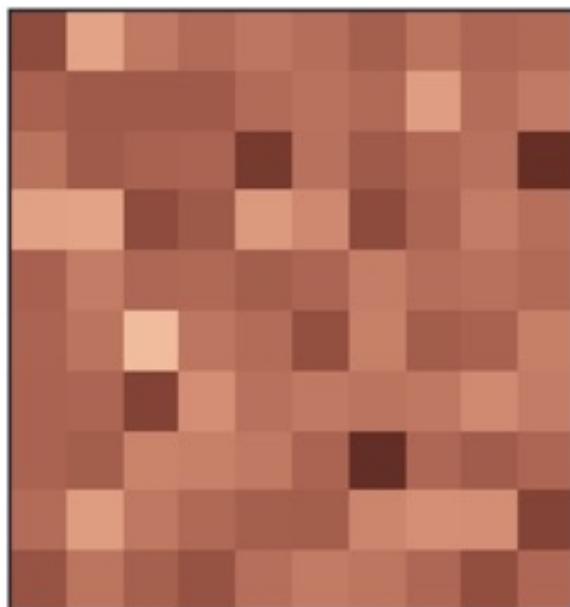
---

$$\frac{\partial \mathcal{F}(\mathbf{X}; \mathbf{w})}{\partial \mathbf{w}^T} \neq \text{constant}$$

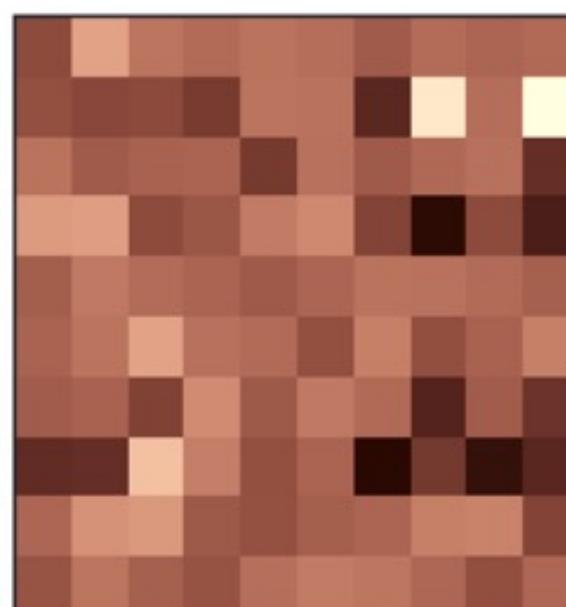


# Solution – Overparametrization?

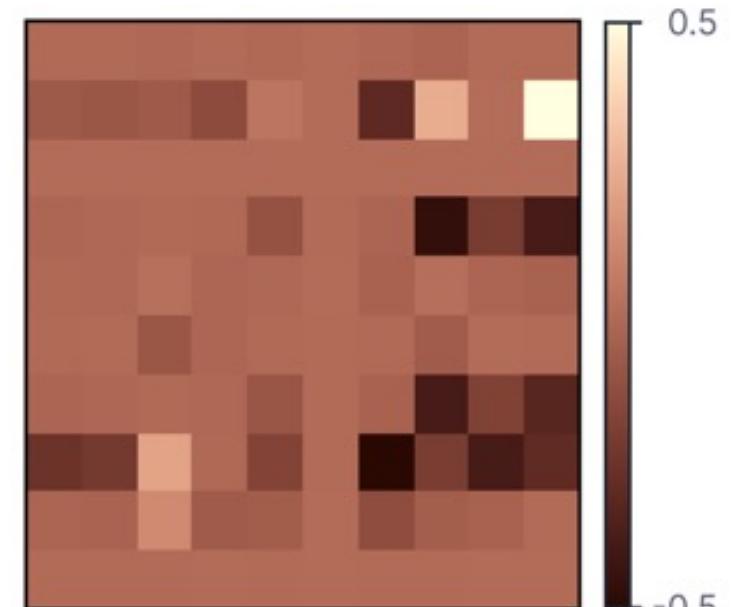
e.g. 2 layer MLP, varying width



$\mathbf{W}_0$



$\mathbf{W}_\infty$



$\mathbf{W}_0 - \mathbf{W}_\infty$

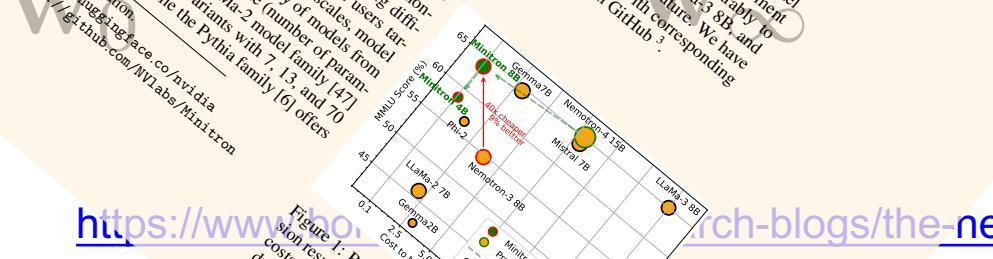
<https://www.borealisai.com/research-blogs/the-neural-tangent-kernel/>

# Solution – Overparameterization?

e.g. 2 layer MLP, varying width



$$\mathbf{w}_0 - \mathbf{w}_\infty$$



<https://www.csail.mit.edu/~davidduvenaud/research-blogs/the-neural-tangent-kernel/>

and have demonstrated real-world applications. To aid users, model sizes, and scales, from the entire size family (number of parameters) to the different variants with 7, 13, and 70 million parameters, different while the Pythia family [6] offers qual contribution  
<https://huggingface.co/nVidia/Pythia>  
<https://github.com/nVidia/Minit>

open  
supplie  
duction

Large language models (LLMs) now dominate real-world natural language processing in understanding different contexts [7, 40, 50, 47, 46]. To aid users getting different train an entire size family of models providers, each with a different variant(s). For instance, three different while the the includes parameters, billion contribution. <https://github.com/microsoft/MINTRON>

**Introduction**

Large language models (LLMs) now dominate real-world natural language processing in understanding different contexts [7, 40, 50, 47, 46]. To aid users getting different train an entire size family of models providers, each with a different variant(s). For instance, three different while the the includes parameters, billion contribution. <https://github.com/microsoft/MINTRON>

Pythia family size	MMLU Score
7	50.5
8	52.5
9	54.5
10	56.5
11	57.5
12	58.5
13	58.5

<https://www.ho.com>

<https://www.ho...>

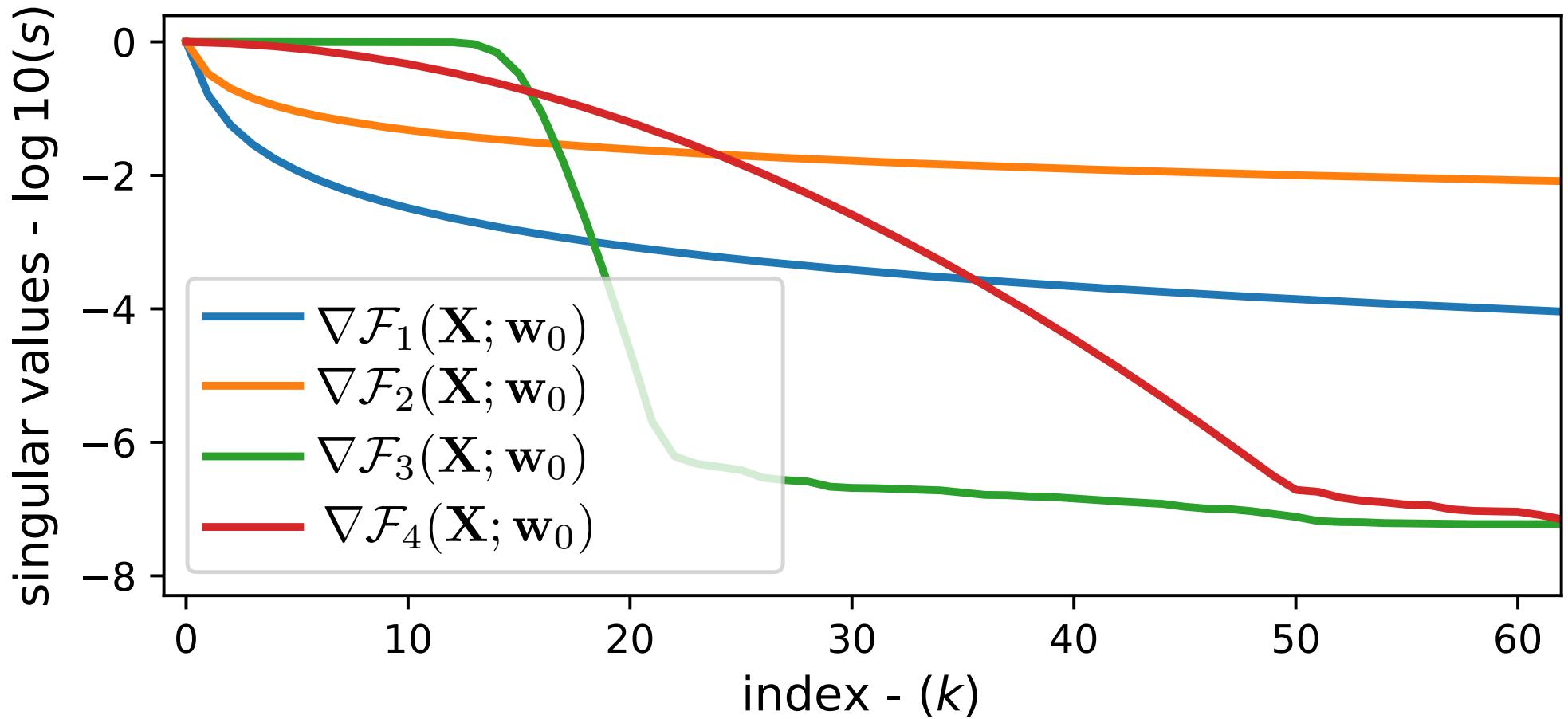
# Neural Tangent Kernel

---

- It has been argued by Jacot et al. that a deep network's Jacobian will become fixed as the width of the network tends towards infinity.
- Allows for the closed form analysis, and better motivates numerous modern interventions including:-
  - Skip connections
  - Batch-norm
  - Layer-norm
  - Adam optimization.

[https://proceedings.neurips.cc/paper\\_files/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf)

# Problem - III-Conditioned Jacobian



<https://www.borealisai.com/research-blogs/the-neural-tangent-kernel/>

## More to read...

---

