

# CS 229br: Advanced Topics in the theory of machine learning

Boaz Barak



Ankur Moitra  
MIT 18.408



Yamini Bansal  
**Official TF**



Dimitris Kalimeris  
Unofficial TF



Gal Kaplun  
Unofficial TF

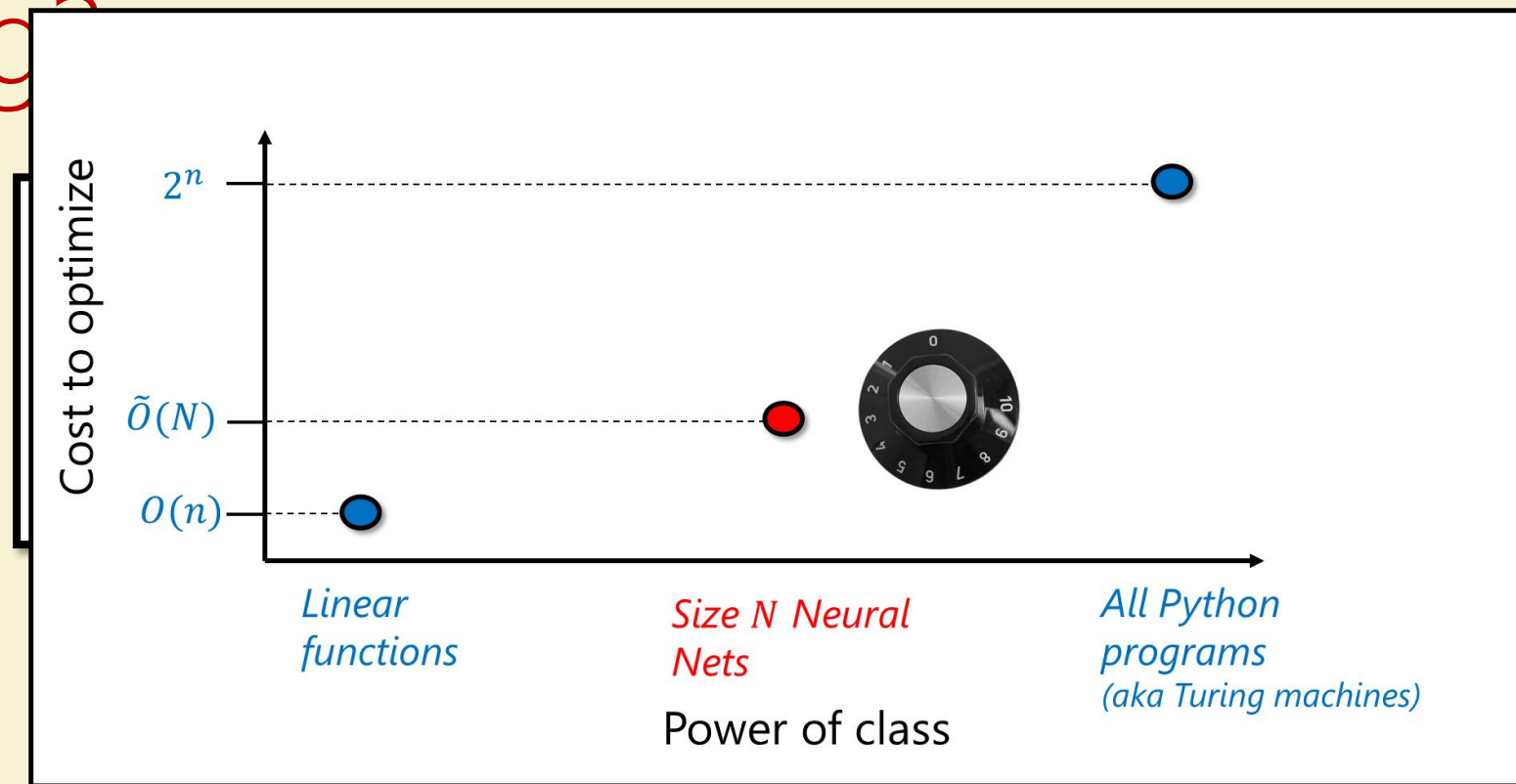
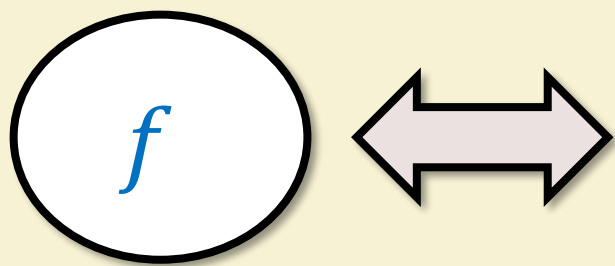


Preetum Nakkiran  
Unofficial TF



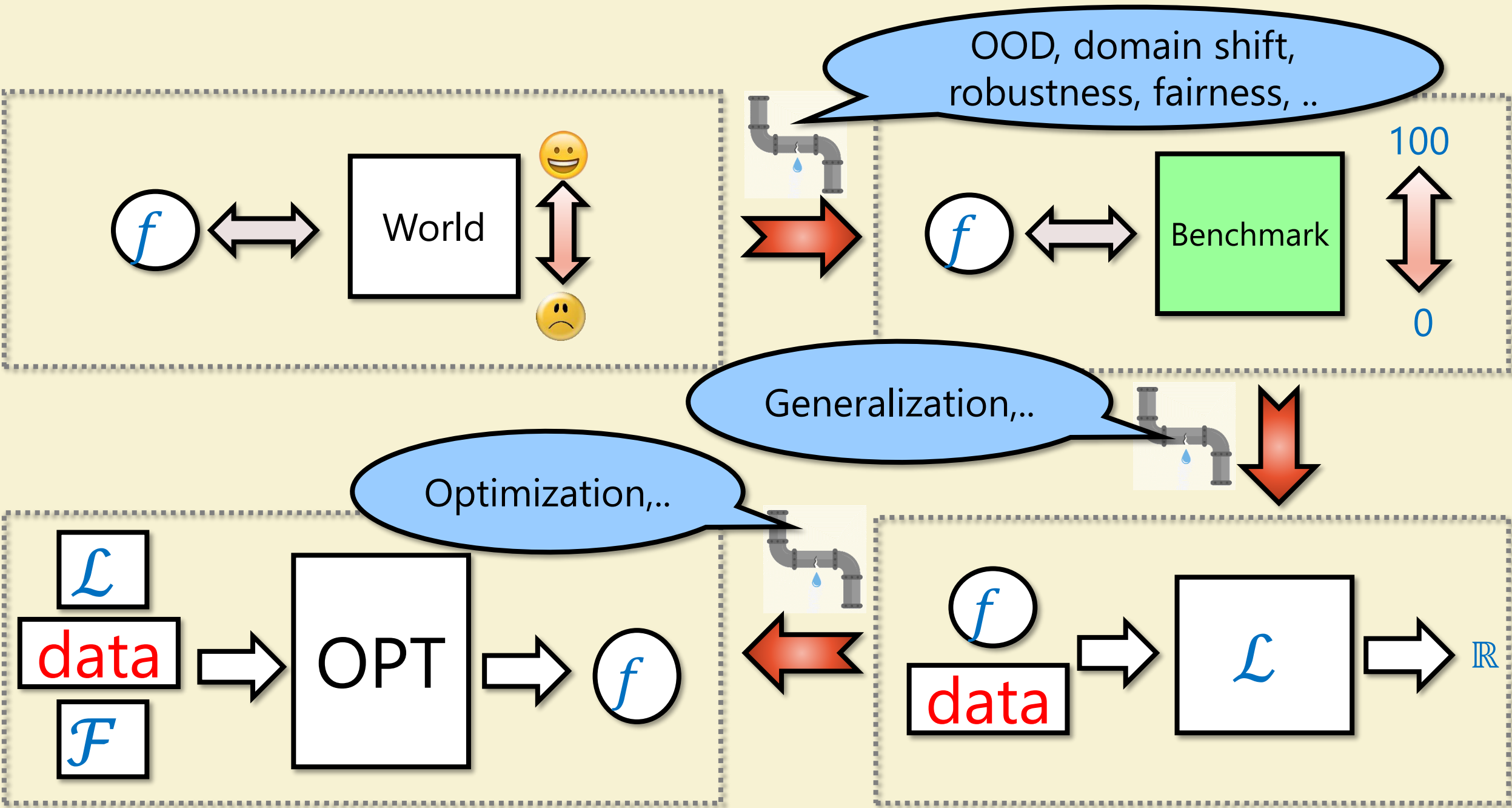
Join **slack team** and sign in on your devices!

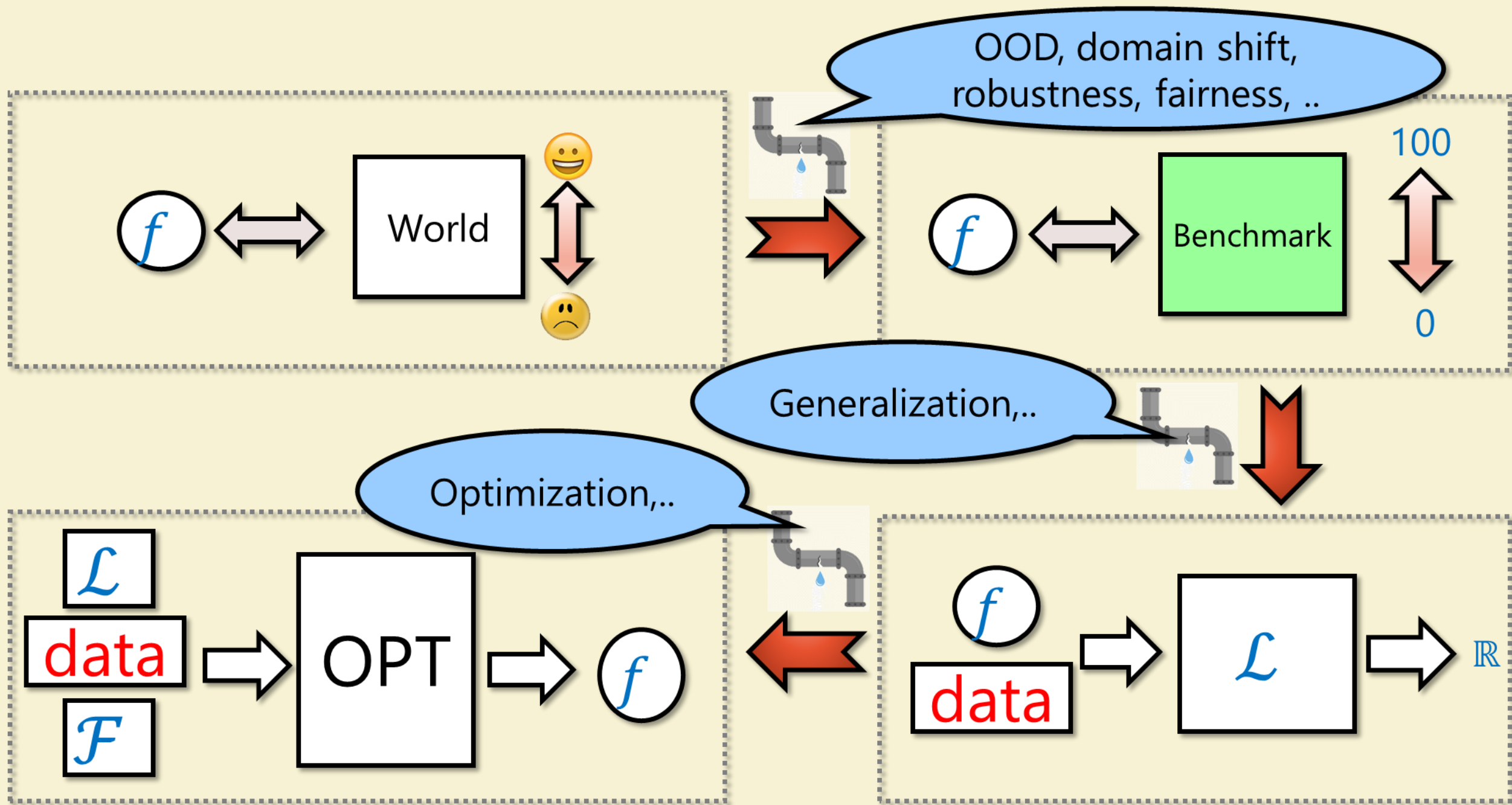
# What is learning?

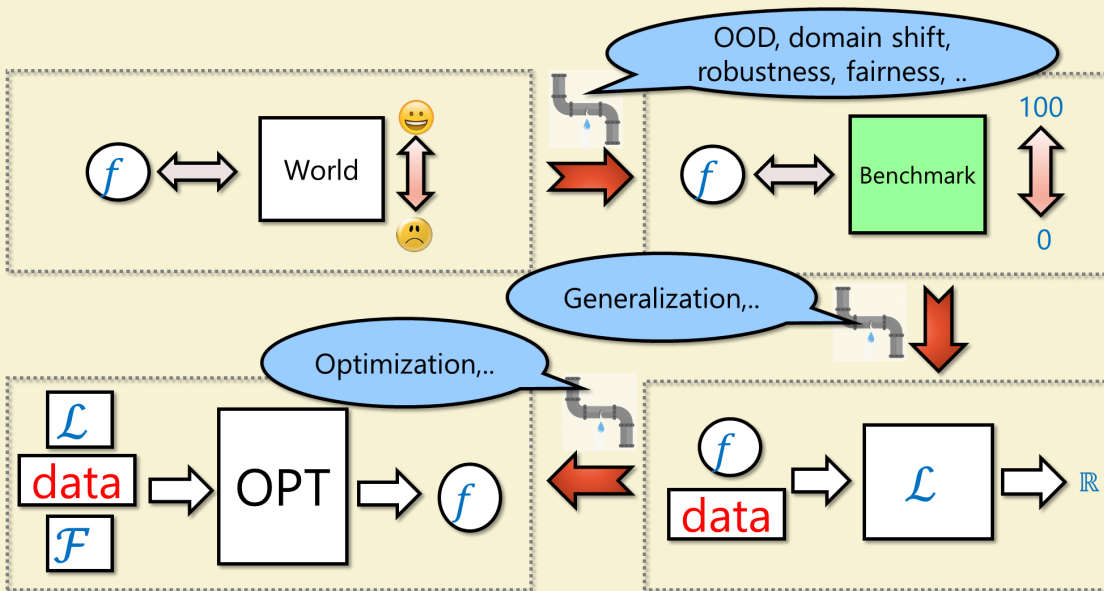


# What is deep learning?







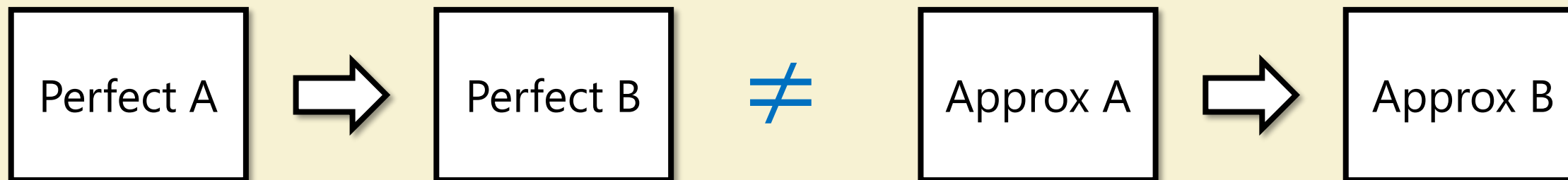


God's view:

$$f = \operatorname{argmax}_f \mathbb{E} [ \text{😊}(f) \mid \text{data} ]$$

$$f = \operatorname{argmax}_f \mathbb{E} [ \text{😊}(f \leftrightarrow w) ]$$

$w \sim W \mid \text{data}$



# This seminar

- Taste of **research results**, **questions**, **experiments**, and more
- **Goal:** Get to state of art research:
  - Background and language
  - Reading and reproducing papers
  - Trying out extensions
- Very experimental and **"rough around the edges"**
- A lot of learning on your own and **from each other**
- Hope: **Very interactive** – in lectures and on slack



**CS 229br:** Survey of very recent research directions, emphasis on experiments



**MIT 18.408:** Deeper coverage of foundations, emphasis on theory

# Student expectations

Not set in stone but will include:

- **Pre-reading** before lectures
- **Scribe notes** / blog posts (adding proofs and details)
- **Applied problem sets** (replicating papers or mini versions thereof, exploring)  
*Note: Lectures **will not** teach practical skills – rely on students to pick up using suggested tutorials, other resources, and each other.*  
*Unofficial TFs happy to answer questions!*
- Some **theoretical problem sets**.
- Might have you grade each other's work
- **Projects** – self chosen and directed.

HW0 on  
slack

**Grading:** We'll figure out some grade – hope that's not your loss function 😊

Rest of today's lecture:

## Blitz through classical learning theory

- Special cases, mostly one dimension
- "Proofs by picture"

PATTERNS, PREDICTIONS, AND ACTIONS

A story about machine learning

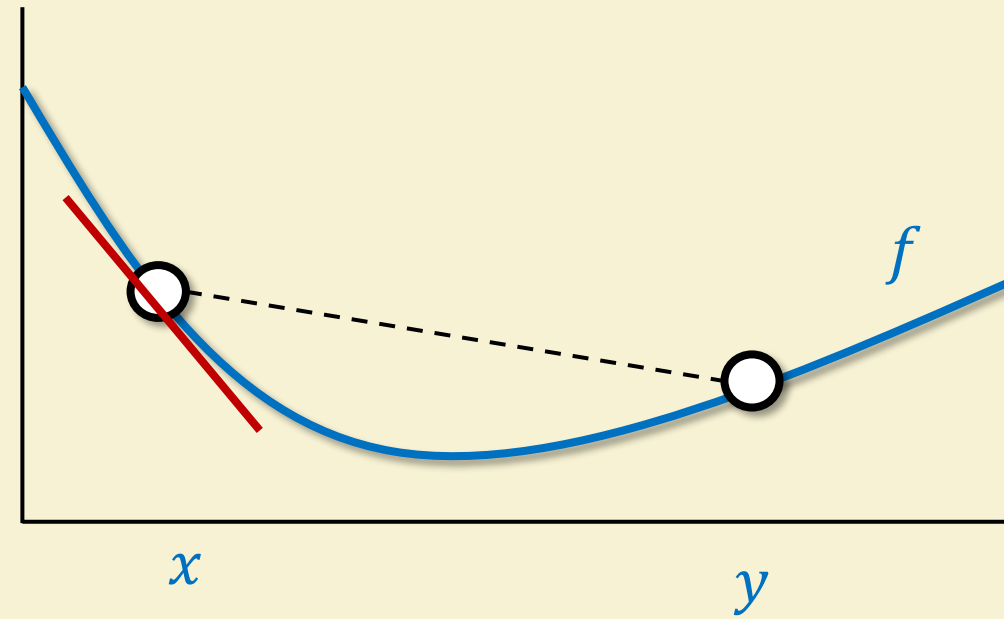
*Moritz Hardt and Benjamin Recht*



# Part I: Convexity & Optimization

# Convexity

1. Line between  $(x, f(x))$  and  $(y, f(y))$  above  $f$ .
2. Tangent line at  $x$  (slope  $f'(x)$ ) below  $f$ .
3.  $f''(x) > 0$  for all  $x$

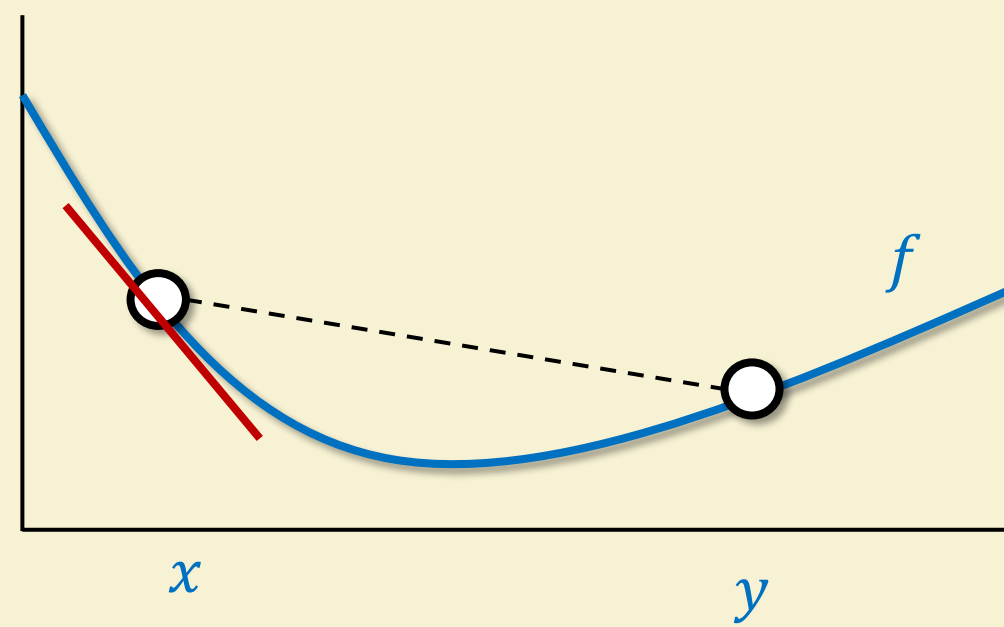


**CLAIM:**  $f''(x) < 0$  implies tangent line at  $x$  above  $f$

**Proof:**  $f(x + \delta) = f(x) + \delta f'(x) + \delta^2 f''(x) + O(\delta^3)$  implies  $f$  below line

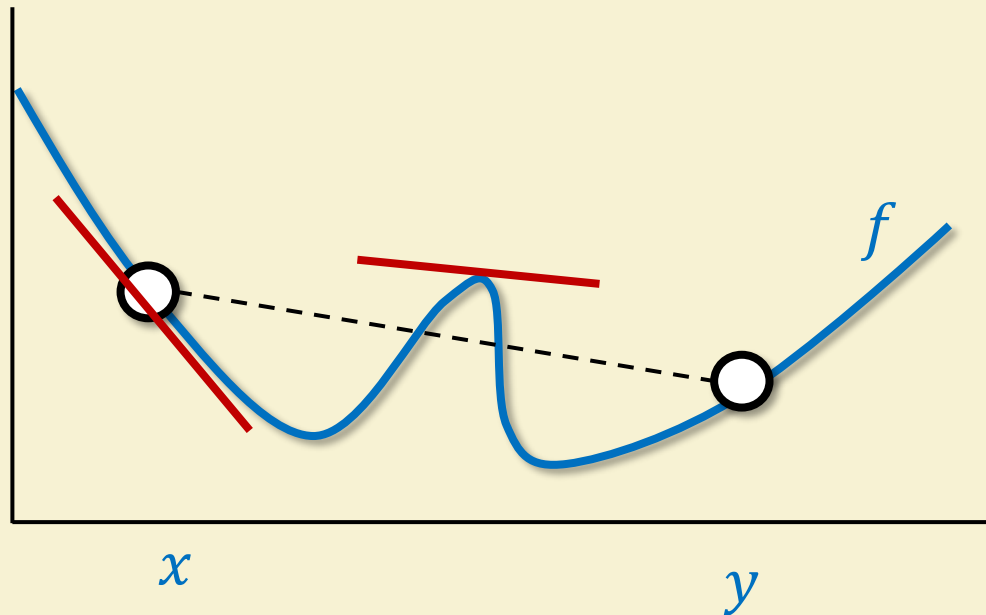
# Convexity

- ↑
1. Line between  $(x, f(x))$  and  $(y, f(y))$  above  $f$ .
  2. Tangent line at  $x$  (slope  $f'(x)$ ) below  $f$ .
  3.  $f''(x) > 0$  for all  $x$



**CLAIM:** If  $f$  above  $(x, f(x)) - (y, f(y))$  then exists  $z$  with tangent line at  $z$  above  $f$

Proof by picture:



# Gradient Descent

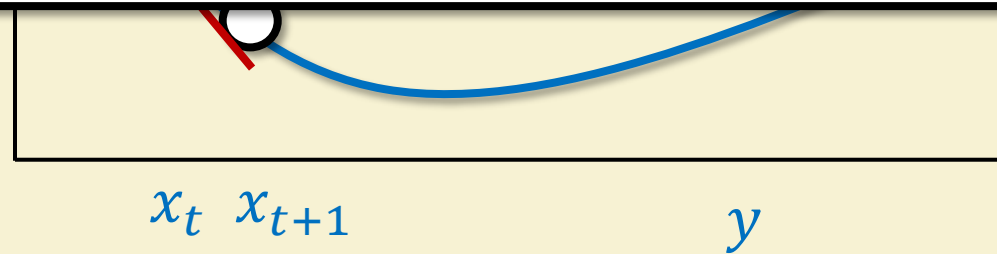
$$x_{t+1} = x_t - \eta f'(x_t)$$

## Dimension $d$ :

$$f'(x) \rightarrow \nabla f(x) \in \mathbb{R}^d$$

$$f''(x) \rightarrow H_f(x) = \nabla_2^2 f(x) \in \mathbb{R}^{d \times d} \text{ (psd)}$$

$$\eta \sim 1/\lambda_d \text{ drop by } \sim \frac{\lambda_1}{\lambda_d} \|\nabla\|^2$$



$$f(x_t + \delta) \approx f(x_t) + \delta f'(x_t) + \frac{\delta^2}{2} f''(x_t)$$

$$f(x_{t+1}) \approx f(x_t) - \eta f'(x_t)^2 + \frac{\eta^2 f'(x_t)^2}{2} f''(x_t) = f(x_t) - \eta f'(x_t)^2 \left(1 - \frac{\eta f''(x_t)}{2}\right)$$

- If  $\eta < 2/f''(x_t)$  then make progress
- If  $\eta \sim 2/f''(x_t)$  then drop by  $\sim \eta f'(x_t)^2$

# Stochastic Gradient Descent

In Machine Learning:

$$f(x) = \frac{1}{n} \sum_{i=1}^n L_i(x)$$

$$\hat{f}'(x_t) = L_i'(x) \text{ for } i \sim [n]$$

Mean  $f'(x_t)$

Variance  $\sigma^2$

Independent

$$x_{t+1} = x_t - \eta \hat{f}'(x_t)$$

$$\mathbb{E}[\hat{f}'(x)] = f'(x_t), V[\hat{f}'(x)] = \sigma^2$$

Assume  $\hat{f}'(x) = f'(x) + N$

$$f(x_t + \delta) \approx f(x_t) + \delta f'(x_t) + \frac{\delta^2}{2} f''(x_t)$$

$$f(x_{t+1}) \approx f(x_t) - \eta f'(x_t)^2 \left(1 - \frac{\eta f''(x_t)}{2}\right) + \eta^2 \sigma^2 f''(x_t)$$

- If  $\eta < 2/f''(x_t)$  and (\*)  $\eta \sigma^2 \ll f'(x_t)^2 / f''(x)$  then make progress
- If  $\eta \sim 2/f''(x_t)$  and (\*) then drop by  $\sim \eta f'(x_t)^2$

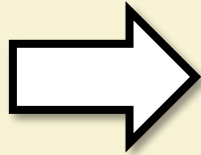
# Part II: Generalization

# Supervised learning

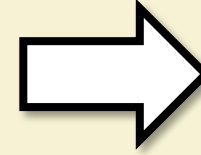
Distribution  $(X, Y)$  over  $\mathcal{X} \times \{\pm 1\}$

$$S = (x_i, y_i)_{i=1..n}$$

**data**



Learning  
Algorithm  
 $A$



$f \in \mathcal{F}$

**Empirical Risk Minimization (ERM):**

$$A(S) = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}_S(f)$$

$$\mathcal{L}(f) = \Pr[f(X) \neq Y]$$

Population 0-1 loss

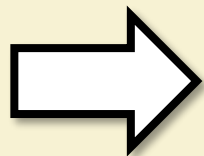
$$\hat{\mathcal{L}}_S(f) = \frac{1}{n} \sum_{i=1}^n 1_{f(x_i) \neq y_i}$$

Empirical 0-1 loss

Generalization gap:  $\mathcal{L}(f) - \hat{\mathcal{L}}_S(f)$  for  $f = A(S)$

# Bias Variance Tradeoff

$$S = (x_i, y_i)_{i=1..n}$$



Learning  
Algorithm  
 $A$

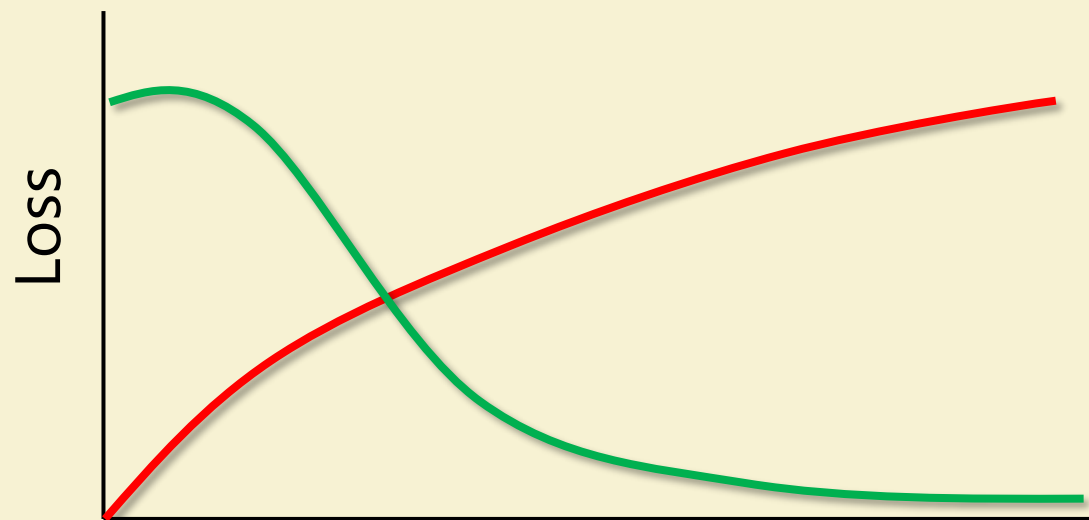
**Empirical Risk Minimization (ERM):**

$$A(S) = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}_S(f)$$

$$f \in \mathcal{F}$$

Population 0-1 loss  $\mathcal{L}(f) = \Pr[f(X) \neq Y]$  **bias** **variance**  $\hat{\mathcal{L}}_S(f) = \frac{1}{n} \sum_{i=1}^n 1_{f(x_i) \neq y_i}$

Assume  $\mathcal{F}_K = \{f_1, \dots, f_K\}$ ,  $\hat{\mathcal{L}}(f_i) = \mathcal{L}(f_i) + N(0, 1/n)$



$$\sim \sqrt{\frac{\log K}{n}} ?$$

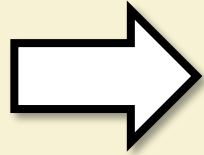
$$\sim (\log K)^{-\alpha} ??$$

$K$  (log scale)



# Bias Variance Tradeoff

$$S = (x_i, y_i)_{i=1..n}$$



Learning  
Algorithm  
 $A$

**Empirical Risk Minimization (ERM):**

$$A(S) = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}_S(f)$$

$$f \in \mathcal{F}$$

Population 0-1 loss  $\mathcal{L}(f) = \Pr[f(X) \neq Y]$

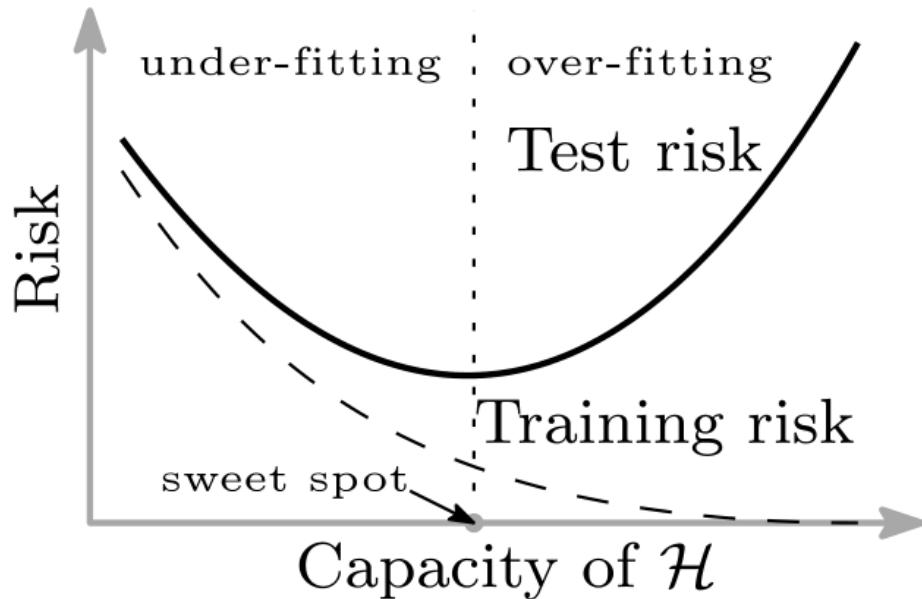
bias

mp

variance

$$\hat{\mathcal{L}}_S(f) = \sum_{i=1}^n 1_{f(x_i) \neq y_i}$$

Assume  $\mathcal{F} = \{f_1, \dots, f_K\}$   $\hat{\mathcal{L}}(f_i) = \mathcal{L}(f_i) + N(0, 1/n)$



$$\sim \sqrt{\frac{\log K}{n}} ?$$

$$\sim (\log K)^{-\alpha} ??$$

Can prove:

$$\text{GAP} \leq O\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$$

Population 0-1 loss  $\mathcal{L}(f) = \Pr[f(X) \neq Y]$

Empirical 0-1 loss  $\hat{\mathcal{L}}_S(f) = \sum_{i=1}^n 1_{f(x_i) \neq y_i}$

Can prove:

$$\text{GAP} \leq o\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$$

**THM:** with high prob over  $S \sim (X, Y)^n$ ,  $\max_{f \in \mathcal{F}} |\mathcal{L}(f) - \hat{\mathcal{L}}_S(f)| \leq o\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$

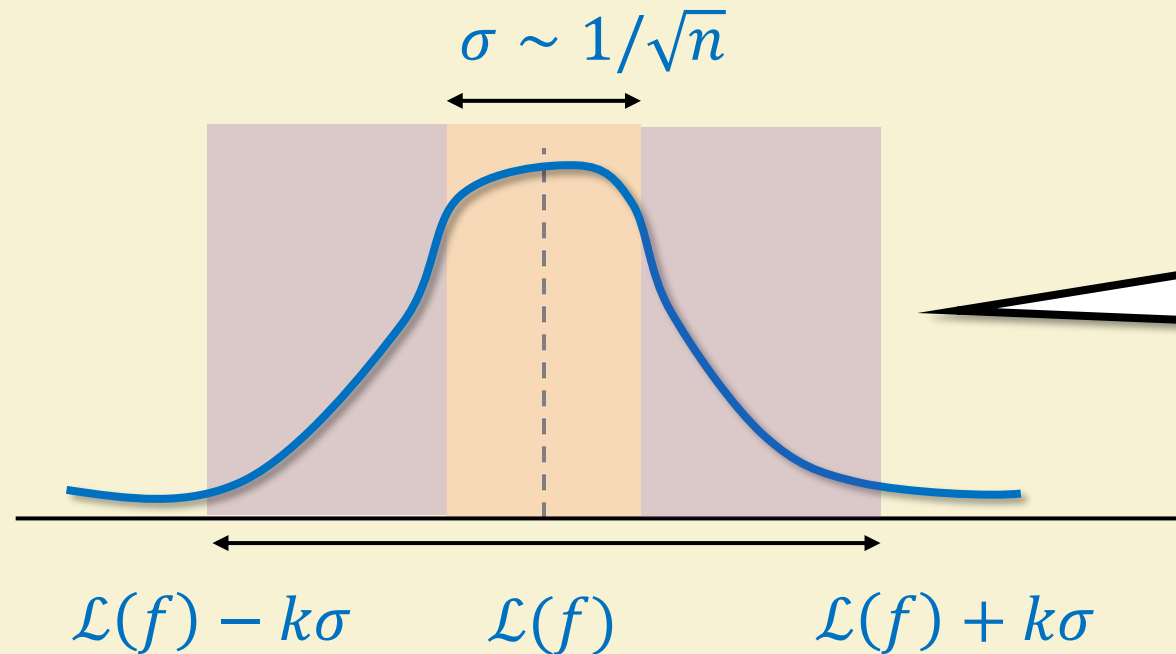
Population 0-1 loss  $\mathcal{L}(f) = \Pr[f(X) \neq Y]$

Empirical 0-1 loss  $\hat{\mathcal{L}}_S(f) = \sum_{i=1}^n 1_{f(x_i) \neq y_i}$

THM: with high prob over  $S \sim (X, Y)^n$ ,  $\max_{f \in \mathcal{F}} |\mathcal{L}(f) - \hat{\mathcal{L}}_S(f)| \leq O\left(\sqrt{\frac{\log|\mathcal{F}|}{n}}\right)$

LEM: For every fixed  $f$ ,  $\Pr_S \left[ |\mathcal{L}(f) - \hat{\mathcal{L}}_S(f)| \geq \frac{k}{\sqrt{n}} \right] \leq \exp(-c \cdot k^2)$

PF: Central limit thm / Chernoff / Hoeffding / ....



$$\Pr[|X - \mathbb{E}X| \geq k\sigma] \leq \exp(-c \cdot k^2)$$

LEM $\Rightarrow$ THM: Set  $k = 10 \cdot \sqrt{\log|\mathcal{F}|} / c$

# Generalization bounds

$\log |\mathcal{F}|$  can be replaced with dimensionality / capacity of  $\mathcal{F}$

Can prove:

$$\text{GAP} \leq O\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$$

VC dimension:

$\max d$  s.t.  $\mathcal{F}$  can fit every labels  $y \in \{\pm 1\}^d$  on every samples  $x_1, \dots, x_d$

Rademacher complexity:

$\max d$  s.t.  $\mathcal{F}$  can fit random labels  $y \sim \{\pm 1\}^d$  on random  $x_1, \dots, x_d \sim X$

PAC Bayes:

$d = I(A(S); S)$  where  $S$  is training set.

Margin bounds:

$\max d$  s.t. linear classifier satisfies  $\langle w_i, x_i \rangle y_i > \|w_i\| \cdot \|x_i\| / \sqrt{d}$  for correct predictions

General form: If  $C(f) \ll n$  then  $\hat{\mathcal{L}}(f) - \mathcal{L}(f) \approx 0$

Intuition: Can't "overfit" – if you do well on  $n$  samples, must do well on population

Challenge: Many modern deep nets and natural  $C$  ,  $C(f) \gg n$

Hope: Find better  $C$ ?

## UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

**Chiyuan Zhang\***

Massachusetts Institute of Technology  
chiyuan@mit.edu

**Samy Bengio**

Google Brain  
bengio@google.com

**Moritz Hardt**

Google Brain  
mrtz@google.com

**Benjamin Recht†**

University of California, Berkeley  
brecht@berkeley.edu

**Oriol Vinyals**

Google DeepMind  
vinyals@google.com

ICLR 2017

General form: If  $C(f) \ll n$  then  $\hat{\mathcal{L}}(f) - \mathcal{L}(f) \approx 0$

Intuition: Can't "overfit" – if you do well on  $n$  samples, must do well on population

Challenge: Many modern deep nets and natural  $C$  ,  $C(f) \gg n$

Hope: Find better  $C$ ?

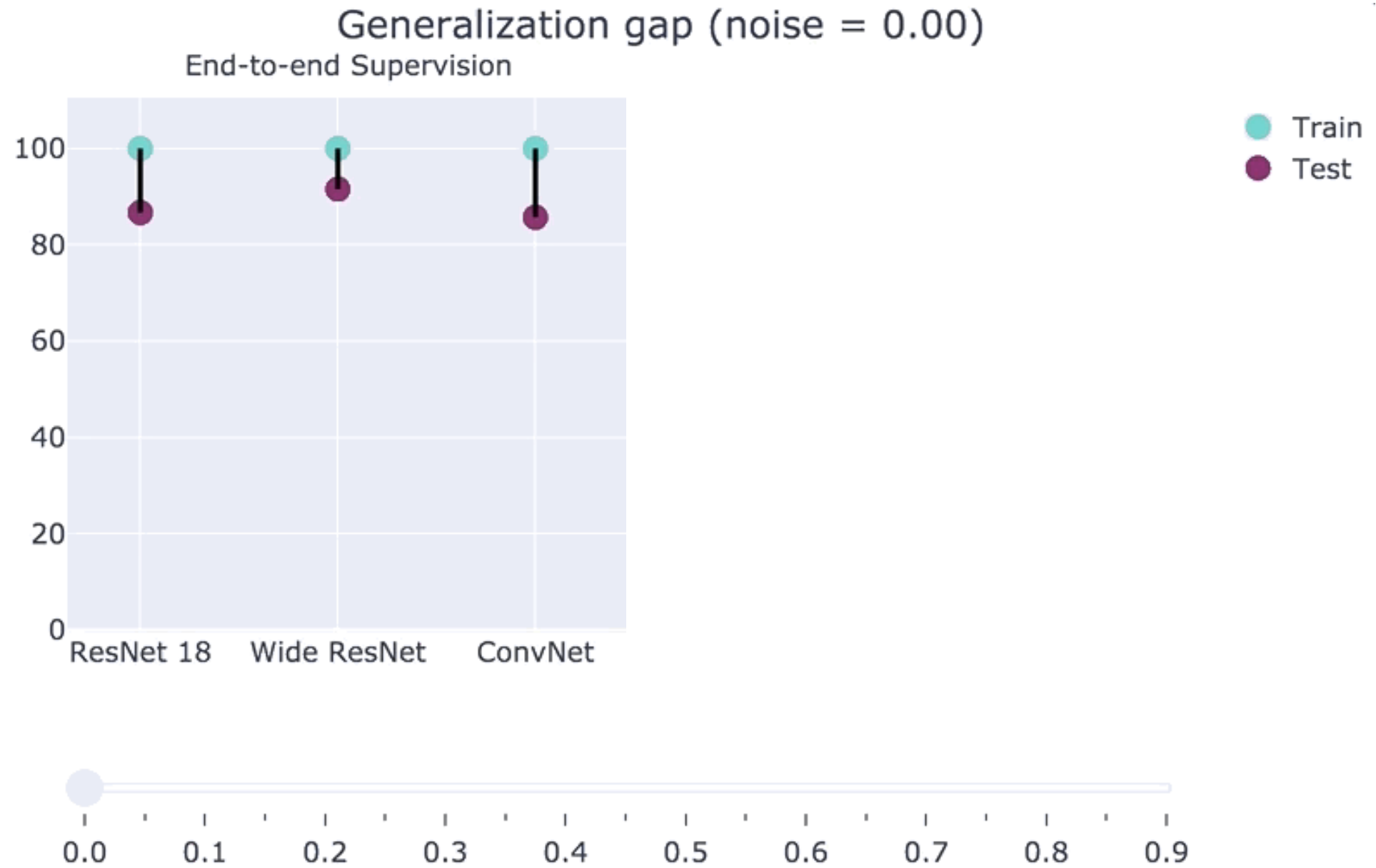
Performance on CIFAR-10

Classifier	Train	Test
Wide ResNet <sup>1</sup>	100%	92%
ResNet 18 <sup>2</sup>	100%	87%
Myrtle <sup>3</sup>	100%	86%

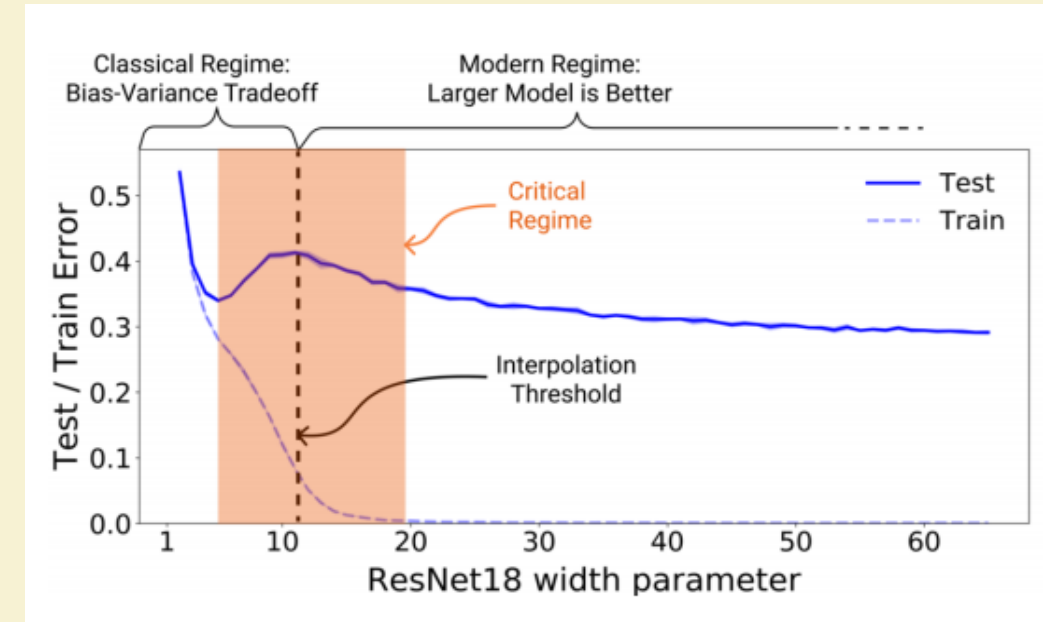
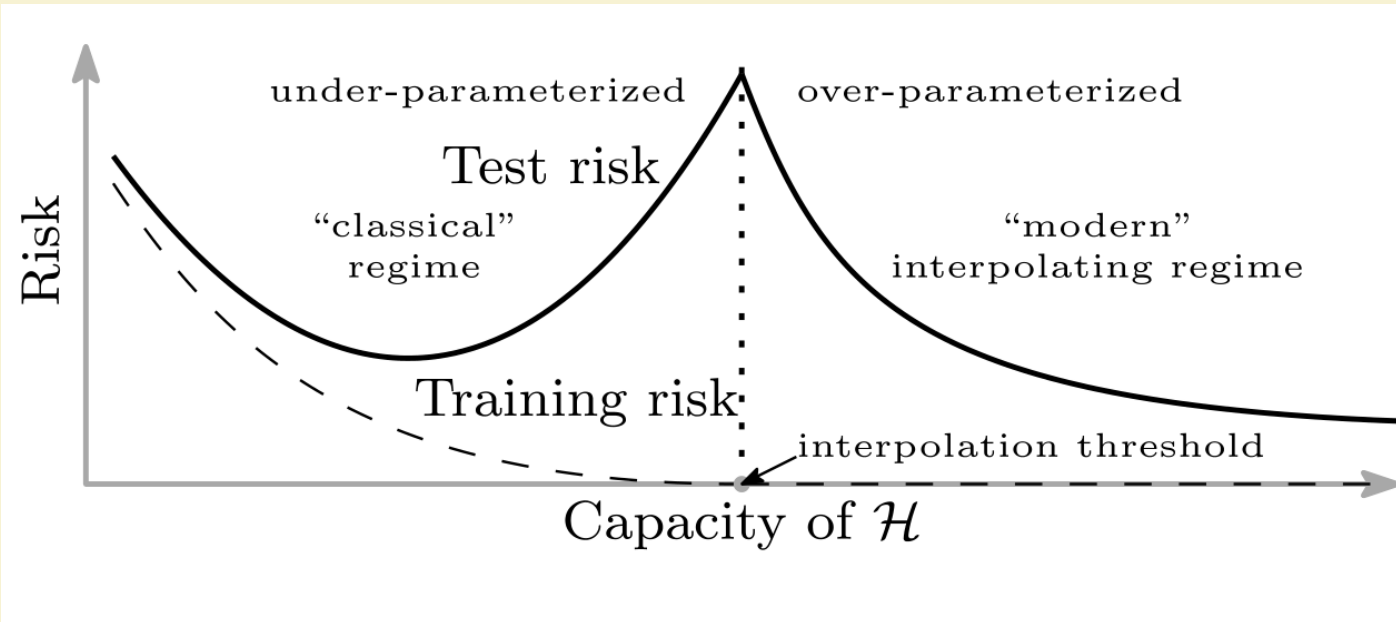
Performance on noisy CIFAR-10:<sup>\*</sup>

Train	Test
100%	10%
100%	10%
100%	10%

<sup>1</sup> Zagoruyko-Komodakis'16 <sup>2</sup> He-Zhang-Ren-Sun'15 <sup>3</sup> Page '19



# "Double descent"



Belkin, Hsu, Ma, Mandal, PNAS 2019.

Nakkiran-Kaplun-Bansal-Yang-B-Sutskever, ICLR 2020

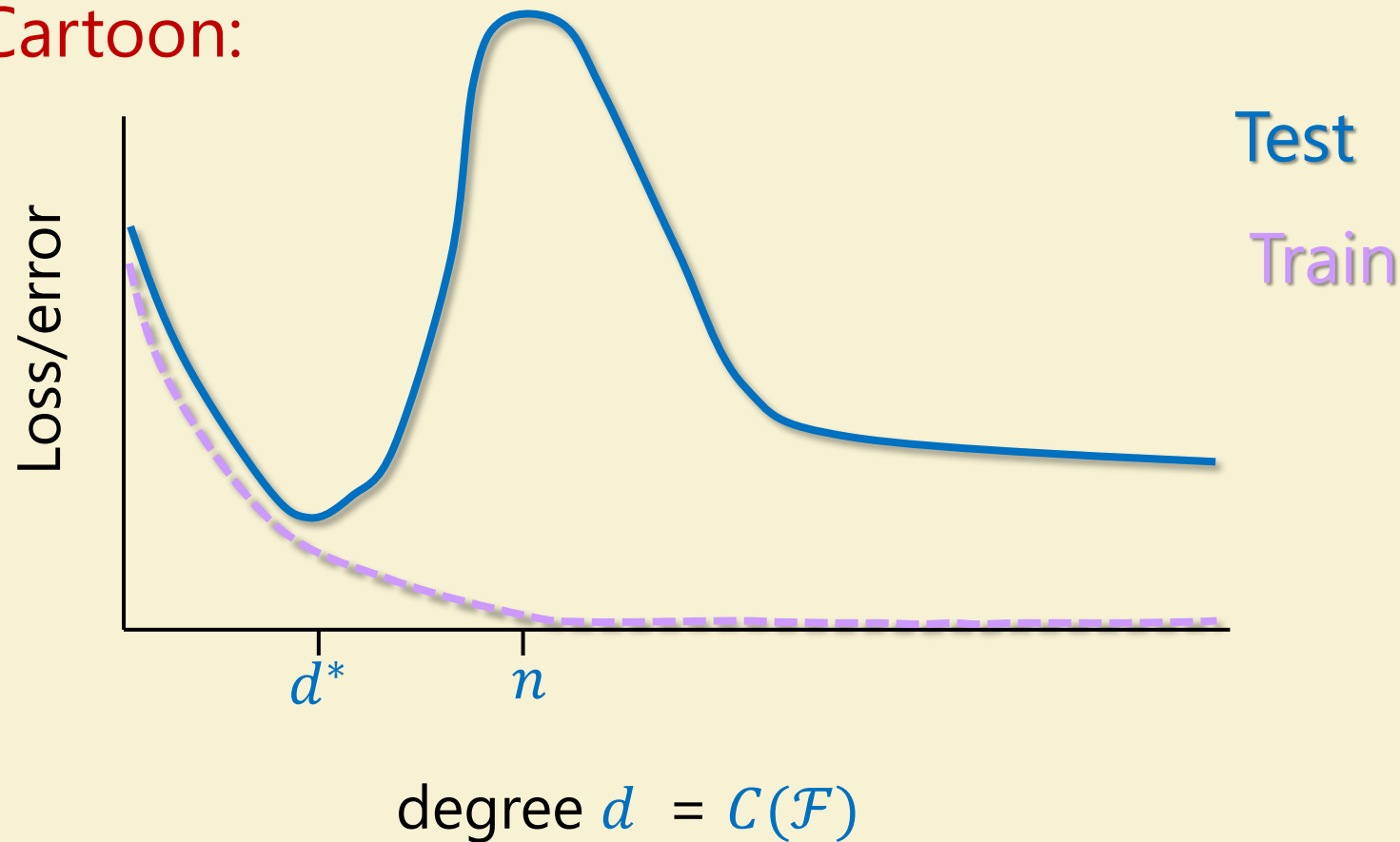


# Example:

**Data:**  $(x_i, f^*(x_i) + N)$  where  $f^*$  degree  $d^*$  polynomial

**Algorithm:** SGD to minimize  $\sum (f(x_i) - y_i)^2$  where  $f \in \mathcal{F}_d$  = degree  $d$  polys

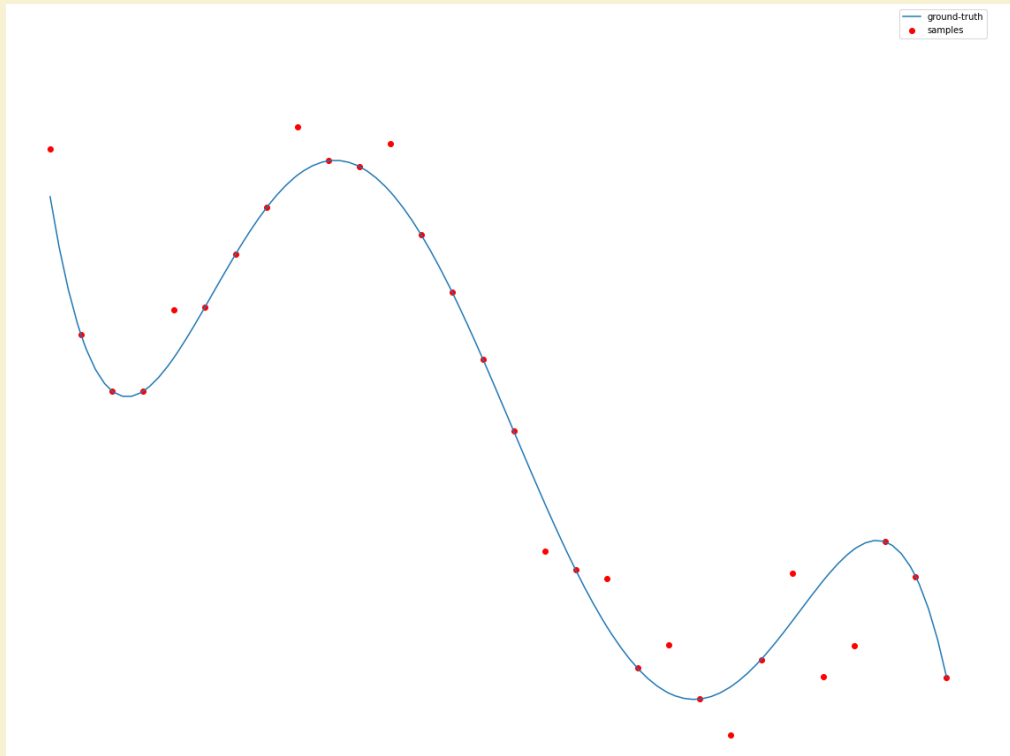
**Cartoon:**



# Example:

**Data:**  $(x_i, f^*(x_i) + N)$  where  $f^*$  degree  $d^*$  polynomial

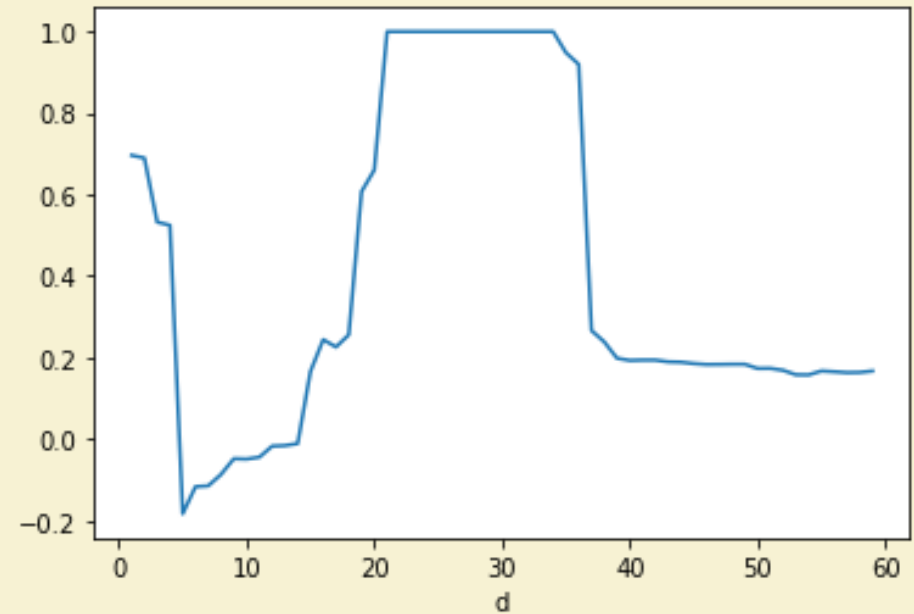
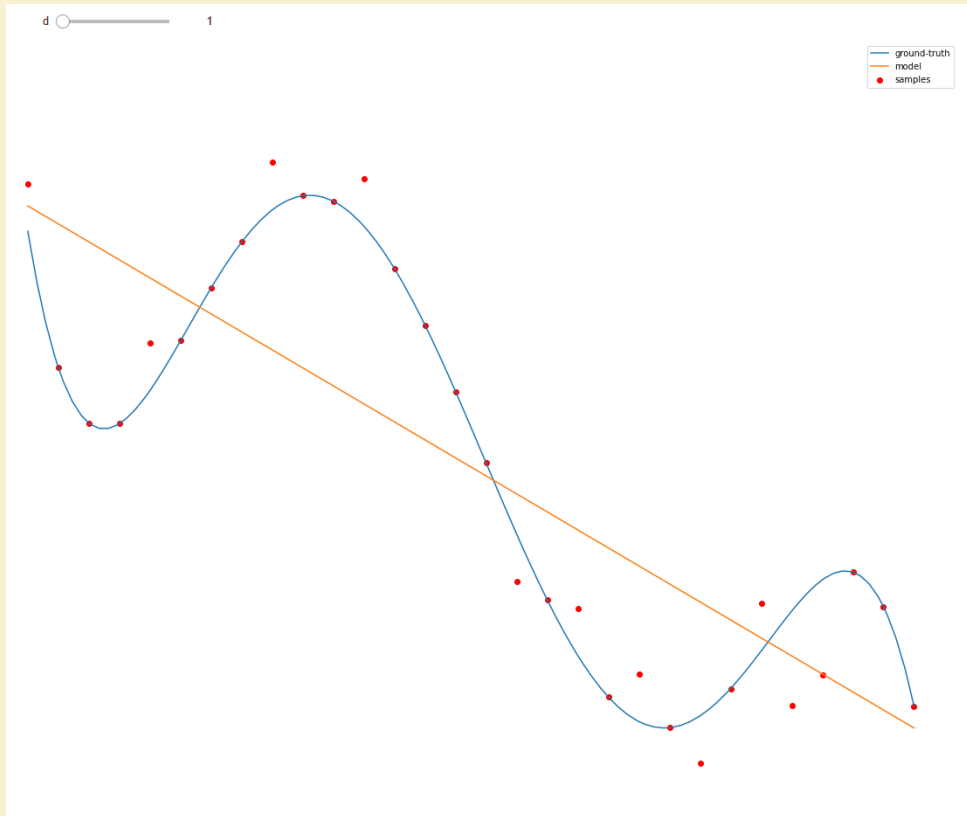
**Algorithm:** SGD to minimize  $\sum (f(x_i) - y_i)^2$  where  $f \in \mathcal{F}_d$  = degree  $d$  polys



# Example:

**Data:**  $(x_i, f^*(x_i) + N)$  where  $f^*$  degree  $d^*$  polynomial

**Algorithm:** SGD to minimize  $\sum (f(x_i) - y_i)^2$  where  $f \in \mathcal{F}_d = \text{degree } d \text{ polys}$

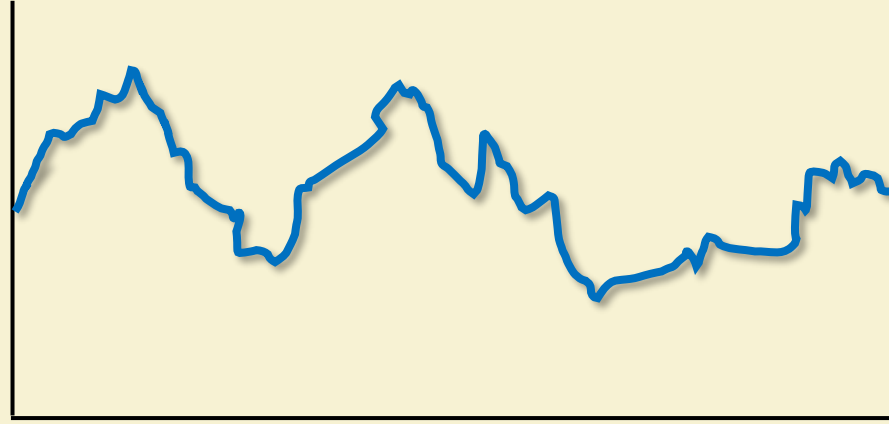


# Part III: Approximation and Representation

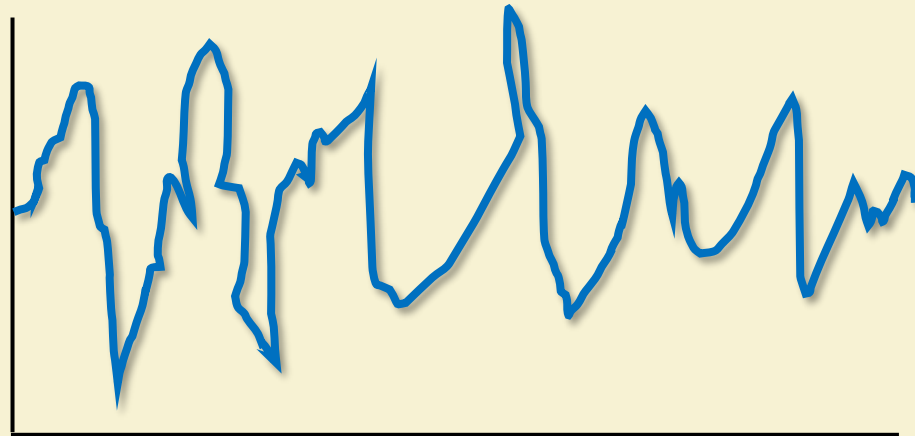
# Example:



Hello



Hello

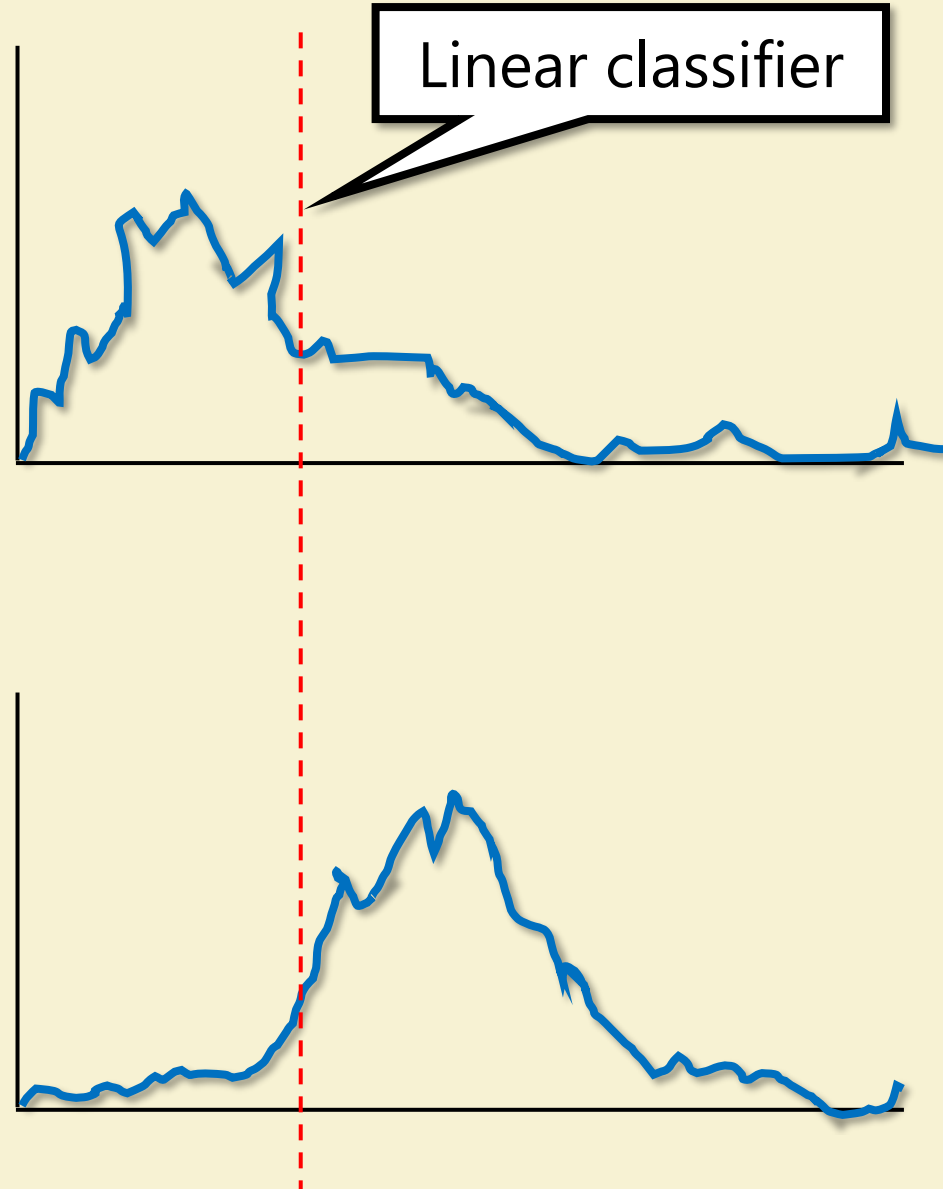


# Example:

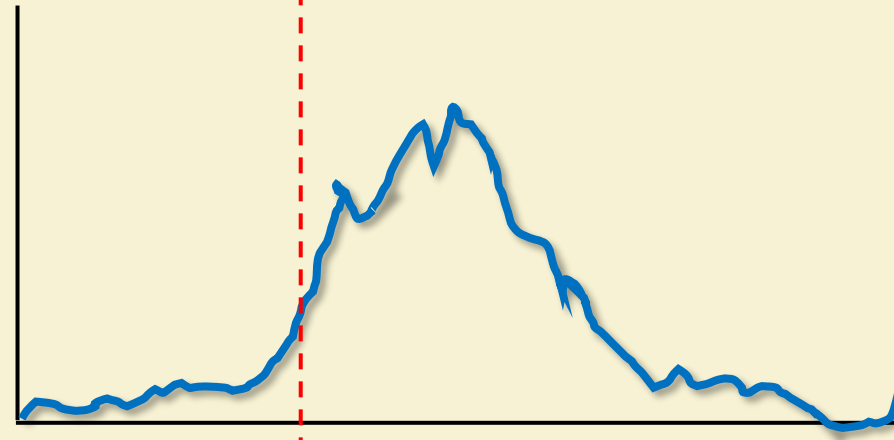
## After Fourier Transform



Hello



Hello



# Fourier Transform

$$\int_0^1 |f - g| < \epsilon$$

Every continuous  $f: [0,1] \rightarrow \mathbb{R}$  can be arbitrarily approximated by  $g$  of form

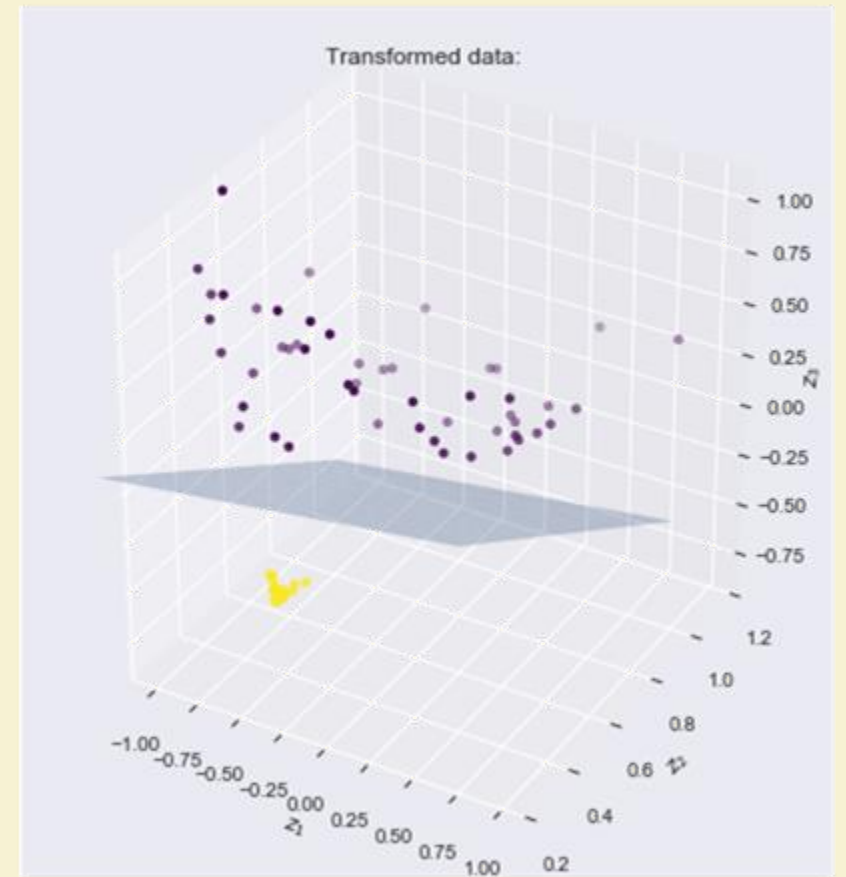
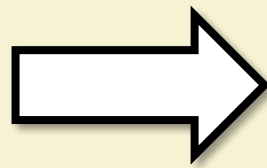
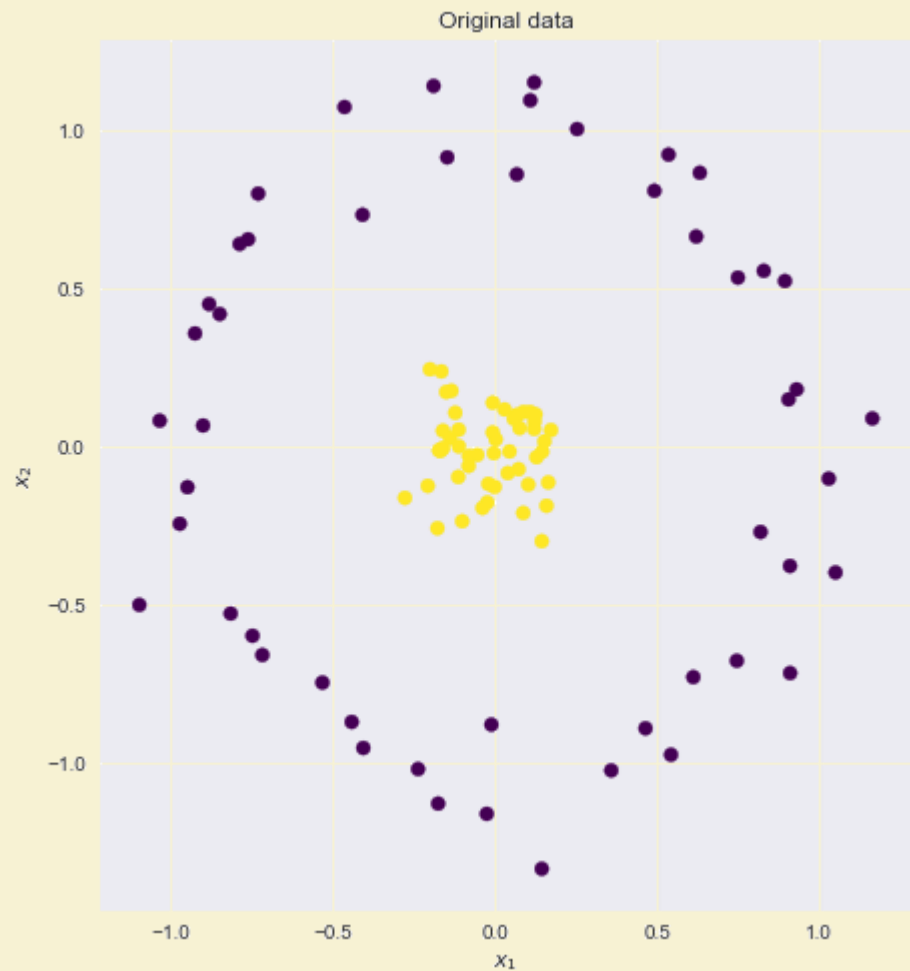
$$g(x) = \sum \alpha_j e^{-2\pi i \beta_j x}$$

i.e.,  $f \approx \text{linear}$  in  $\varphi_1(x), \dots, \varphi_N(x)$  where  $\varphi_j(x) = e^{-2\pi i \beta_j x}$

$\varphi: \mathbb{R} \rightarrow \mathbb{R}^N$  embedding

For some natural data, representation  $\varphi$  is "nice"  
(e.g. sparse, linearly separable,...)

Tasks become simpler after transformation  $x \rightarrow \varphi(x)$



$$\varphi(x, y) = (xy, x^2, y^2)$$



# More approximation

$$\text{ReLU}(x) = \max(x, 0)$$

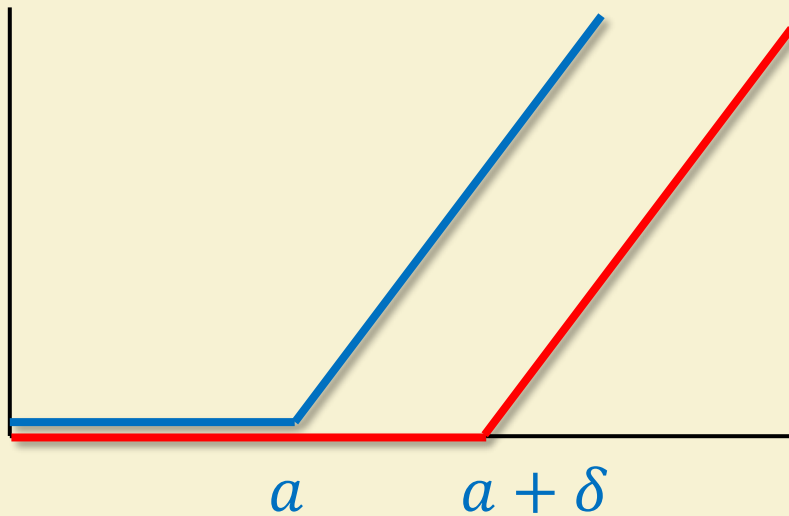
Every continuous  $f: [0,1] \rightarrow \mathbb{R}$  can be arbitrarily approximated by  $g$  of form

$$g(x) = \sum \alpha_i \text{ReLU}(\beta_i x + \gamma_i)$$

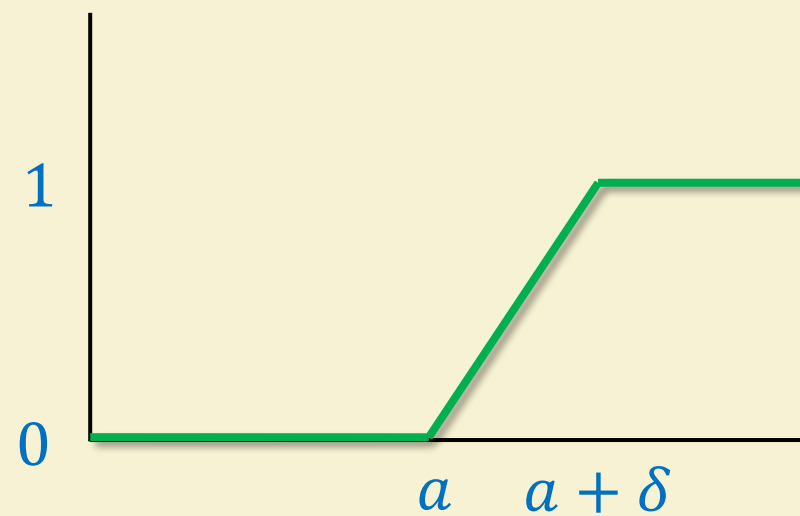
Proof by picture:

$$g_a(x) = \text{ReLU}(x/\delta - a)$$

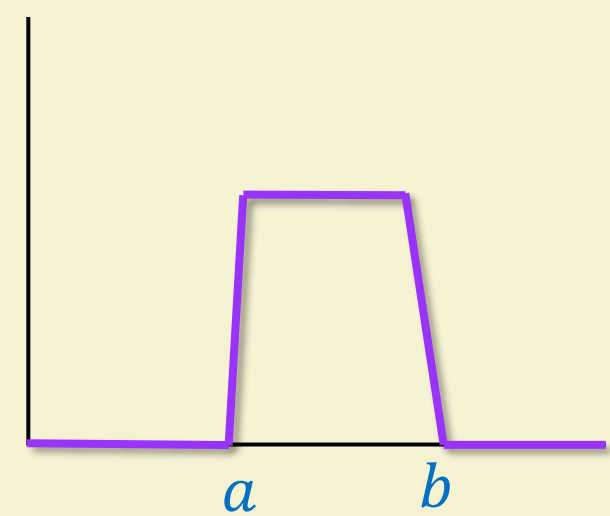
$$g_{a+\delta}(x) = \text{ReLU}(x/\delta - a - \delta)$$



$$h_a = g_a - g_{a+\delta}$$



$$I_{a,b} = h_a - h_b$$



# More approximation

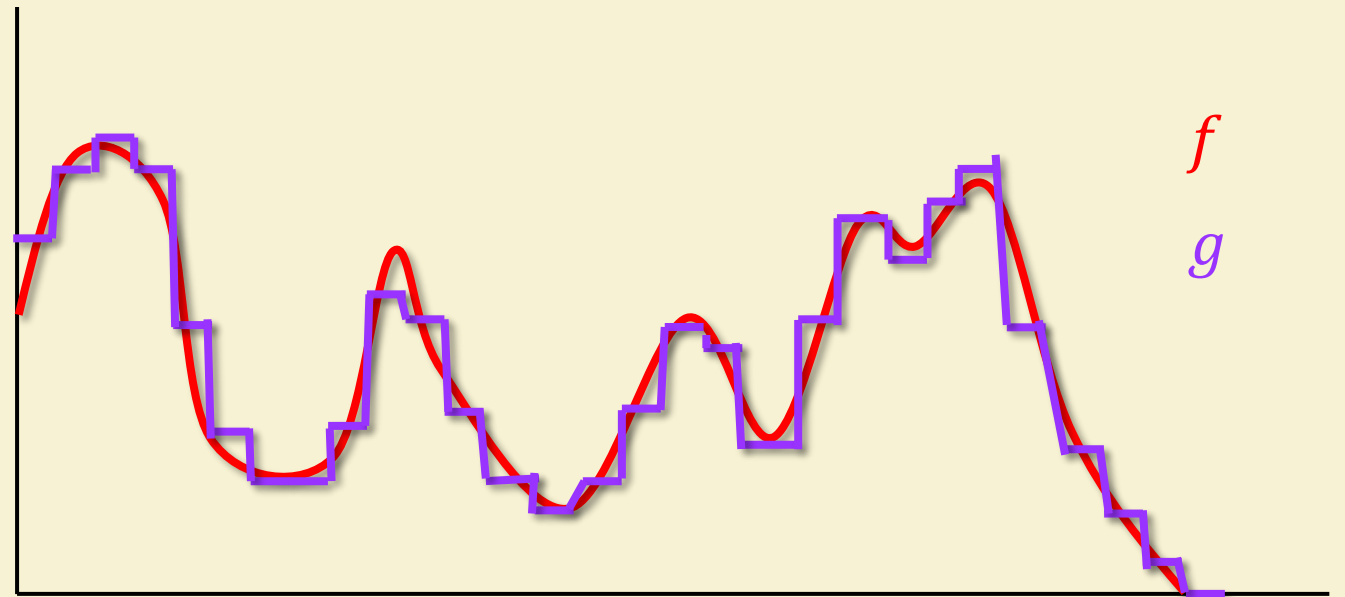
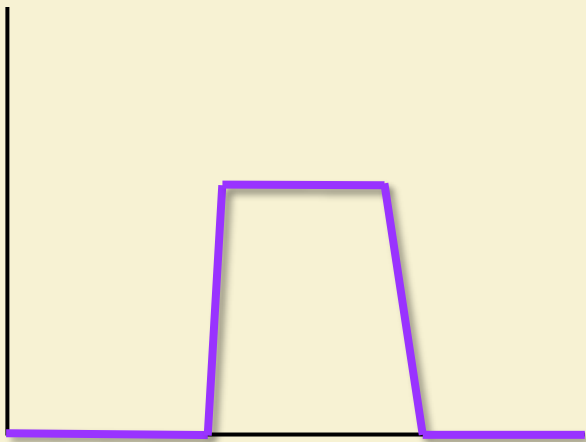
$$\text{ReLU}(x) = \max(x, 0)$$

Every continuous  $f: [0,1] \rightarrow \mathbb{R}$  can be arbitrarily approximated by  $g$  of form

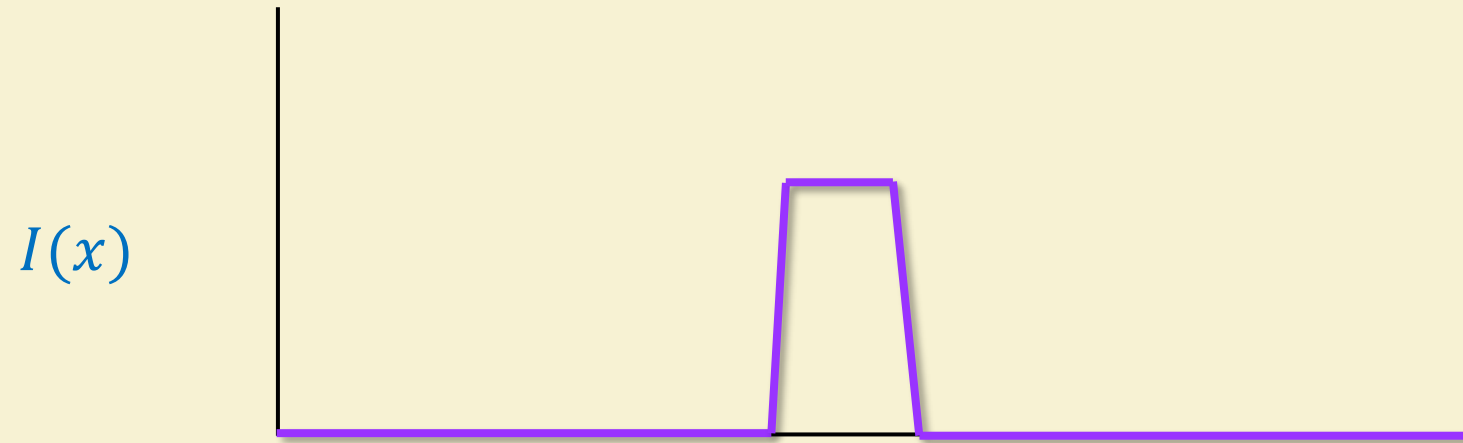
$$g(x) = \sum \alpha_i \text{ReLU}(\beta_i x + \gamma_i)$$

Proof by picture:

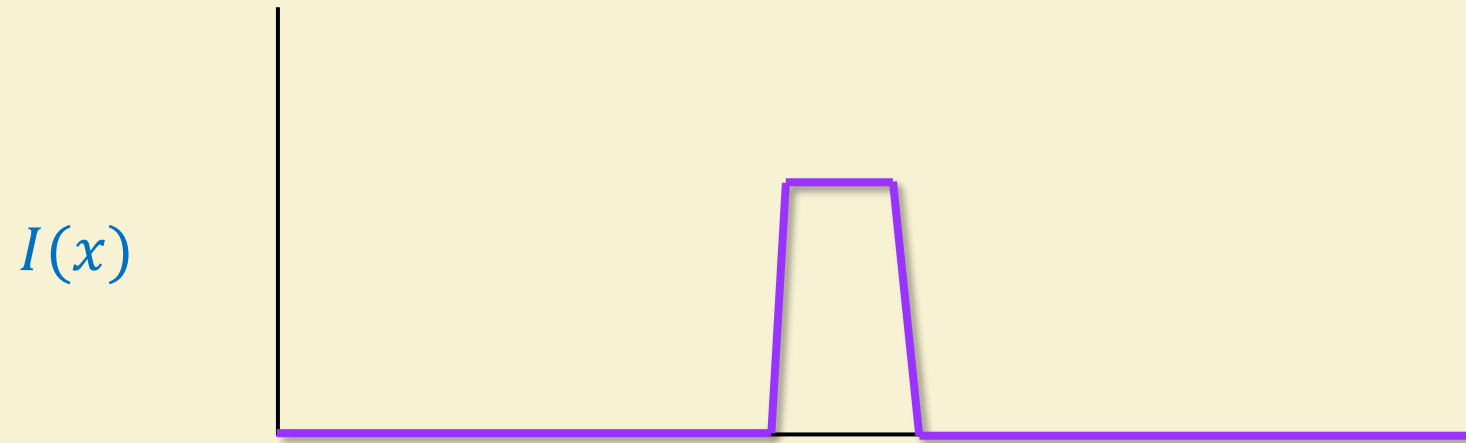
$$I_{a,b} = h_a - h_b$$



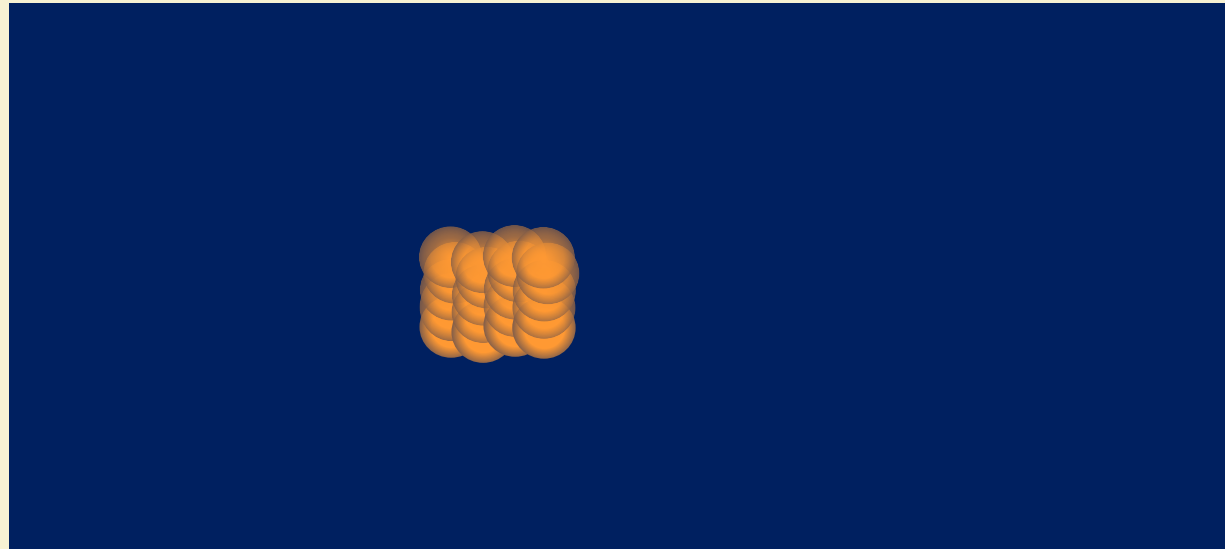
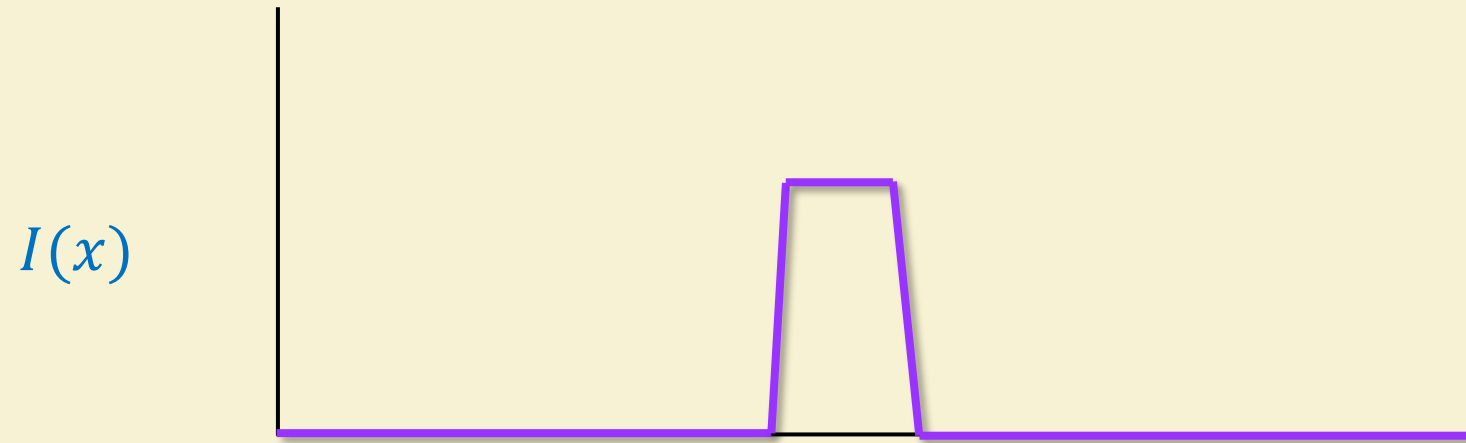
Higher dimension  $r(x) = \max(\langle \alpha, x \rangle + \beta, 0)$  ,  $d + 1$  parameters



Higher dimension  $r(x) = \max(\langle \alpha, x \rangle + \beta, 0)$  ,  $d + 1$  parameters



Higher dimension  $r(x) = \max(\langle \alpha, x \rangle + \beta, 0)$  ,  $d + 1$  parameters



# How many ReLU's

Every  $f: [0,1]^d \rightarrow \mathbb{R}$  can be approximated as sum of  $r_{\alpha,\beta}(x) = \max(\langle \alpha, x \rangle + \beta)$

Can discretize  $\mathbb{R}^{d+1}$  to  $O(1)^d$  points – need at most  $O(1)^d = 2^{O(d)}$  ReLUs

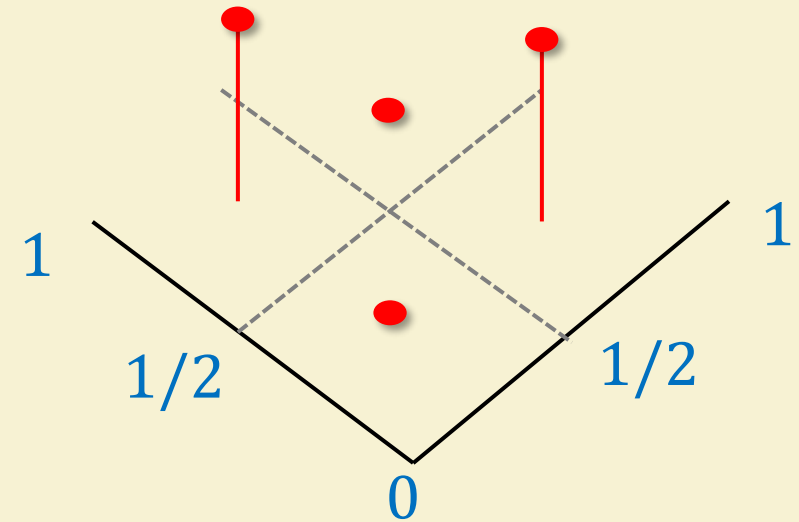
(For formal proofs, see Cybenko '89, Hornik '91, Barron '93; MIT 18.408)

**THM:** For some  $f$ 's,  $2^{c \cdot d}$  ReLU's are **necessary**

Random  $f: [0,1]^d \rightarrow \mathbb{R}$

Choose  $f(x) \in \{0,1\}$  independently at random  
for  $x \in \{1/4, 3/4\}^d$

Extend linearly



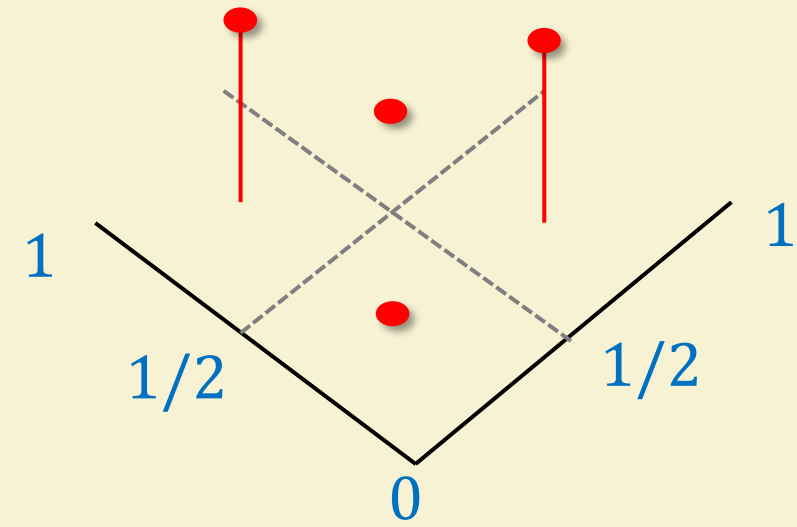
**LEMMA:** For every fixed  $g: [0,1]^d \rightarrow \mathbb{R}$ ,  $\Pr[ \int |f - g| < 0.1 ] \leq 2^{-\epsilon \cdot 2^d}$

**THM:** For some  $f$ 's,  $2^{c \cdot d}$  ReLU's are **necessary**

Random  $f: [0,1]^d \rightarrow \mathbb{R}$

Choose  $f(x) \in \{0,1\}$  independently at random  
for  $x \in \{1/4, 3/4\}^d$

Extend linearly



**LEMMA:** For every fixed  $g: [0,1]^d \rightarrow \mathbb{R}$ ,  $\Pr\left[\int |f - g| < 0.1\right] \leq 2^{-\epsilon \cdot 2^d}$

**LEMMA  $\Rightarrow$  THEOREM:** Need  $2^{\epsilon \cdot 2^d}$  distinct  $g$ 's  $\Rightarrow$  need  $\approx 2^d$  distinct ReLUs

**PROOF OF LEMMA:**

For  $x \in \{1/4, 3/4\}^d$  let  $Z_x = 1$  iff  $\int |g(y) - f(y)| dy > 0.2$  over  $y \in [-1/4, +1/4]^d + x$

Each  $Z_x$  is independent with expectation at least  $1/2$

$$\Rightarrow \Pr\left[\frac{1}{2^d} \sum Z_x < 1/4\right] \leq 2^{-0.1 \cdot 2^d}$$

$$\frac{1}{2^d} \sum Z_x \geq 1/4 \Rightarrow \int |g - f| \geq 1/20$$



# Bottom line

For every\* function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  many choices of  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^N$  s.t.  
$$f \approx \text{Linear}_f \circ \varphi$$

Want to find  $\varphi$  that is **useful** for many interesting functions  $f$

**useful:**

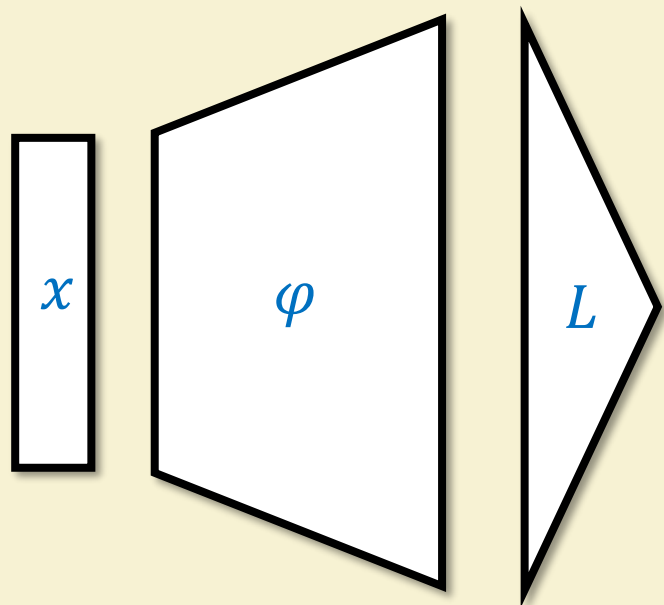
- $N$  is not too big
- Efficiently compute  $\varphi(x)$  or  $\langle \varphi(x), \varphi(y) \rangle$
- $\langle \varphi(x), \varphi(y) \rangle$  large  $\Leftrightarrow x$  and  $y$  are “semantically similar”
- For interesting  $f$ 's, coefficients of  $\text{Linear}_f$  are *structured*.
- ...

Or  $f$ 's such that  $\psi \circ f$  interesting for some non-linear  $\psi: \mathbb{R} \rightarrow \mathbb{R}$

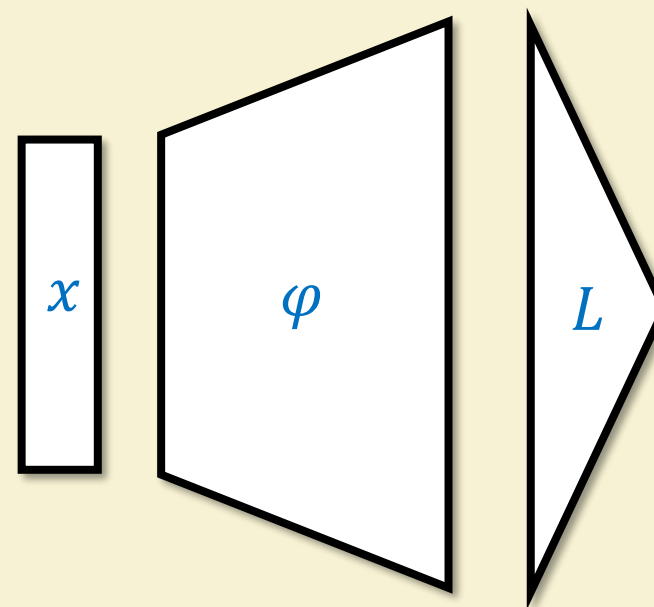


# Part IV: Kernels

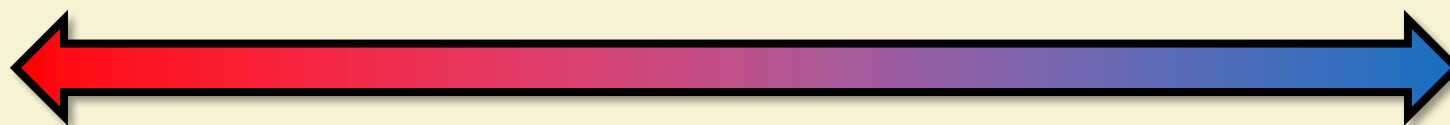
# Neural Net



# Kernel



## Definitely NN



## Definitely Kernel

$\varphi$  trained  
end-to-end  
for task at  
hand

$\varphi$  pretrained  
with same  
data

$\varphi$  pretrained  
on different  
data

$\varphi$  pretrained  
on synthetic  
data

$\varphi$  hand-  
crafted based  
on neural  
architectures

$\varphi$  hand-  
crafted  
before 1990s

# Kernel methods (intuitively)

Distance measure  $K(x, x')$

**Input:**  $(x_1, y_1), \dots, (x_n, y_n)$  s.t.  $f^*(x_i) \approx y_i$

To approx  $f^*(x)$  output  $y_i$  for  $x_i$  closest to  $x$

or output  $\sum \alpha_i y_i$  with  $\alpha_i$  depending on  $K(x, x_i)$

# Hilbert Space

Linear space:  $v + u, c \cdot v$

Dot product:  $\langle u, v + c \cdot w \rangle = \langle u, v \rangle + c \langle u, w \rangle, \langle u, v \rangle = \langle v, u \rangle, \langle v, v \rangle \geq 0$

Can solve linear equations of form  $\{ \langle v_i, x \rangle = b_i \}$  knowing only  $\langle v_i, v_j \rangle$

Also do least-square minimization  $\min_x \sum (\langle v_i, x \rangle - b_i)^2$  knowing only  $\langle v_i, v_j \rangle$

**DEF:** Sym matrix  $K \in \mathbb{R}^{n \times n}$  is p.s.d. if  $v^T K v \geq 0$  for all  $v \in \mathbb{R}^n$

Equivalently  $\lambda_1(K), \dots, \lambda_n(K) \geq 0$

**CLAIM:**  $K$  is p.s.d. iff\*  $u^T K v$  is inner product

**PROOF** ( $\Rightarrow$ ):  $u^T K v = \psi(u) \cdot \psi(v)$  where  $\psi(u_1, \dots, u_n) = (\sqrt{\lambda_1} u_1, \dots, \sqrt{\lambda_n} u_n)$ , expressed in eigenbasis

# Kernel Methods

**Goal:** Solve linear equations, least-square minimization, etc. under  $\varphi$

**Observation:** Enough to compute  $K(x, x') = \langle \varphi(x), \varphi(x') \rangle$

Let  $\hat{x} = \varphi(x)$ , given  $\{\hat{x}_i, y_i\}_{i=1..n}$  want to find  $\hat{w} \in \mathbb{R}^n$  s.t.  $\langle \hat{x}_i, \hat{w} \rangle \approx y_i \quad \forall i$

can compute  $\hat{w} = \sum \alpha_i \hat{x}_i$  using  $K(x_i, x_j)$

To compute prediction on new  $x$  using  $\hat{w}$ , can compute  $\hat{x} = \sum \beta_i \hat{x}_i$  using  $K(x, x_i)$

