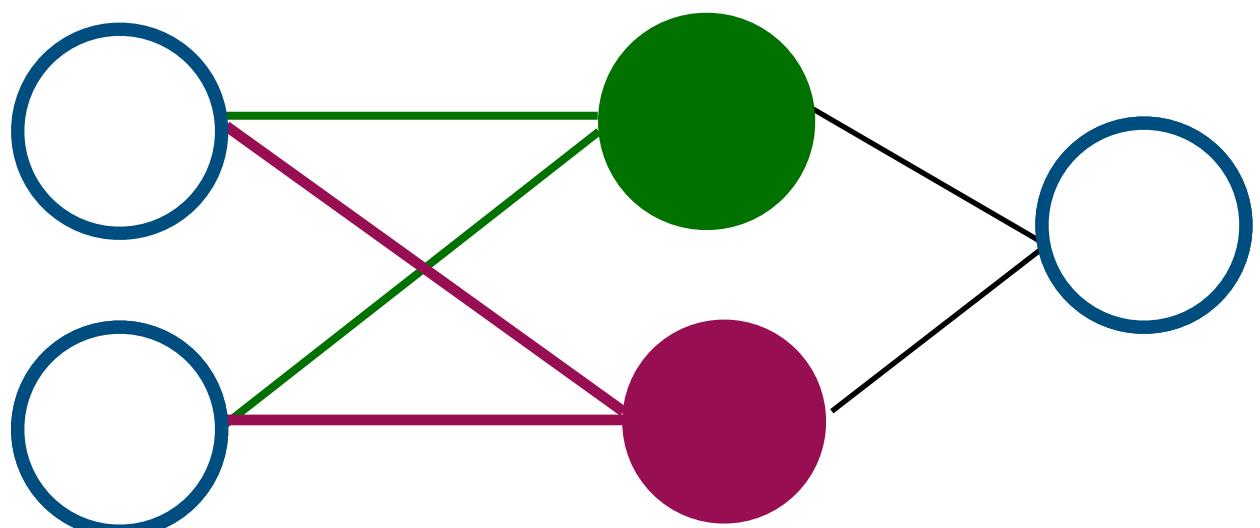


# Neural Networks: A Theory Lab

## Shallow neural networks



# Regression problem

---

2

- ▶ Given  $(x_1, y_1), \dots, (x_n, y_n) \Rightarrow (x_{n+1}, ?)$
- ▶ Where  $y_i = f(x_i)$  where function  $f$  is unknown

# Neural Networks vs Random Features

3

## Random Features

$$\min_{\alpha_i} \sum_j (y_j - \sum_i \alpha_i \cos(\langle v_i, x_j \rangle + b_i))^2$$

$v_i \sim p(w)$  depending on  $k(x, y)$

## Neural Networks

$$\min_{\alpha_i, v_i, b_i} \sum_j (y_j - \sum_i \alpha_i \cos(\langle v_i, x_j \rangle + b_i))^2$$

$v_i$  are optimized since  $p(w)$  is unknown

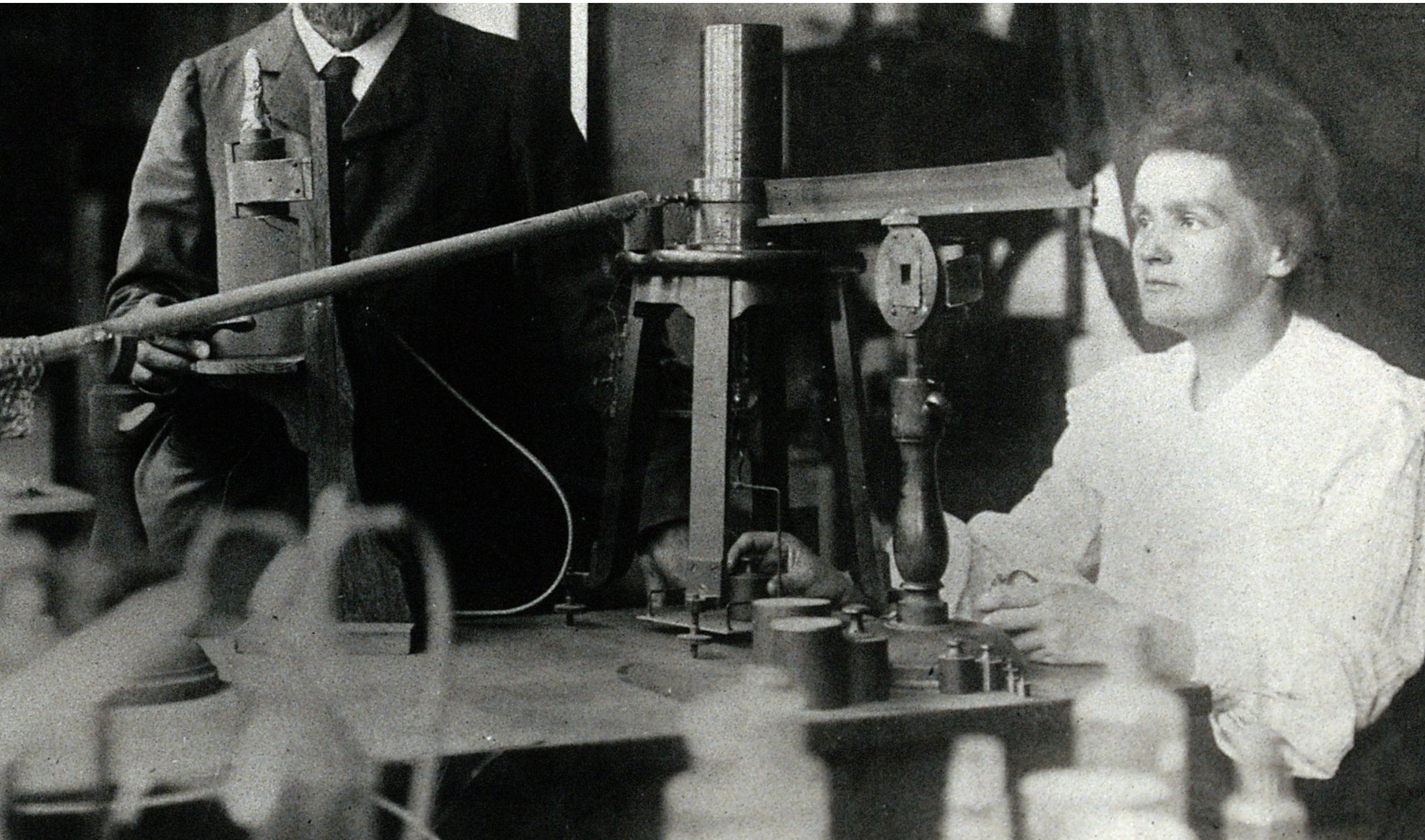
# Goal and outline

---

4

- ▶ Experiments in groups
  - Observing issues with random features
- ▶ Theoretical background
  - Neural networks do not suffer from the curse of dimensionality

- ▶ A brief recap
- ▶ Experiments
- ▶ Theory



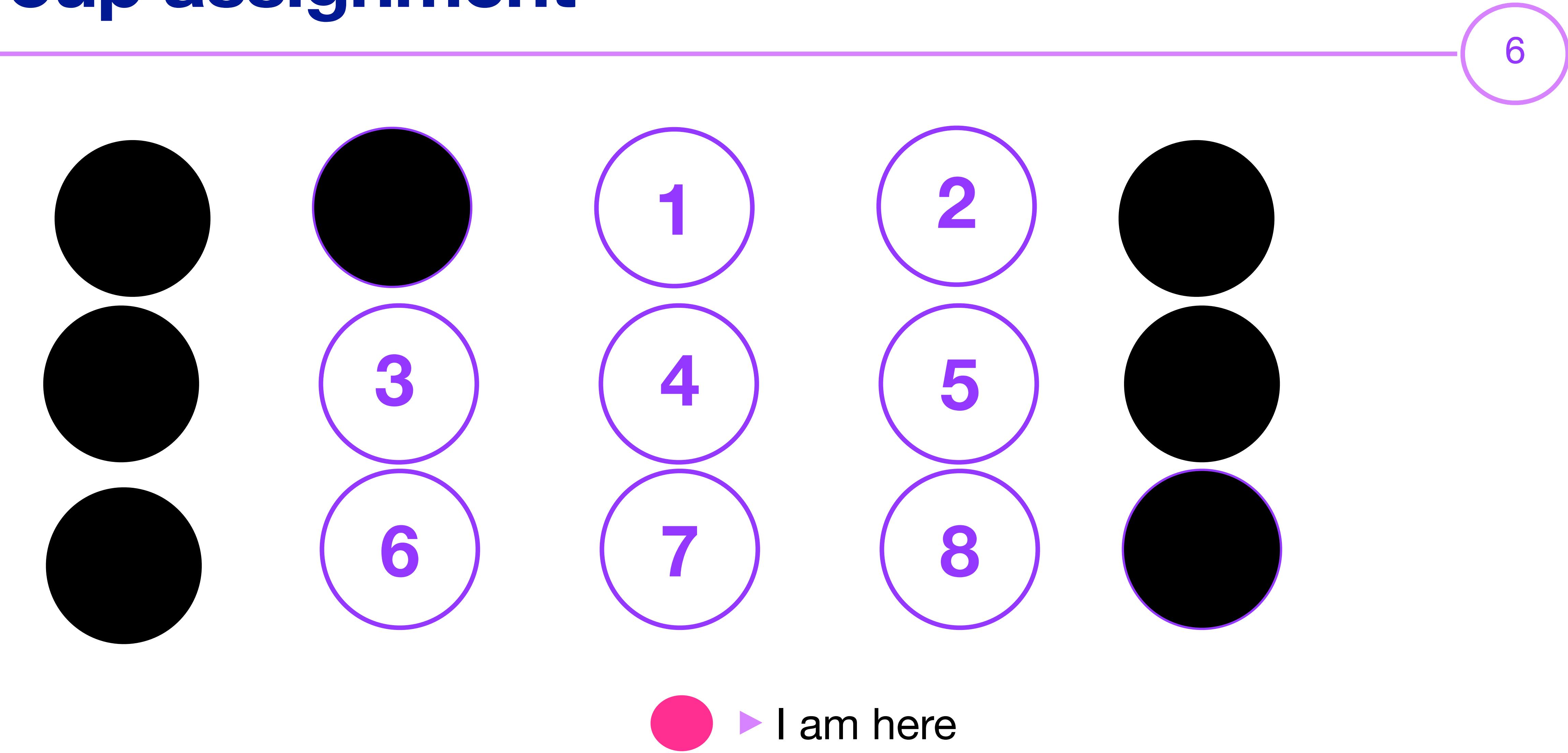
wikipedia: Marie Curie

# Experiments

Get ready for hands-on group activity

# Group assignment

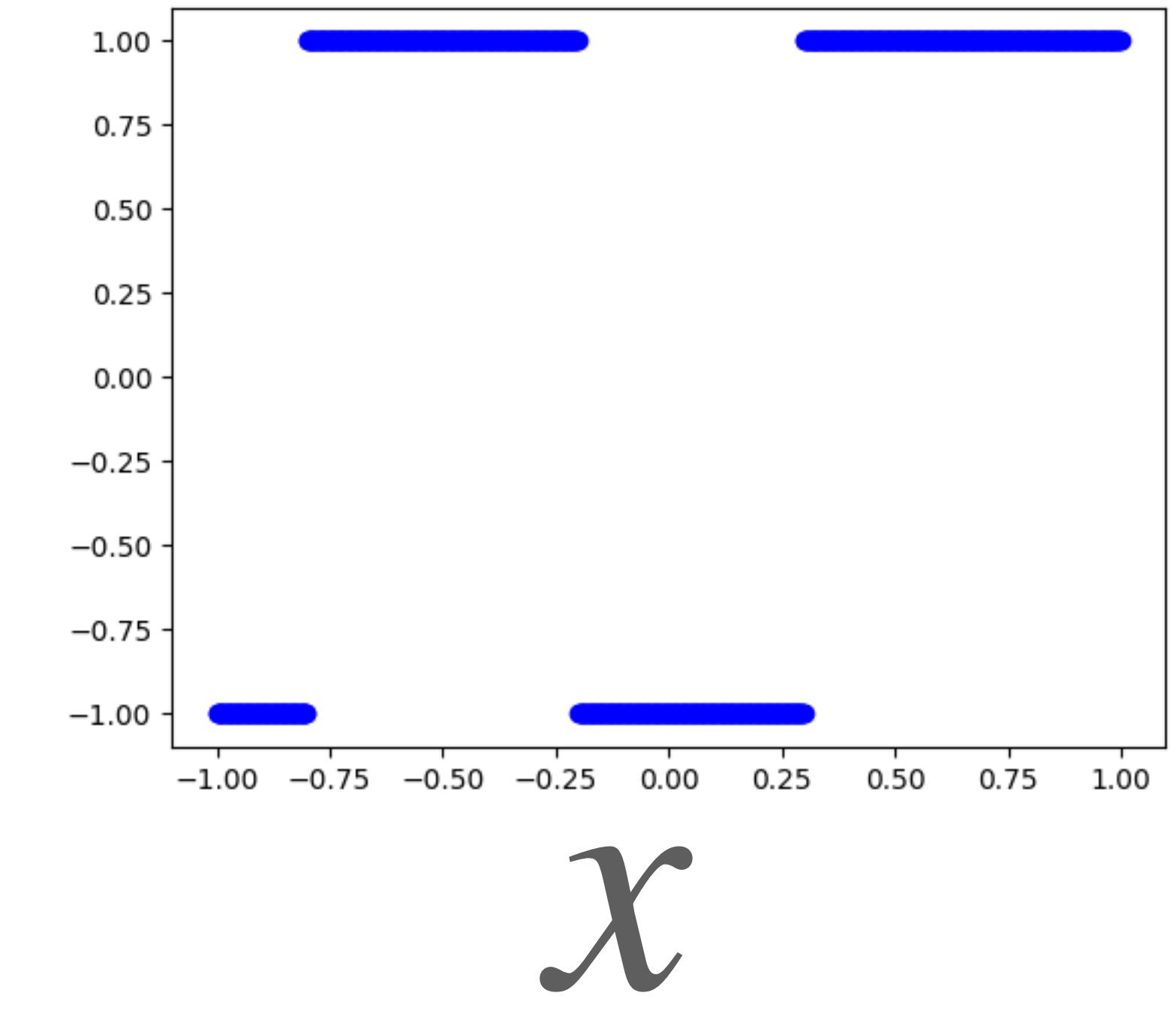
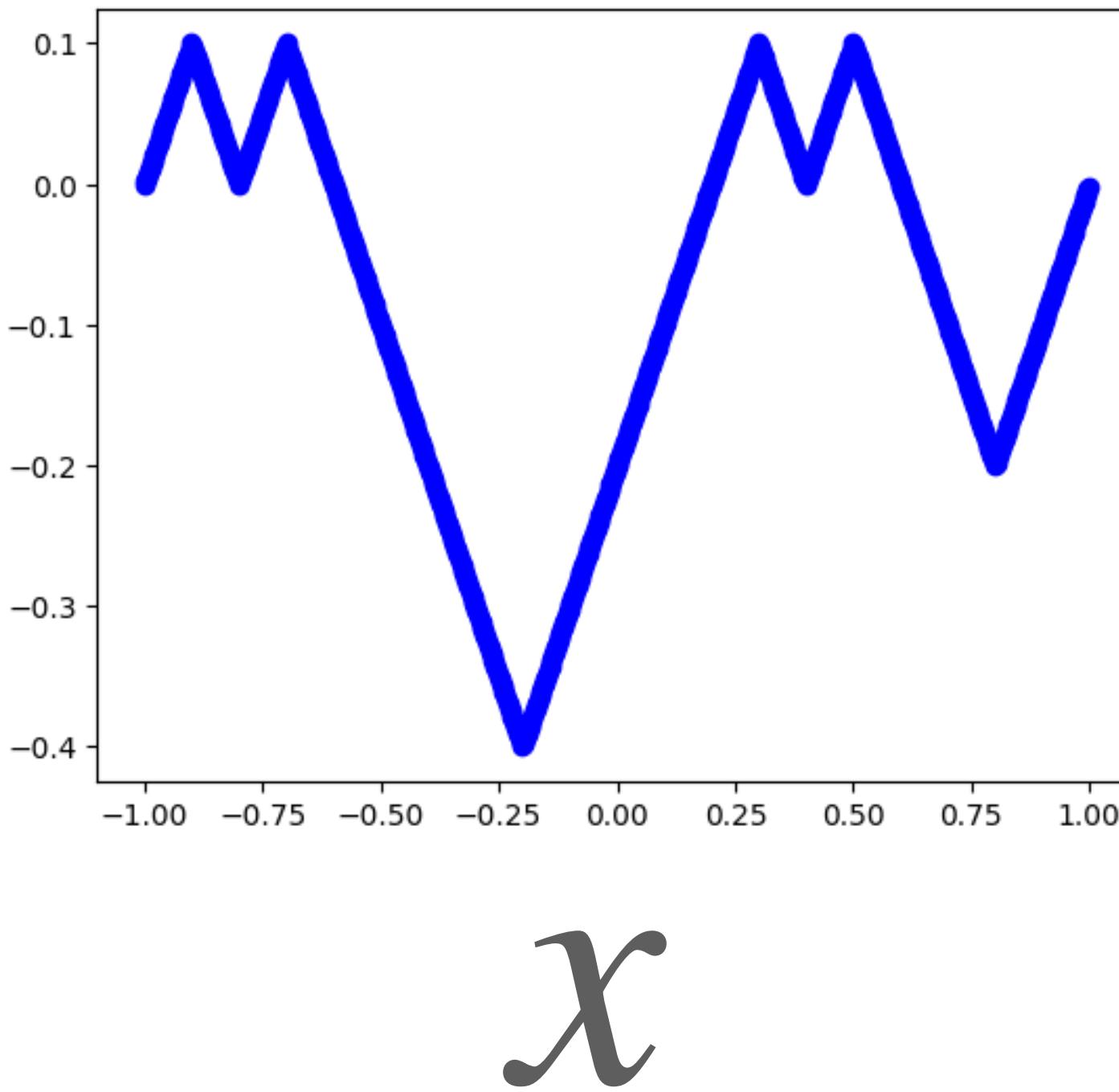
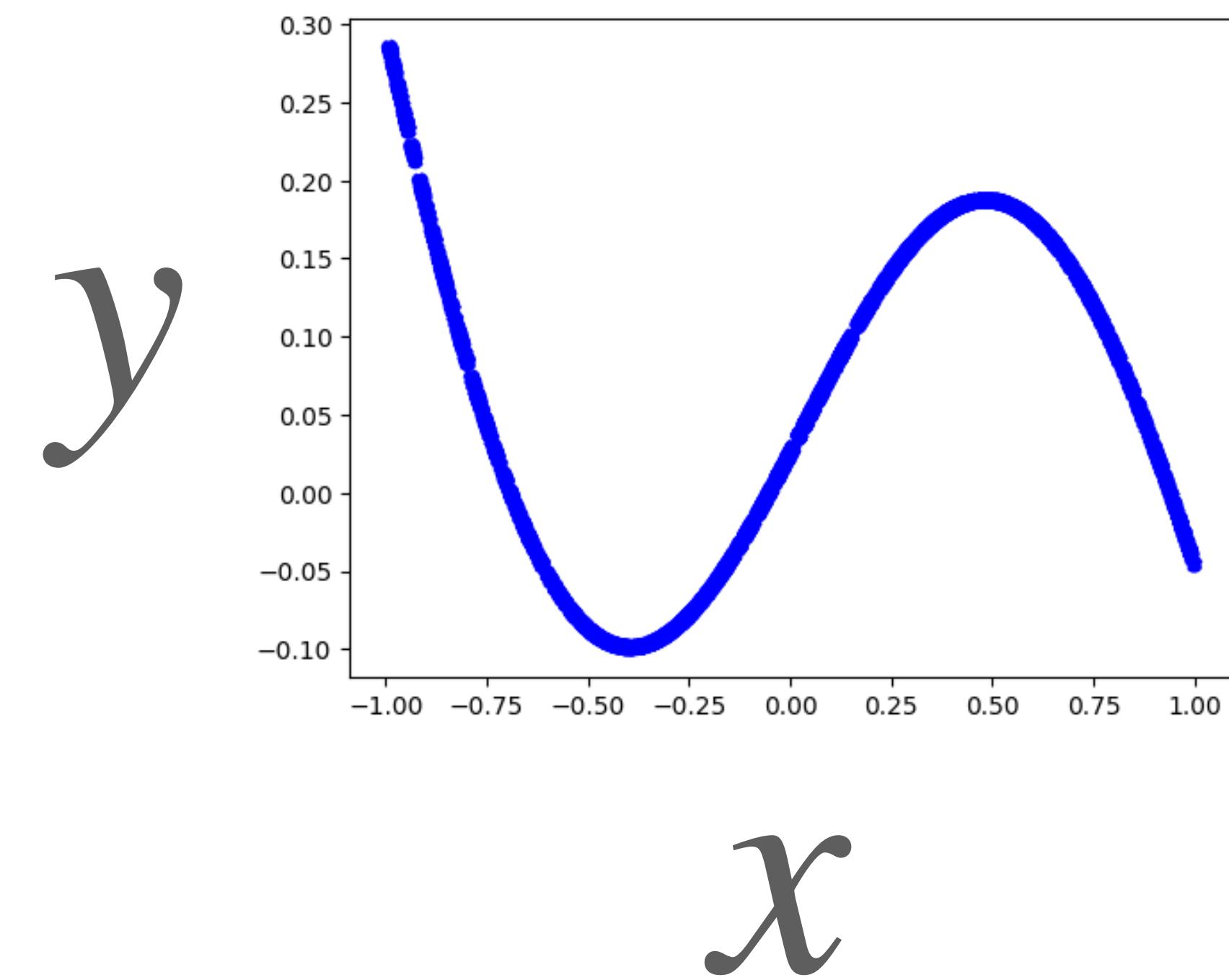
---



# Recap

7

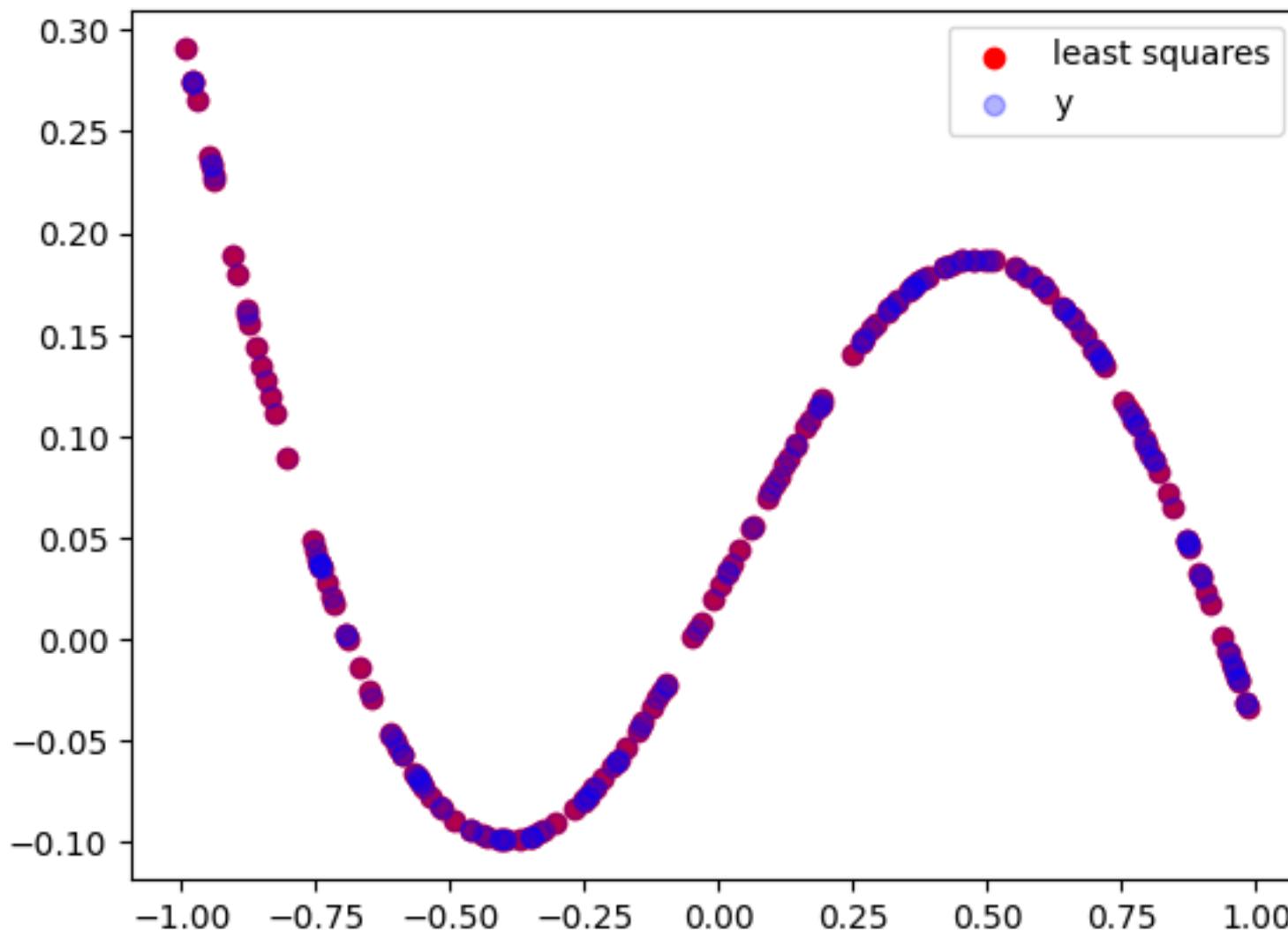
- Goal: We want to design **universal** non-linear features for 3 functions



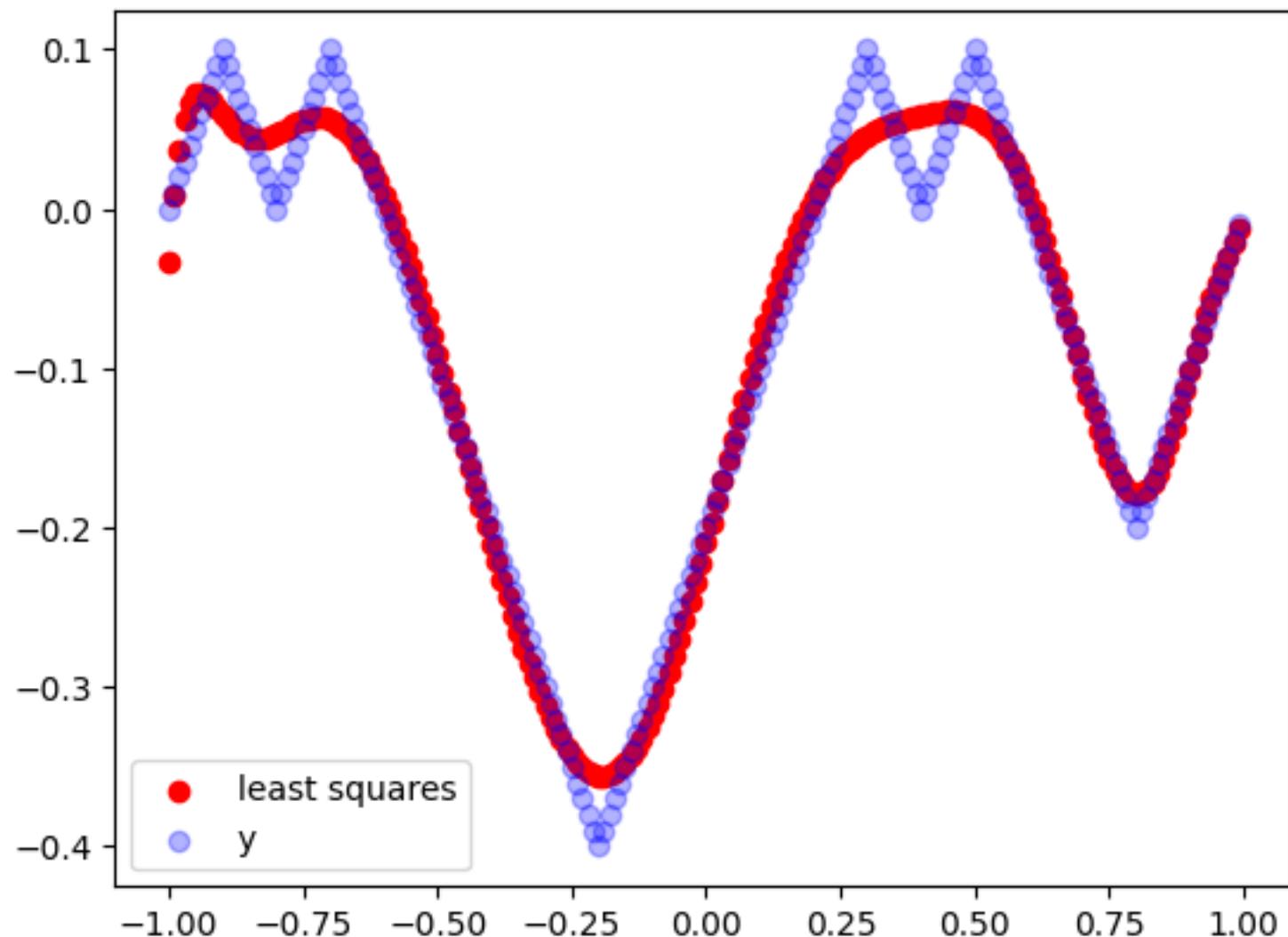
# Performance of random gaussian features

8

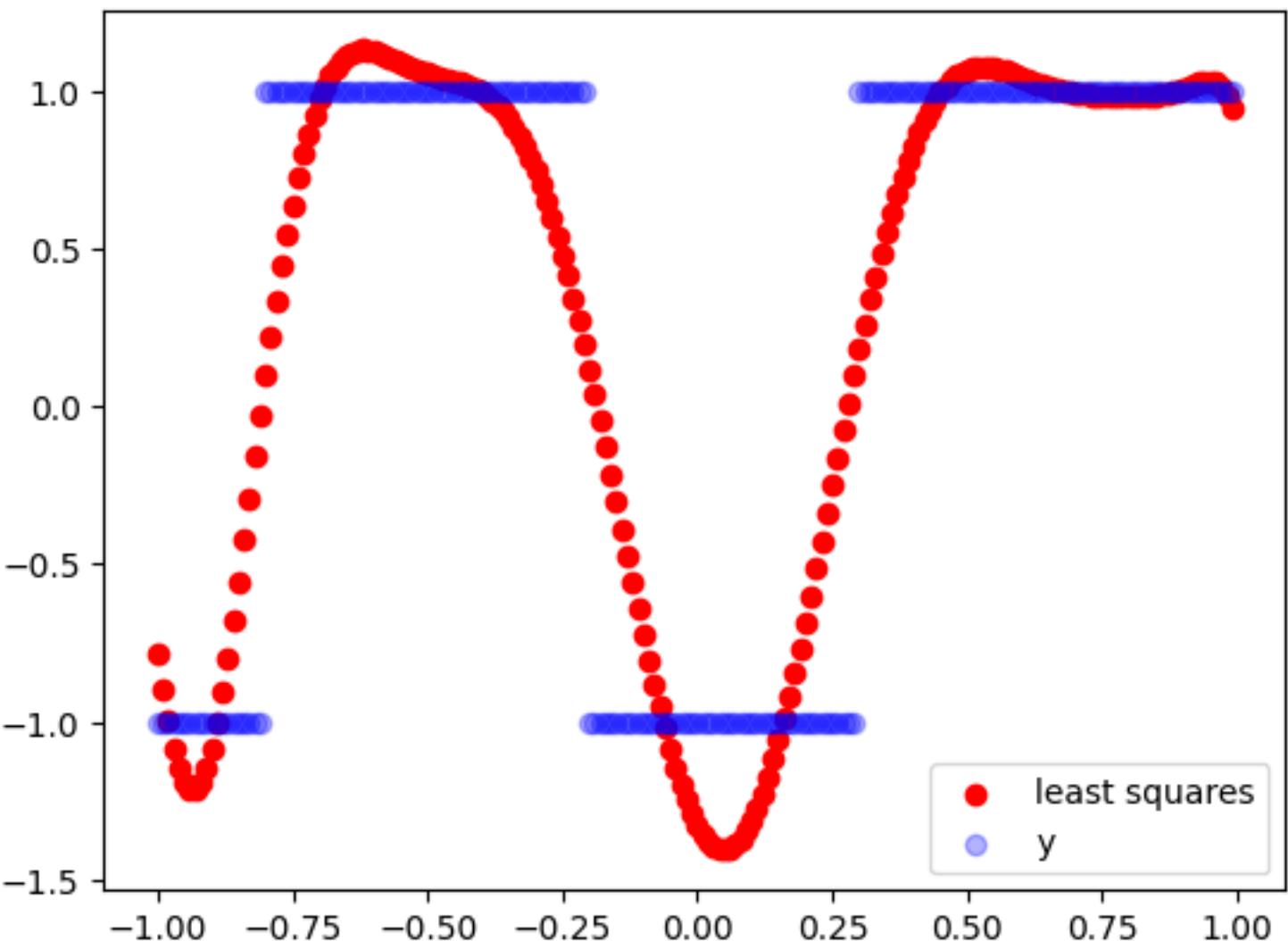
- Blue:  $y$  and red is the solution of regression on random features



Good approximation



Reasonable approximation

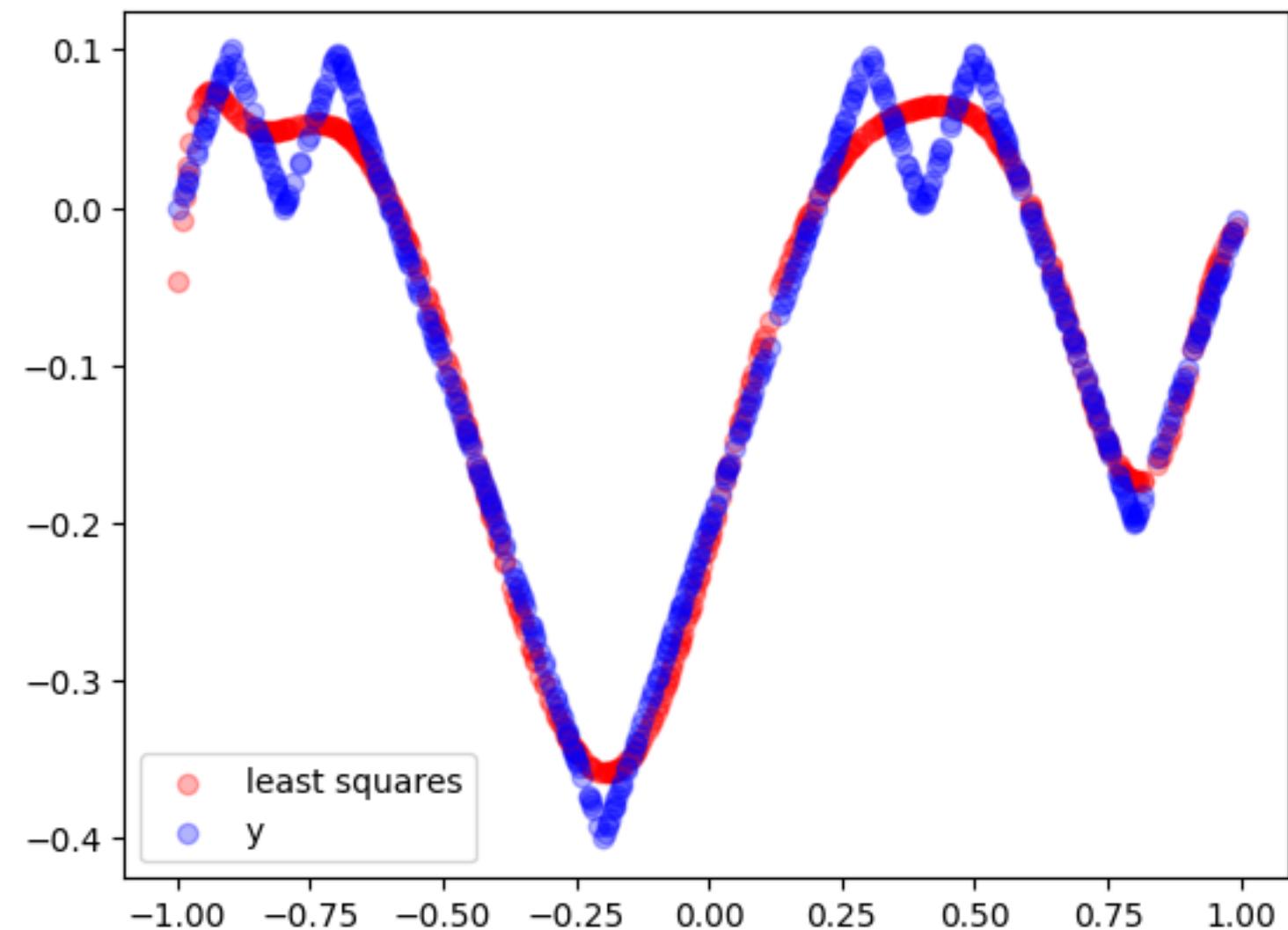


Poor approximation

# Discuss: How to improve random features?

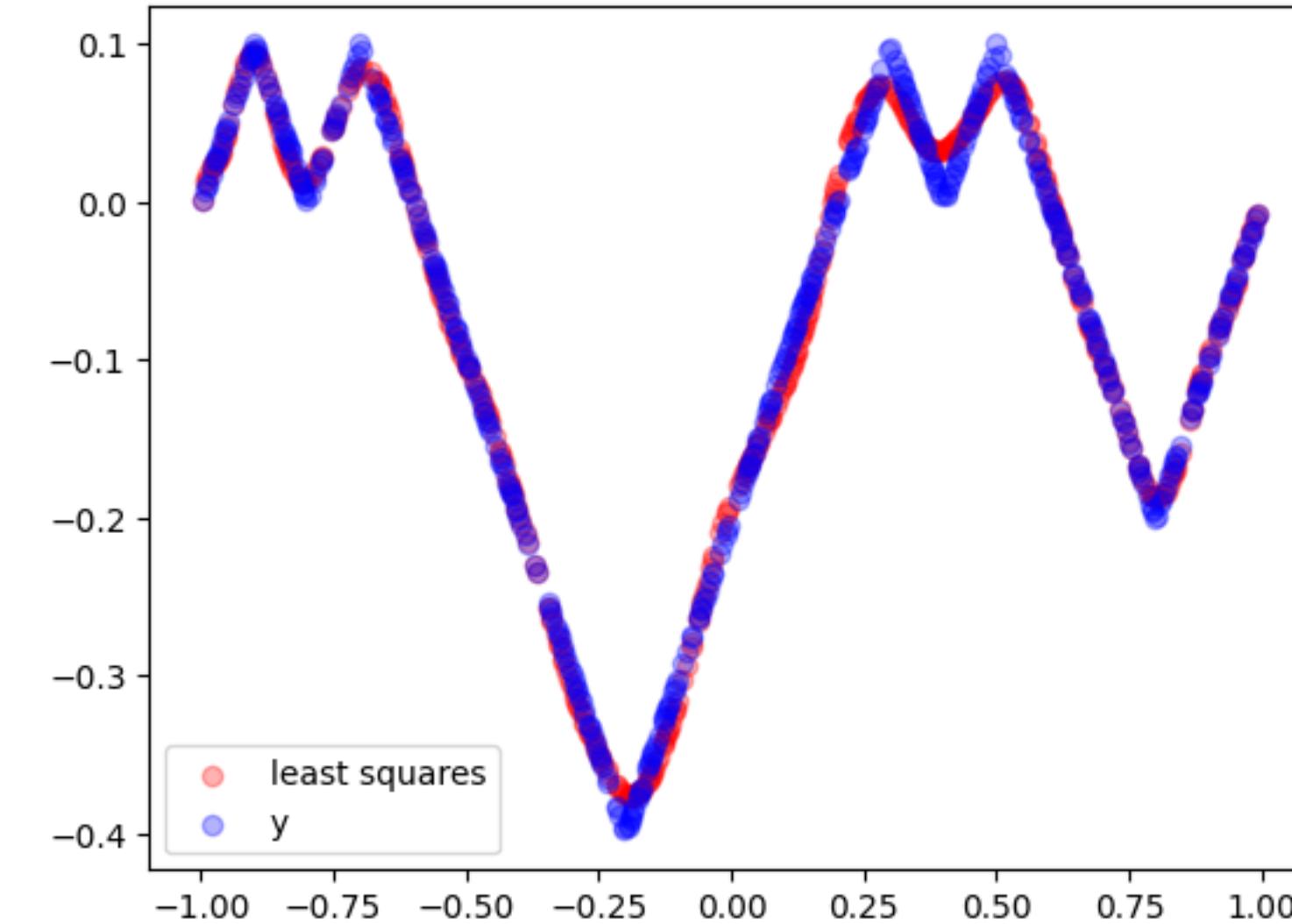
9

I used Gaussian features:  $y \approx \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i), v_i \sim N(0, I_d)$



Reasonable approximation

How?



Ideal approximation

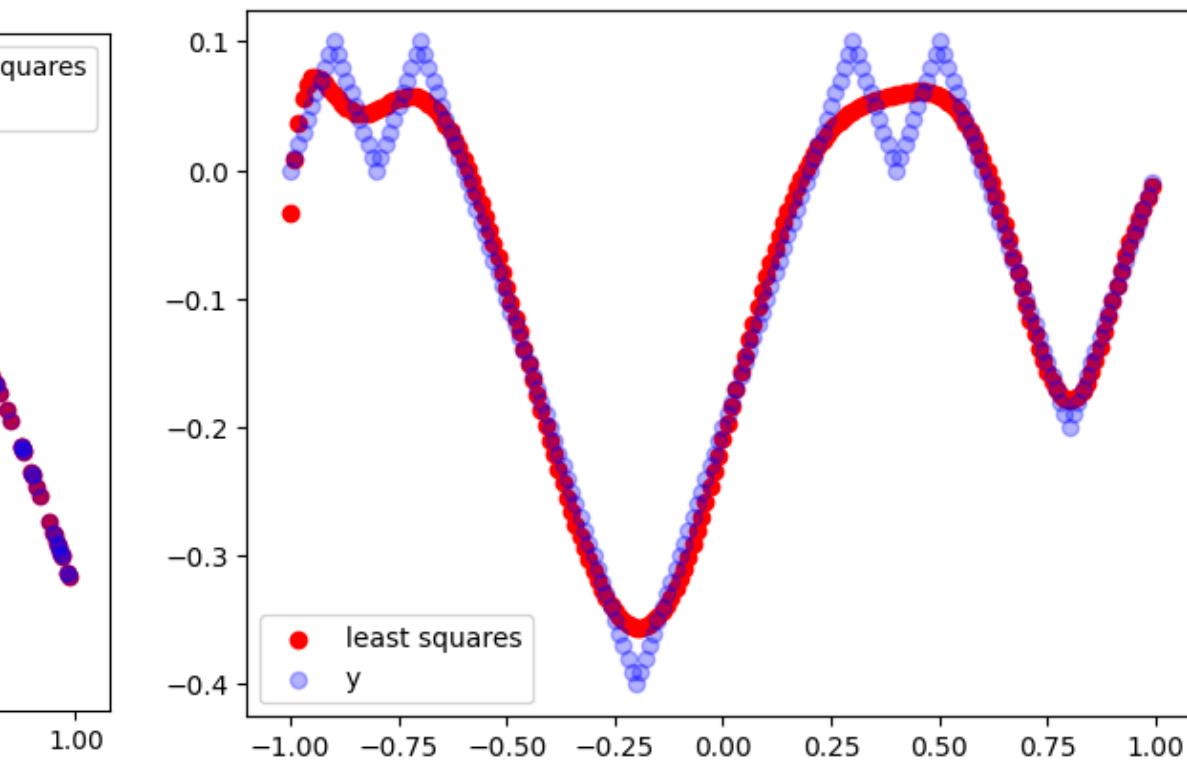
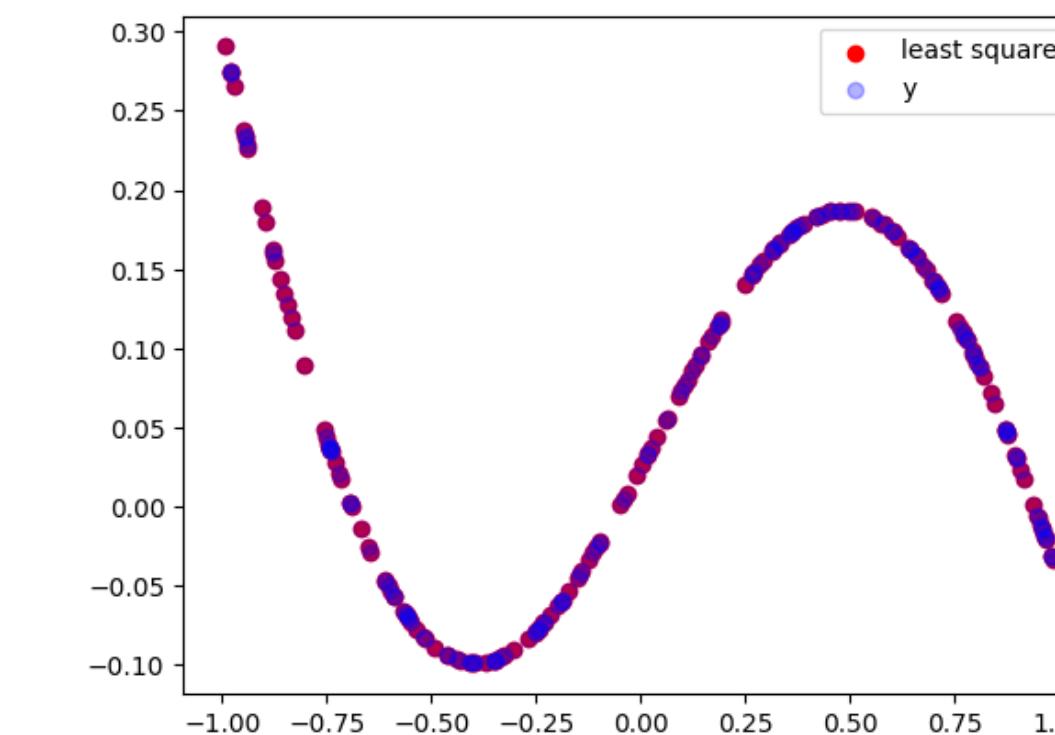
# Hint: use theoretical insights from last lecture

10

Theory

$$k(x - y) = \mathbb{E} [\cos(\langle w, x \rangle + b) \cos(\langle w, y \rangle + b)]$$

Observation



Random features approximate a kernel

Random features can fit various functions

# Distribution choice for random features

11

$$k(x - y) = 2\mathbb{E}_{w \sim p} [\cos(\langle w, x \rangle + b)\cos(\langle w, y \rangle + b)]$$

Kernel name	$k(\Delta)$	$p(w)$
Gaussian	$e^{-\frac{1}{2}\ \Delta\ _2^2}$	$(2\pi)^{-d/2}e^{-\frac{1}{2}\ w\ _2^2}$
Cauchy	$e^{-\ \Delta\ _1}$	$\prod_i \frac{1}{\pi(1 + w_i^2)}$
Laplacian	$\prod_i \frac{1}{\pi(1 + \Delta_i^2)}$	$e^{-\ w\ _1}$

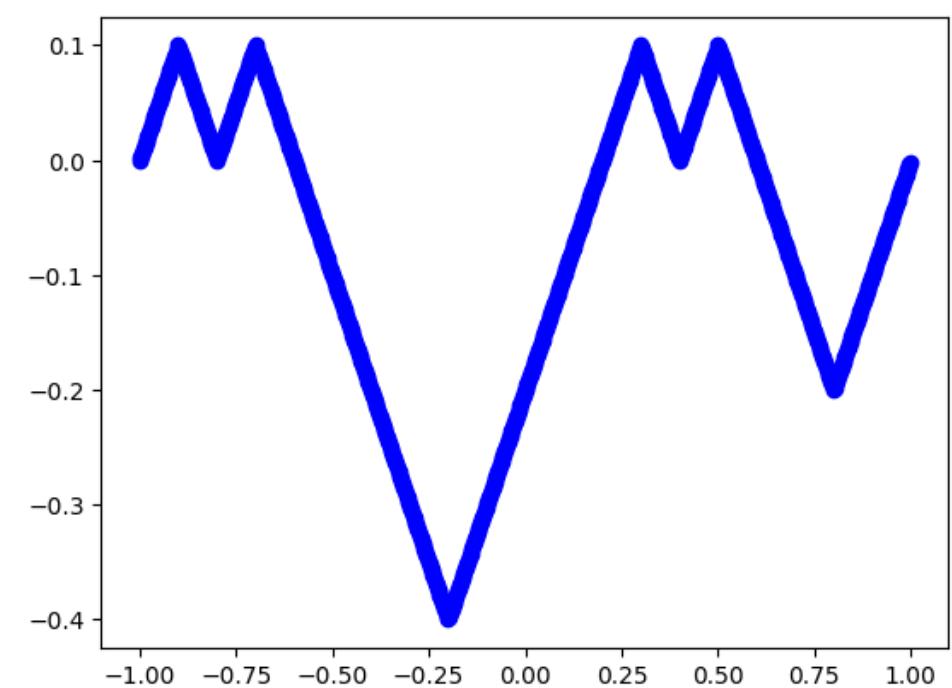
[Rahimi &amp; Recht 2007]

# A solution

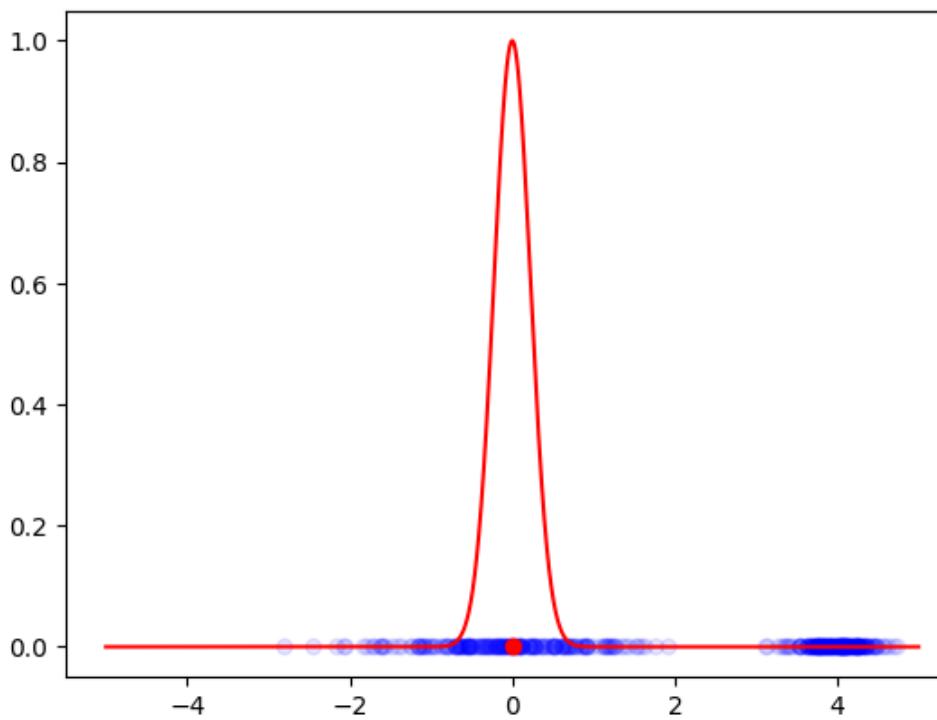
12

- Kernel estimation works well for functions obeying **reproducing property**

Function  $f(x)$



Kernels  $k(x, y)$



Random feature  $\phi_i(x)$

Kernel name	$k(\Delta)$	$p(w)$
Gaussian	$e^{-\frac{1}{2}\ \Delta\ _2^2}$	$(2\pi)^{-d/2} e^{-\frac{1}{2}\ \Delta\ _2^2}$
Cauchy	$e^{-\ \Delta\ _1}$	$\prod_i \frac{1}{\pi(1 + w_i^2)}$
Laplacian	$\prod_i \frac{1}{\pi(1 + \Delta_i^2)}$	$e^{-\ w\ _1}$

$$\phi_i(x) = \cos(\langle w^{(i)}, x \rangle + b_i), w_i \sim p(w)$$

# A solution

- ▶ Kernel estimation works well for functions obeying **reproducing property**
- ▶ Idea: adapting the distribution of random features

Random feature  $\phi_i(x)$

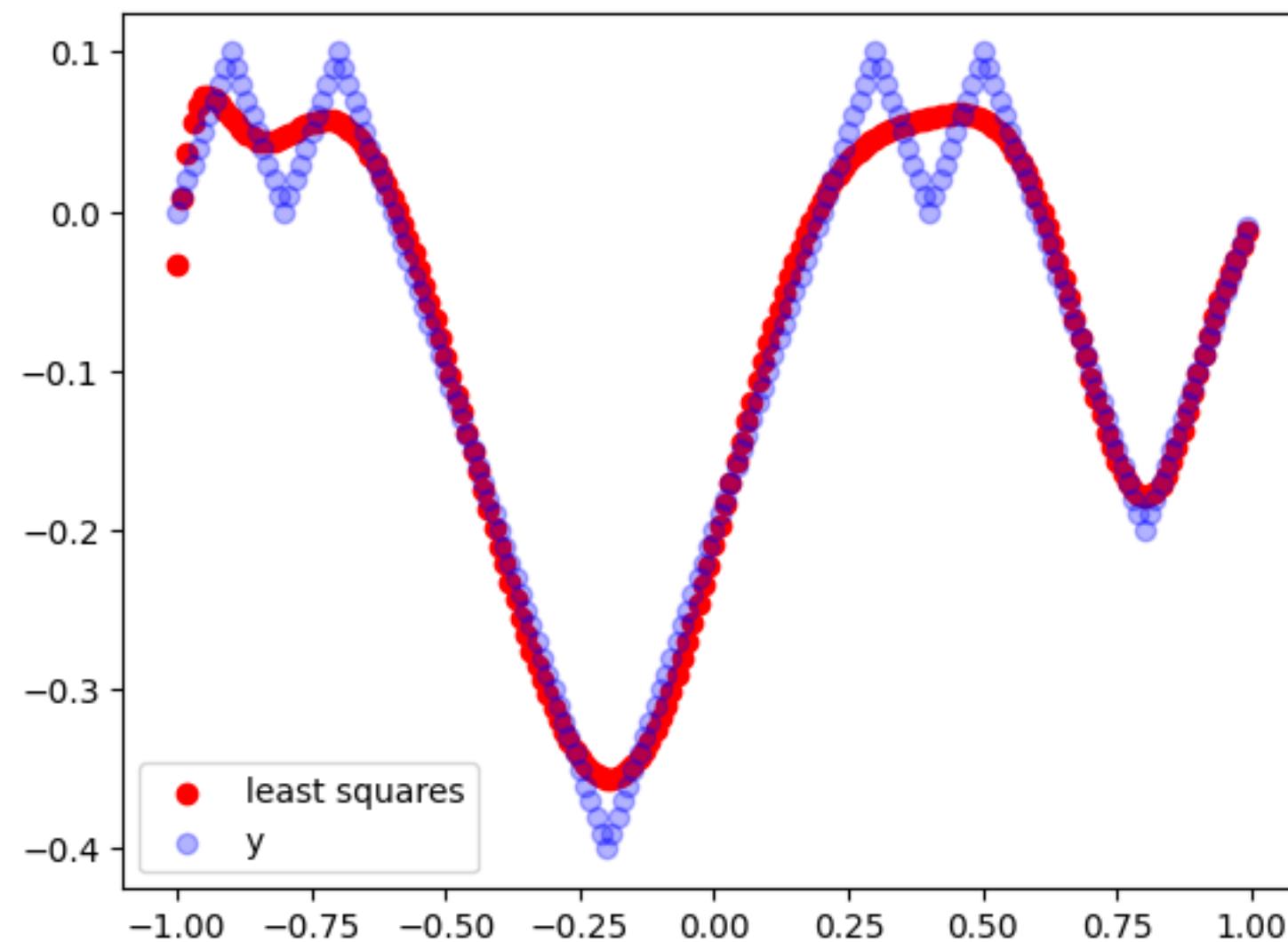
Kernel name	$k(\Delta)$	$p(w)$
Gaussian	$e^{-\frac{1}{2}\ \Delta\ _2^2}$	$(2\pi)^{-d/2}e^{-\frac{1}{2}\ w\ _2^2}$
Cauchy	$e^{-\ \Delta\ _1}$	$\prod_i \frac{1}{\pi(1 + w_i^2)}$
Laplacian [Rahimi &	$\prod_i \frac{1}{\pi(1 + \Delta_i^2)}$	$e^{-\ w\ _1}$

$$\phi_i(x) = \cos(\langle w^{(i)}, x \rangle + b_i), w_i \sim p(w)$$

# Group activity: implement Cauchy random features

14

Gaussian random feature



Cauchy random feature



$$p(w) = (2\pi)^{-d/2} e^{-\frac{1}{2}\|\Delta\|_2^2}$$

$$p(w) = \prod_i \frac{1}{\pi(1 + w_i^2)}$$

# Group activity: implementation

---

15

$$\triangleright y \approx \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i), v_i \sim \text{cauchy}$$

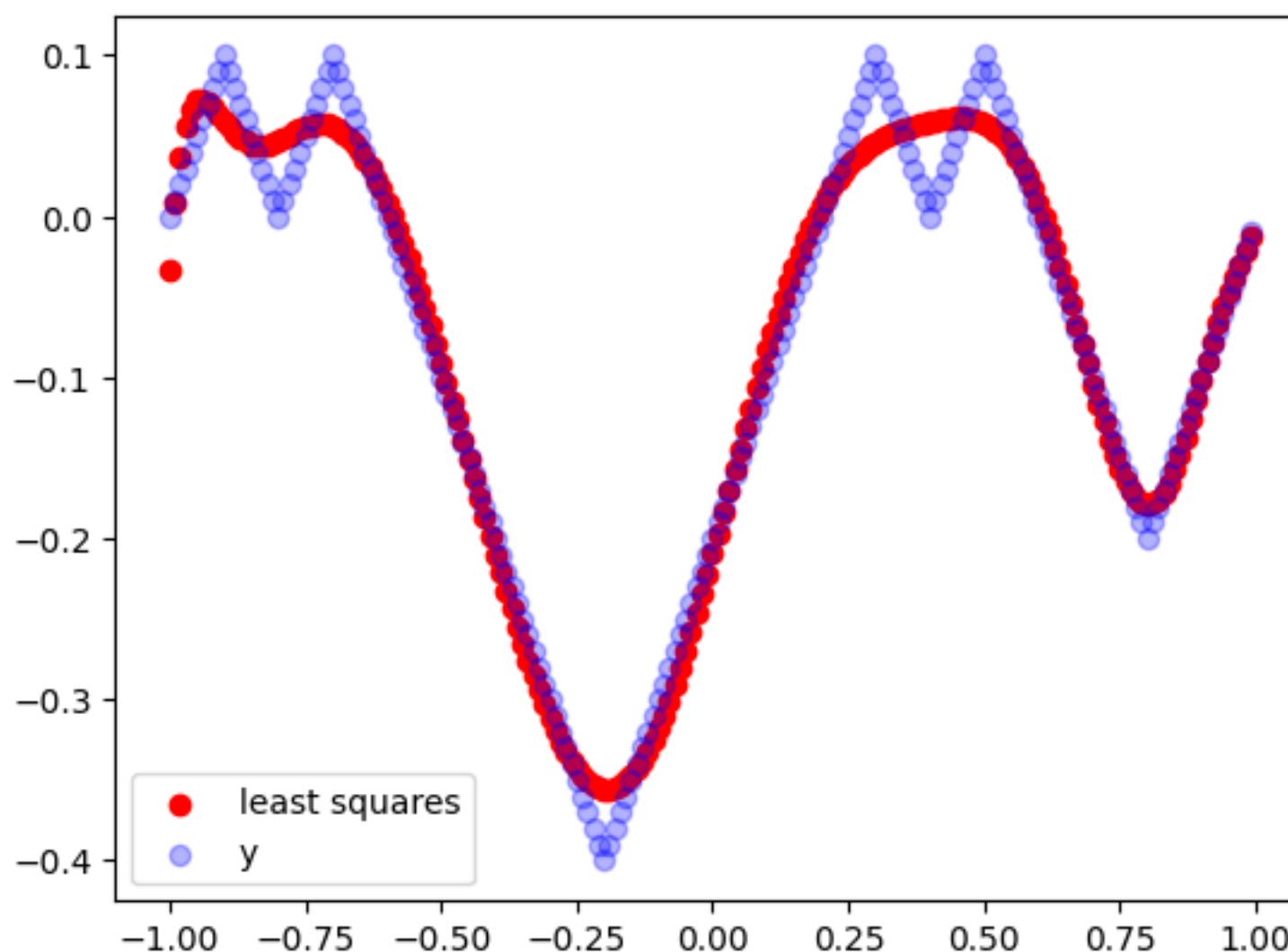
**Task 3 in colab:** <https://shorturl.at/ABb6H>

<https://colab.research.google.com/drive/1M-nRBdhg1XJiV8sy4wMkUsfPJKKKSCB0?usp=sharing>

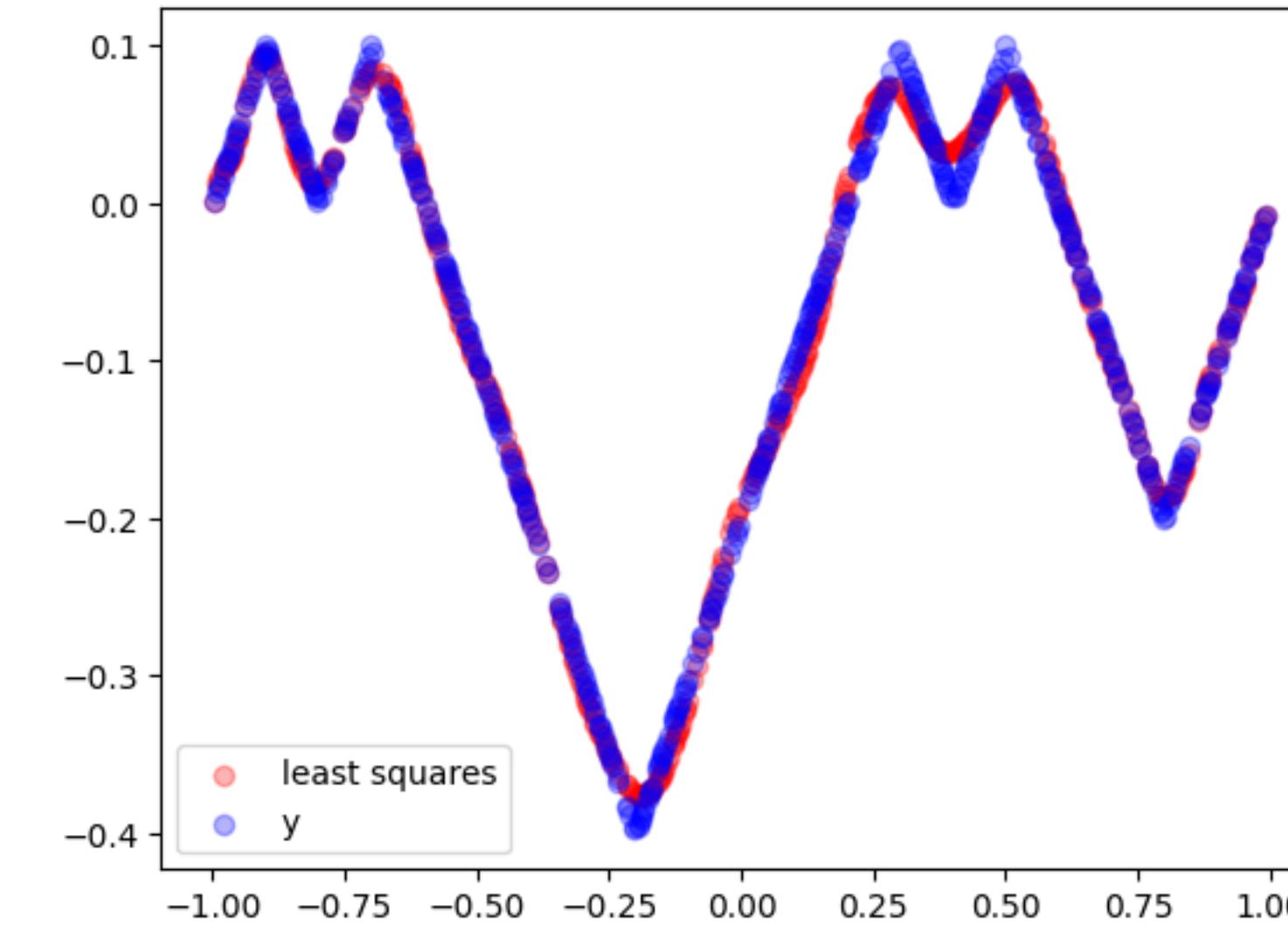
# Comparing random features

16

Gaussian random feature



Cauchy random feature



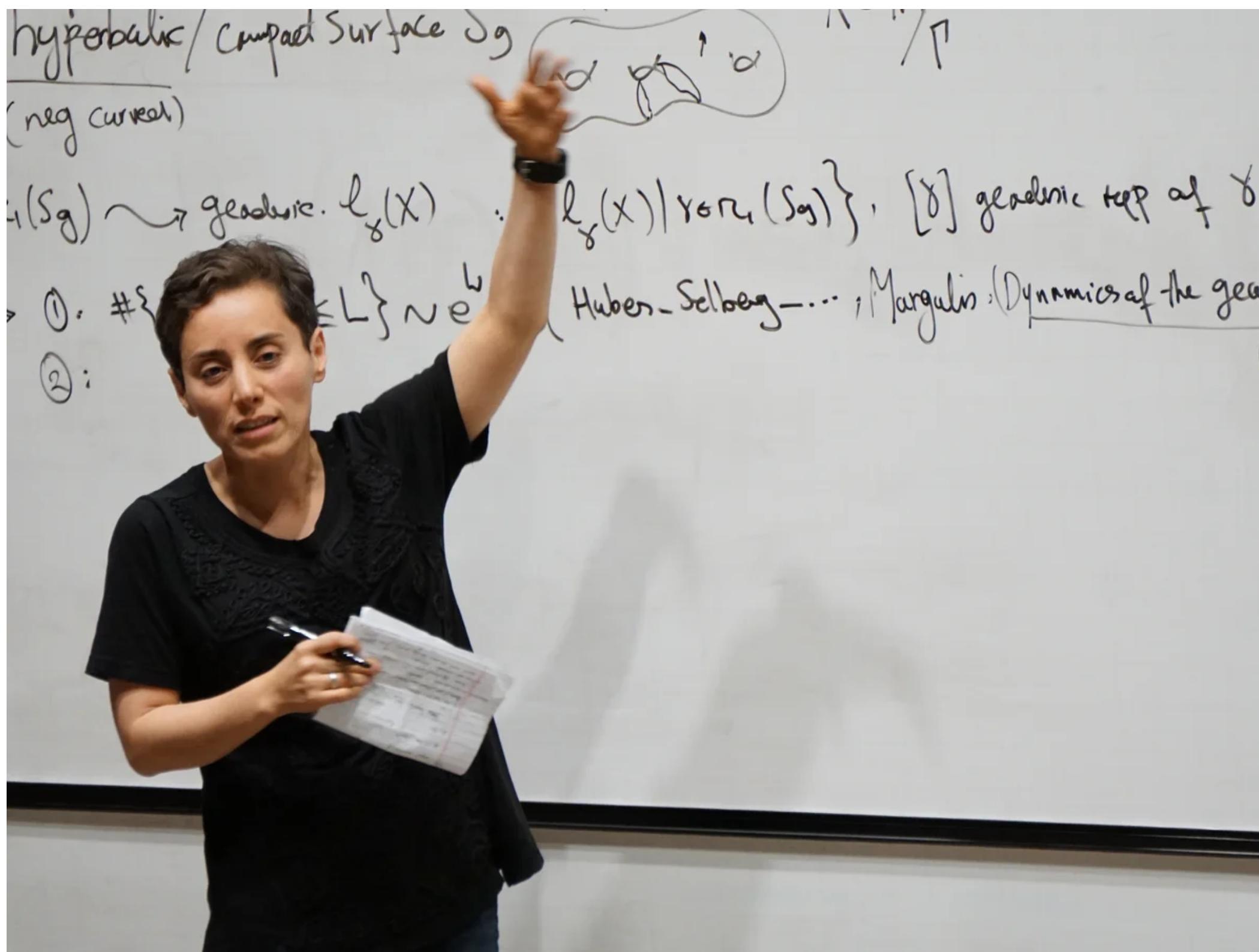
$$p(w) = (2\pi)^{-d/2} e^{-\frac{1}{2}\|\Delta\|_2^2}$$

$$p(w) = \prod_i \frac{1}{\pi(1 + w_i^2)}$$

► Recap

► Experiments

► Theory



theguardian: Maryam Mirzakhani

# Theory

Get ready for math

# Optimization over distributions

---

- Random features approximate  $y$  by a linear transformation of non-linear features:

$$\min_{\alpha_i} \mathbb{E} \left( \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i) - y \right)^2, v_i \sim p(w), b_i \sim \text{uniform}$$

- For an unknown  $p$ :  $\min_{\alpha,p,b} \mathbb{E} \left( \int \alpha(v) \cos(\langle v, x \rangle + b(v)) p(v) dv - y \right)^2$
- $p, \alpha, b : \mathbb{R}^d \rightarrow \mathbb{R}$  are functions.

# Particle optimization

---

Abstract problem:  $\arg \min_p \mathcal{L} \left( \int \phi(v)p(v) \right)$ ,  $p$  is a probability measure

How can we optimize with respect to probability measure?

Particle method:  $\arg \min_{v_1, \dots, v_m} \mathcal{L} \left( \frac{1}{m} \sum_{i=1}^m \phi(v_i) \right)$

# An example of particle optimization

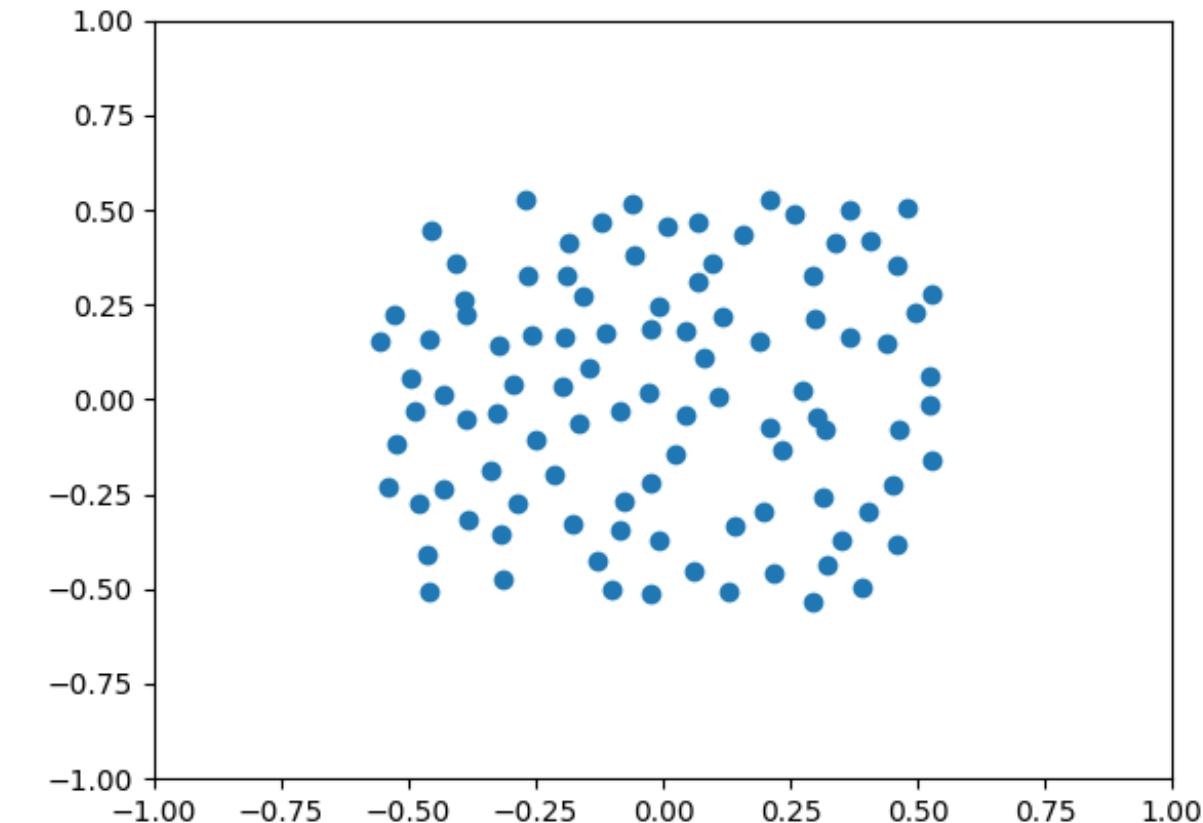
20

$$\mathcal{L}(p) := \frac{1}{2} \min_p \int \mathcal{K}(x - y) dp(x) dp(y), \quad \mathcal{K}(\Delta) = -\log(\|\Delta\|) + \|\Delta\|^2$$

Surprise 1: the solution in 2D is the uniform measure over a circle

[See [Global Minimizers of a Large Class of Anisotropic Attractive-Repulsive Interaction Energies in 2D](#)]

Particle Method:  $\min_{v_1, \dots, v_m} \sum_{i=1}^m \sum_{j \neq i, 1}^m \mathcal{K}(v_i - v_j)$



Particle gradient descent is extensively studied in physics

[see works of [José A. Carrillo](#)]

# Neural Networks as a particle method

- Random features with an unknown distribution  $p$

$$\min_{\alpha, p, b} \mathbb{E} \left( \underbrace{\int \alpha(v) \cos(\langle v, x \rangle + b(v)) p(v) dw - y}_{f(x)} \right)^2$$

- $p, a, b : \mathbb{R}^d \rightarrow \mathbb{R}$  are functions.

► Neural Networks:  $\min_{\alpha_i, v_i, b_i} \mathbb{E} \left( \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i) - y \right)^2$

# Barron's theorem

A. R. Barron, "Universal approximation bounds for superpositions of a sigmoidal function," in *IEEE Transactions on Information Theory*, 1993

22

- ▶ Suppose  $y = f(x)$  has the Fourier transform

$$\int e^{i\langle w, x \rangle} \hat{f}(w) dw \text{ such that } C_f = \int |w| |\hat{f}(w)| dw \text{ is finite.}$$

- ▶ Then, there are  $v_1, \dots, v_n \in \mathbb{R}^d$  and  $b_1, \dots, b_n \in \mathbb{R}$  such that

$$\int \left( y - \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i) \right)^2 p(x) dx \leq \frac{(2rC_f)^2}{n}$$

holds where  $p(x)$  is a probability measure over the ball of radius  $r$ .

# Barrons' comparison

A. R. Barron, "Universal approximation bounds for superpositions of a sigmoidal function," in *IEEE Transactions on Information Theory*, 1993

23

$$f(x) \approx \underbrace{\sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i)}_{f_n(x)}$$

Random Features

$v_i$  independent of  $y$

$$\mathbb{E}(f(x) - f_n(x))^2 \geq c \frac{C_f}{d} \left( \frac{1}{n} \right)^{1/d}$$

$n = \Omega((\frac{1}{\epsilon})^d)$  for  $\epsilon$ -accuracy

Neural Networks

$v_i$  are optimized for  $y$

$$\mathbb{E}(f(x) - f_n(x))^2 = O\left(\frac{1}{n}\right)$$

$n = O(\frac{1}{\epsilon})$  for  $\epsilon$ -accuracy

# Curse of dimensionality for function approximation

24

## Random Features

$v_i$  independent of  $y$

$$\mathbb{E}(f(x) - f_n(x))^2 \geq c \frac{C_f}{d} \left(\frac{1}{n}\right)^{1/d}$$

$n = \Omega((\frac{1}{\epsilon})^d)$  for  $\text{Neural networks can break curse of dimensionality}$

# Neural Networks vs Random Features

25

## Random Features

$$y \approx \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i)$$

$v_i \sim p(w)$  depending on  $k(x, y)$

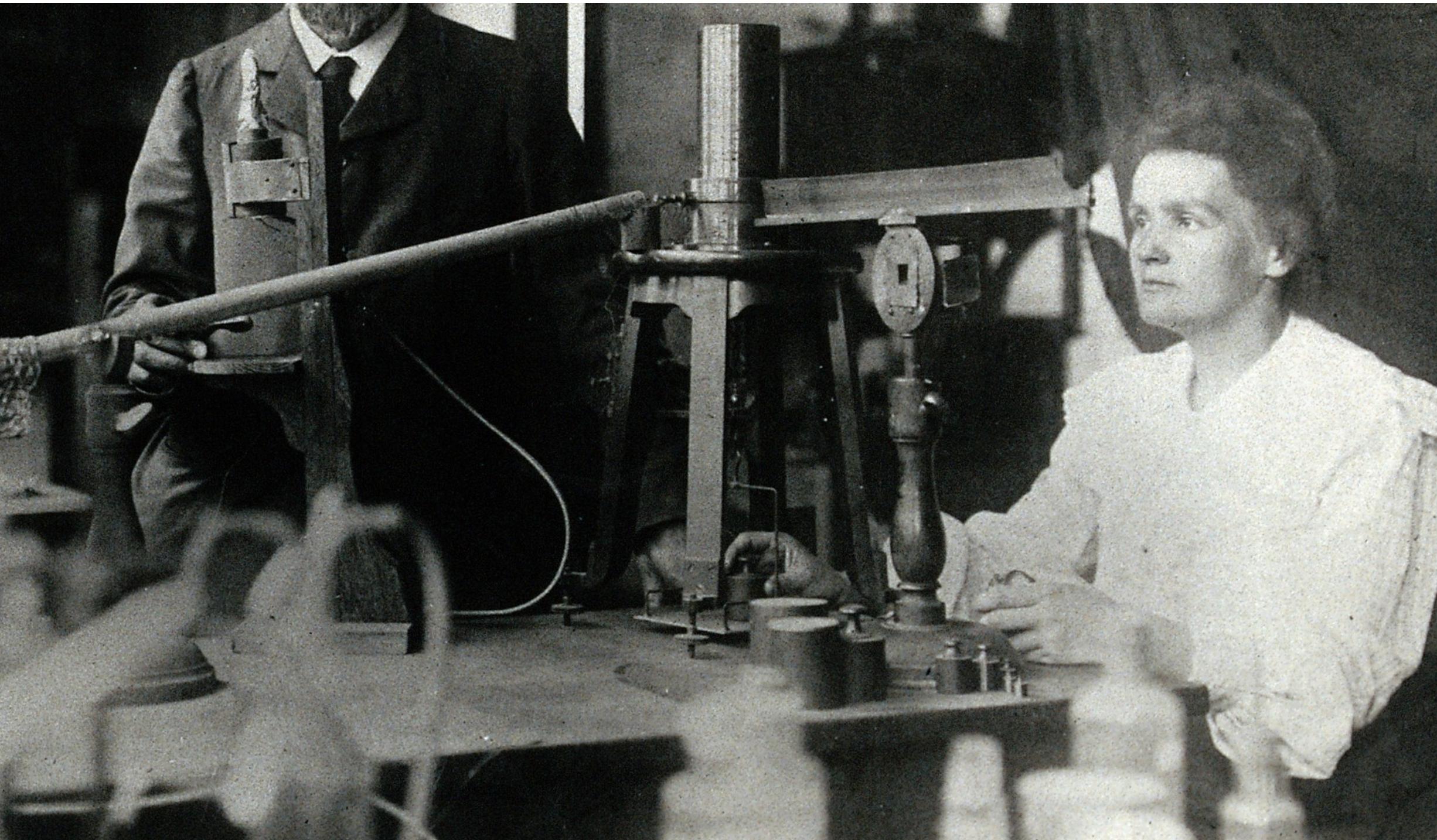
## Neural Networks

$$y \approx \sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i)$$

$v_i$  are optimized since  $p(w)$  is unknown

Two challenges of random features:

{  
The choice of  $p$   
Curse of dimensionality}



wikipedia: Marie Curie

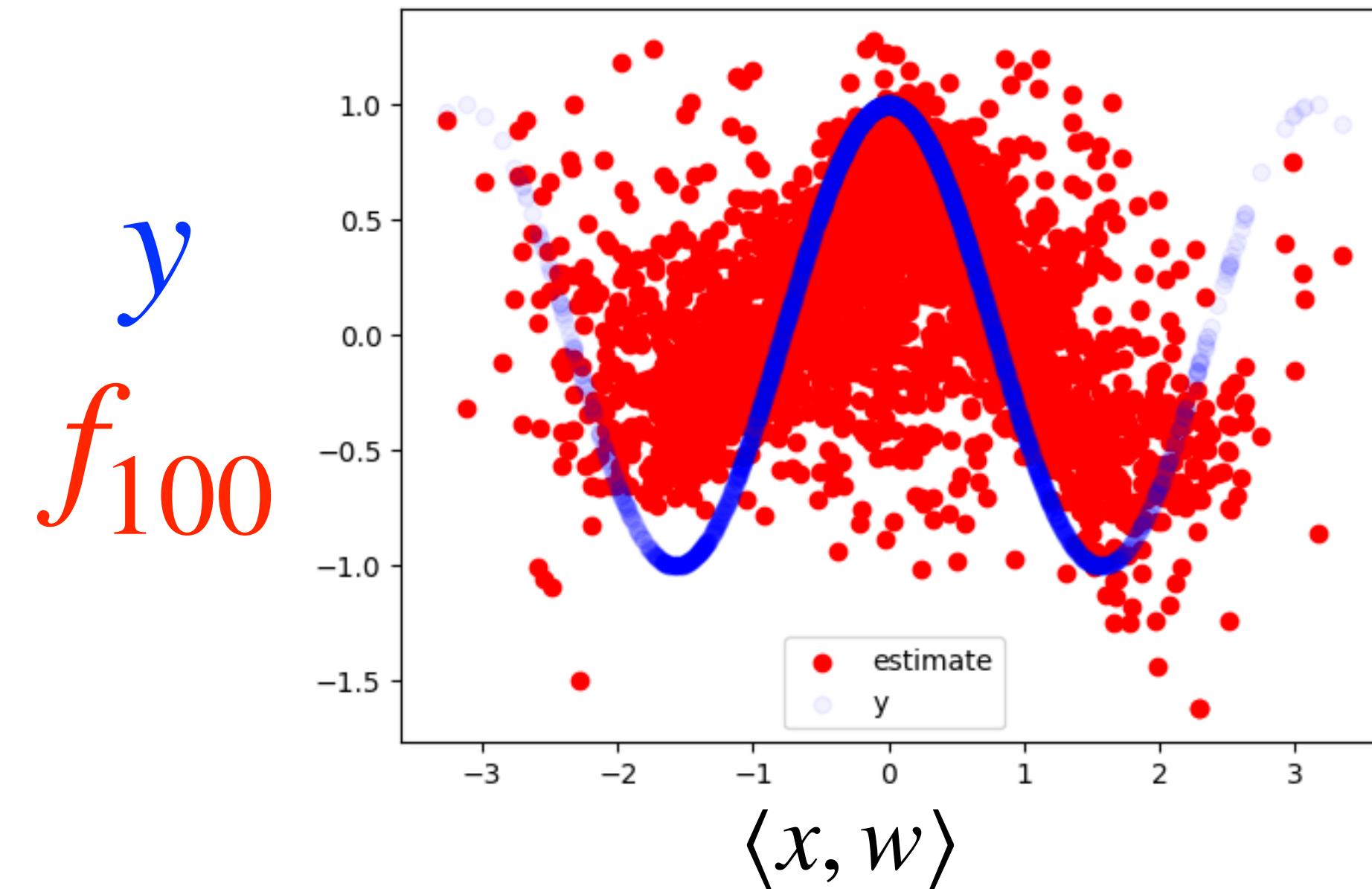
# Experiments

Get ready for hands-on group activity

# 4-dimensional case

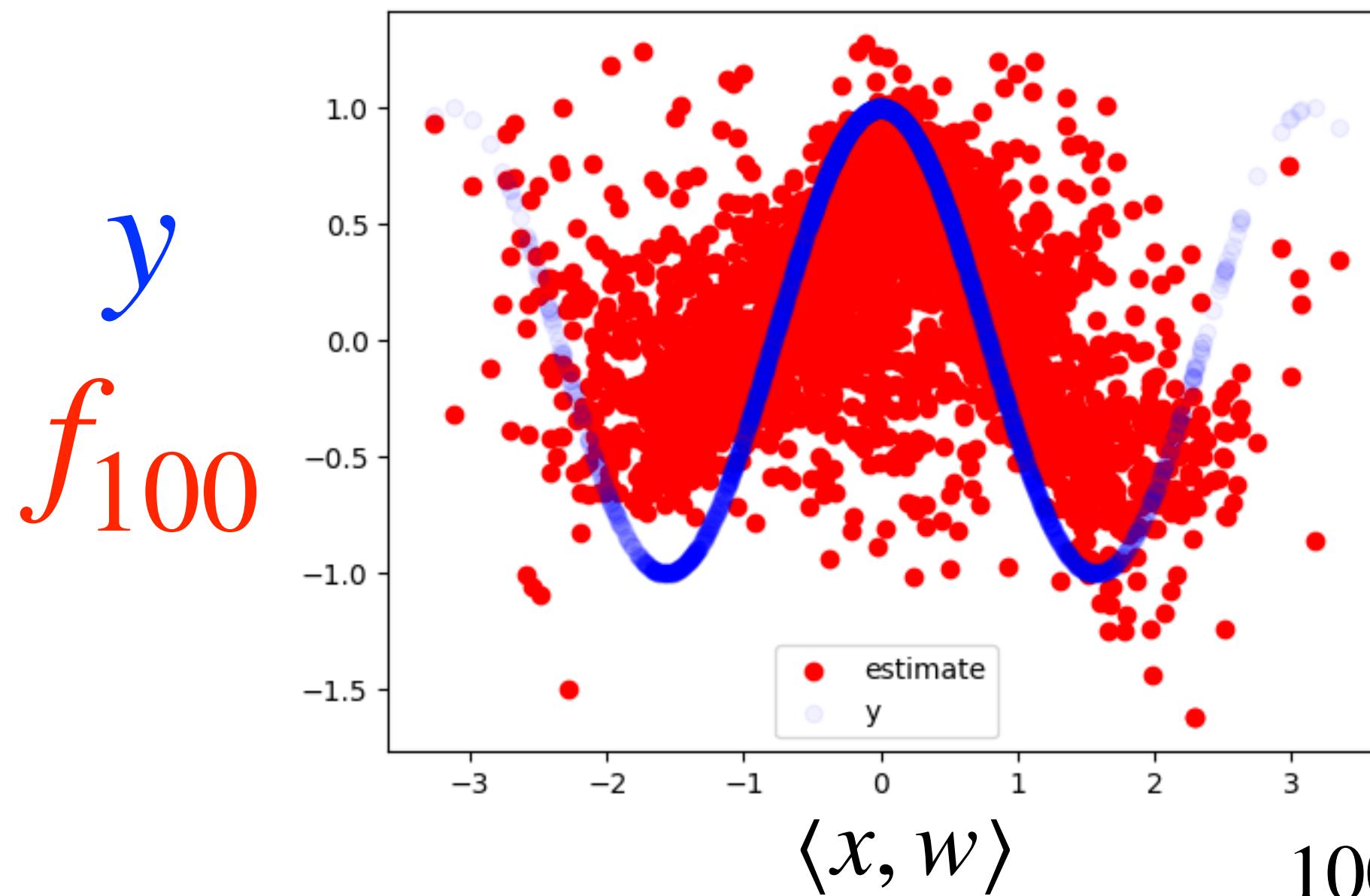
27

- ▶  $y = \cos(\langle w, x \rangle), x, w \in \mathbb{R}^4$
- ▶ Estimate with 100 random feature

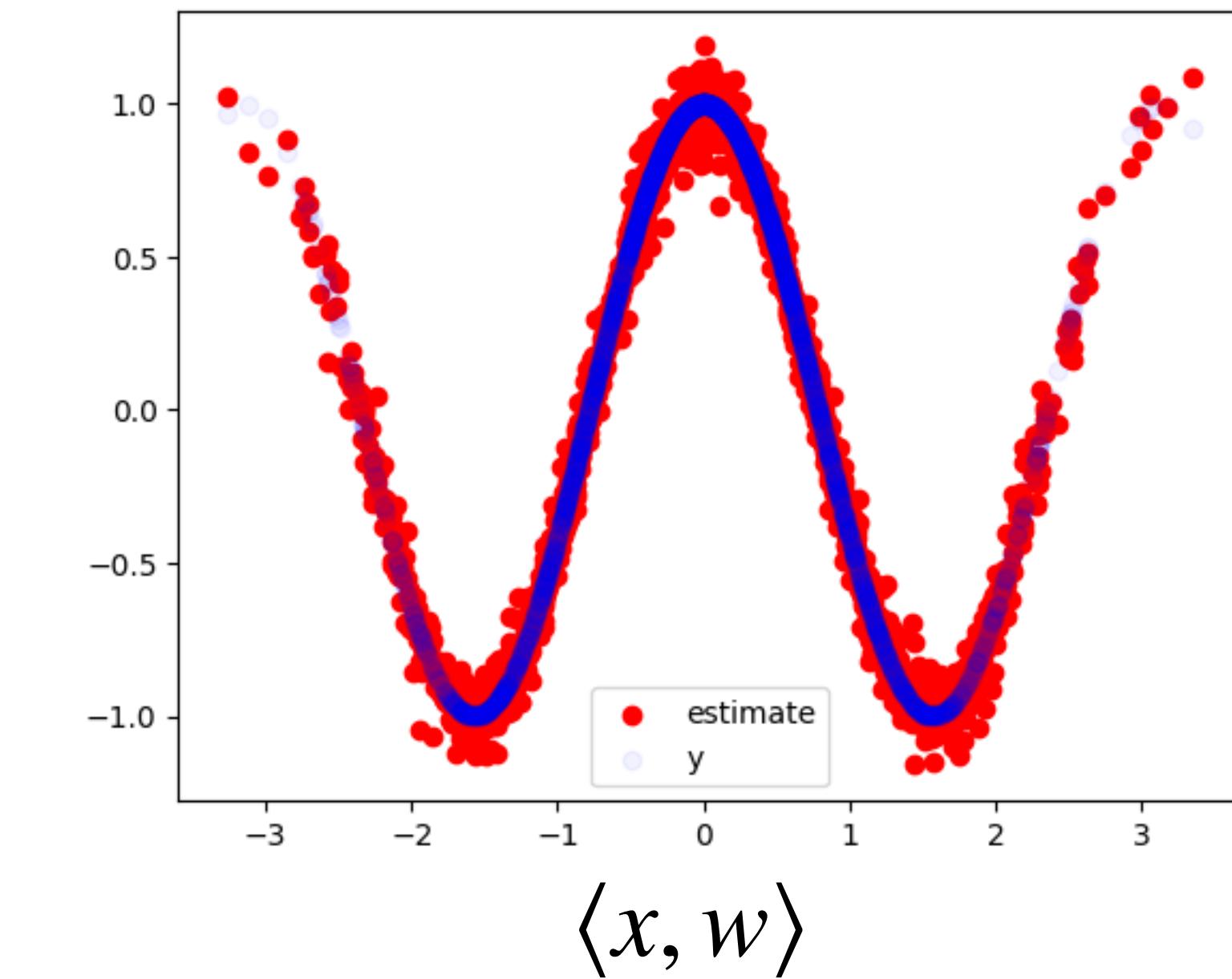


# Task: improve random feature estimate

28



How?



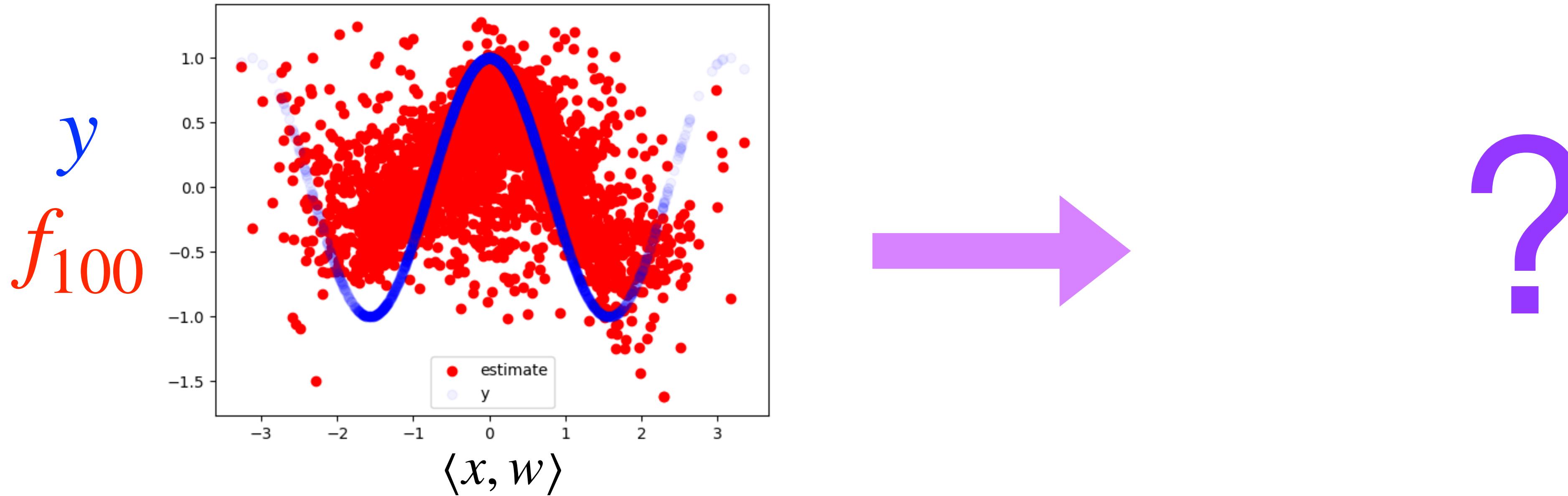
$$y \approx \sum_{i=1}^{100} a_i \cos(\langle v_i, x \rangle + b_i), v_i, x \in \mathbf{R}^d$$

$f_{100}(x)$

# Task: increase the number of features to 1000

29

Colab: <https://shorturl.at/mWTLt>



<https://colab.research.google.com/drive/1SzLEF2EEMhjAXZkx86GuRxIwSkeTKbvs#scrollTo=V5xanwu8fDBt>

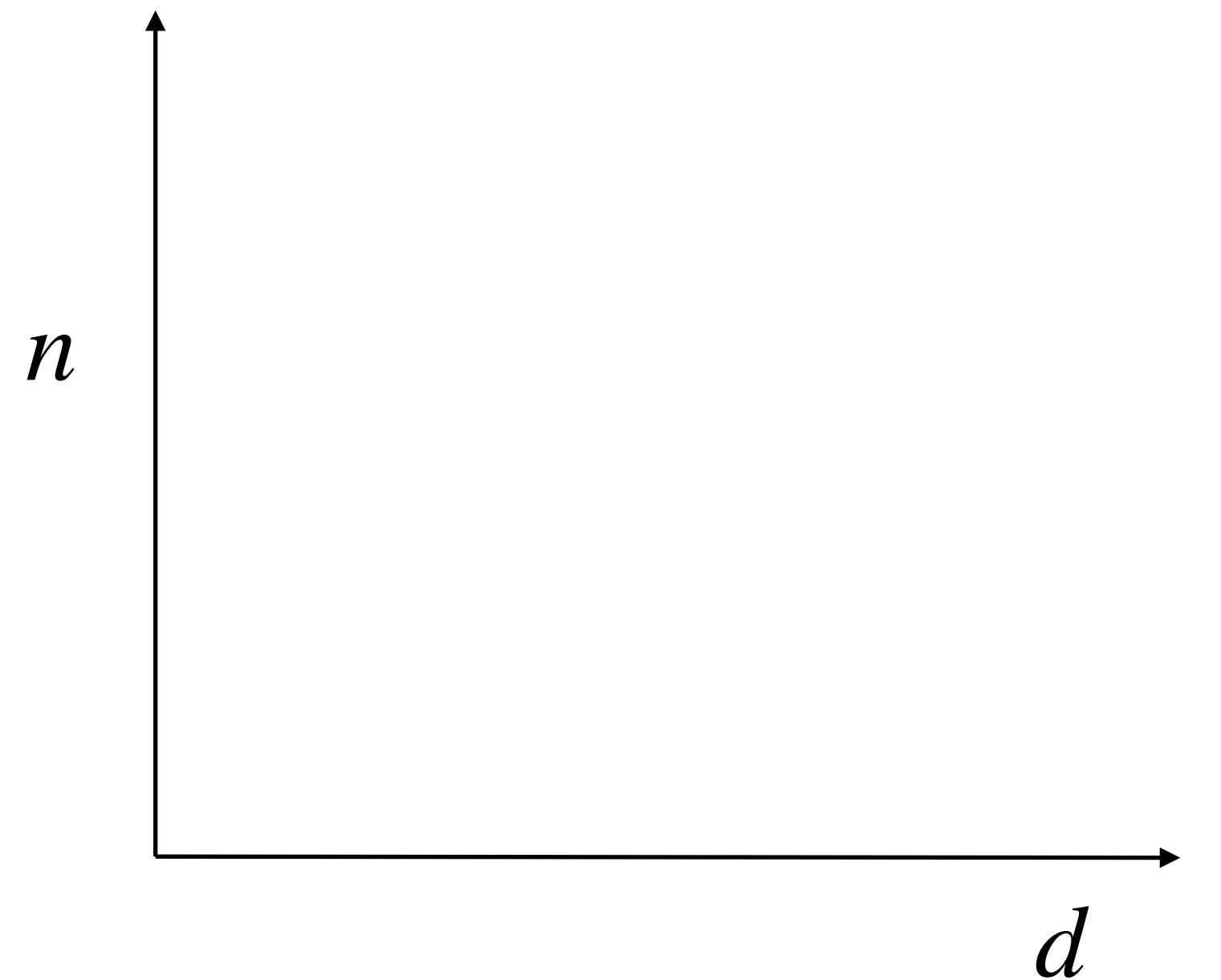
# Task 2: observing the curse of dimensionality

30

$$f(x) \approx \underbrace{\sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i)}_{f_n(x)}, v_i, x \in \mathbf{R}^d$$

- ▶  $\mathbb{E} [(f(x) - f_n(x))^2] \leq 0.0005$
- ▶ How large  $n$  for  $d = 1, 2, 3, 4, 5$

Colab: <https://shorturl.at/mWTLt>

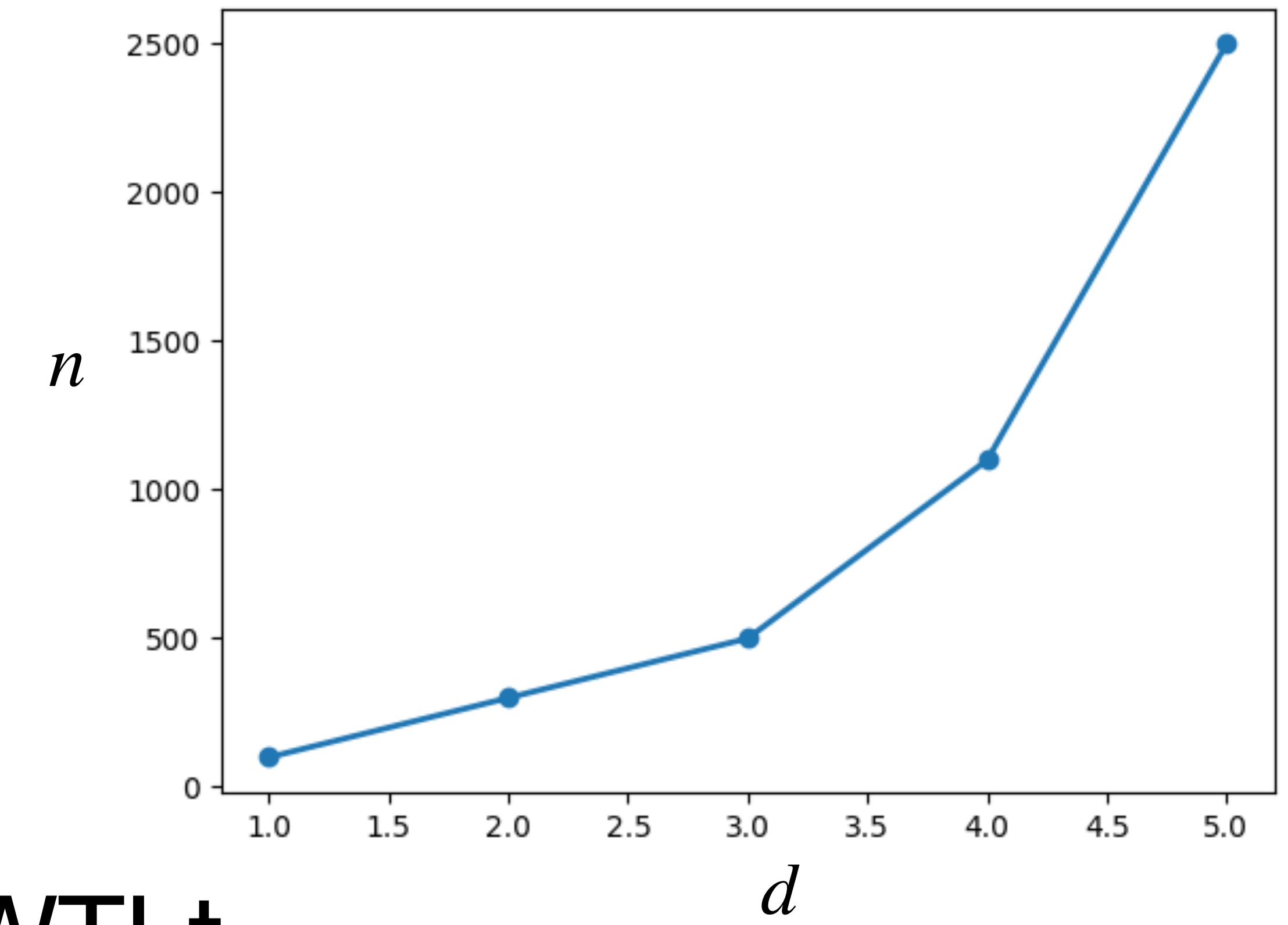


# Solution

31

$$f(x) \approx \underbrace{\sum_i \alpha_i \cos(\langle v_i, x \rangle + b_i)}_{f_n(x)}, v_i, x \in \mathbf{R}^d$$

- $\mathbb{E} [(f(x) - f_n(x))^2] \leq 0.0005$
- How large  $n$  for  $d = 1, 2, 3, 4, 5$



Colab: <https://shorturl.at/mWTLt>

# Theory vs Lab

32

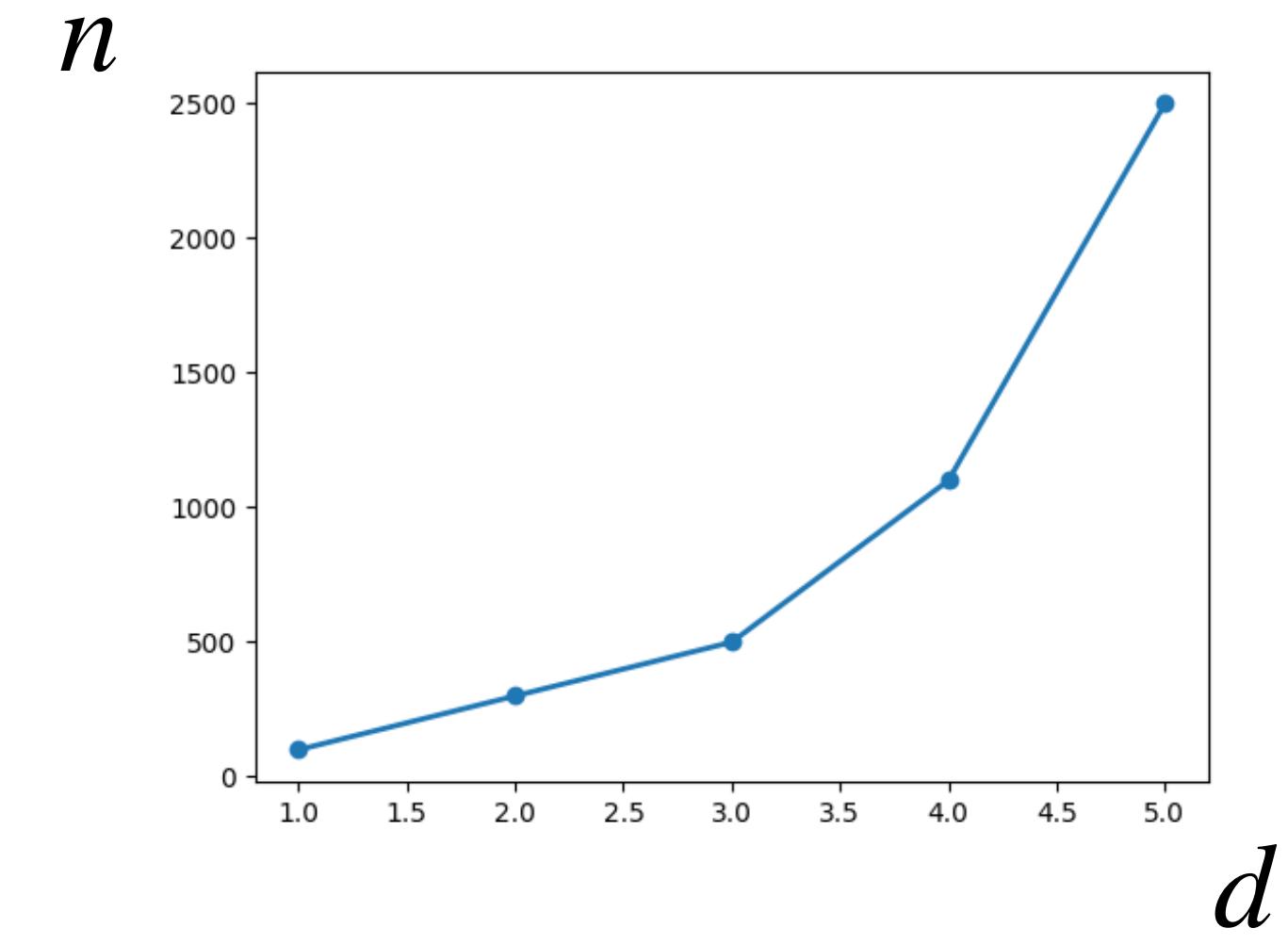
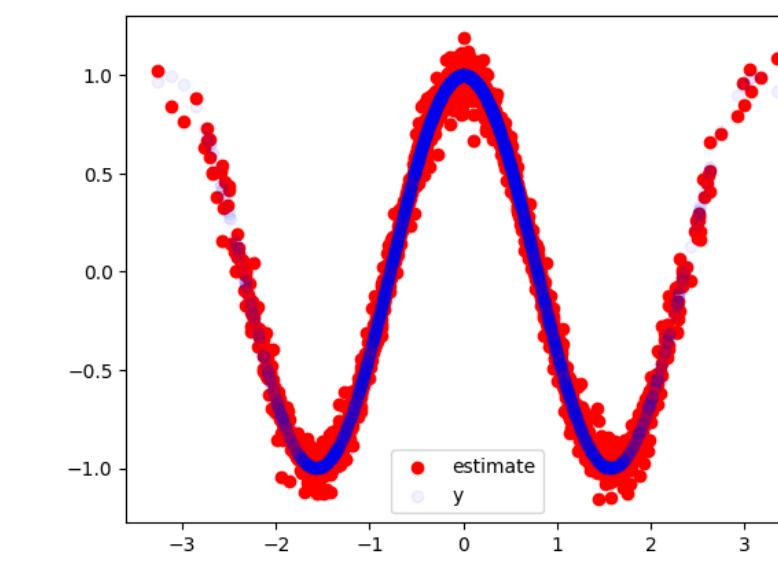
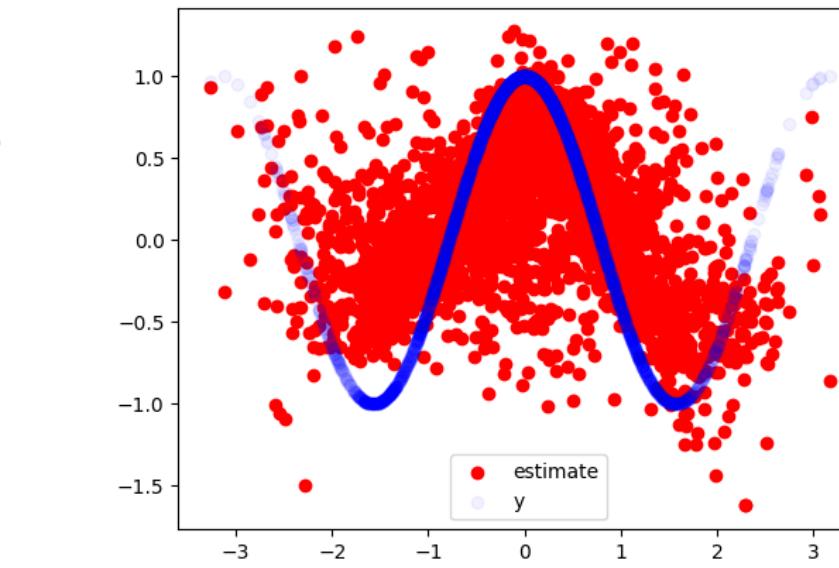
Theory (last lecture)

$$\mathbb{E} \left[ \max_{i \leq n} \langle v_i, w \rangle^2 \right] \approx \frac{\log(n)}{d}$$

- ▶ Order statistics
- ▶ Memoryless property
- ▶ Integral bound for  $\sum_i \frac{1}{i}$

$$y = \cos(2 * \langle w, x \rangle)$$

Lab



# Connecting theory and lab

---

33

- ▶  $y = \cos(\langle w, x \rangle) \approx \sum_i \alpha_i \cos(\langle v_i, x \rangle)$
- ▶ To approximate  $y$ , there must exist a  $v_i$  such that  $\|w - v_i\| \leq \epsilon$  Why?
  - $\cos(\langle v_i, x \rangle) - \cos(\langle w, x \rangle) \leq -2 \sin\left(\frac{\langle v_i + w, x \rangle}{2}\right) \sin\left(\frac{\langle v_i - w, x \rangle}{2}\right)$
  - $|\cos(\langle v_i, x \rangle) - \cos(\langle w, x \rangle)| \leq \|x\| \max\{\|w - v_i\|, \|w + v_i\|\}$
- ▶  $|\langle w, v_i \rangle| \geq 1 - \epsilon$

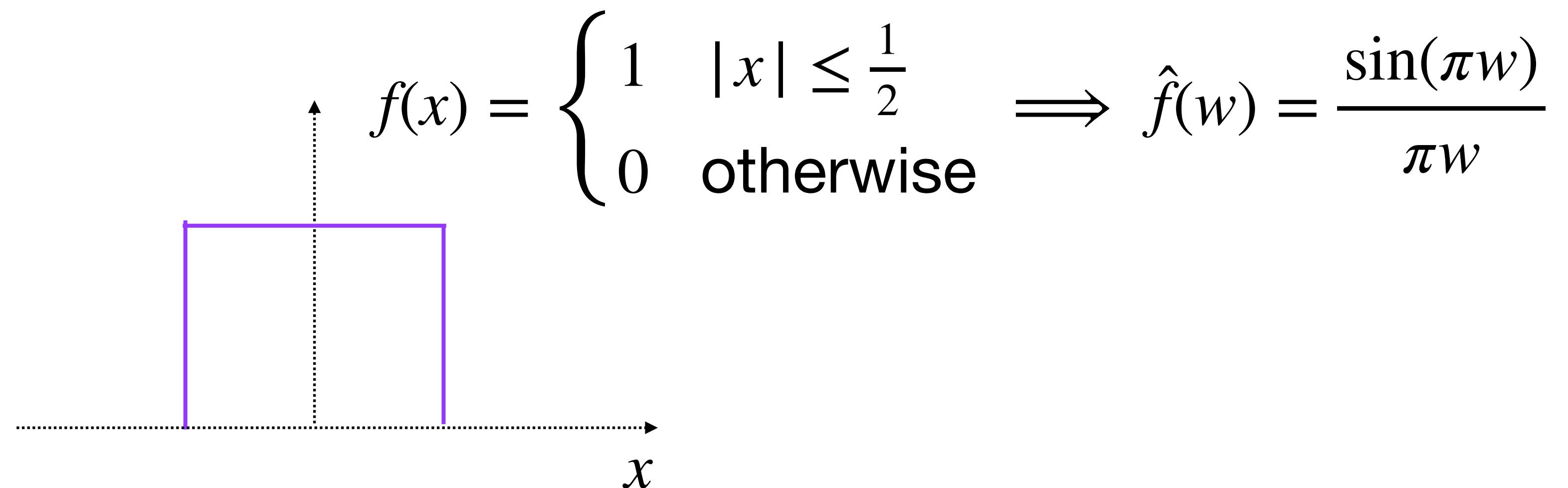
# Barron's function class

A. R. Barron, "Universal approximation bounds for superpositions of a sigmoidal function," in *IEEE Transactions on Information Theory*, 1993

34

**Barrons' function:**  $f(x) = \int e^{i\langle w, x \rangle} \hat{f}(w) dw$  such that  $C_f = \int |w| |\hat{f}(w)| dw < \infty$

- ▶ For Barrons' functions, neural networks break curse of dimensionality
- ▶ **Question:** Is the step function (see bellow) a Barrons' function?



# Thank you very much!

---

35

