

Neural Networks: A Theory Lab

Goal: Introducing the notion of universality

► Intro

Outline: ► Lab

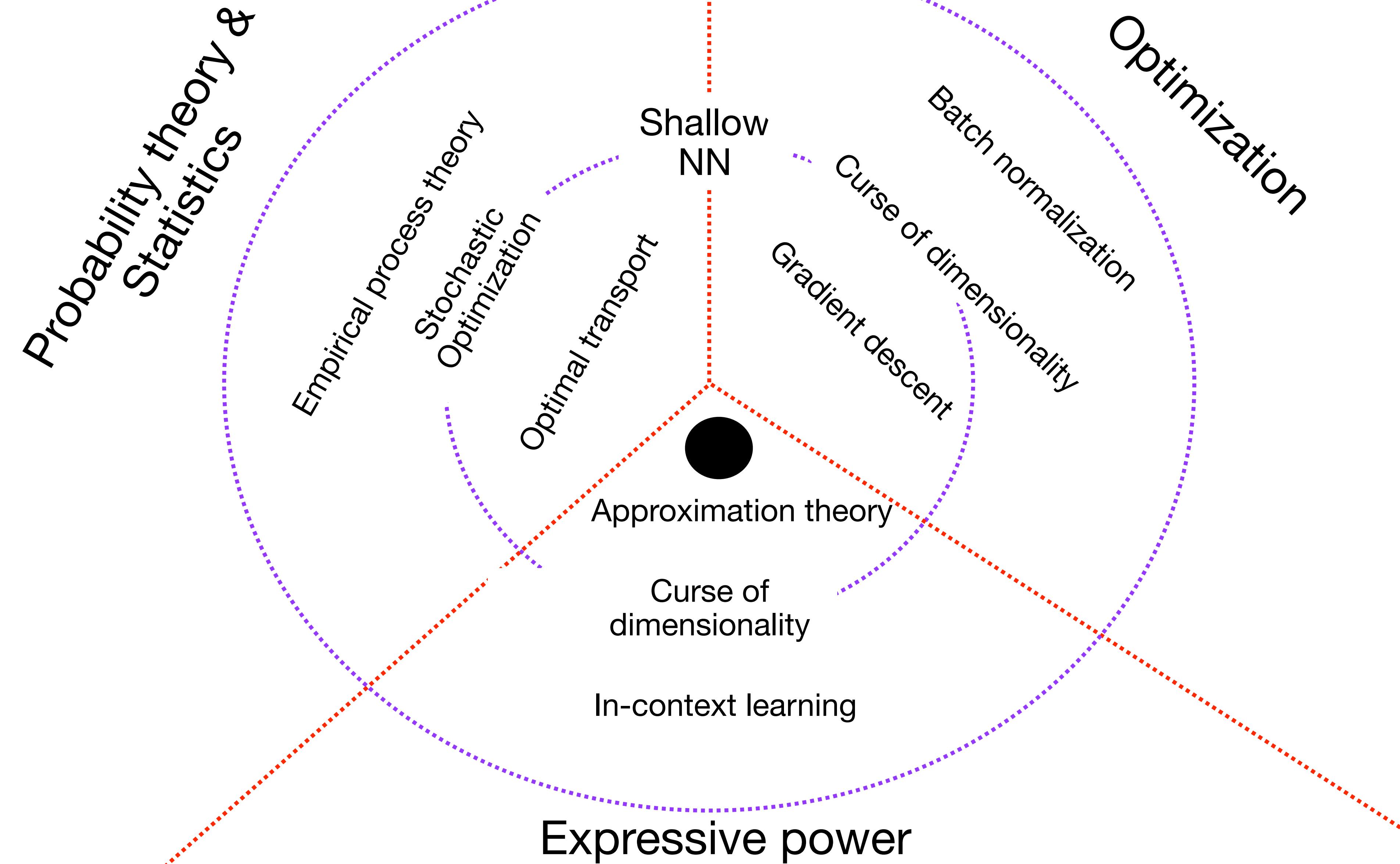
► Theory

Before starting the lecture

2

- ▶ Canvas is published
 - You can submit assignments and project reports in Canvas.
- ▶ I will record lectures from today on and post it on Canvas
- ▶ Please interrupt me for **questions**

Big picture



Why neural networks?

4

- ▶ Neural networks are designed to replace feature engineering
- ▶ What is feature engineering?
- ▶ We explain feature engineering via an example

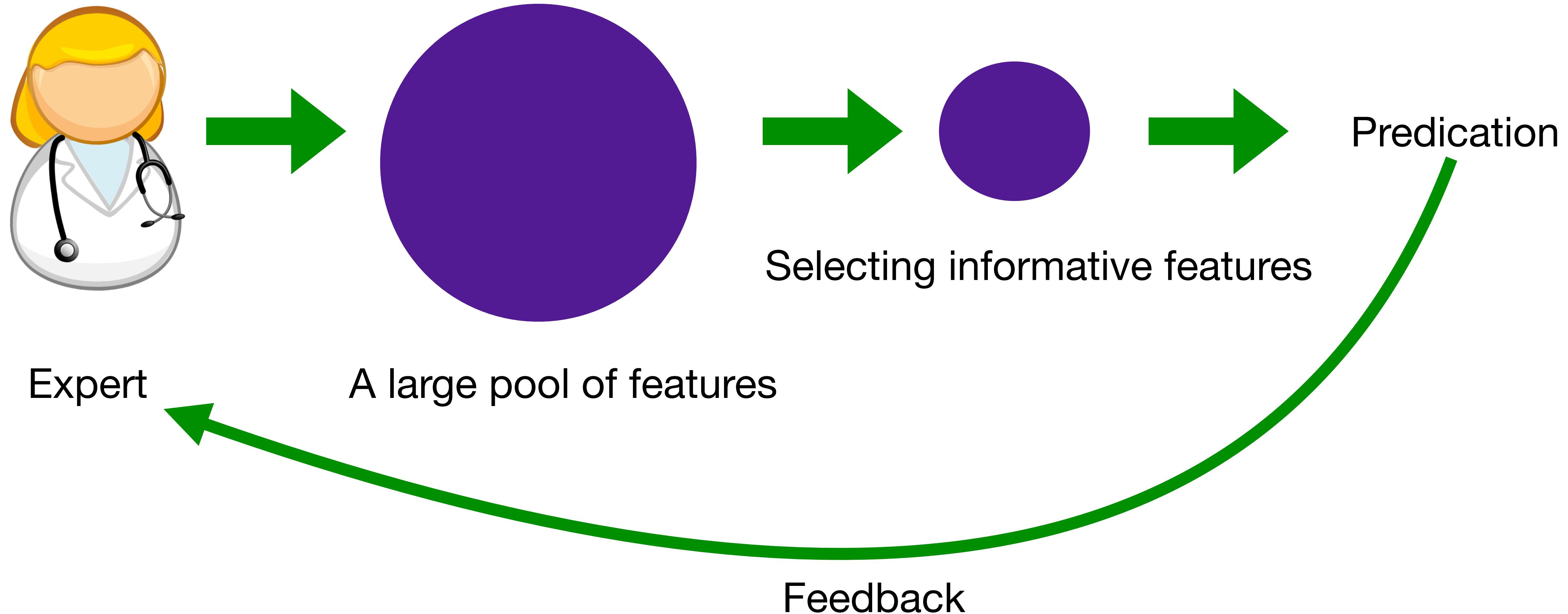
Heart disease predication

5

- ▶ Problem: we want to predict whether a person in future have a heart disease
- ▶ Question: what information about the person will help for the predication?
- ▶ Ideas?
- ▶ Can I use the above information to predict another disease?

Conventional feature engineering

6



Challenges of feature engineering

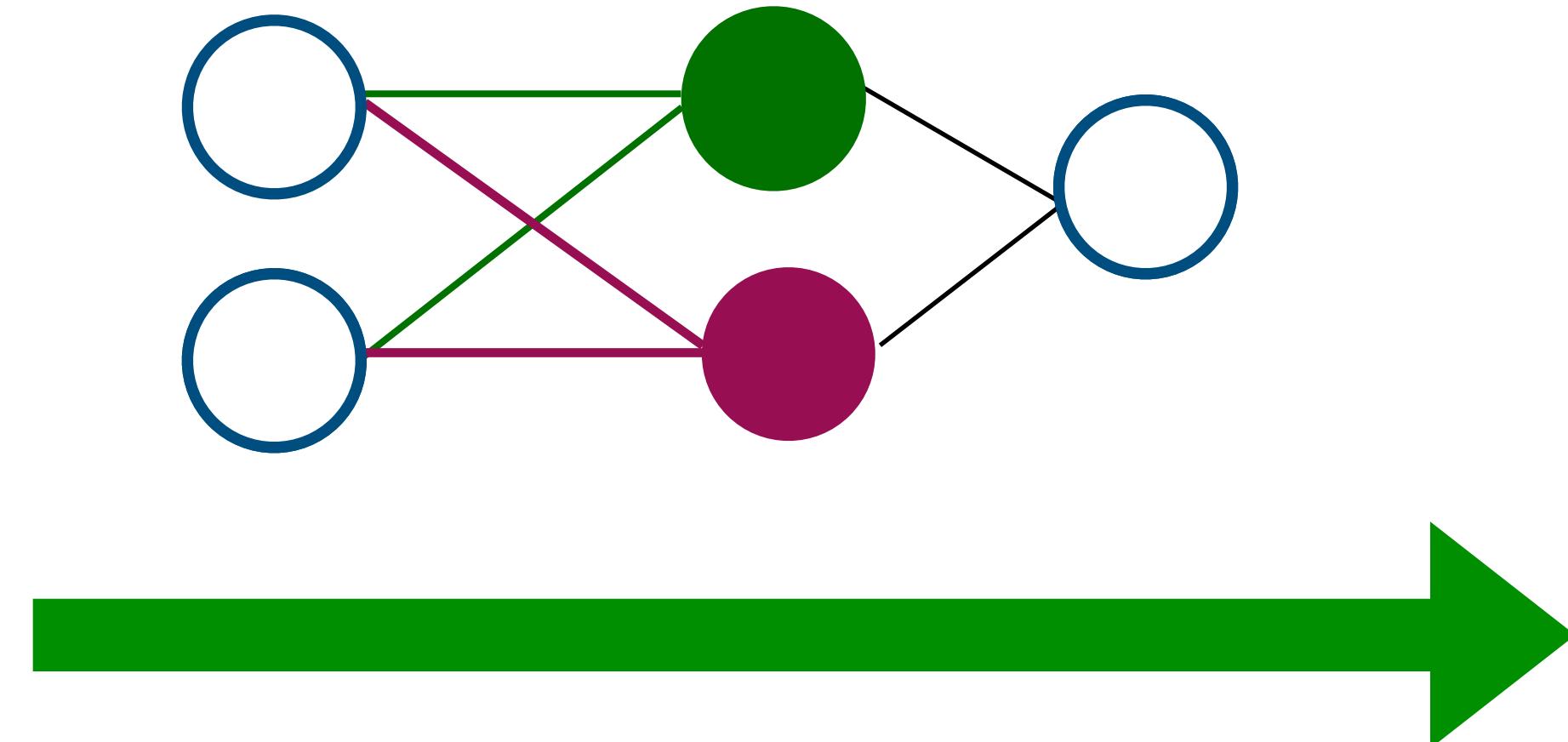
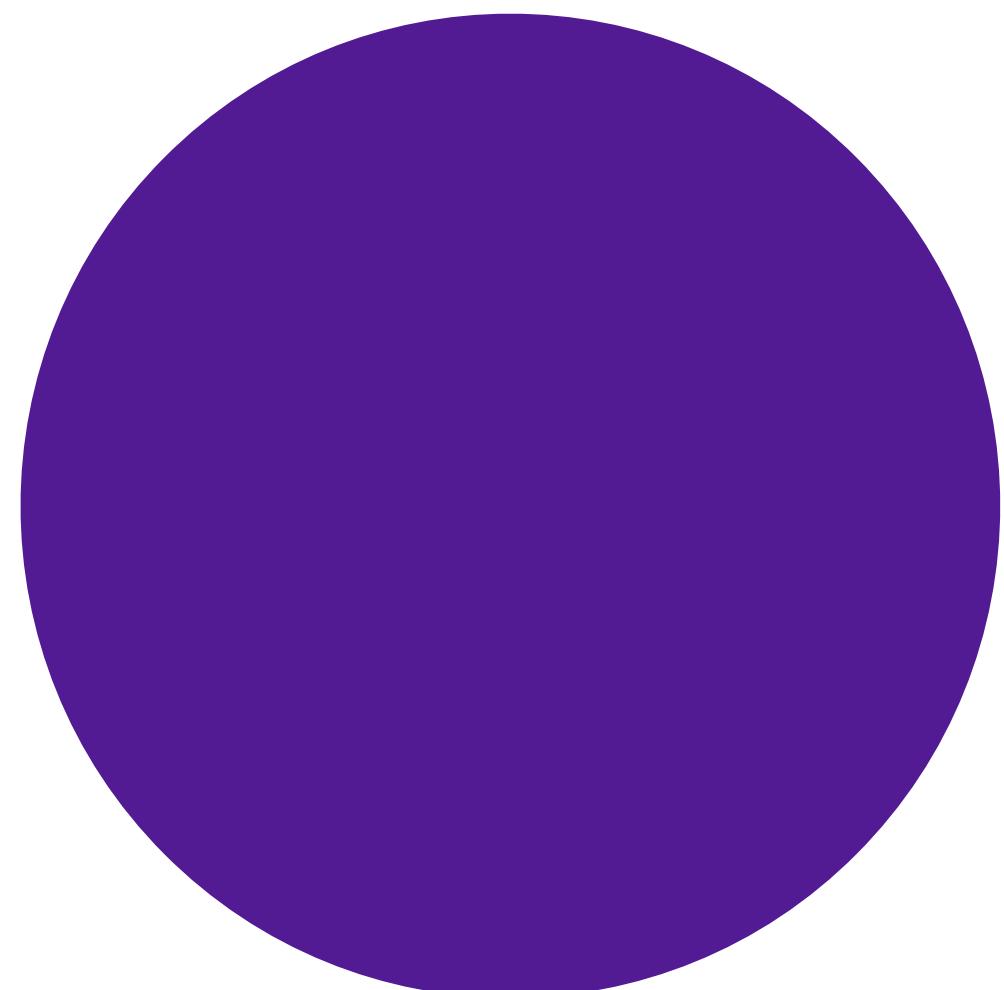
- ▶ Labor-intensive
- ▶ Requires domain knowledge
- ▶ It is task dependent

Deep Learning Goal

8

- ▶ Deep learning automates feature extraction (representation learning)

neural networks or non-parametric methods

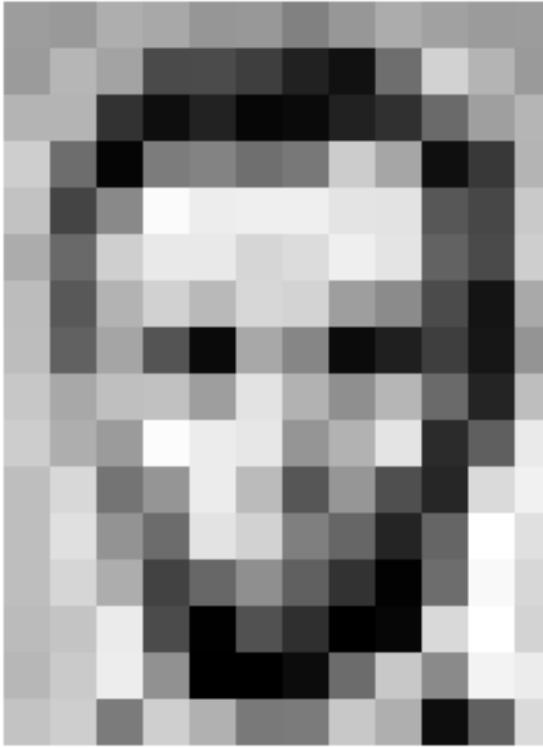


A VERY large pool of features

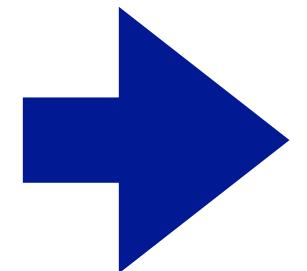
Abstraction of predication

9

- ▶ Image processing example



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	84	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	257	259	259	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	164	134	11	31	62	22	148
199	168	191	193	158	227	17	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	9	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218



A face

Large Pool of Features: χ

Output: $y = f(x)$

<https://ai.stanford.edu/~syyeung/cvweb/tutorial1.html>

Regression problem

10

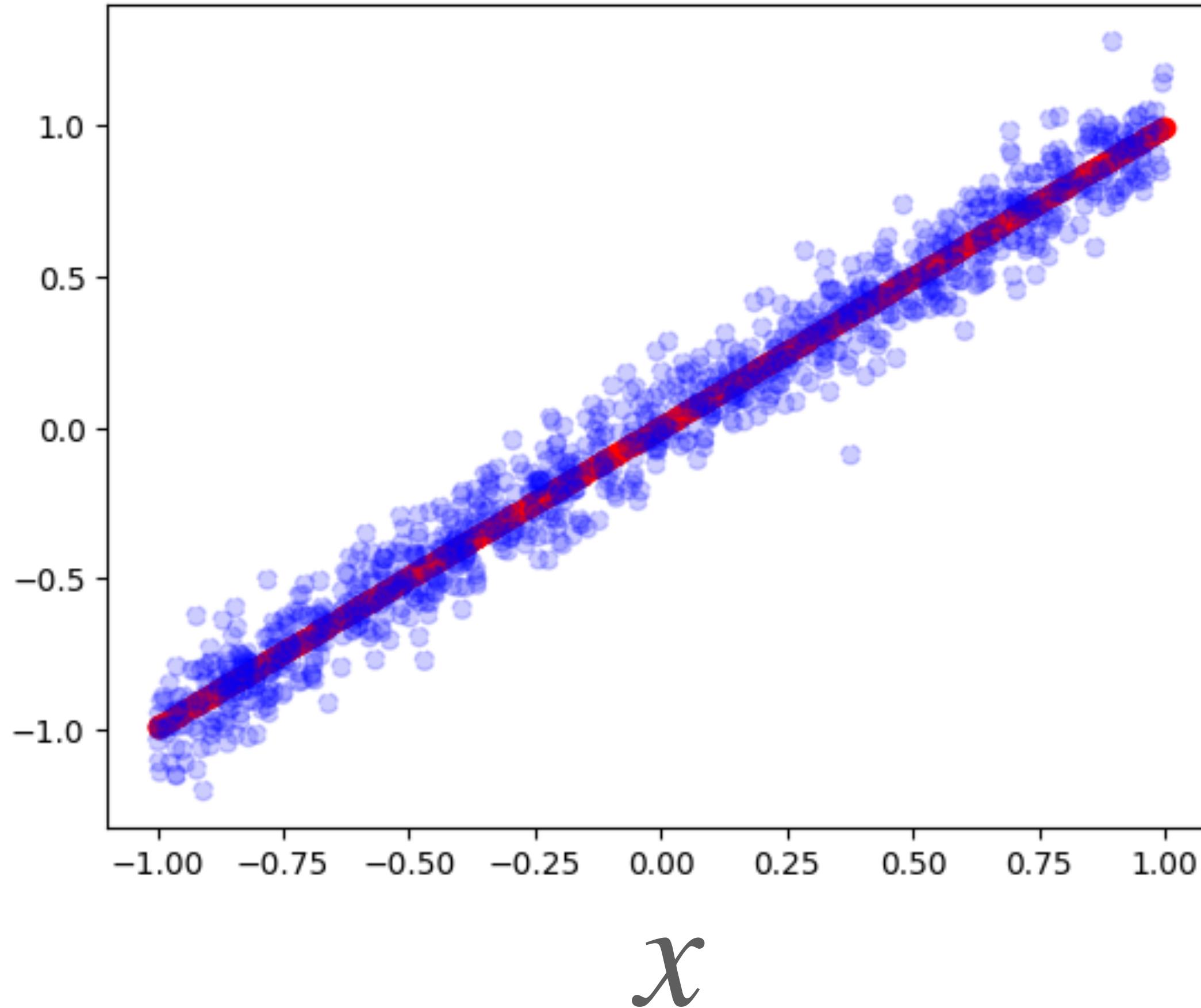
- ▶ Given $(x_1, y_1), \dots, (x_n, y_n) \Rightarrow (x_{n+1}, ?)$
- ▶ Where $y_i = f(x_i)$ where function f is unknown

Recall linear regression

11

- ▶ Recall the example of linear regression that you learned about in Intro to ML course

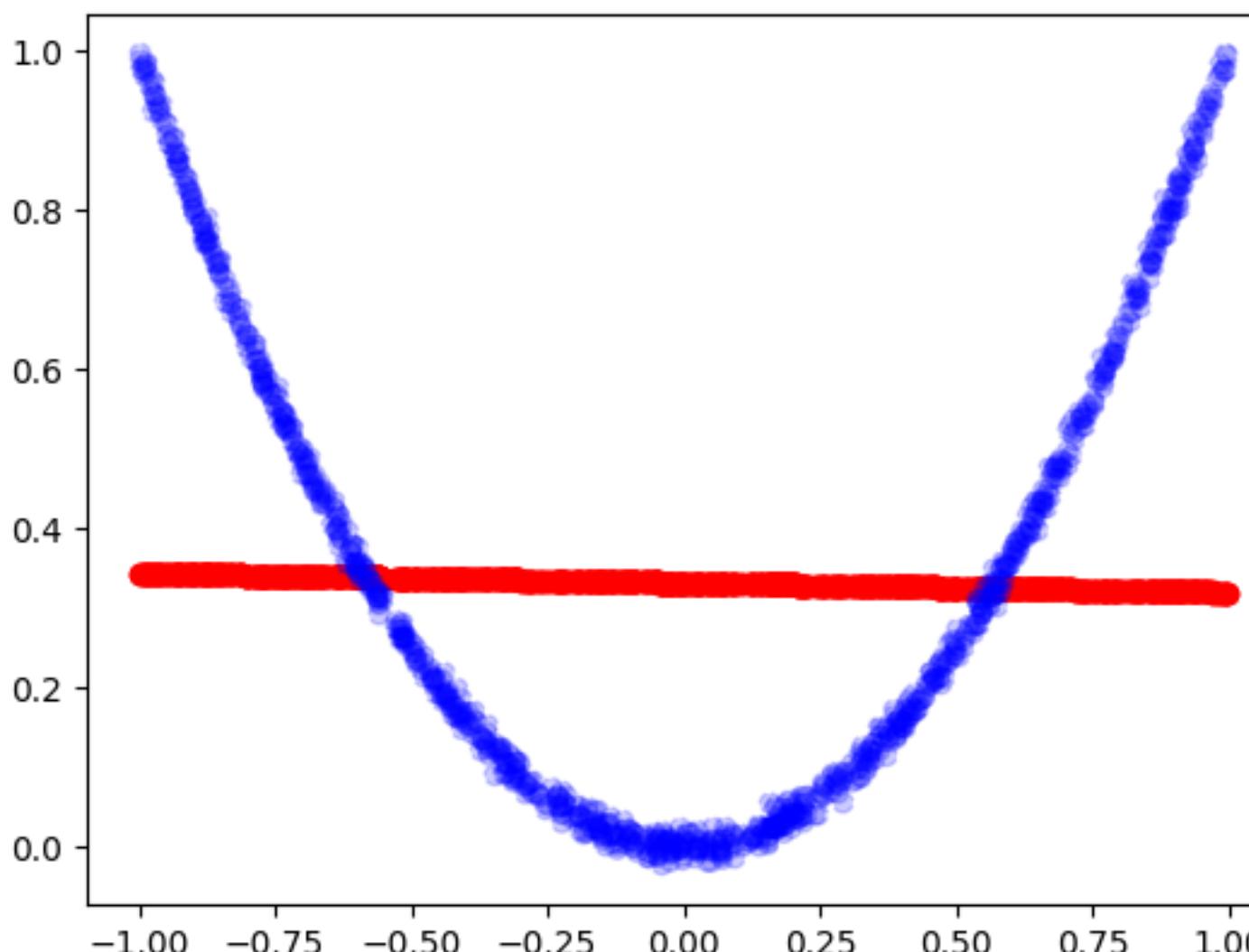
$$y = f(x)$$



Recap: linear regression

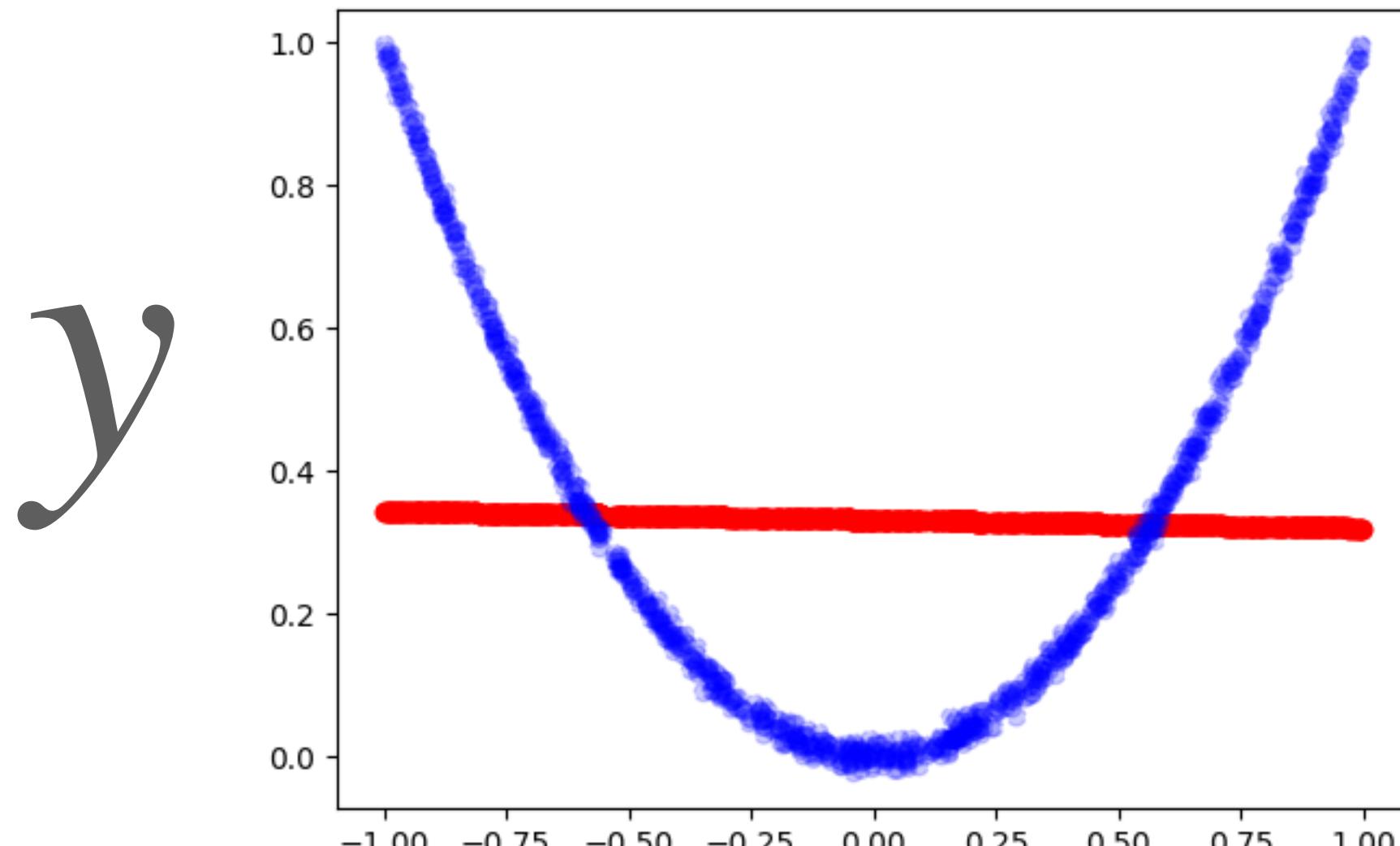
12

- ▶ Given: $x_1, x_2, \dots, x_n \in \mathbb{R}^d, y_1, y_2, \dots, y_n \in \mathbb{R}$
- ▶ Find $w \in \mathbb{R}^d$ such that $\sum_{i=1}^n (x_i^\top w - y_i)^2$ is minimized.
- ▶ Challenge of linear regression: approximating non-linear functions



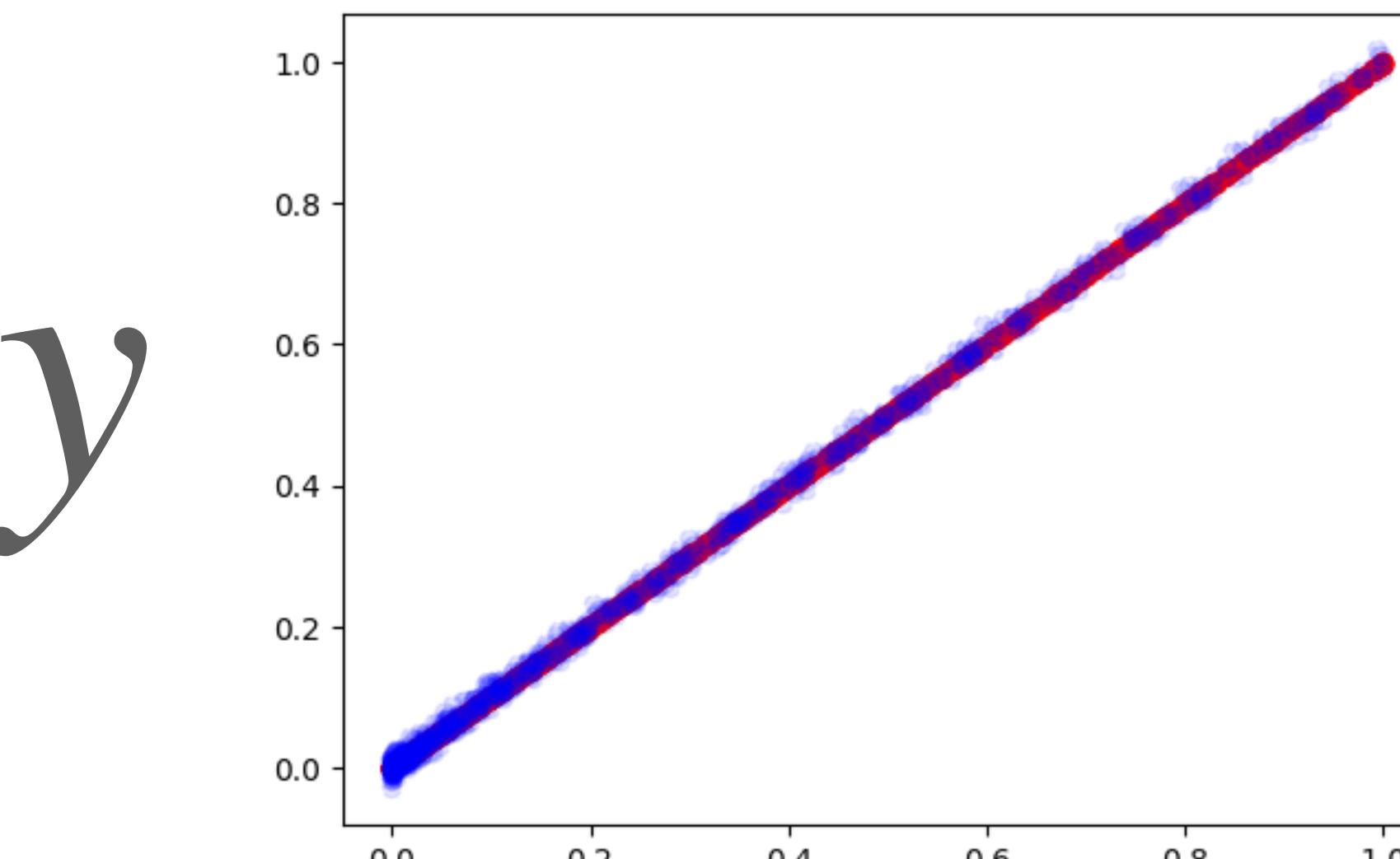
Question: How to solve the non-linearity challenge

13



y

x



y

x^2

$$\phi(x) = x^2$$



A conventional method in statistics

14

- ▶ You won't believe: It is the standard approach in **statistics** to ~~manually~~ design a non-linear transformation of x

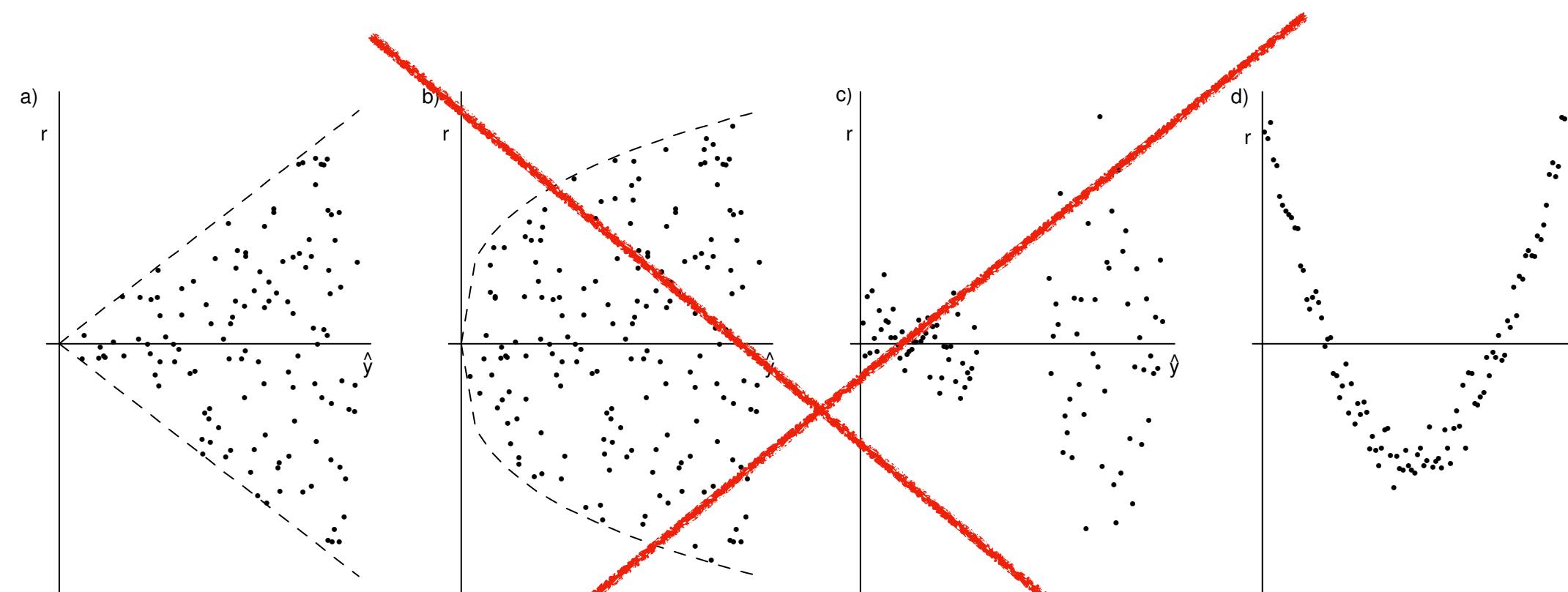


Figure 1.5: a) linear increase of standard deviation, b) nonlinear increase of standard deviation, c) 2 groups with different variances, d) missing quadratic term in the model.

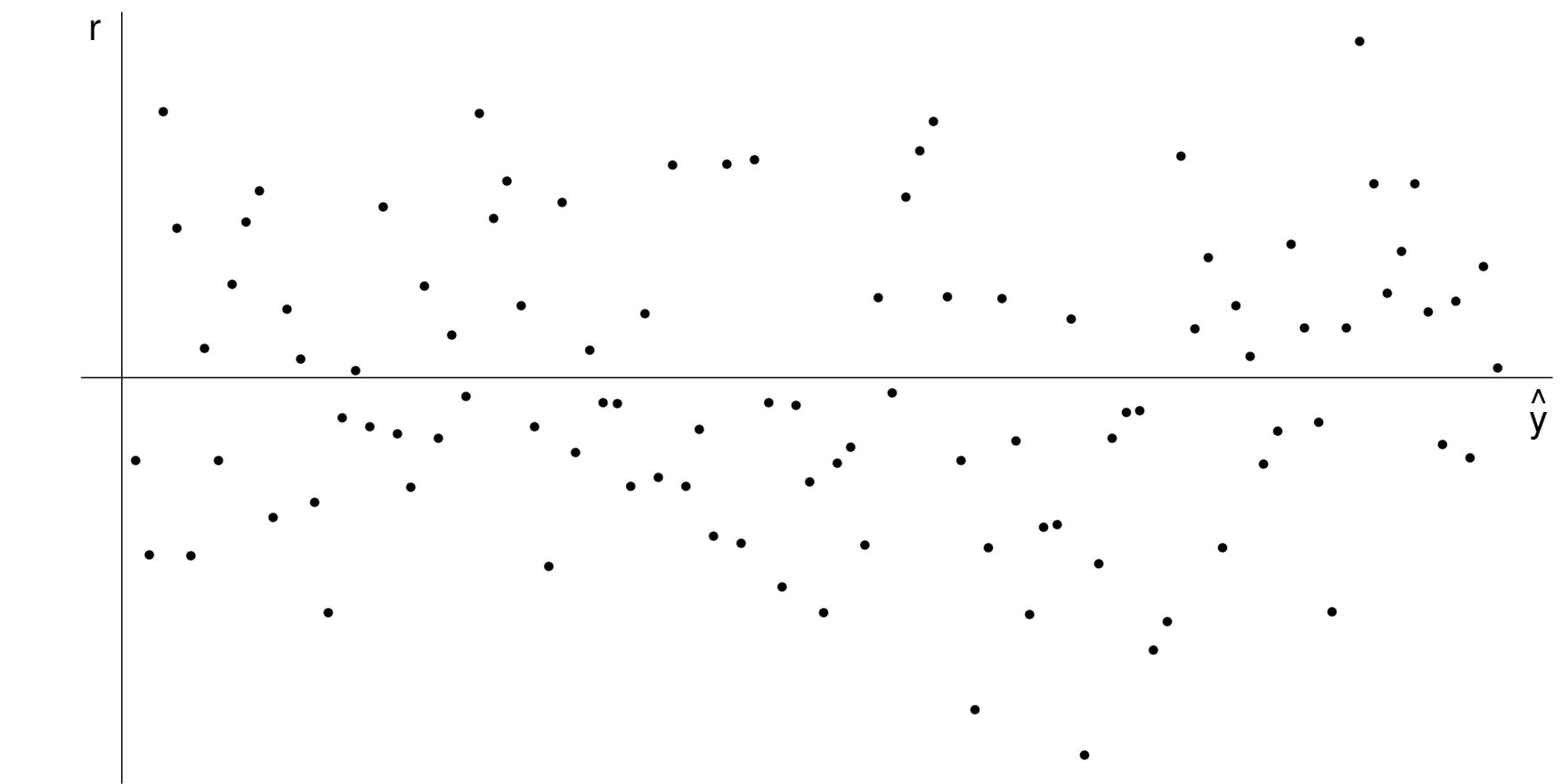
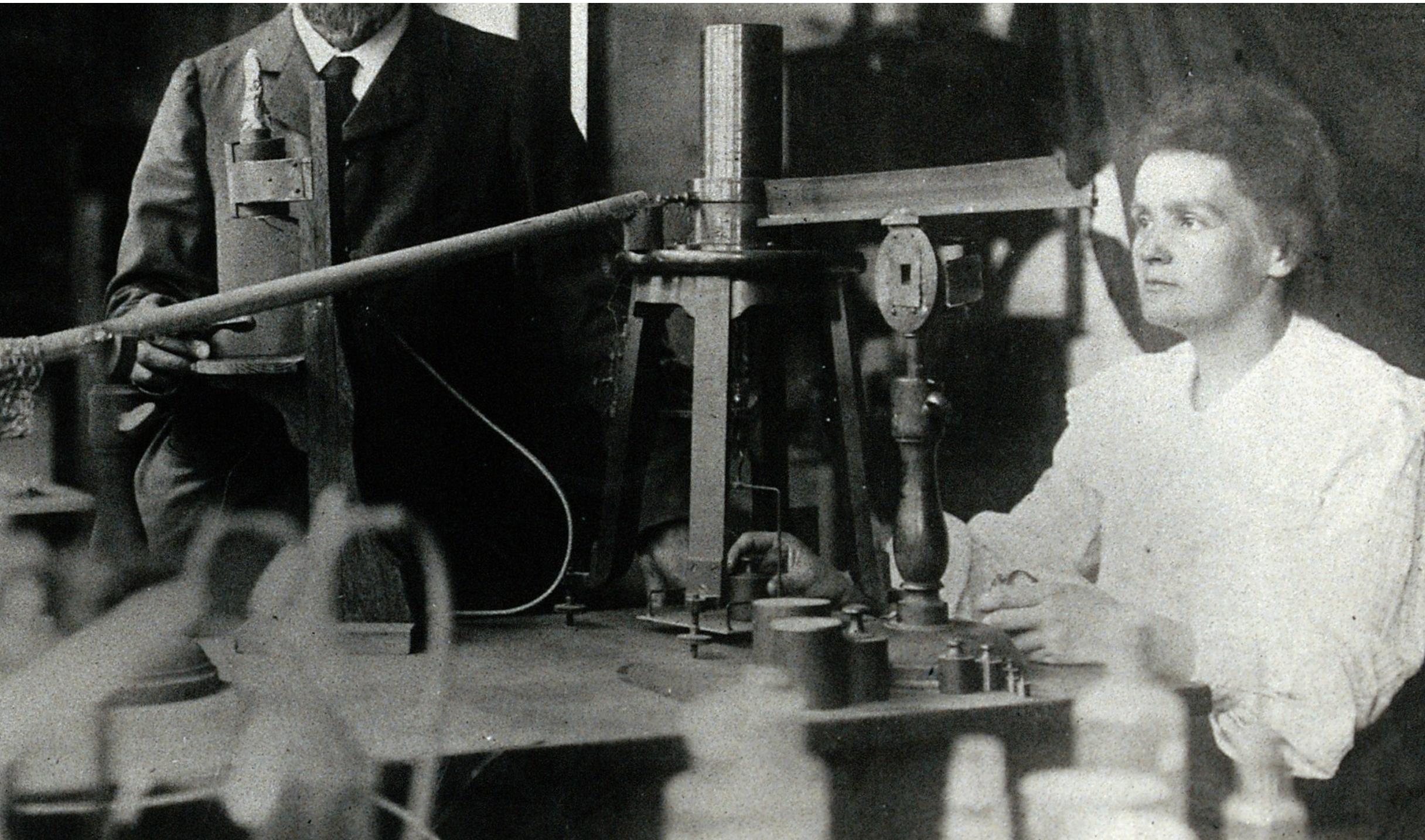


Figure 1.4: Ideal Tukey-Anscombe plot: no violations of model assumptions.

- ▶ Taken from my lecture notes in computational statistics course

- ▶ A brief intro
- ▶ Lab
- ▶ Theory



wikipedia: Marie Curie

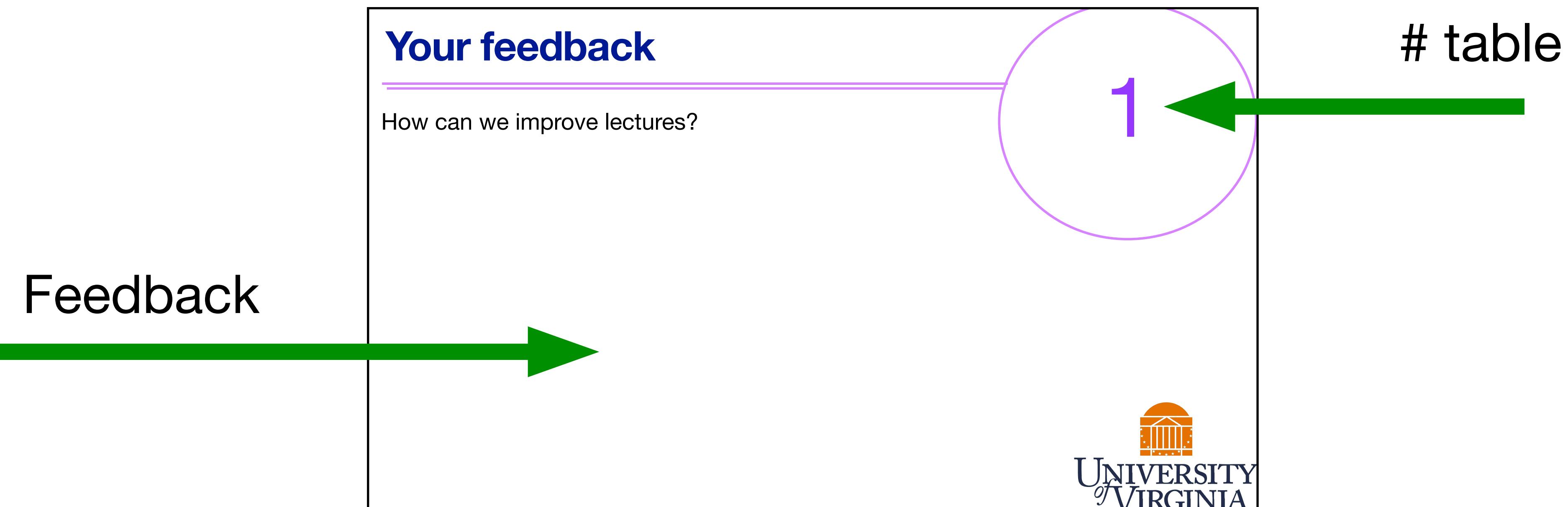
Lab

Get ready for hands-on group activity

Group assignment

16

- ▶ Please join the table with the number on the paper



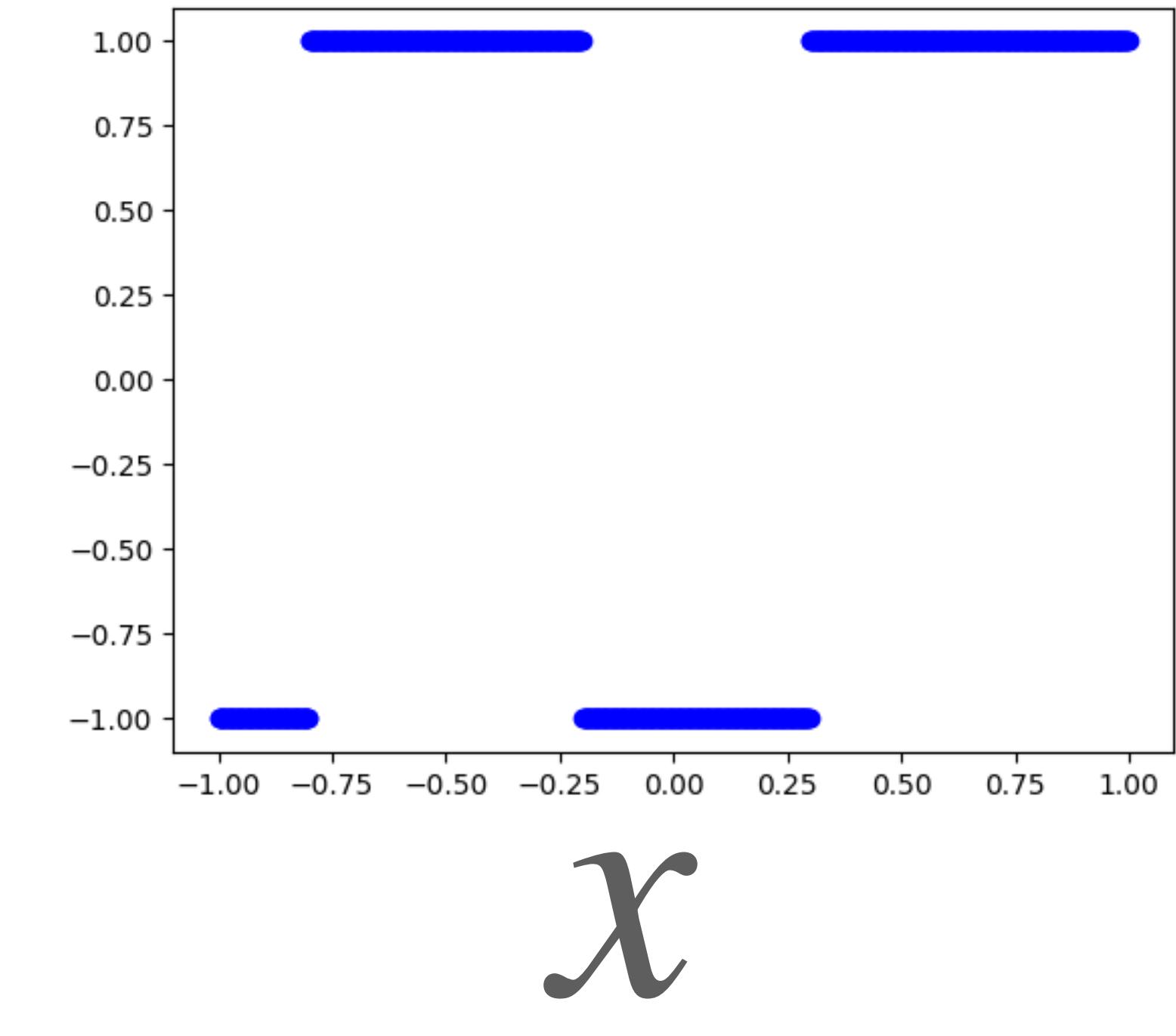
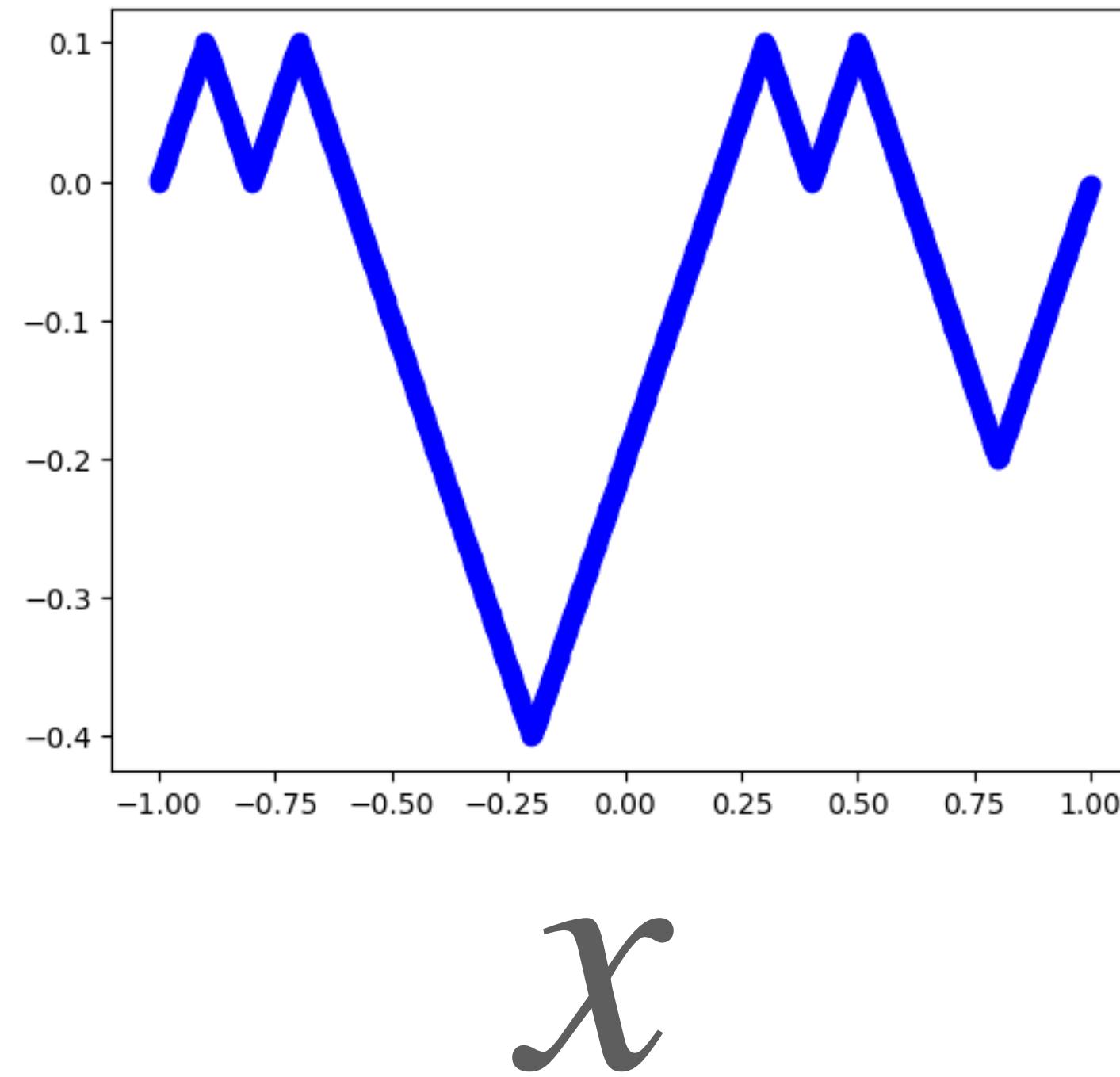
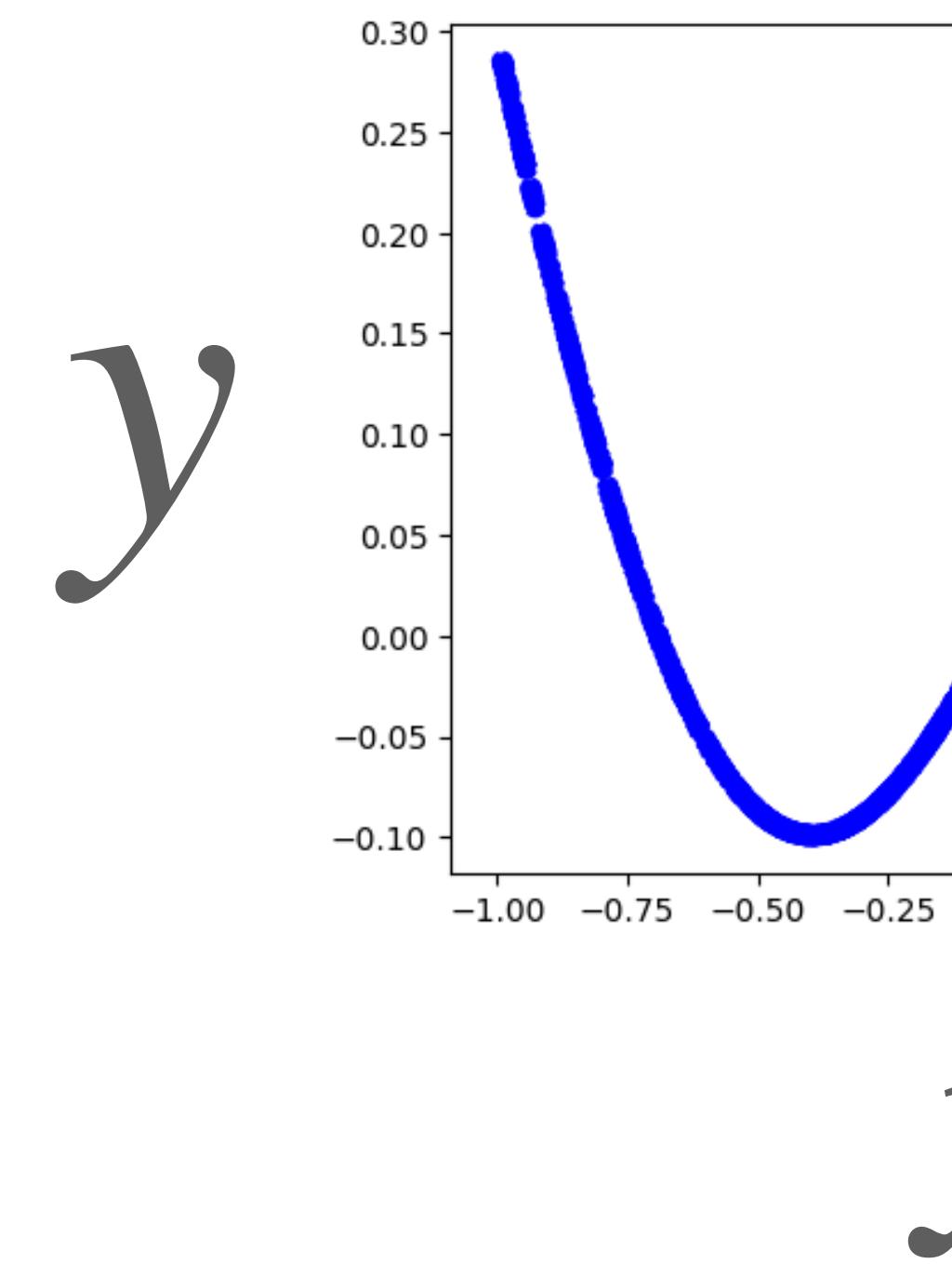
Work in groups

17

- ▶ Goal: We want to design **universal** non-linear features
- ▶ **Universal** features can approximate various y 's



image: Flaticon.com¹



Your task (8 mins)



- ▶ Given x , design features $\phi_1(x), \dots, \phi_n(x)$ such that $y \approx \sum_{i=1}^n w_i \phi_i(x)$
- ▶ $\phi_i : \mathbb{R} \rightarrow \mathbb{R}, x \in \mathbb{R}$



Scan to Start

<https://shorturl.at/ABb6H>

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

- ▶ You can search or use ChatGPT

Ideas?

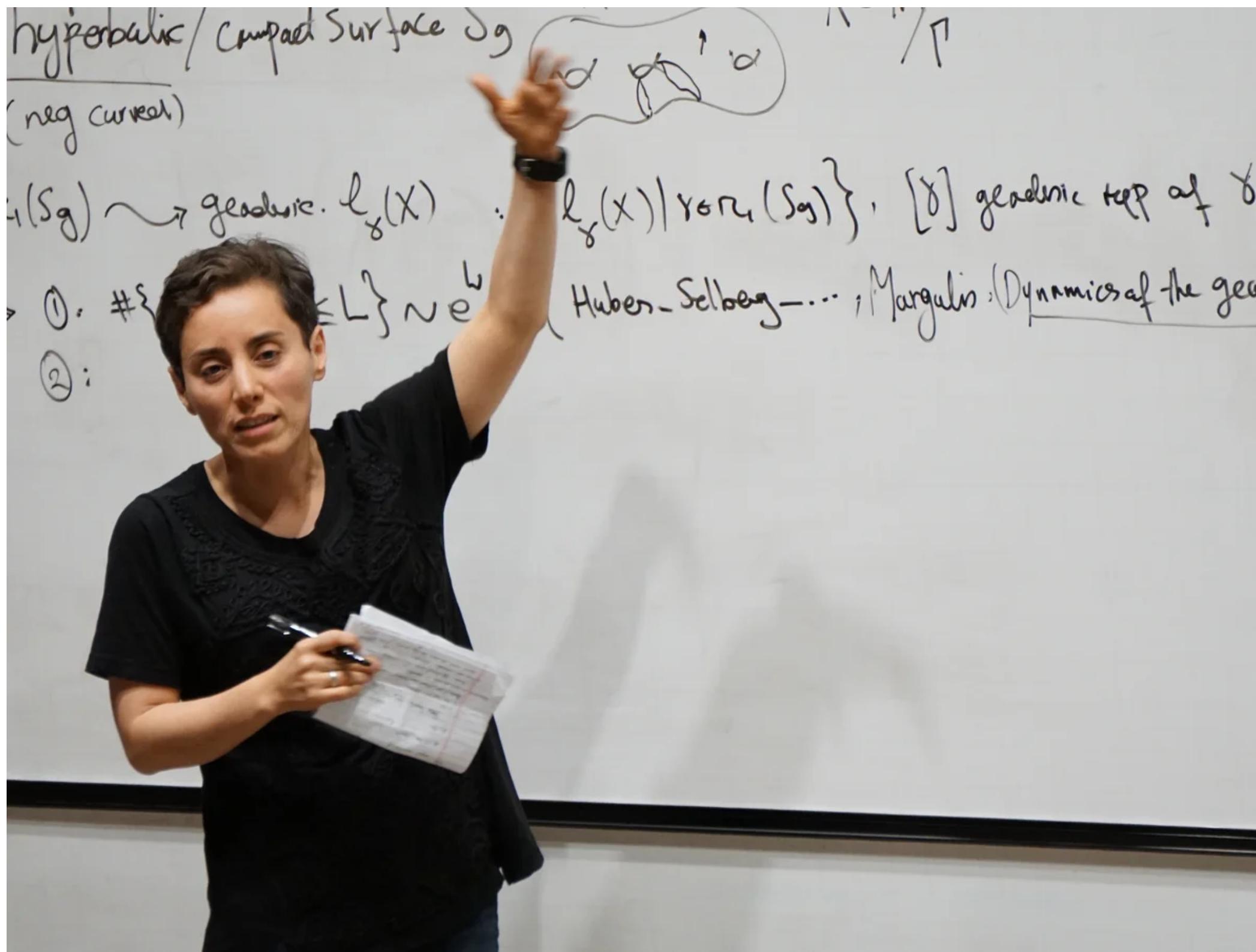
19



► Intro

► Lab

► Theory



theguardian: Maryam Mirzakhani

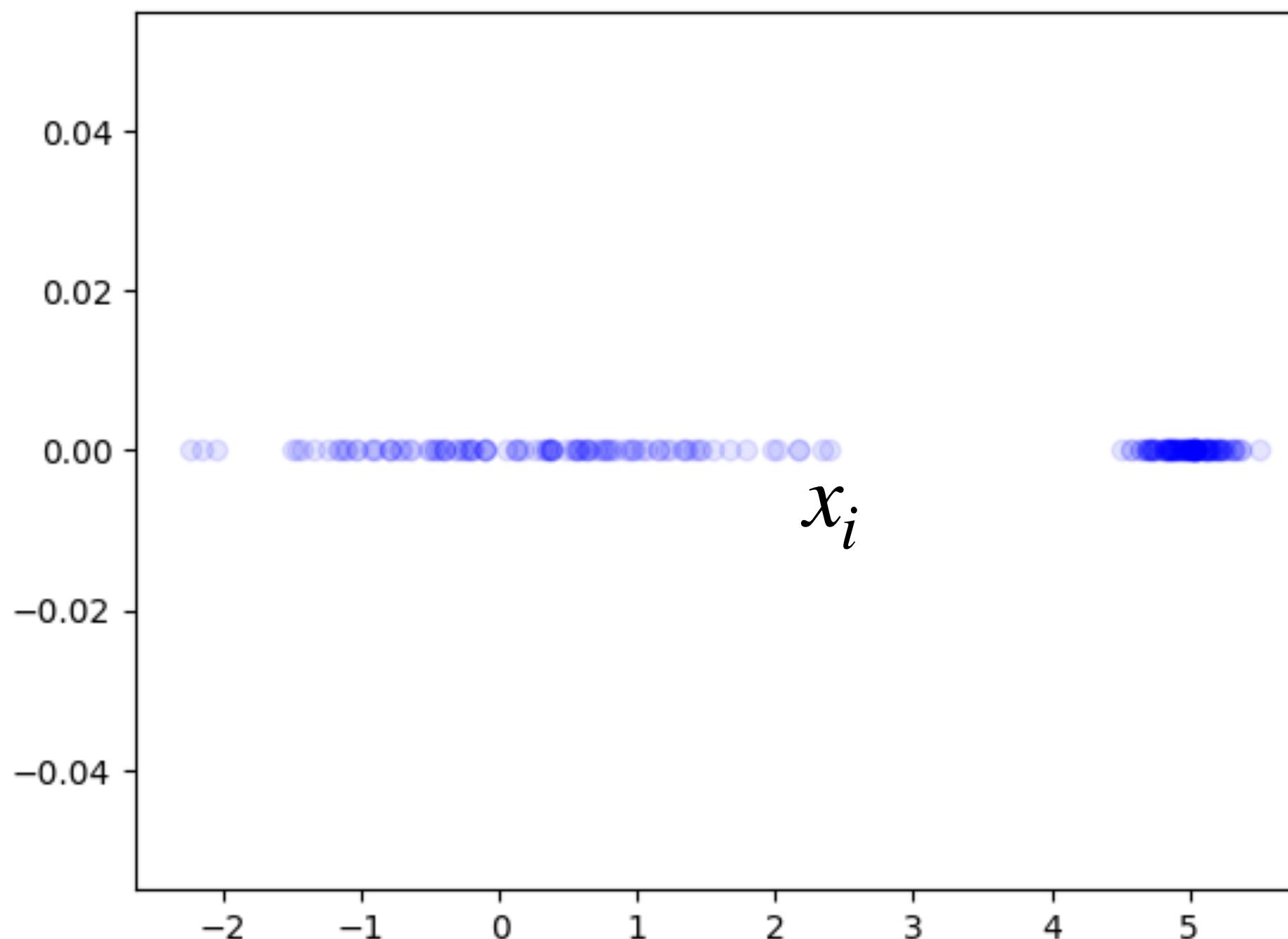
Theory

Get ready for math

Density estimation problem

21

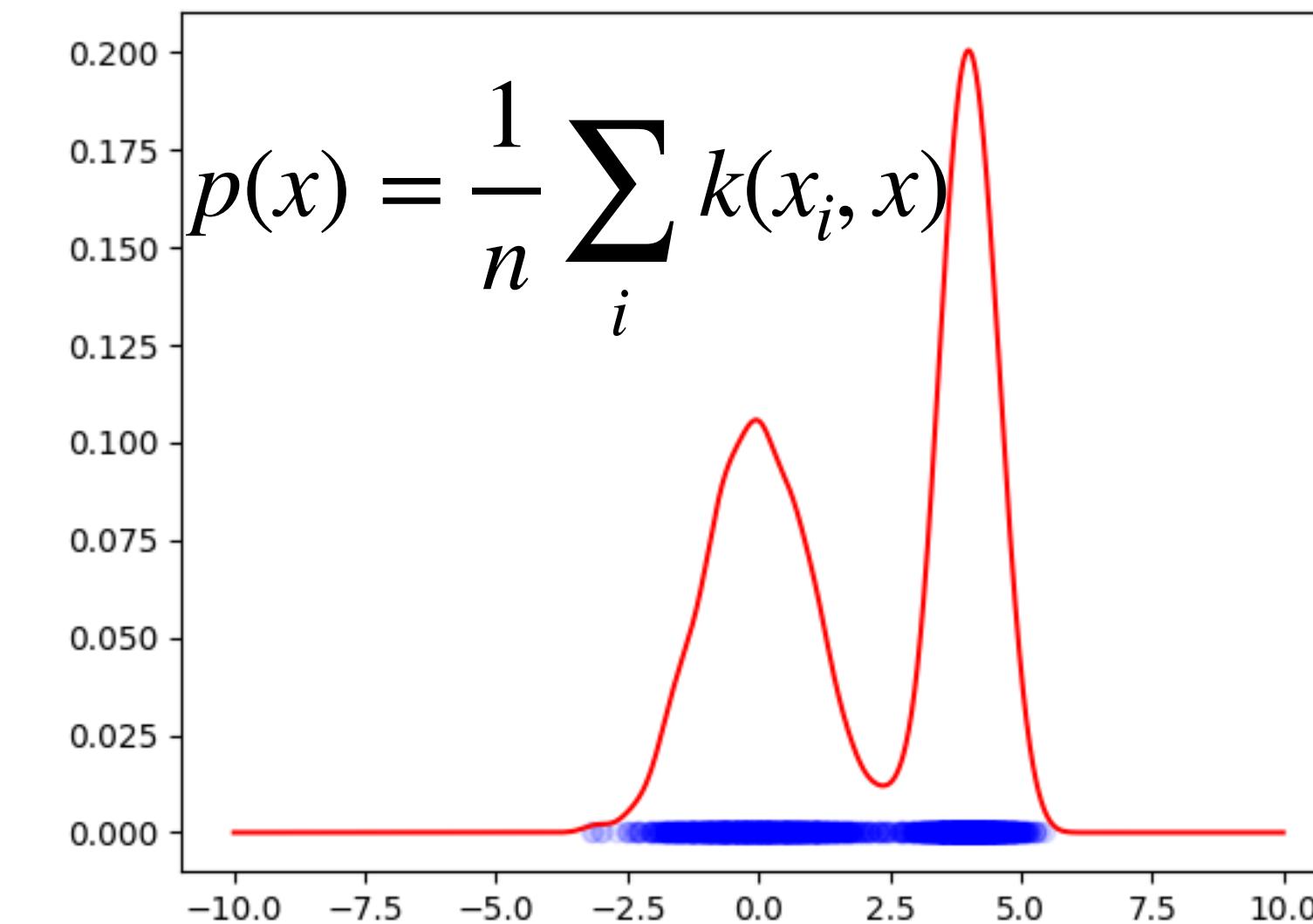
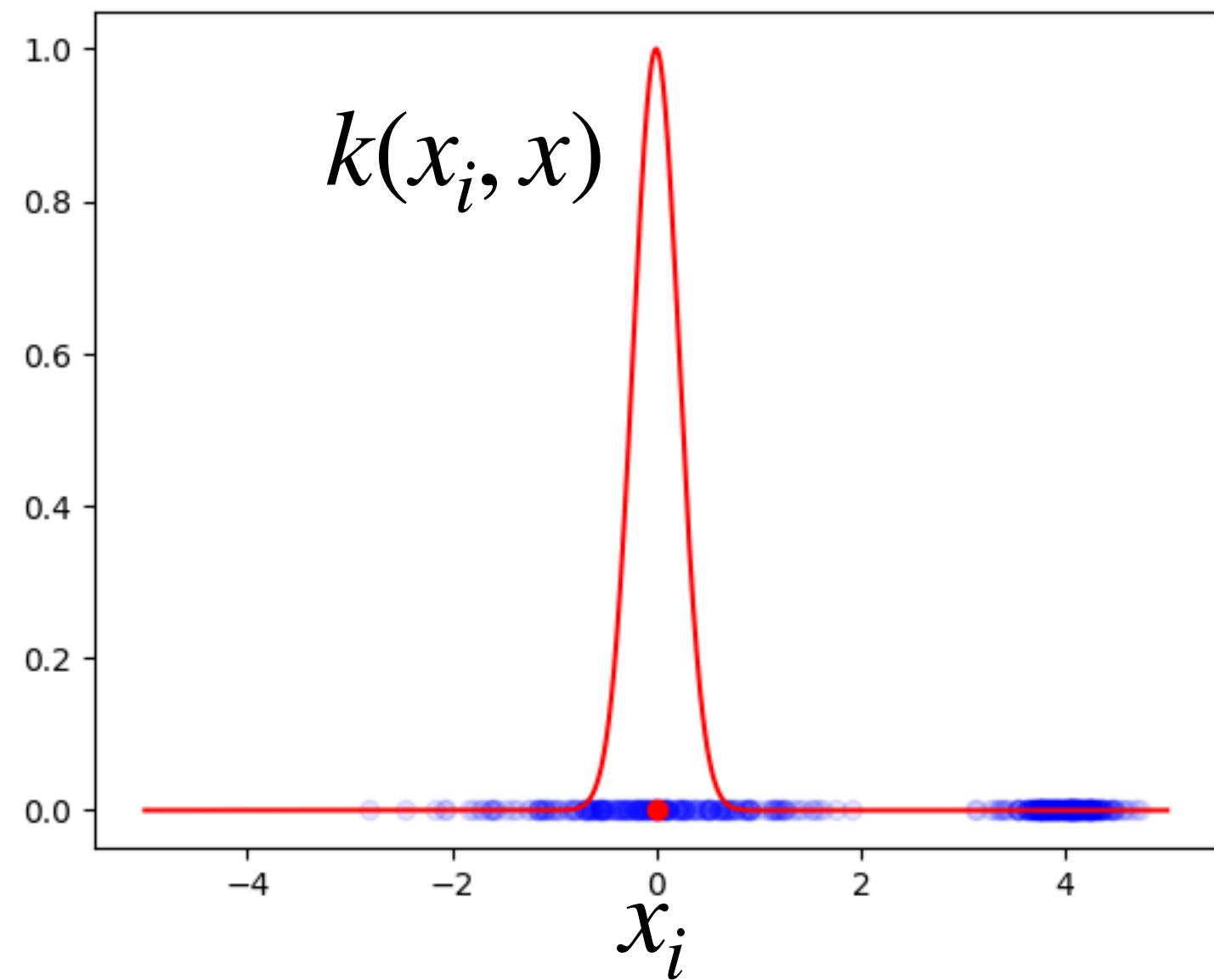
- Given i.i.d. samples x_1, \dots, x_n , what is the probability density of the underlying distribution? (Pioneered by Tukey in 1977)



Kernel density estimation

22

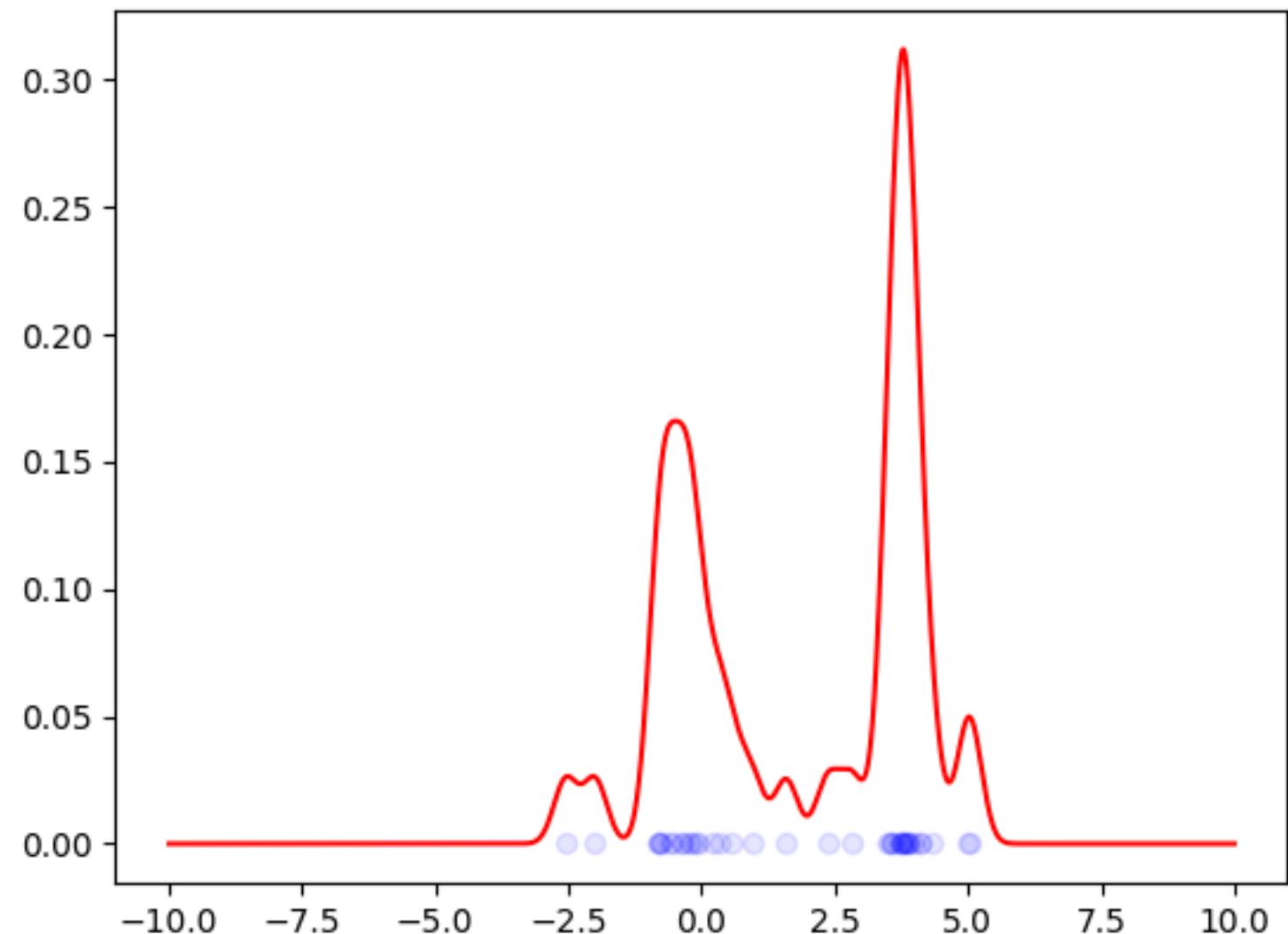
- ▶ Main idea is putting **bumps** on each point and compute the average
- ▶ **Bumps** are kernels $k(x, x_i) = \exp(-0.5\|x - x_i\|^2/\sigma^2)$



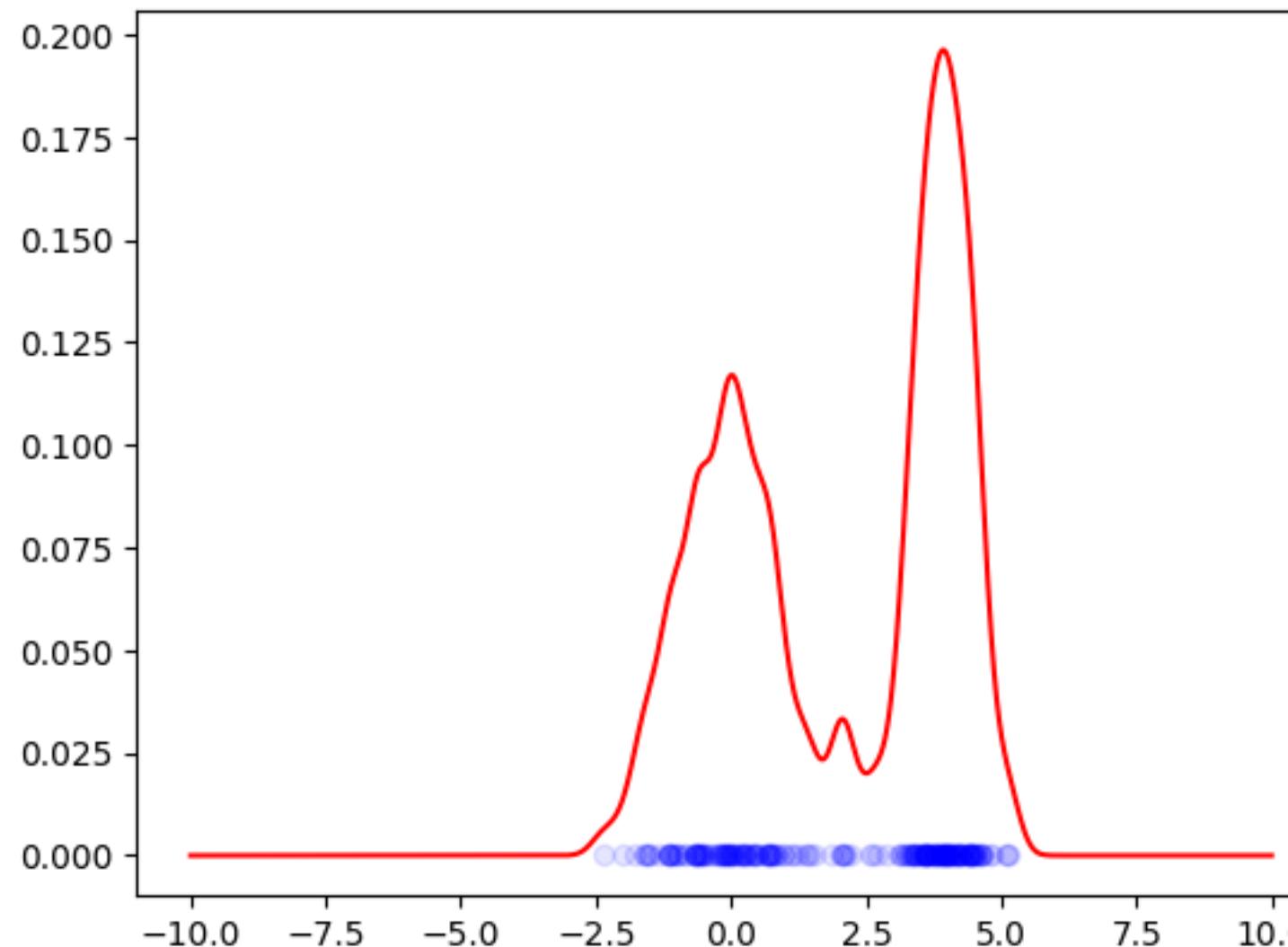
Asymptotic density estimation

23

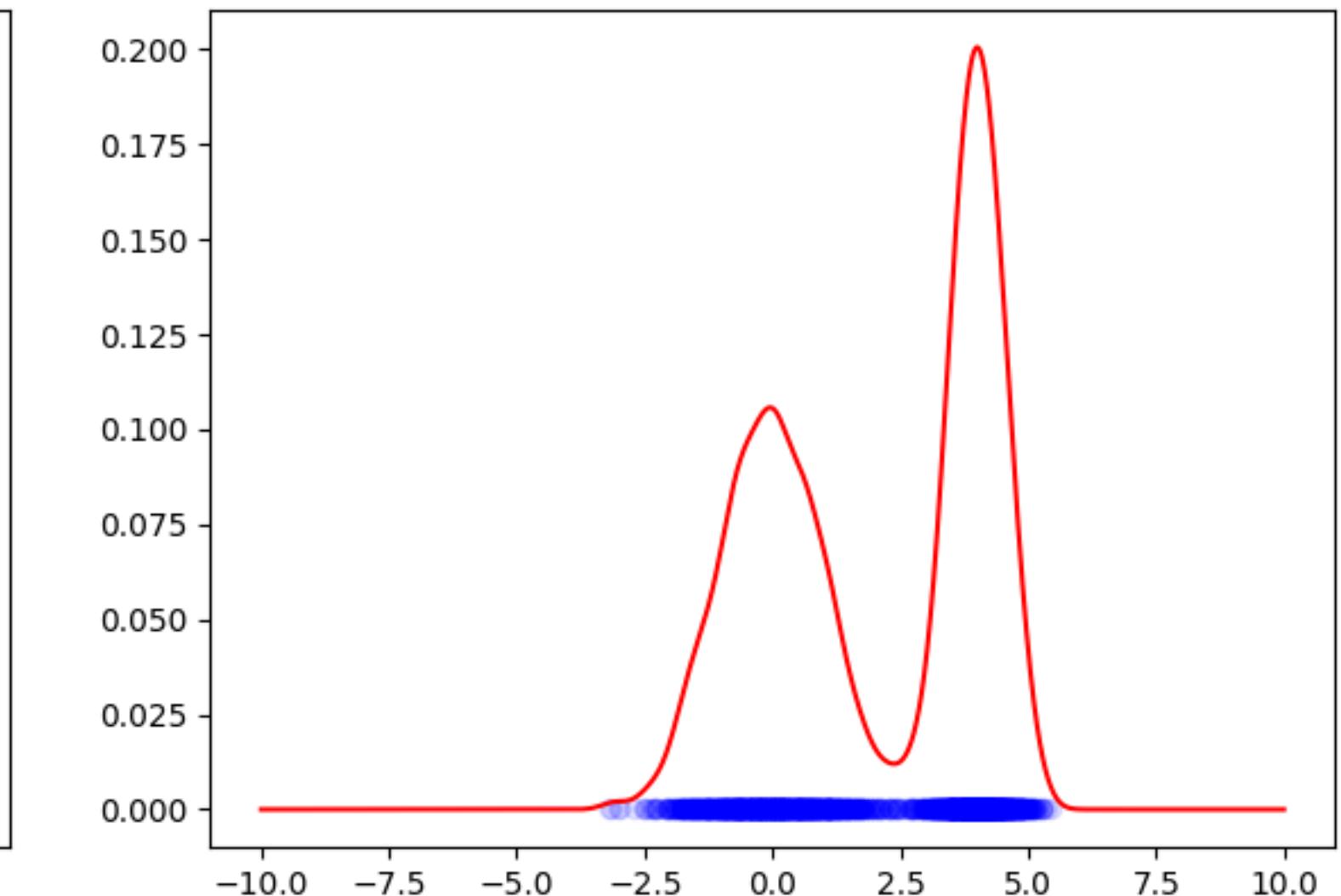
- ▶ For many distributions, the kernel density estimation converges to the true density as the number of samples goes to infinity



20 samples



100 samples



1000 samples

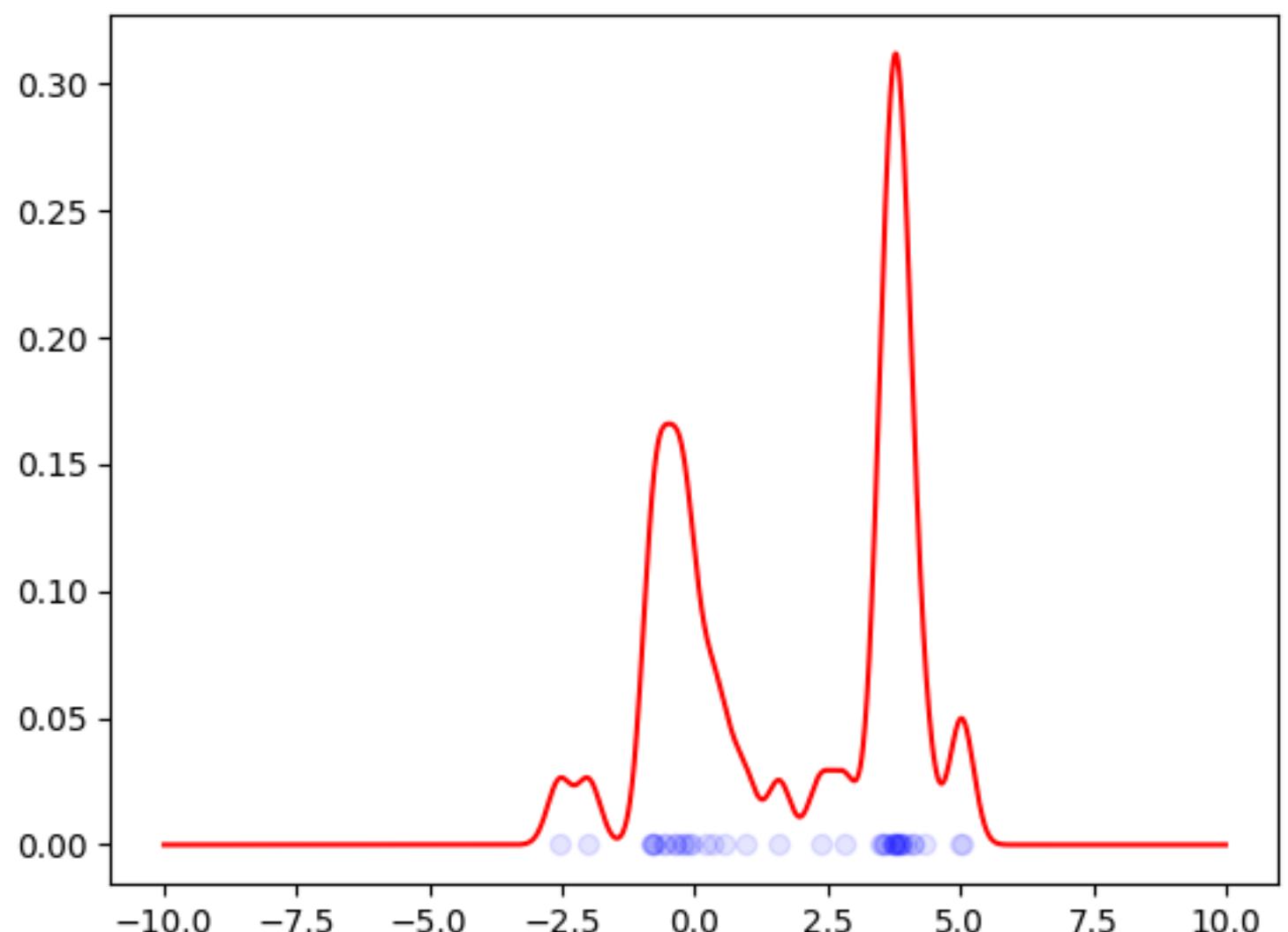
Group activity (3mins)



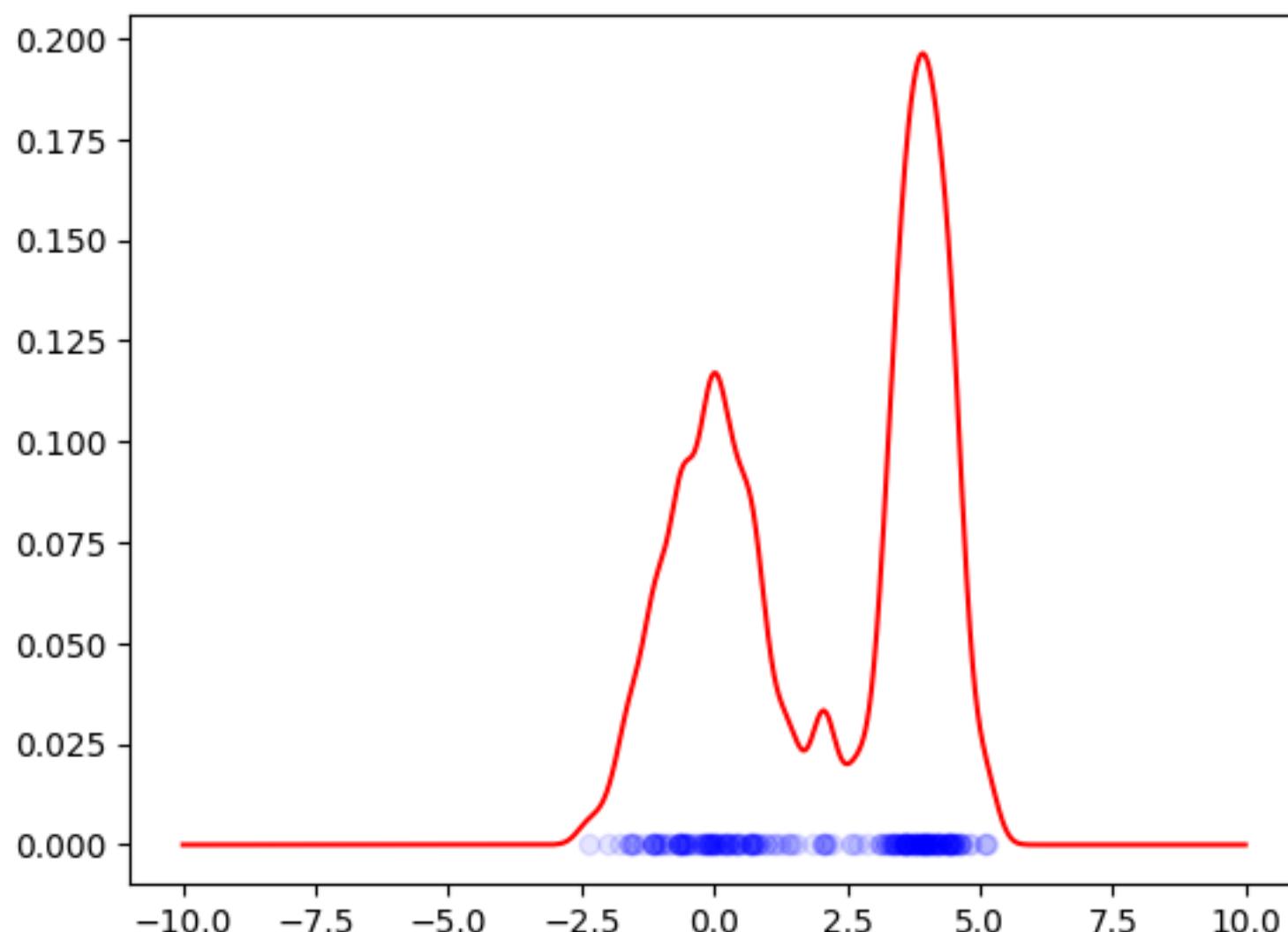
image: Flaticon.com¹

24

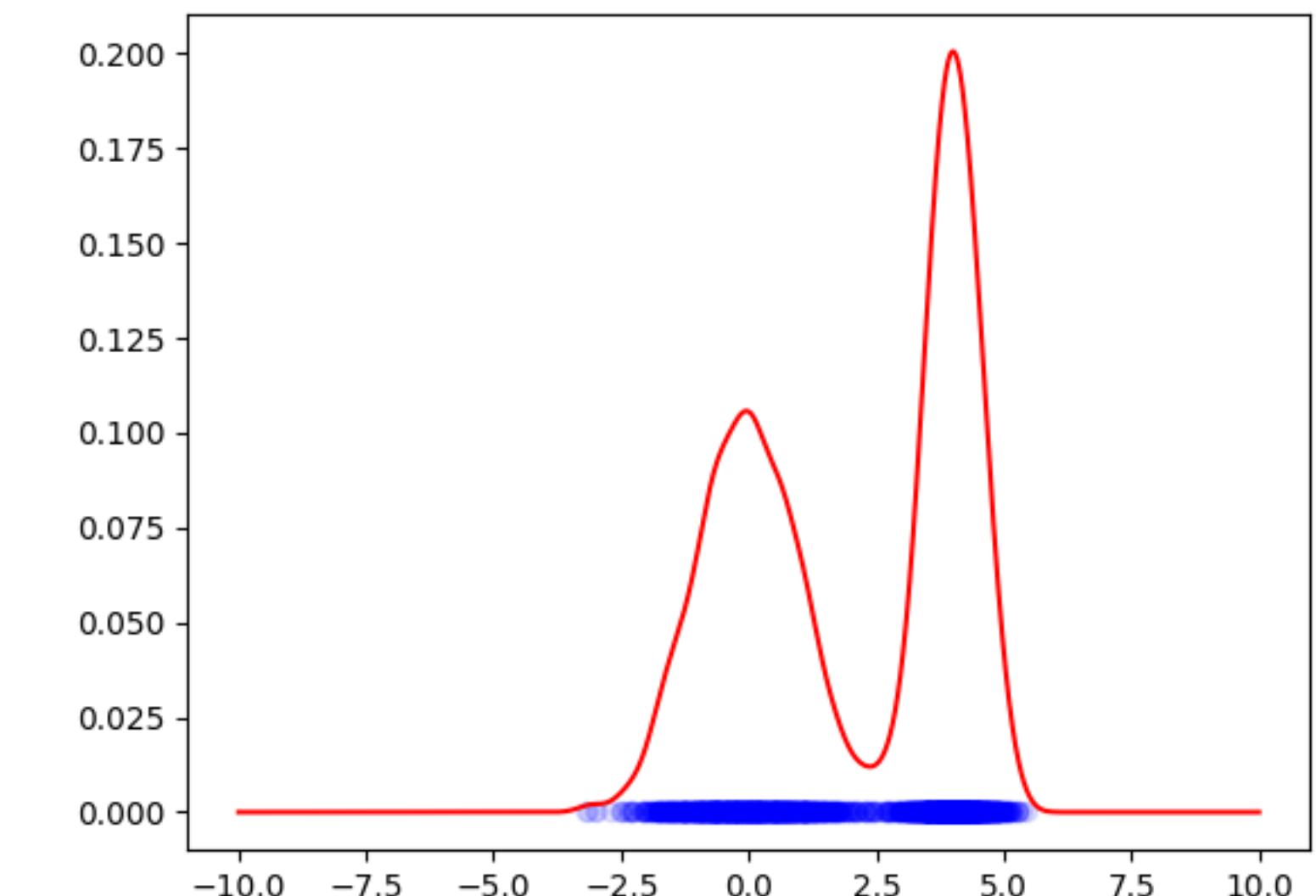
- When kernel density estimation become exact with an infinite sample size?



20 samples



100 samples



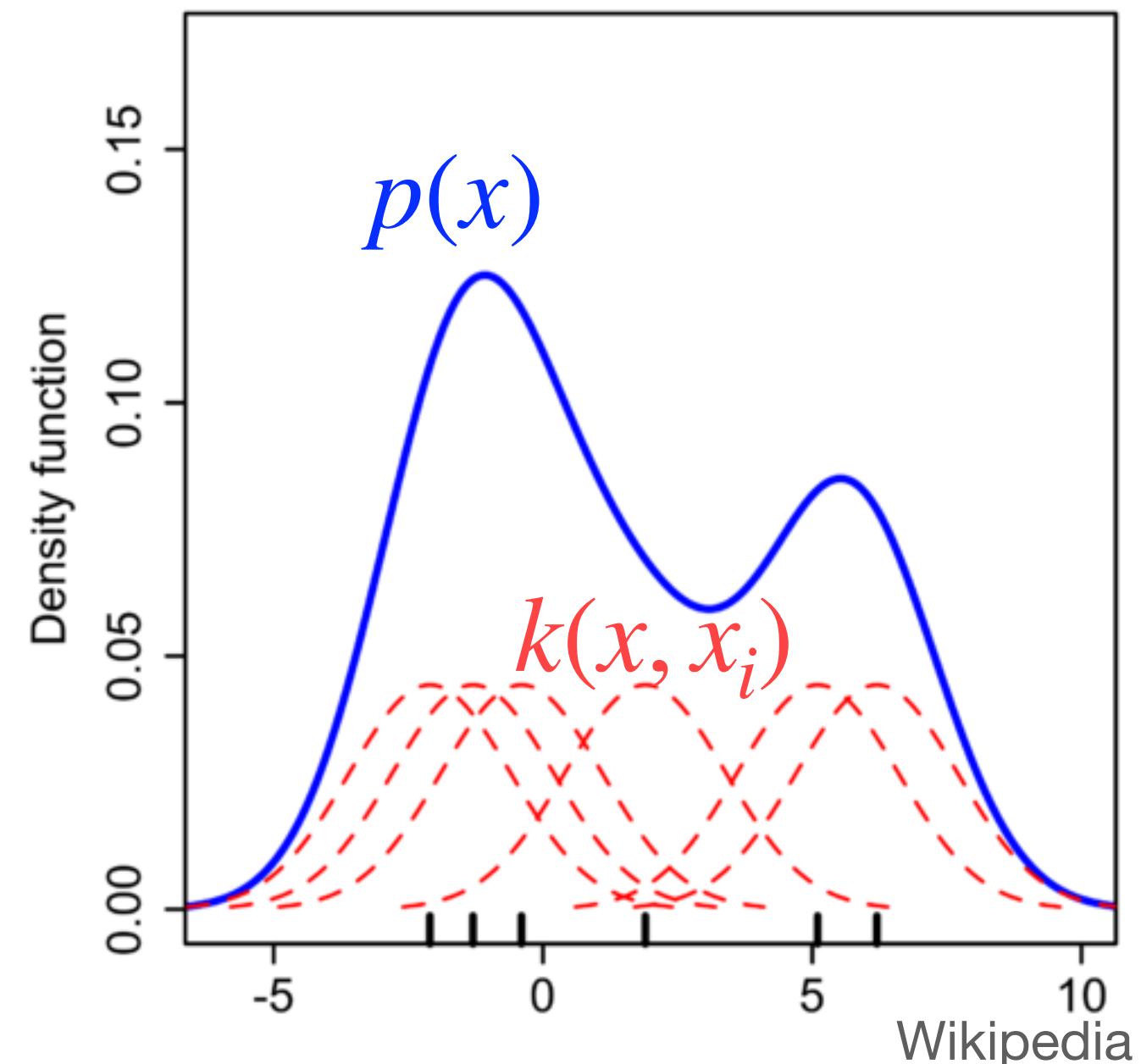
1000 samples

$$p(x) = \frac{1}{n} \sum_i k(x_i, x)$$

When does error vanishes with sample size?

25

- ▶ $p(x) = \lim_{n \rightarrow \infty, y_i \sim p(y)} \frac{1}{n} \sum_i k(y_i, x) = \int k(y, x)p(y)dy$: reproducing property
- ▶ Under the reproducing property, the error of kernel estimation vanishes as $n \rightarrow \infty$



Kernel regression and representer theorem

26

▶ Kernel regression: $f^* = \arg \min_f \sum_i (f(x_i) - \underbrace{y_i}_{p(x_i)})^2 + \lambda \|f\|^2$

subject to f obeys the reproducing property $f(x) = \int_{-\infty}^{\infty} k(y, x)f(y)dy$

▶ **Representer theorem [Schölkopf et al. (2001)]:** $f^* = \sum_i \alpha_i k(x, x_i)$ for $\alpha_1, \dots, \alpha_n \in \mathbb{R}$

▶ As long as: k is positive semi-definite, namely $\sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) \geq 0$ holds for all $x_i, x_j \in \mathbb{R}^d, \alpha_i \in \mathbb{R}$

Revisiting the task (7 mins)

27

- Given x , design features $\phi_1(x), \dots, \phi_n(x)$ such that $y \approx \sum_{i=1}^n w_i \phi_i(x)$



Scan to Start

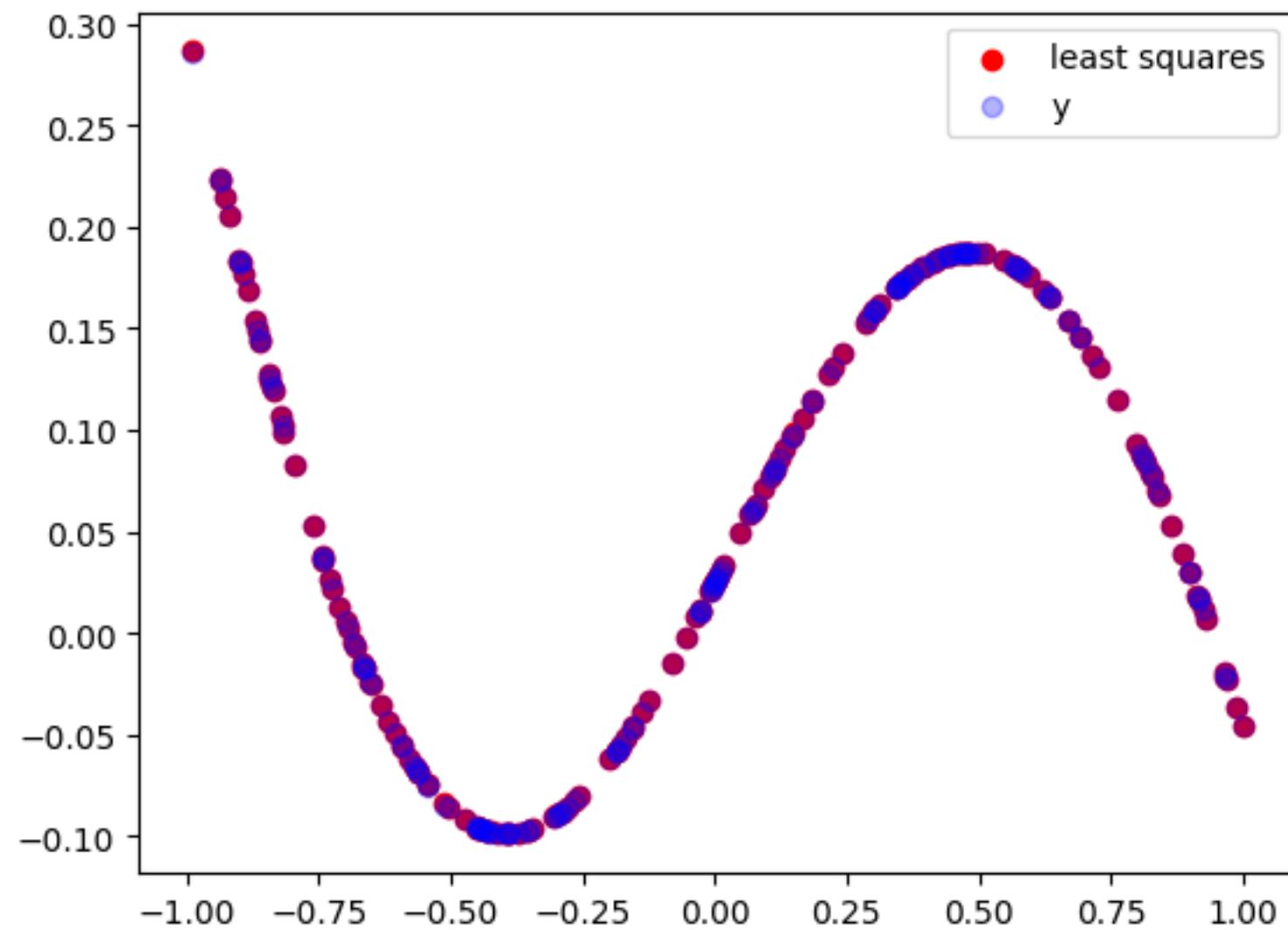
<https://shorturl.at/ABb6H>

<https://colab.research.google.com/drive/1M-nRBdhg1XJiV8sy4wMkUsfPJKKSCB0?usp=sharing>

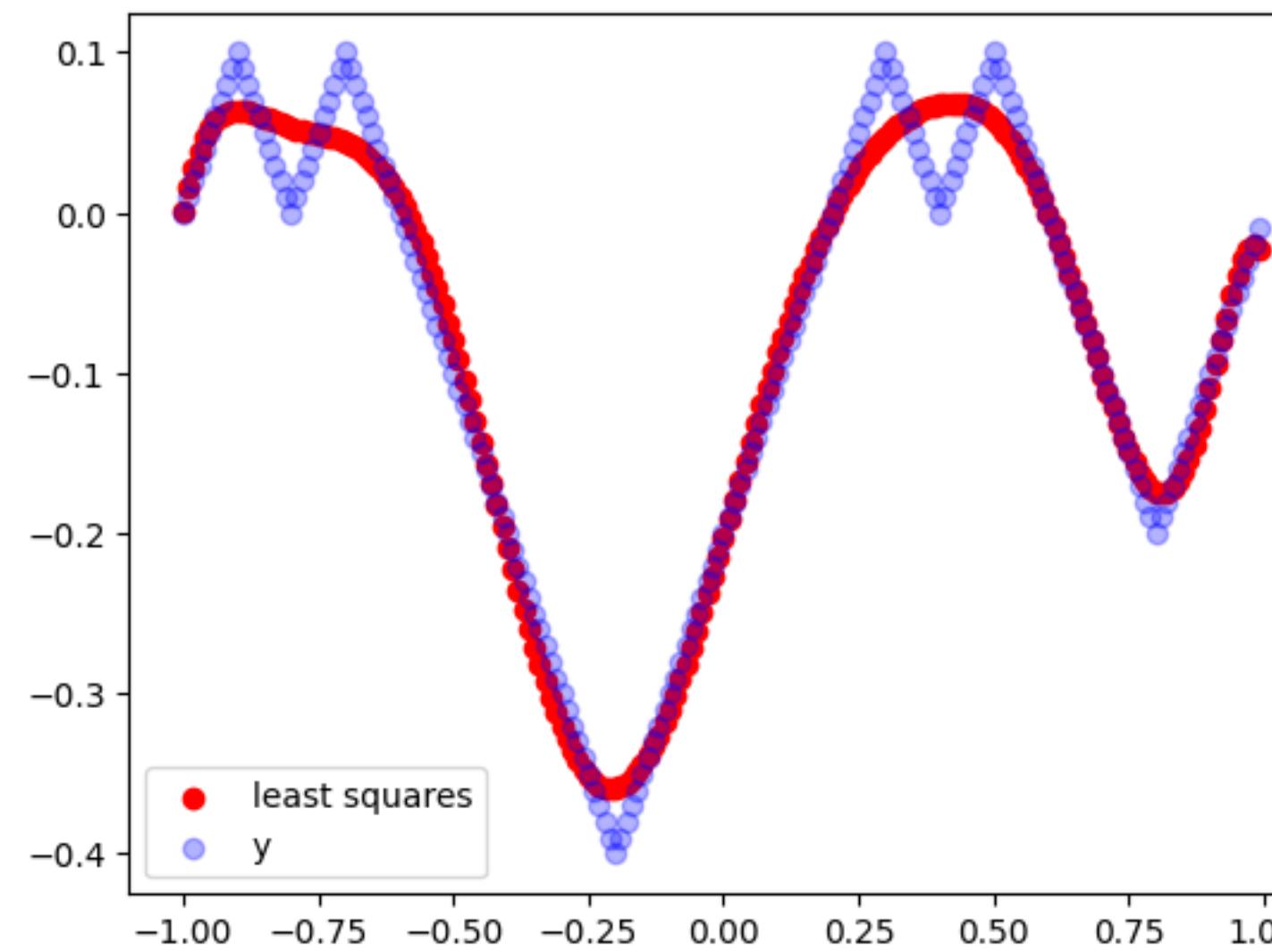
$$\phi_i(x) = e^{-0.5\|x-x_i\|_2^2}$$

My results

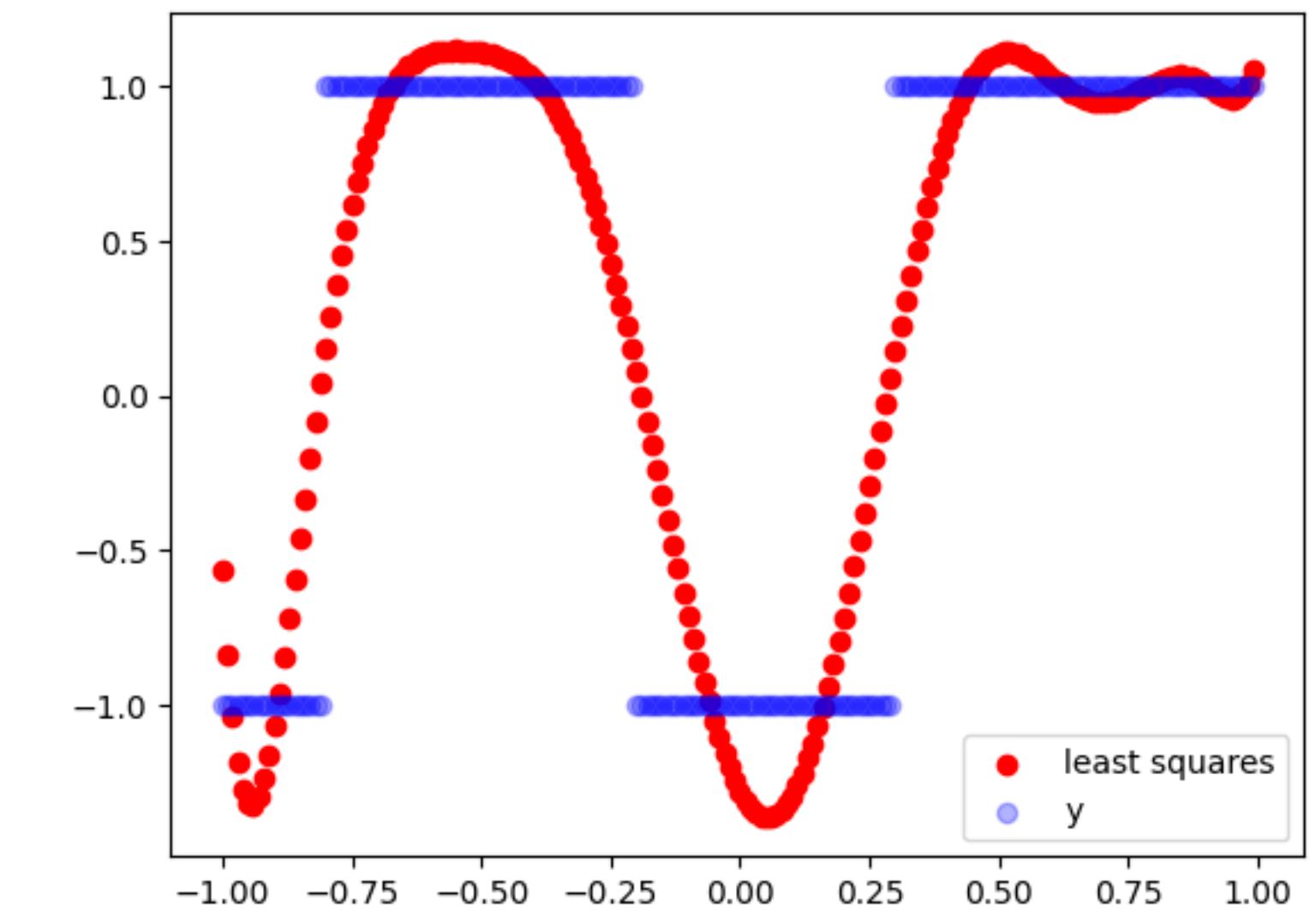
28



Good approximation



Reasonable approximation



Poor approximation

Summary

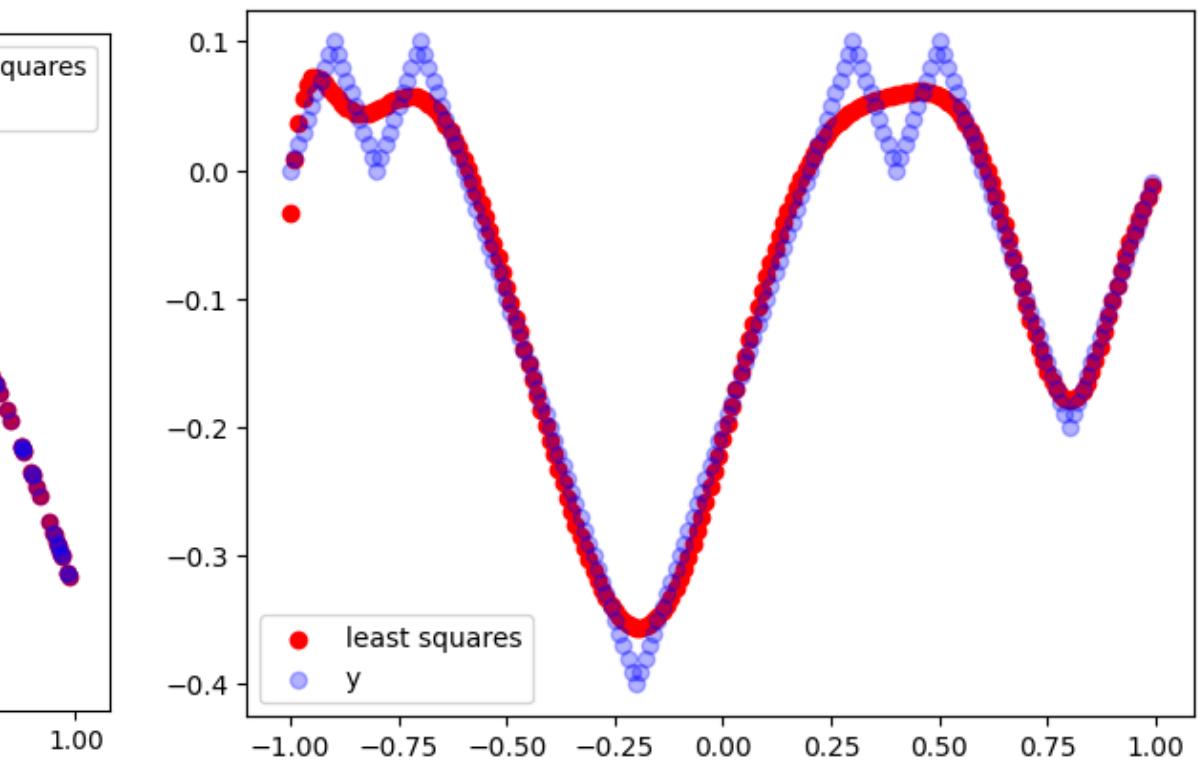
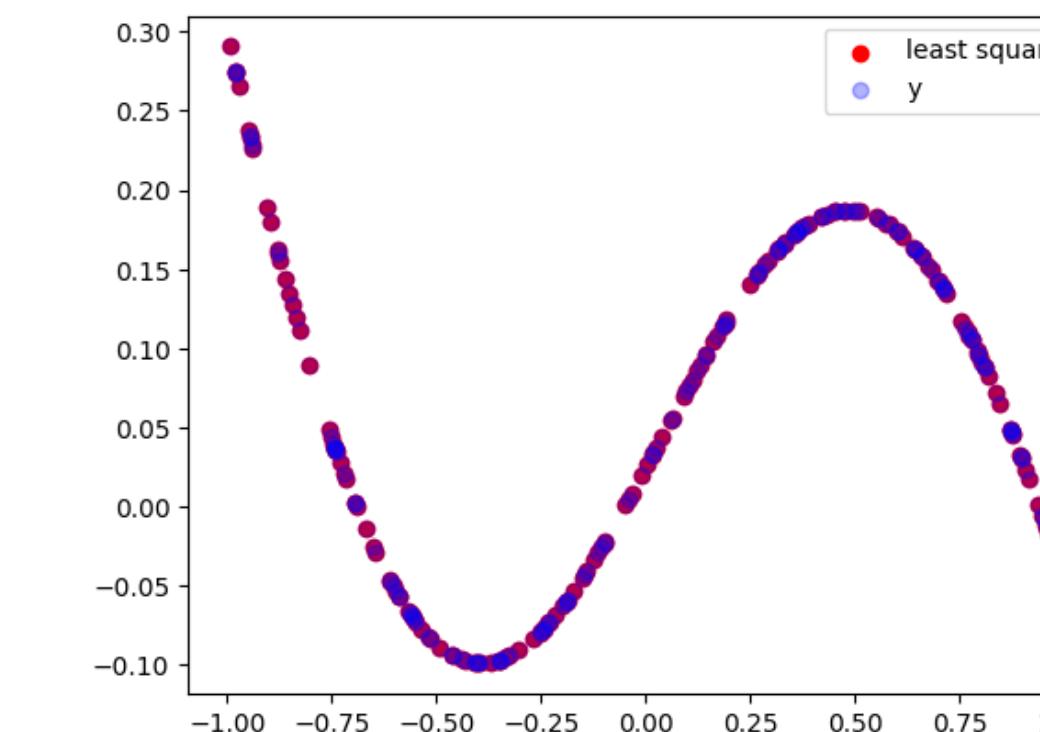
29

Theory

$$f^* = \sum_i \alpha_i k(x, x_i)$$

Representer theorem

Observation



Kernel methods provide universal features

Group activity (2mins)

30

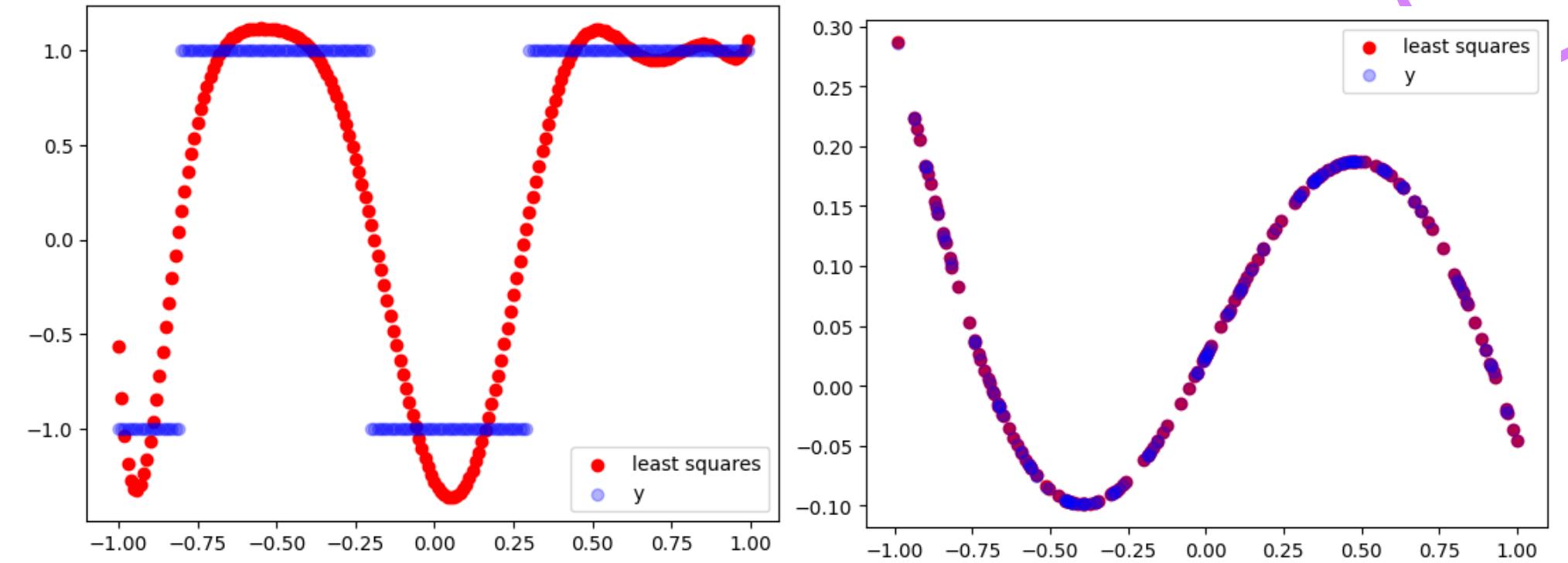
- ▶ For which functions, the estimation is good?

Recall

Kernel regression: $f^* = \arg \min_f (f(x_i) - \underbrace{y_i}_{p(x_i)})^2 + \lambda \|f\|^2$

subject to f obeys the reproducing property $f(x) = \int_{-\infty}^{\infty} k(y, x)p(y)dy$

▶ Representer theorem: $f^* = \sum_i \alpha_i k(x, x_i)$ for $\alpha_1, \dots, \alpha_n \in \mathbb{R}$



Group activity (3mins)

- What is the accuracy of kernel methods on ImageNet dataset?

3 issue with kernels

32

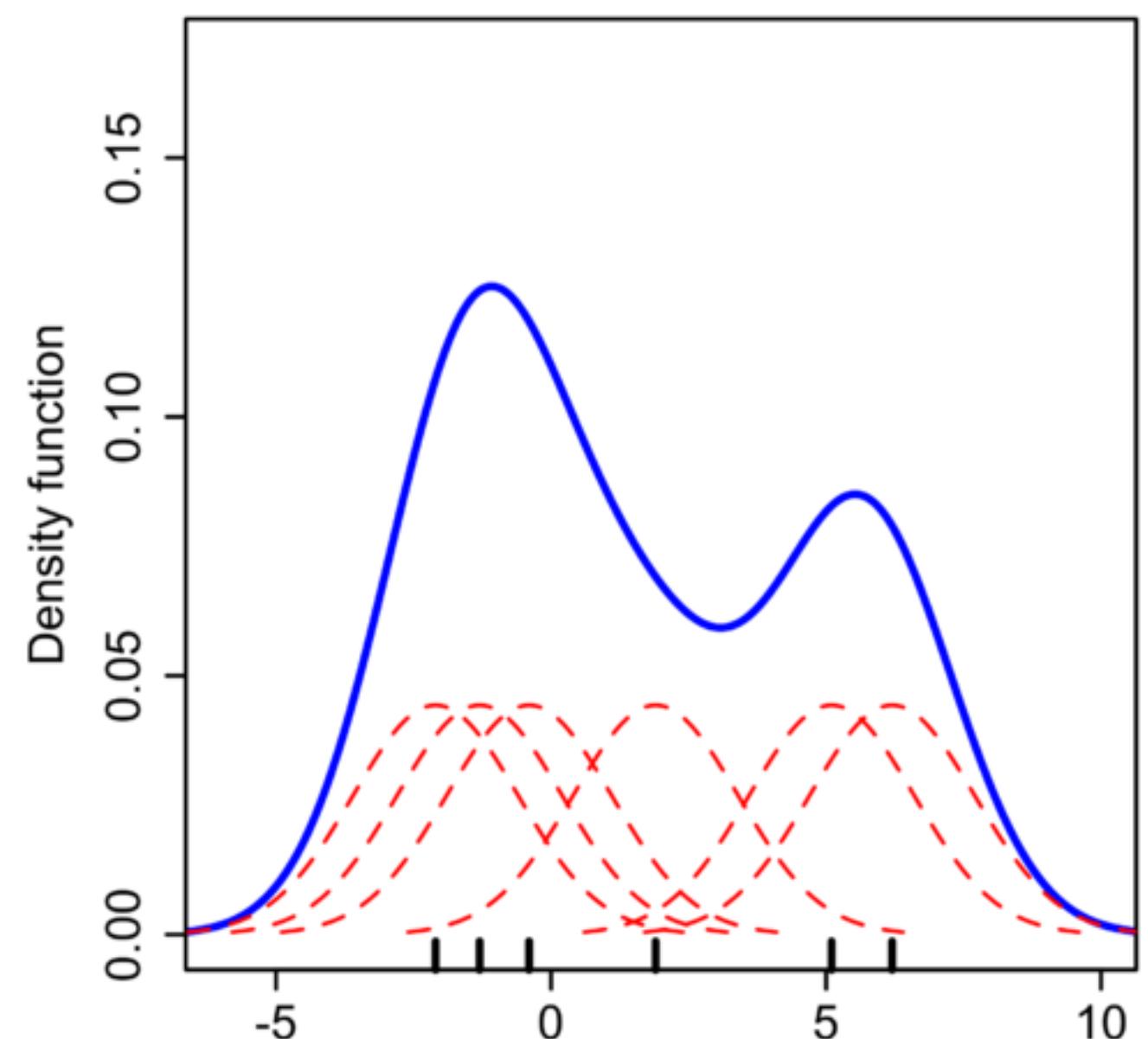
- ▶ Computational issue: Computational complexity of $O(\text{samplesize}^2)$
- ▶ Function approximation issue: **Curse of dimensionality**
- ▶ Choosing and designing proper kernels requires domain knowledge

Machine learning universal tools

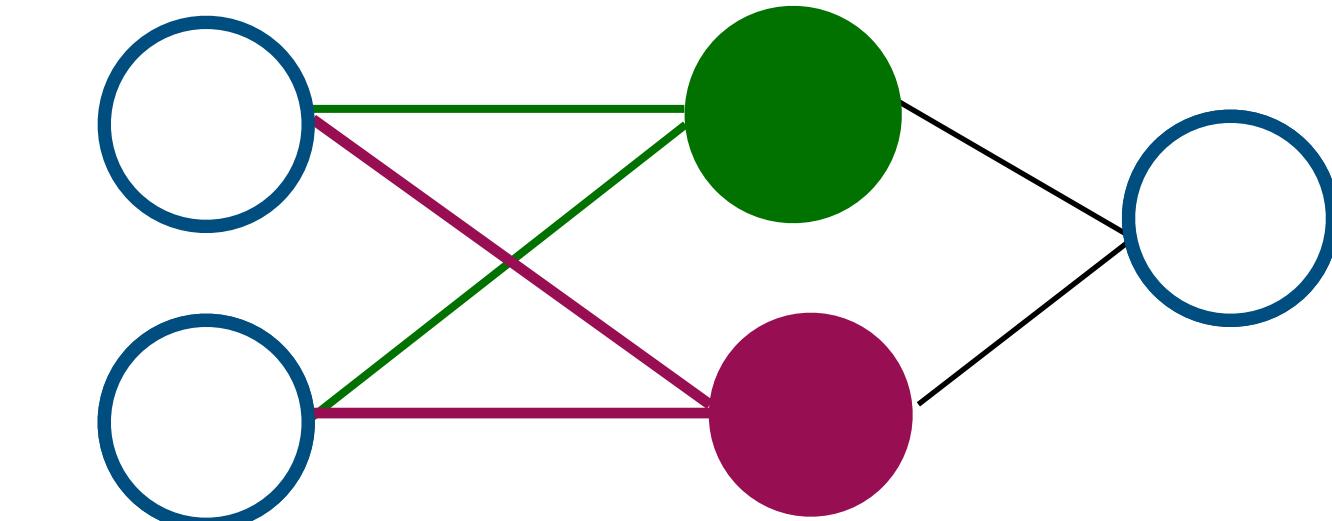
33

Kernel methods

Neural Nets



There is a bridge between these two



Thank you very much!

34

