

# Lecture 9

## 17.11.2025

# Today's plan and announcements

- Hour 1:
  - Recurrent neural networks training
  - Probability review
- Hour 2: Naive Bayes classification
  
- Announcements
  - Problem set 4 , Problem set 5
  - You can also ask questions on Problem set ~~7~~ and the python homework  
↙  
previous
  - Quiz 2 ... next Wed 26.11.2025

Goal: given  $\{x_\tau\}_{\tau=0}^t$ , predict  $x_{t+1}$ .

- Step 1: encode the input data  $a_\tau = c(x_\tau)$ ,  $y_\tau = c(x_{\tau+1})$ , where  $c(\cdot)$  is an encoding of the signal
- Step 2: fix the state dimension and the activation functions  $g_s, g_o$ .

$s_{t+1} = g_s(W_{ss}s_t + W_{sa}a_t + b_s)$ , hidden state of the recurrent neural network

$\hat{y}_t = g_o(W_{os}s_t + b_o)$ , output of the recurrent neural network

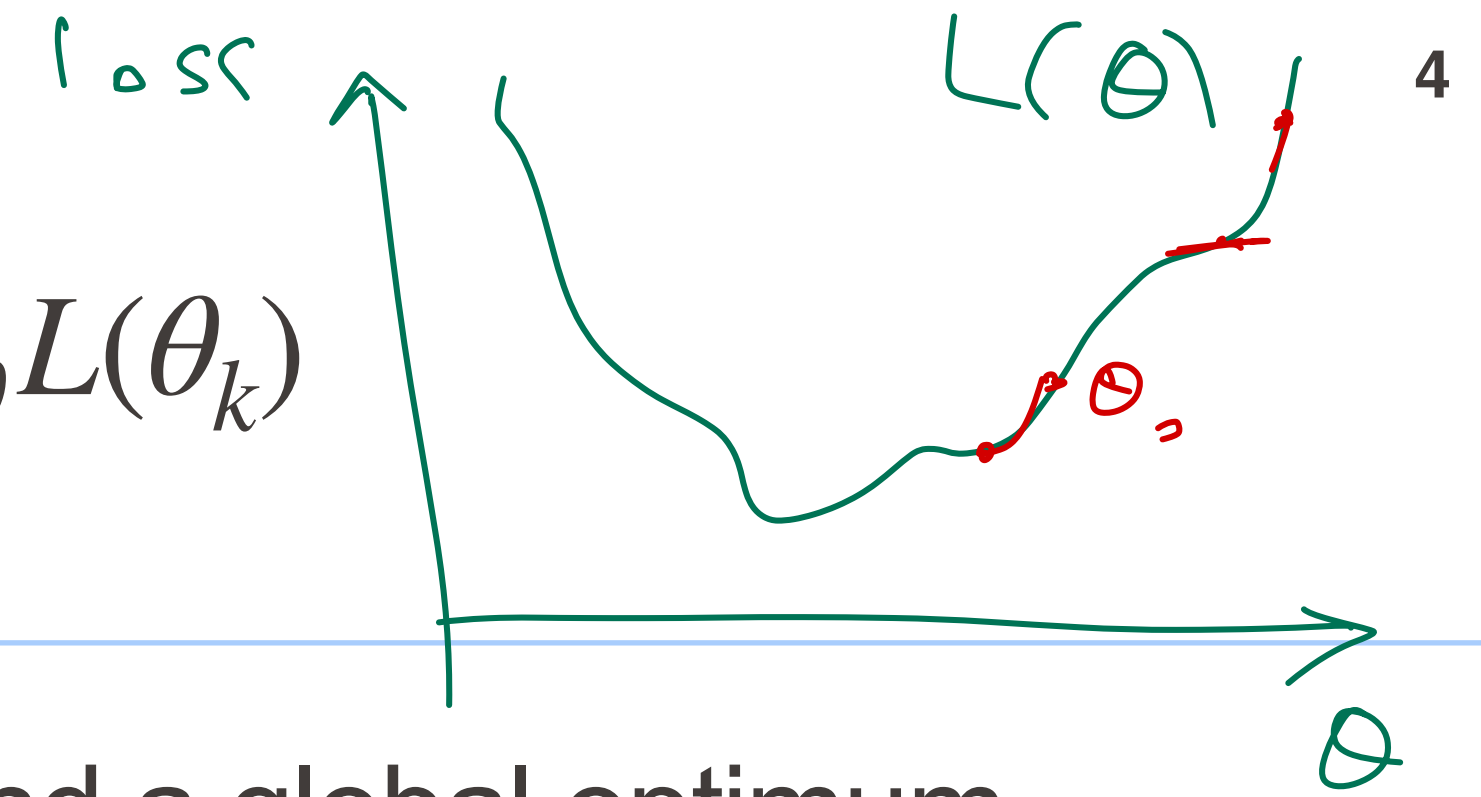
training data

- Step 3: train the parameters  $\theta = \{W_{ss}, W_{sa}, W_{os}, b_s, b_o\}$  by minimizing a loss.  $L(\theta) = \frac{1}{t} \sum_{\tau=1}^t l(y_\tau, \hat{y}_\tau)$

our model

- Step 4: given input  $a_t = x_t$  predict the next signal value  $\hat{y}_t \approx x_{t+1}$ .

# Training challenges



- Training approach: gradient descent on the loss:  $\theta_{k+1} = \theta_k - \alpha \nabla_{\theta} L(\theta_k)$
- The loss function is generally non-convex in  $\theta$ : we don't expect to find a global optimum
- Training RNNs suffers from gradient vanishing/exploding. Let's see through a simple example
  - Same issues observed for more complex examples - the growth/decay dependent on eigenvalues of  $W_{ss}$

- Consider a simple linear RNN (scalar discrete-time linear dynamical system)

$$s_{t+1} = ws_t + a_t, \quad w \in \mathbb{R}, \quad s_t \in \mathbb{R}, \quad \underbrace{w \in \mathbb{R}}_{\text{parameter}}$$

$$\hat{y}_t = s_t \quad \text{predicted output}$$

this is a recurrent neural net with 1 hidden state, & identity activation function

- Let  $L_T(\theta) = \frac{1}{T} \sum_{t=1}^T \frac{1}{2} (y_t - \hat{y}_t)^2$ , with  $\theta = w$ . (consider loss over  $T$  time steps)

## Simple example - back propagation

$$s_{t+1} = w s_t + a_t$$

$$\hat{y}_t = s_t$$

Define  $\delta_t = \frac{\partial L}{\partial s_t}$ . Observe that  $s_t$  affects the loss through  $(y_t - \hat{y}_t)^2$  term (local) and also through  $(y_\tau - \hat{y}_\tau)^2$ , for  $\tau > t$ .

Now,  $\delta_t$  satisfies the following backward recursion

$$\delta_T = y_T - \hat{y}_T \quad L(\theta) = \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{2} (y_t - \hat{y}_t)^2 \quad \Rightarrow \quad \frac{\partial L}{\partial s_T} = -2 \frac{1}{2} (y_T - \hat{y}_T)$$

$$\delta_t = (y_t - s_t) + w \delta_{t+1}, \quad \text{for } t = 0, 1, \dots, T-1$$

$$\frac{\partial L}{\partial w} = \sum_{t=0}^{T-1} \delta_{t+1} \frac{\partial s_{t+1}}{\partial w} = \sum_{t=0}^{T-1} \delta_{t+1} s_t$$

$\frac{\partial L}{\partial s_{t+1}} = \delta_{t+1}$

Observe that by recursion we get  $\delta_t = (y_t - s_t) + w(y_{t+1} - s_{t+1}) + w^2(y_{t+2} - s_{t+2}) + \dots + w^{T-t}(y_T - s_T)$ .

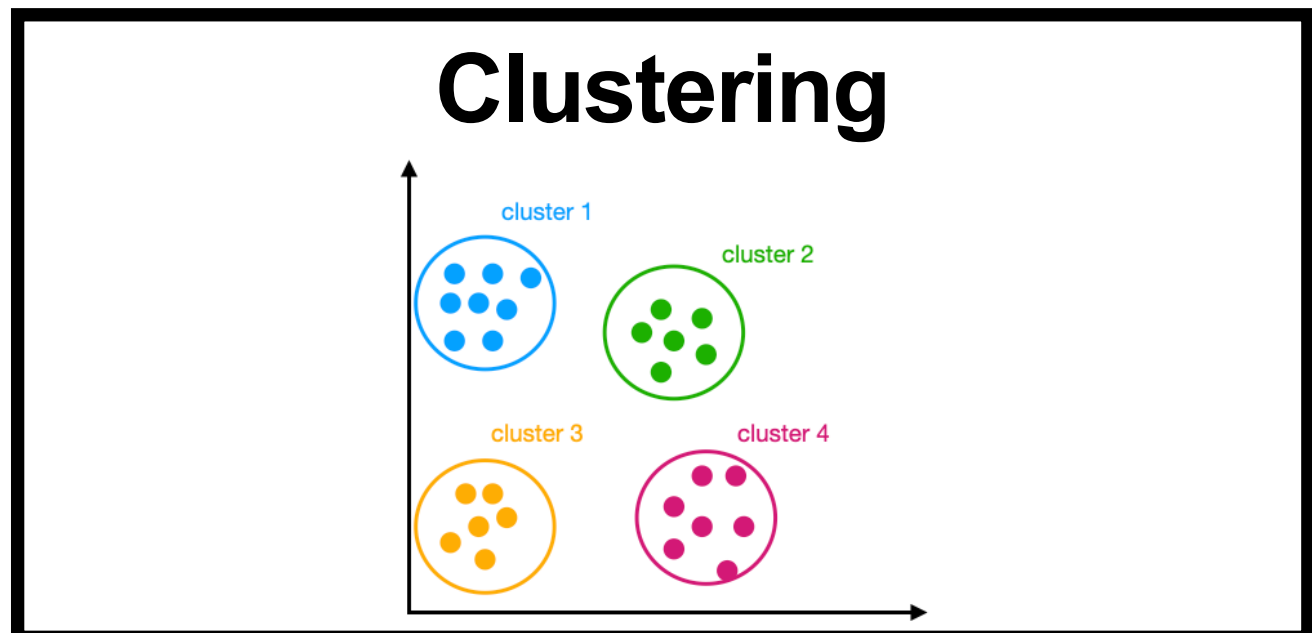
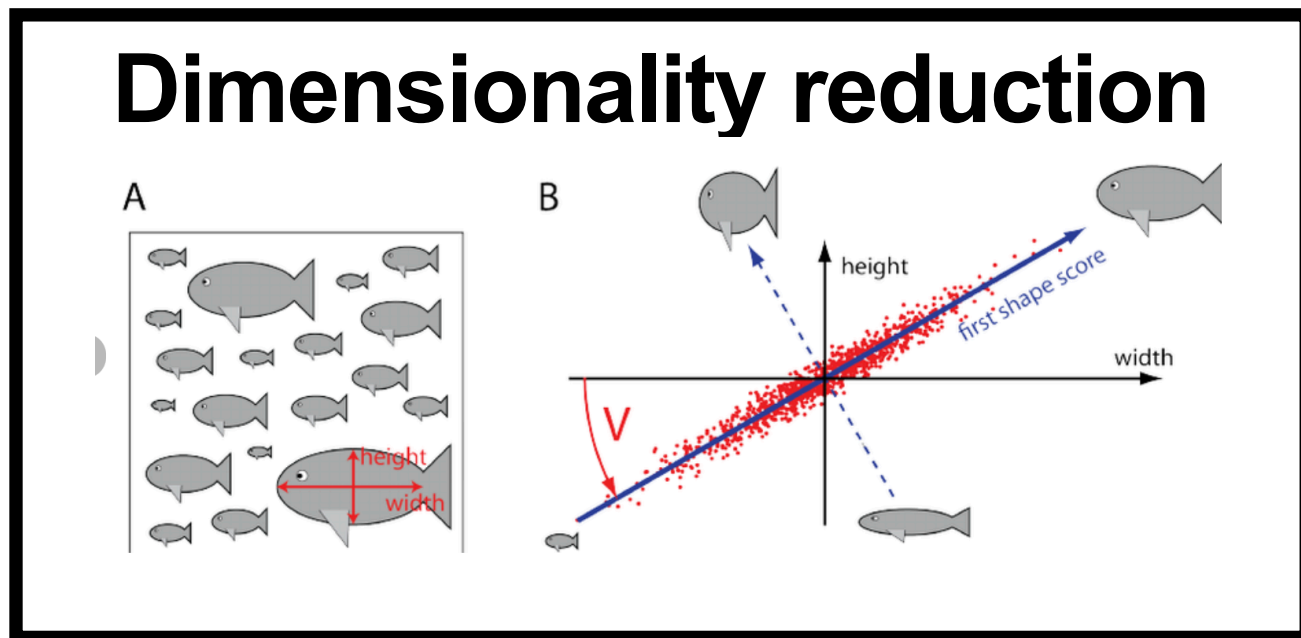
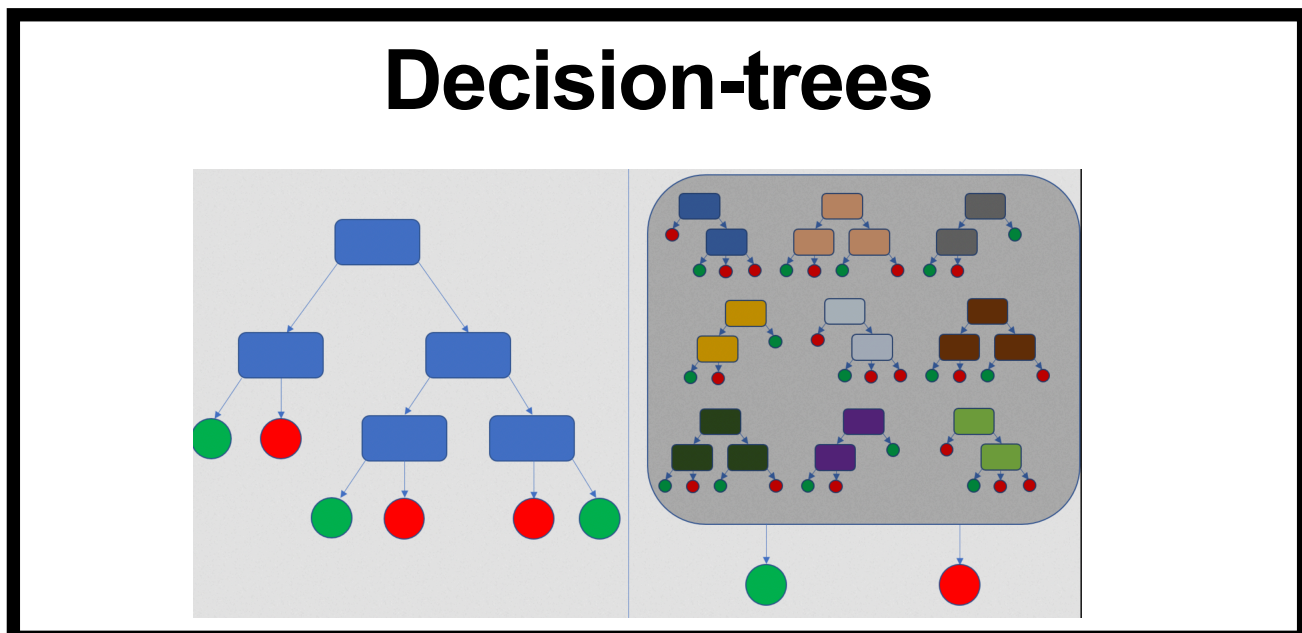
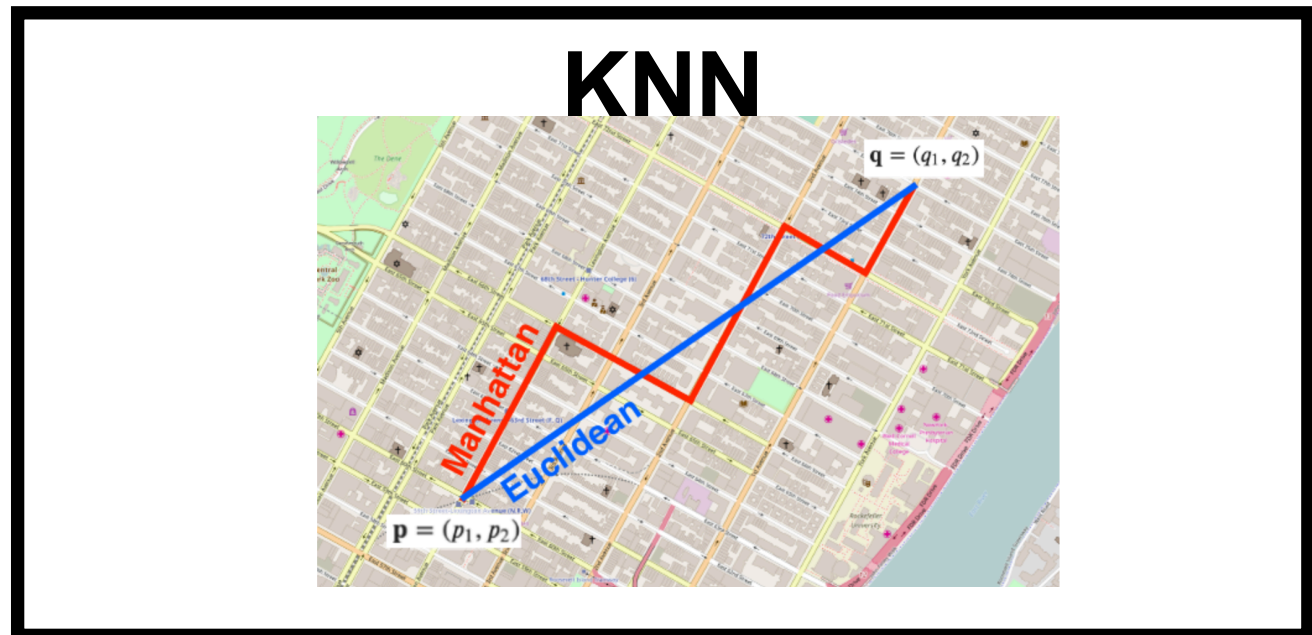
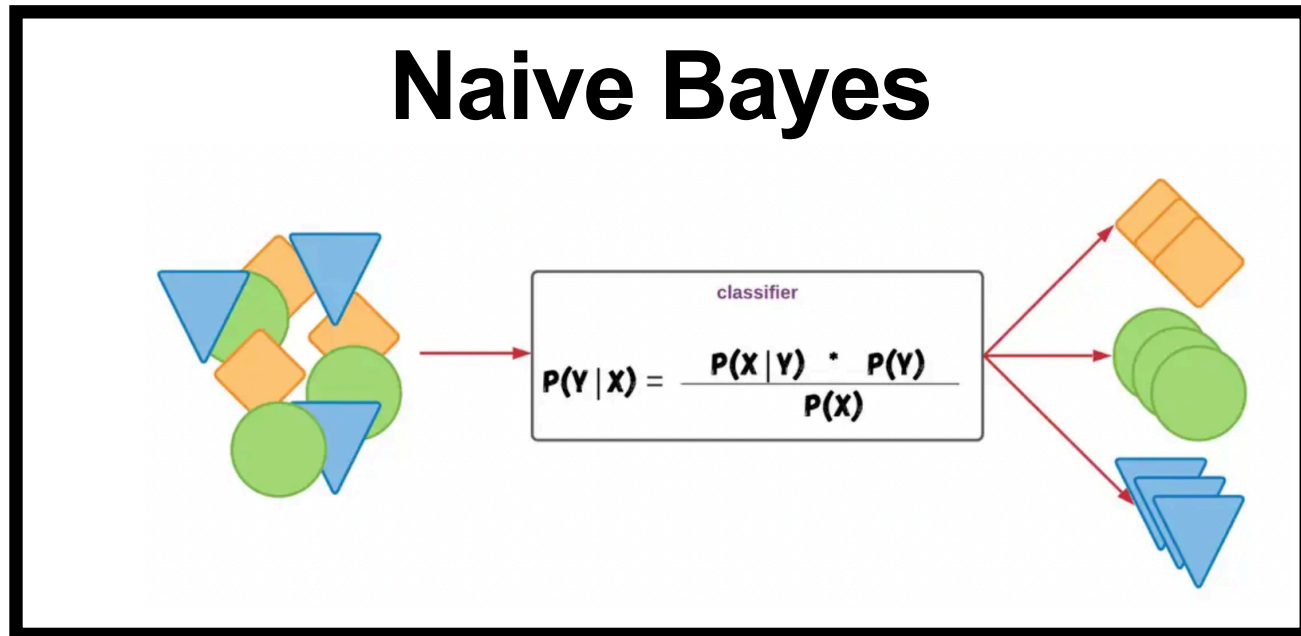
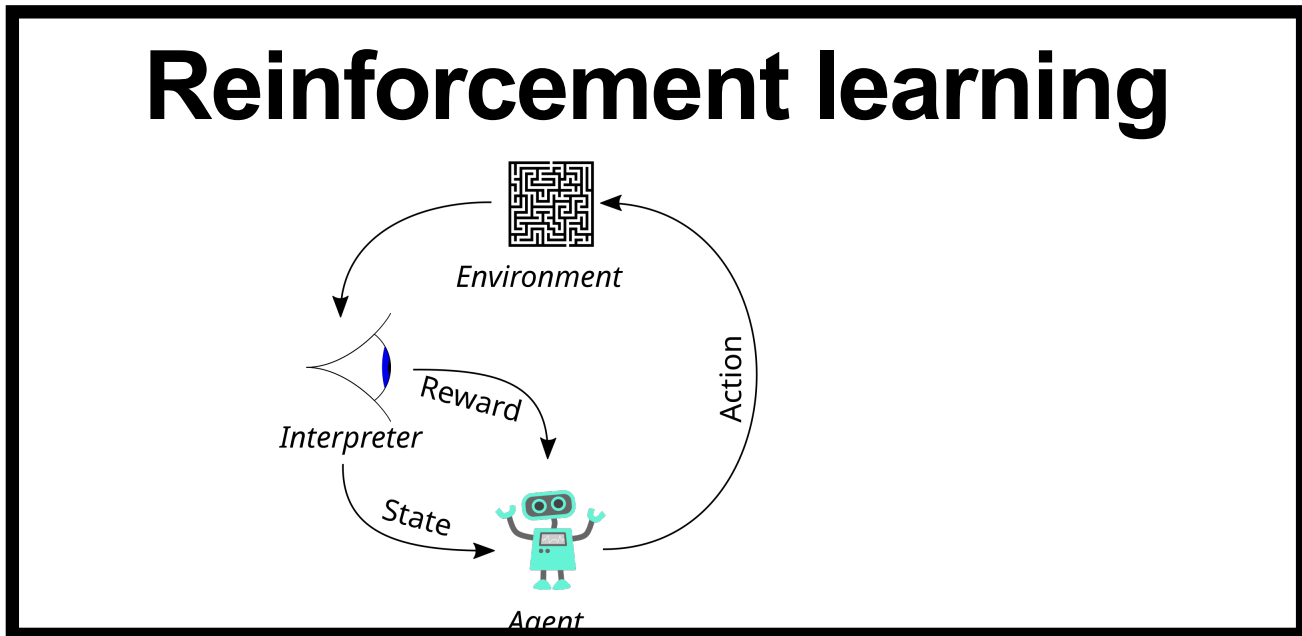
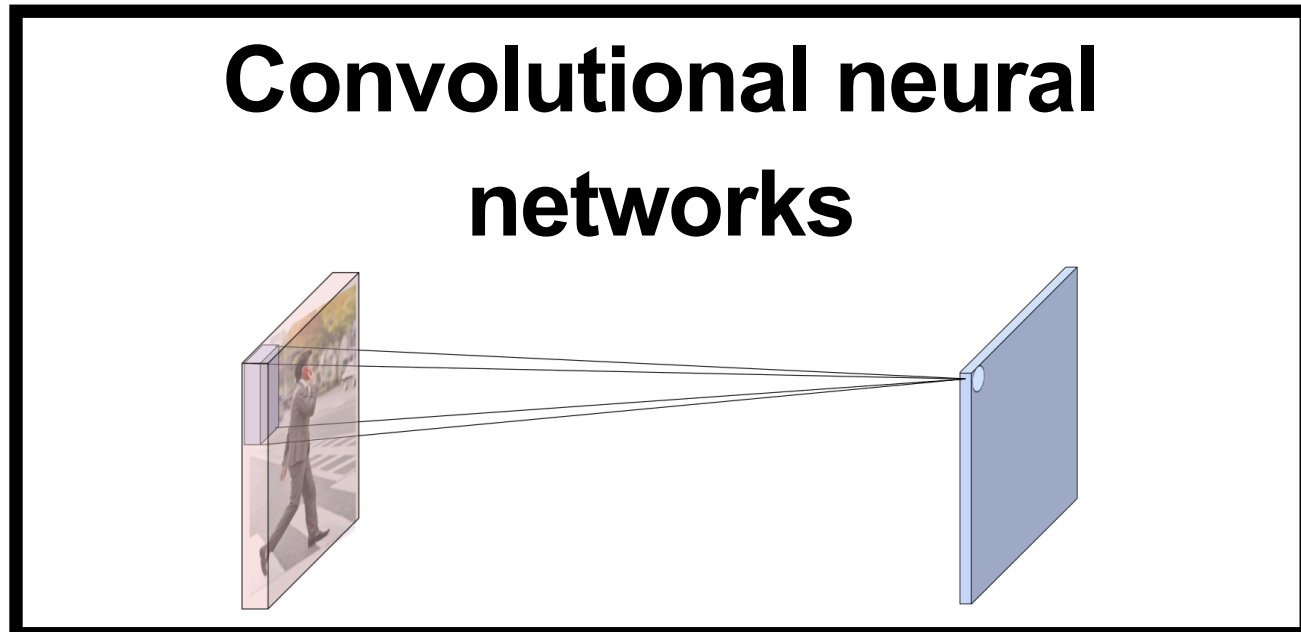
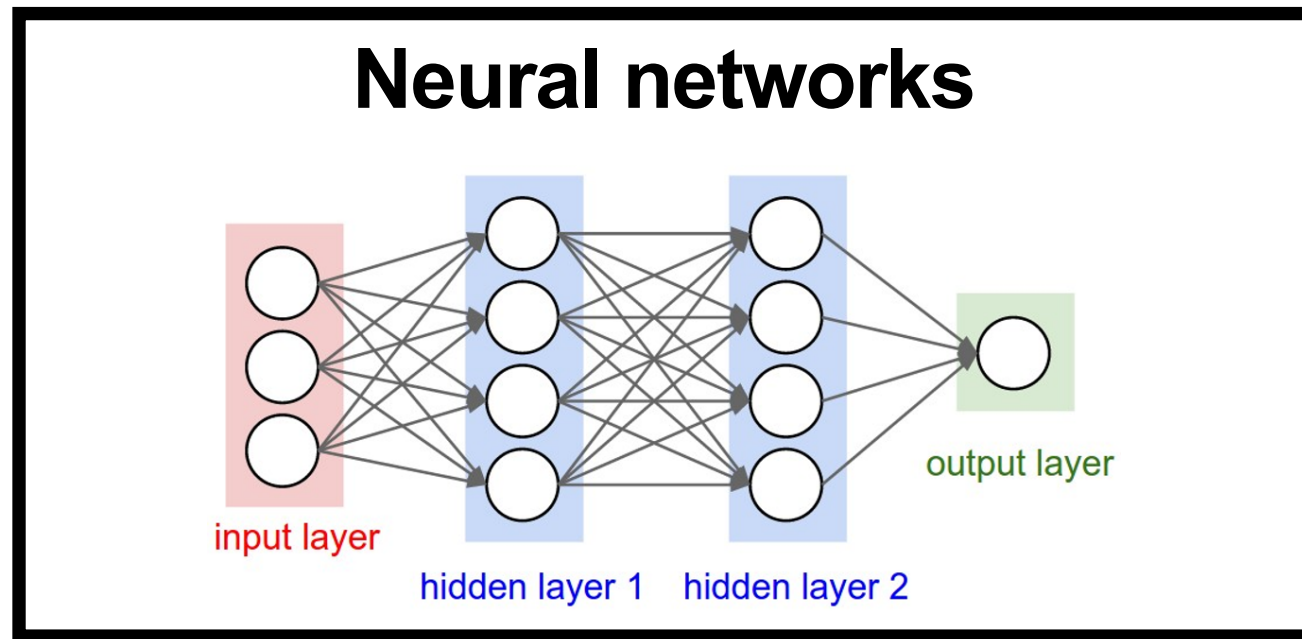
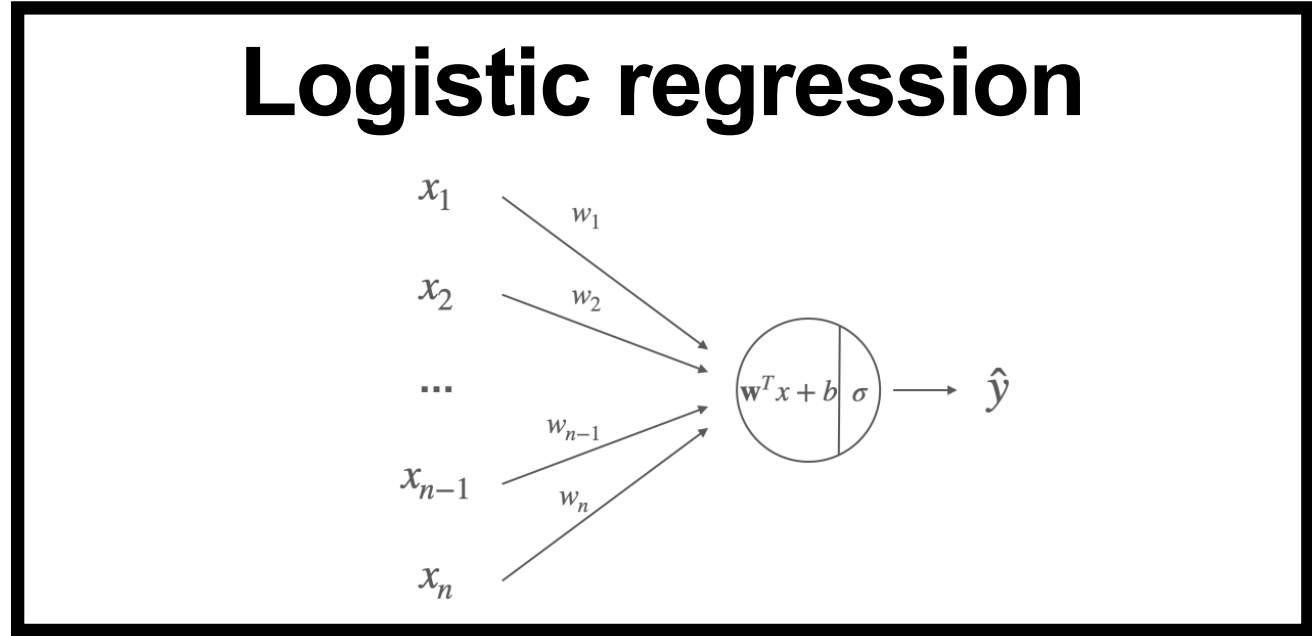
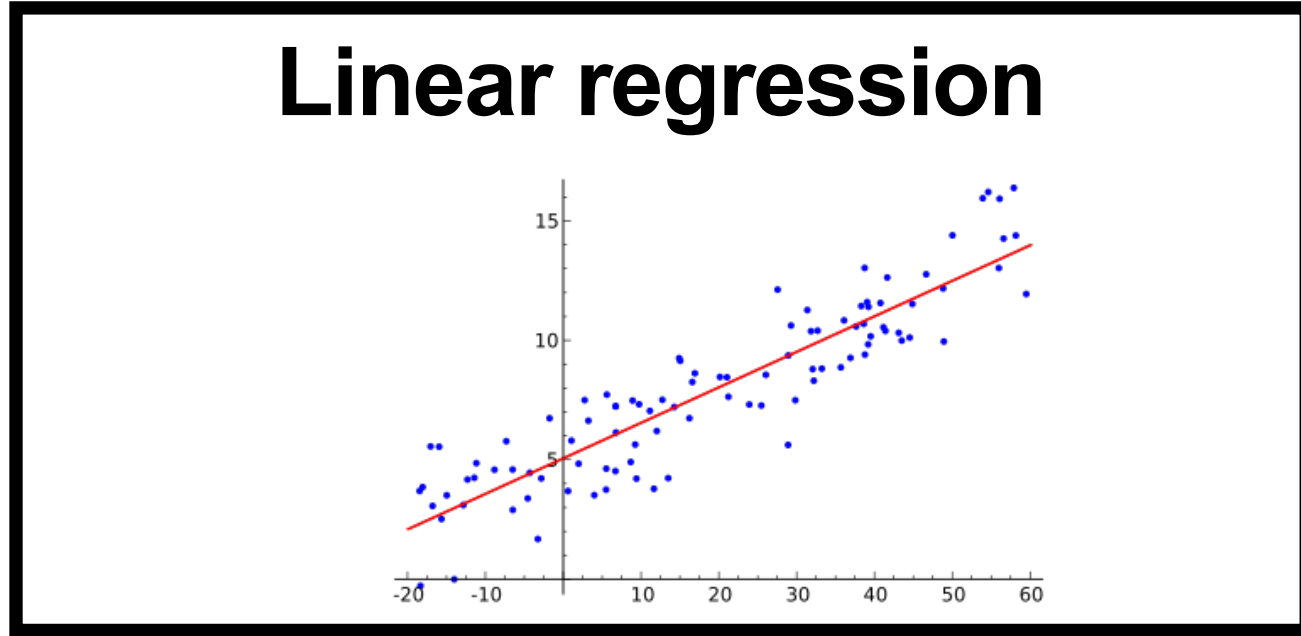
- for  $|w| < 1$  the gradient terms above vanish and for  $|w| > 1$  they explode as  $T-t$  increases.

vanishing gradient

explodes gradient

- Powerful and early models for time-series/sequence prediction
  - Language models, forecasting
- Backpropagation generally suffers from gradient vanishing/exploding as the time horizon increases -> this has led to engineering hacks of the same mathematical idea (dynamical system prediction) that are more successful in training and inference
  - Long Short-Term Memory (LSTM)
  - Transformers: backbone of large language models such as ChatGPT
- Summary: learn the math fundamentals. You can then apply it to understand increasingly complex problems and come up with new approaches.

### Introduction



+ Recurrent Neural Net

# Approaches we look at in this course

## ■ Supervised learning:

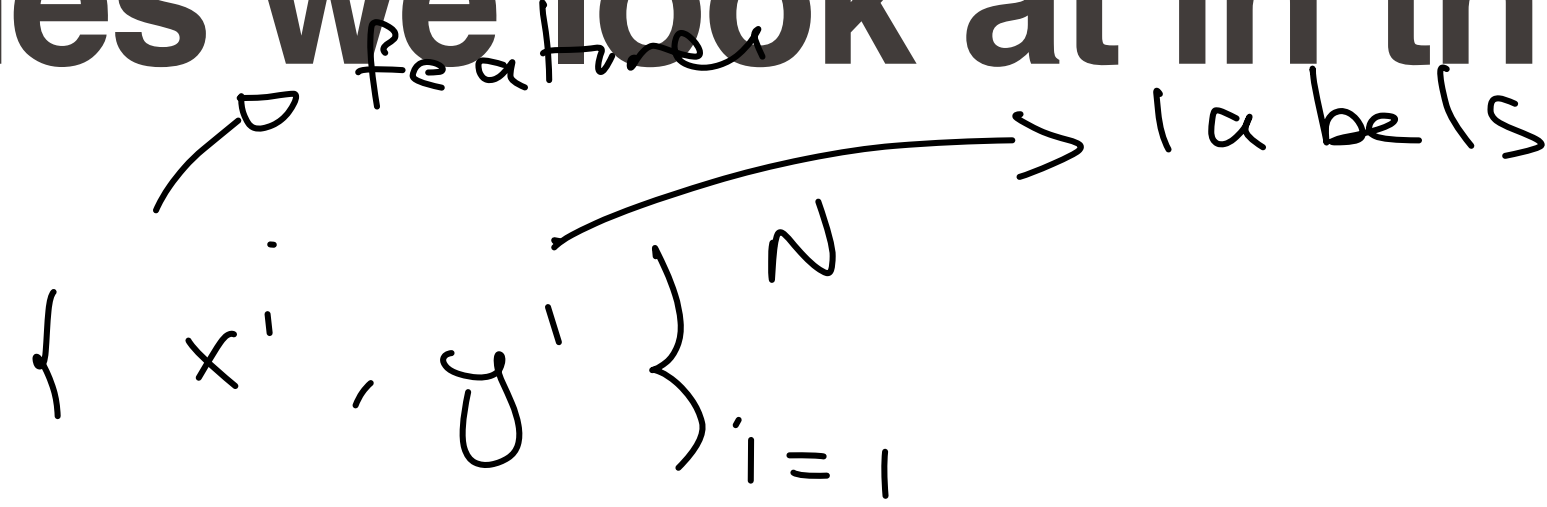
- Regression  $y^i \in \mathbb{R}^m$  : linear regression, neural network,
- Classification  $y^i \in \{1, 2, \dots, k\}$  : logistic regression, neural networks, CNN
- Time-series prediction : recurrent neural network

## ■ Reinforcement learning

- Policy gradient : learning optimal control policy using data

## ■ Unsupervised learning

- Dimensionality reduction
- Clustering



Handwritten text: data :  $\{(s_{0:H+1}^i, a_{0:H}^i), R^i \text{ (reward)}\}_{i=1}^N$

Handwritten text:  $\{x^i\}_{i=1}^N$

# Naive Bayes Classifier

classical approach  $\{x^i, y^i\}_{i=1}^N$   
 $y^i \in \{1, 2, \dots, k\}$

C supervised learning

# Supervised learning - probabilistic classifier

Recall probabilistic interpretation of logistic regression

- Data

$$\{x^i, y^i\}_{i=1}^N, \quad x^i \in \mathbb{R}^d, \quad y^i \in \{1, 2, \dots, k\}$$

- Goal: a probabilistic classifier gives probability of class of a data point  $x^{\text{test}} \Rightarrow$  probability on  $\{1, 2, \dots, k\}$

- Example: multinomial logistic regression

$$P(y = c | x) = \frac{e^{z_c}}{\sum_{j=1}^k e^{z_j}}, \quad z_c = \omega_c^T x + b_c, \quad c \in \{1, 2, \dots, k\}$$

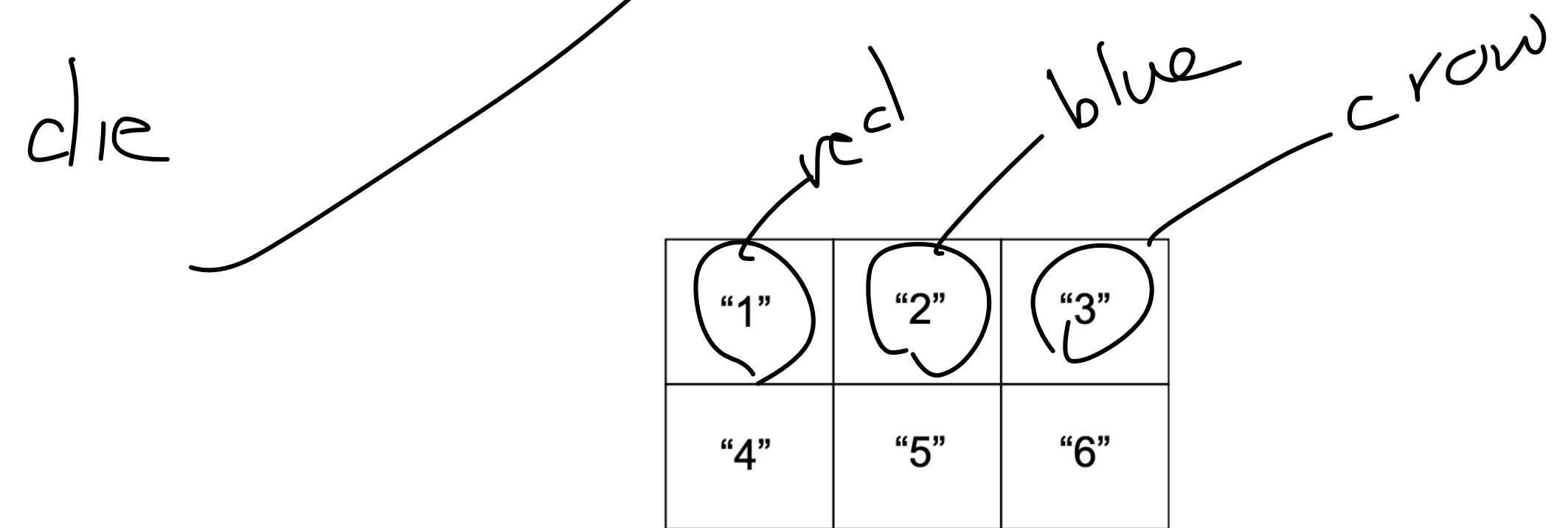
$$z_j = \omega_j^T x + b_j, \quad j \in \{1, 2, \dots, k\}$$

We will now discuss another approach to probabilistic classification

# Probability review

## Probability, random variable

- Event: outcome of a random trial, example: rolling a die
- Probability of an event : how likely an event will happen . . .  
 ex : how likely the dice shows a "crow"
- Random variable :  $X$  assign a real value to each event



# Probability review

- Random variable example

$$X_{\text{all}} = \{1, 2, \dots, 6\}$$

$$P(X = x) = \frac{1}{6} \quad (\text{fair dice})$$

$$\sum_{x \in X_{\text{all}}} P(X = x) = 1,$$

$$P(X = 6)$$

"1"	"2"	"3"
"4"	"5"	"6"

$$1 - P(X = 6) = P(X \neq 6)$$

"1"	"2"	"3"
"4"	"5"	"6"

$$P(X_{\text{all}}) = 1$$

"1"	"2"	"3"
"4"	"5"	"6"

- Note: In this course we use probability notation informally:  $P(A)$  denotes probability of event  $A$ , and  $P(X = x), f(x)$  denotes distribution of a random variable, in discrete and continuous spaces, respectively.

in continuous space

$$\int_{X_{\text{all}}} f(x) dx = 1$$

# Probability review

## Joint probability, independence

probability of both A & B happening

- Joint probability of two events  $P(A \cap B)$

$$P(\text{"odd"} \cap 6) = 0,$$

$$P(\text{"even"} \cap 6) = P(X=6) = \frac{1}{6}$$

$$\{6\} \subset \{2, 4, 6\}$$

- Independence of two events A and B are independent:

$$P(A \cap B) = P(A)P(B). \text{ Similarly, two random variables}$$

X, Y are independent, if  $\forall x, y$  they can take,  $P(X=x, Y=y) = P(X=x)P(Y=y)$

"1"	"2"	"3"
"4"	"5"	"6"

"1"	"2"	"3"
"4"	"5"	"6"

"1"	"2"	"3"
"4"	"5"	"6"

# Probability review

## conditional distribution

- Conditional distribution : probability distribution  $\mathcal{P}$  is defined as

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \text{ assuming } P(B) \neq 0.$$

$$P(X=6 | \text{"even"}) = \frac{1}{3}$$

- Conditional independence : two events  $A_1, A_2$  are conditionally independent given event  $B$  if
 
$$P(A_1 \cap A_2 | B) = P(A_1 | B) P(A_2 | B)$$

- Product rule comes from definition of conditional distribution

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

"1"	"2"	"3"
"4"	"5"	"6"
"1"	"2"	"3"
"4"	"5"	"6"

# Probability review

## Bayes rule

- Recall product rule  $P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$

- Bayes rule  $P(B|A) = \frac{P(A|B)P(B)}{P(A)}$ , whenever  $P(A) \neq 0$

more generally

$$P(B|A_1, A_2, \dots, A_d) = \frac{P(A_1 \cap A_2 \cap \dots \cap A_d | B) P(B)}{P(A_1 \cap A_2 \cap \dots \cap A_d)}$$

# EPFL Bayes rule for classification

$$\left\{ x^i, y^i \right\}_{i=1}^N \quad y^i \in \{1, 2, \dots, K\}$$

Probability to observe this  $x$  knowing the label  $c$

Probability to label with  $c$  knowing  $x$

Probability to label with class  $c$

$$P(y = c | x) = \frac{P(x | y = c)P(y = c)}{P(x)}$$

*(Note: In the original image, the term  $c$  in the numerator of the left-hand side is circled and labeled "class c".)*

Probability of the  $x$

$P(y)$  is called the **prior probability** of the classes  
Describes the probability to encounter the data labeled  $y$

# EPFL Bayes rule for classification

1 compute for each class  $c \in \{1, 2, \dots, K\}$ , the following

$$P(y = c | x) = \frac{P(x | y = c)P(y = c)}{P(x)} ;$$

2 take  $\hat{y} = \underset{c \in \{1, 2, \dots, K\}}{\operatorname{argmax}} P(y = c | x)$

example  $y \in \{0, 1\}$

$$P(y = 0 | x) = \frac{P(x | y = 0)P(y = 0)}{P(x)}, \text{ versus } P(y = 1 | x) = \frac{P(x | y = 1)P(y = 1)}{P(x)}$$

as both probabilities above have the same denominator, we just compare the numerators. But how do we use data to compute numerators?

# Spam detection example

## Bag of words to encode text to features

List all words encountered in the set of texts you have (say  $d$  words)

Use a Boolean feature for the presence (1) or absence (0) of each word,  $x \in \mathbb{R}^d$ , with  $x_i \in \{0,1\}$

In spam email detection, data: emails, labelled as: spam, or no-spam

Suppose your emails contain “Dear”, “Friend”, “Lunch”, “Money”.

Our bag of word in this simple model has 4 words only

# Spam detection example

Features and label

$$x^i \in \mathbb{R}^4, \quad i: \text{email } i$$

Data	"Dear" $x_1$	"Friend" $x_2$	"Lunch" $x_3$	"Money" $x_4$	class
email 1	1	1	1	0	non-spam
email 2	1	0	0	1	spam
⋮					
email N					

# Spam detection example

$$\text{Dataset } D = \{x^i, y^i\}_{i=1}^N$$

## Use of Bayes' rule for spam detection

Need to compute

$$P(x | y = \text{"spam"}) P(y = \text{"spam"})$$

$$P(x | y = \text{"non-spam"}) P(y = \text{"non-spam"})$$

$$P(y = \text{"spam"}) = \frac{\# \text{ spam emails in } D}{N}$$

$$P(y = \text{"non-spam"}) = \frac{\# \text{ non-spam emails in } D}{N}$$

say  $x = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ . How do we compute

$$P(x | y = 1), P(x | y = 0) ?$$

$$P\left(\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}, y = 1\right), P\left(\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}, y = 0\right)$$

# Naive Bayes assumption

Conditional independence of different features given the label

$$P(x | y) = \prod_{j=1}^d P(x_j | y), \quad x \in \mathbb{R}^d, \quad y \in \{1, 2, \dots, K\}$$

$$P\left(\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} | y = 1\right) = P(x_1 = 1 | y = 1) P(x_2 = 0 | y = 1) P(x_3 = 1 | y = 1) P(x_4 = 0 | y = 1)$$

Now,

$$P(x_j = 1 | y = \text{"spam"}) = \frac{\# \text{ of spam emails containing word } j}{\# \text{ of spam emails}}, \quad j = 1, \dots, d$$

$$P(x_j = 0 | y = \text{"spam"}) = \frac{\# \text{ of spam emails not containing word } j}{\# \text{ spam emails}}, \quad j = 1, \dots, d$$

# Naive Bayes classifier

$$P(y = \text{"spam"} | x) = \frac{\prod_{j=1}^d P(x_j | y = \text{"spam"}) P(y = \text{"spam"})}{P(x)}$$

$P(x)$  — no need to compute  
since this term is the  
same for all classes ..

for any classification  $y \in \{1, 2, \dots, k\}$

$$P(y = c | x) = \frac{P(x | y = c) P(y = c)}{P(x)} = \frac{\prod_{j=1}^d P(x_j | y = c) P(y = c)}{P(x)}$$

from  
Bayes rule

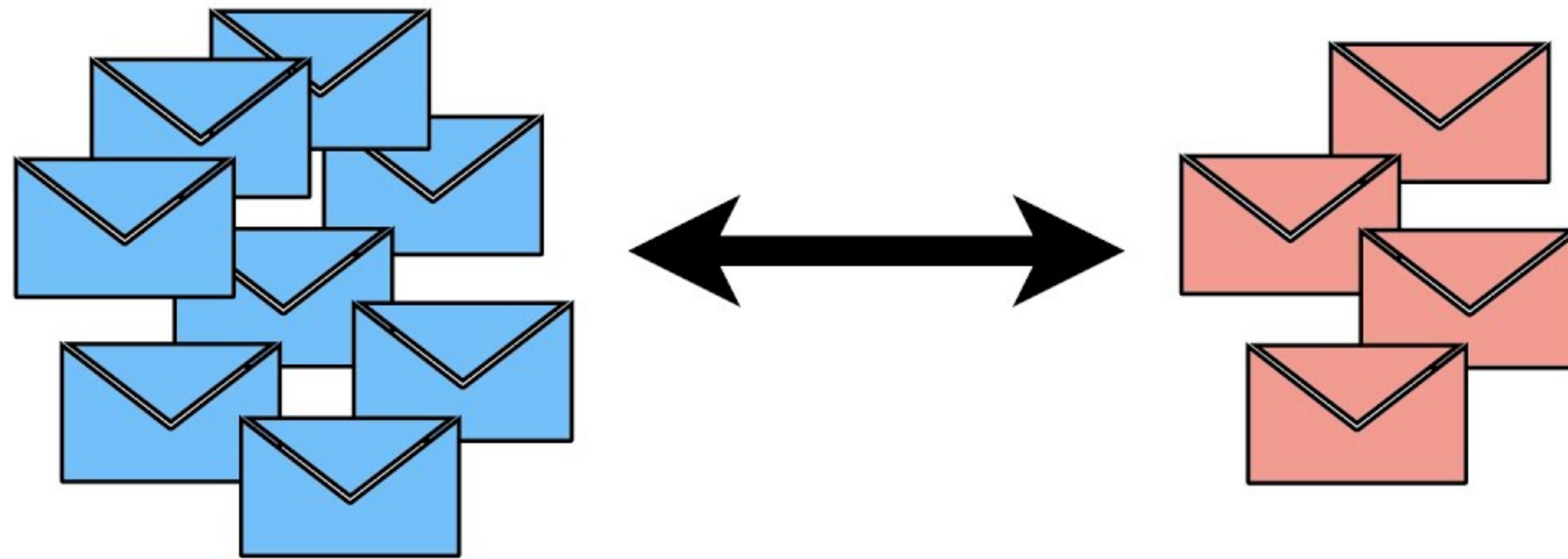
from Naive  
Bayes assumption

Lastly, we estimated  $P(y = c)$  and  $P(x_j | y = c)$  from our training data.

# Naive Bayes classifier for spam detection

## Example with worked out numbers

### Naive Bayes....



**...Clearly Explained!!!**

# Activity

- Watch: <https://www.youtube.com/watch?v=O2L2Uv9pdDA>
  - This takes 15 minutes
- After watching the video, answer the following questions
  1. How are the features defined differently than the binary scores we defined earlier?
  2. How would you evaluate the accuracy of the classifier?
  3. Are there any hyper parameters to tune for this classifier?
  4. Is there any other approach you could imagine using for this classification?
  
- Next class: discuss your answers with your two nearest neighbours
- Pick one person to represent your answers to class

# **Naive Bayes Classifier**

## **Continuous-valued features**

# Bayes rule

## Continuous-valued features

$$P(y = c | x) = \frac{f_{x|y=c}(x)P(y = c)}{f_x(x)} \quad \left( \begin{array}{l} \text{what change of} \\ \frac{P(x|y=c)P(y=c)}{p(x)} \end{array} \right)$$

$f_{x|y=c}(x)$  Probability density function: probability distribution for a continuous random variable

Note:  $f_{x|y=c}(x)$  is not the probability of random variable taking value  $x$ . This probability is generally zero for any  $x \in \mathbb{R}^d$ , namely, for continuous random variables.

$$P(a \leq x \leq b) = \int_a^b f_{x|y=c}(x) dx, \quad x \in \mathbb{R}$$

$$P(x \in S) = \int_S f_{x|y=c}(x) dx, \quad x \in \mathbb{R}^d, \quad S \subseteq \mathbb{R}^d$$

# Naive Bayes assumption

## Continuous features

$$P(y = c | x) = \frac{f_{x|y=c}(x)P(y = c)}{f_x(x)}$$

Naive assumption

$$= \frac{\prod_{j=1}^d f_{x_j|y=c}(x_j) P(y = c)}{f_x(x)}$$

Naive Bayes's assumption: conditioned on class  $c$ , the features in different dimension have independent distributions

$f_{x|y=c}(x)$ : conditional density of  $x \in \mathbb{R}^d$ , given class  $c$ .

# Gaussian Naive Bayes

Assumes that the probability density function for each feature follows a gaussian distribution

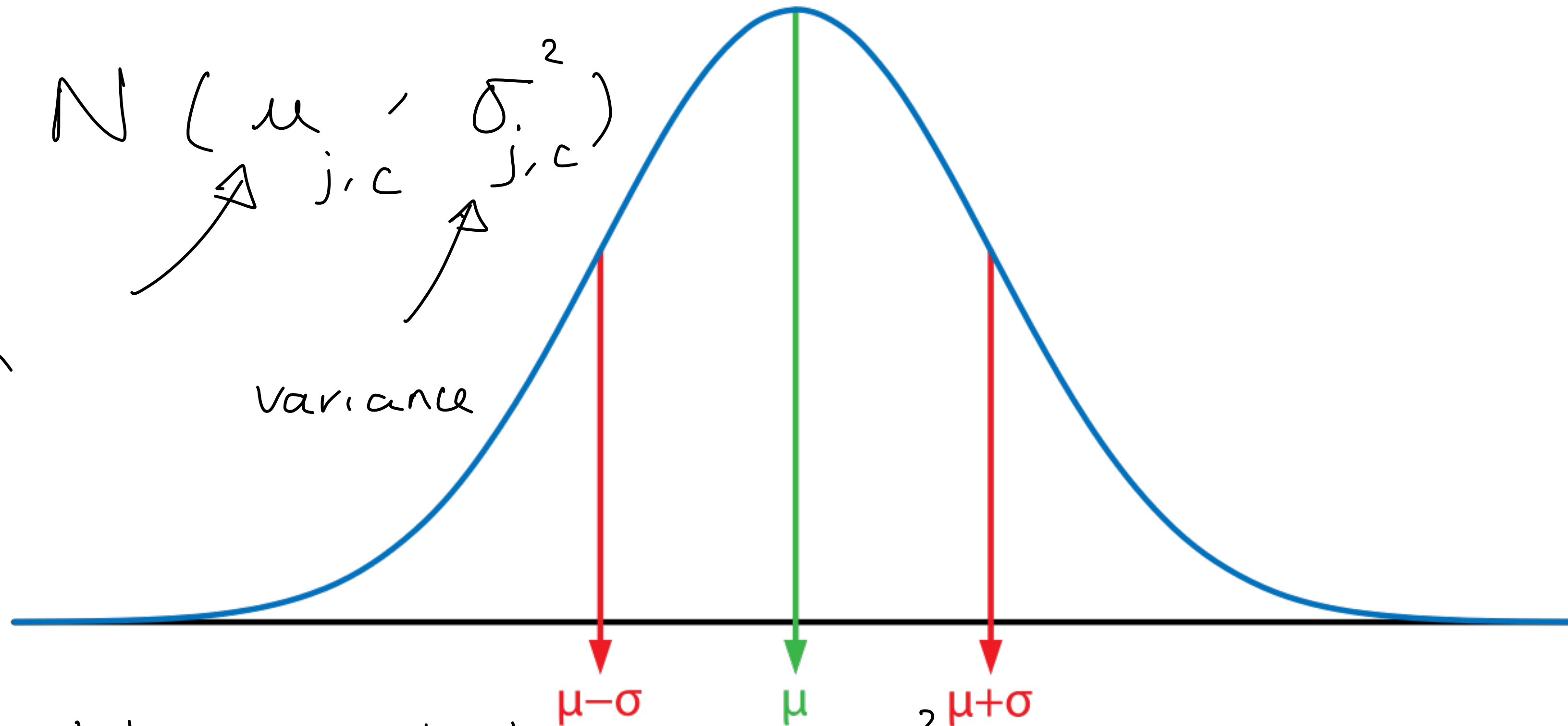
$f_{x_j|c}(x_j)$ : conditional density of feature  $j$  given class  $c$  has as

a Gaussian distribution

$$f_{x_j|c}(x_j) = \mathcal{N}(\mu_{j,c}, \sigma_{j,c}^2)$$

Mean

Variance



use training data to estimate  $\mu_{j,c}, \sigma_{j,c}^2$

# Gaussian Naive Bayes

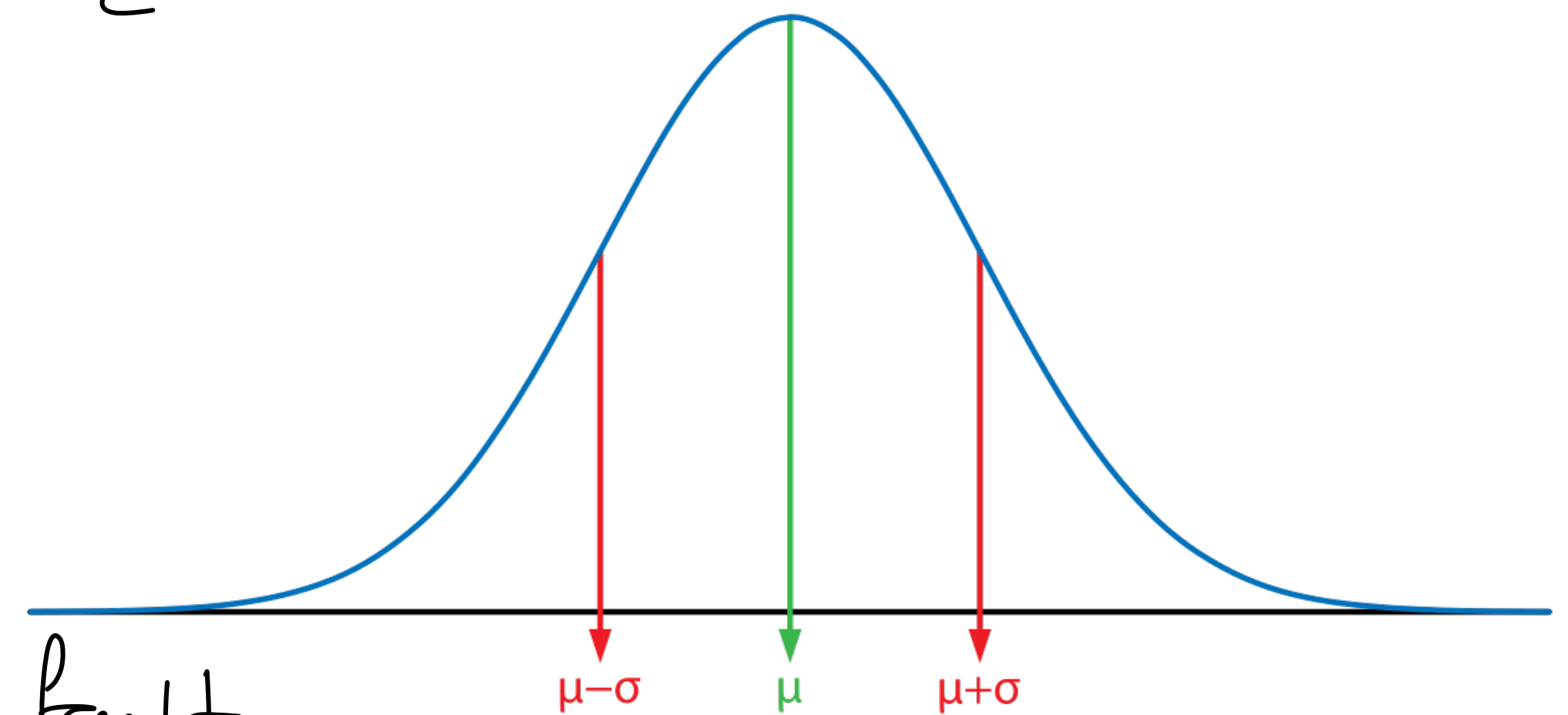
Data set  $\{x^i, y^i\}_{i=1}^n$ ,  $x^i \in \mathbb{R}^d$

Estimating the conditional Gaussian distribution for each feature using sample data

Recall :  $\mu_{j,c} \approx \frac{1}{|I_c|} \sum_{i \in I_c} x_j^i$

$I_c$  : indices of data that have label class  $c$ .

example  $x_1$  weight, class: fault, no fault  
 $x_2$  height, 1, 2  
 $x_3$  width



$\mu_{1,1}$  is average weight of faulty objects  
 $\mu_{2,1}$  is average height of faulty object.

# Gaussian Naive Bayes

## Classifier prediction

$$P(y = c | x) = \frac{f_{x|y=c}(x)P(y = c)}{f_x(x)} = \frac{\prod_{j=1}^d f_{x_j|y=c}(x_j) P(y = c)}{f_x(x)}$$

these are now  
Gaussian with  
mean & variance  
estimated from data

Note:  $f_{x|y=c}(x)$  in statistics is also referred to as likelihood function, when considered as a function of class  $c$ , given a fixed  $x$ .

Suppose we get new  $x^{\text{test}}$ , you want to predict  $y^{\text{test}} \in \{1, 2, \dots, K\}$ .

you compare the numerators above for each class  $c$ .

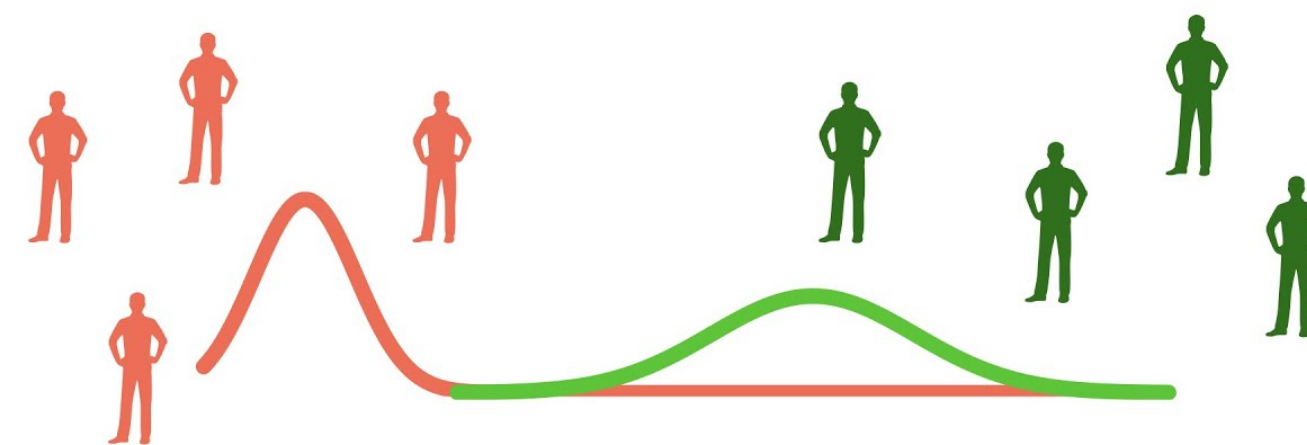
$$\prod_{j=1}^d f_{x_j|y=c}(x_j^{\text{test}}) P(y = c), \quad f_{x_j|y=c}(x_j^{\text{test}}) = \frac{1}{\sqrt{2\pi} \sigma_{j,c}} e^{-\frac{(x_j^{\text{test}} - \mu_{j,c})^2}{2(\sigma_{j,c})^2}}$$

# Gaussian Naive Bayes

Naive Bayes assumption and final classifier

$$P(y = c | x) = \frac{f_{x_1|y=c}(x_1)f_{x_2|y=c}(x_2)\cdots f_{x_d|y=c}(x_d)P(y = c)}{f_x(x)}$$

**Guassian Naive Bayes....**



**...Clearly Explained!!!**

# Naive Bayes

## Summary

- Probabilistic classifier
  - Features can be finite- or continuous-valued
- Uses Bayes rule and empirical probabilities computed based on data
- Assumes different features are independent given the class label
  - In practice hard to verify the assumption but works well on certain problems

- Recurrent neural network for time-series prediction
  - Back propagation for computing gradients
  - Gradient vanishing/exploding -> motivates other approaches (LSTM, transformers)
  
- Back to classification
  - Naive Bayes probabilistic classifier
  - Conditional probability distribution and Naive Bases assumption
  
- Your tasks this week
  - Problem set 5
  - Work on python homework and exercises