

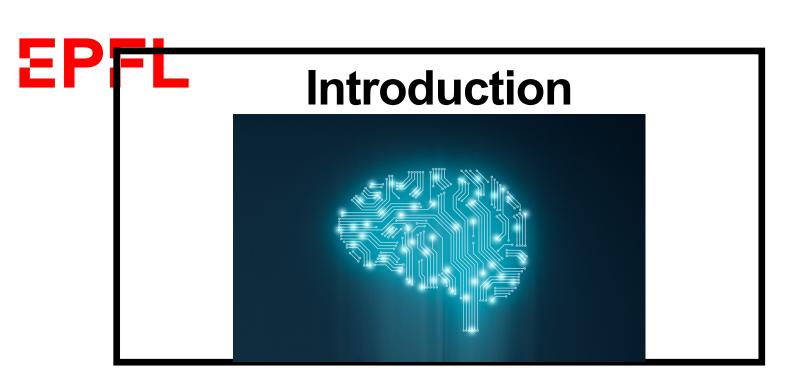
Lecture 05
14.10.2024

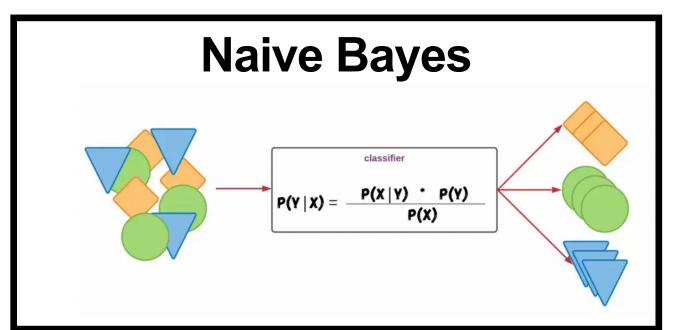
Naive Bayes Classifier

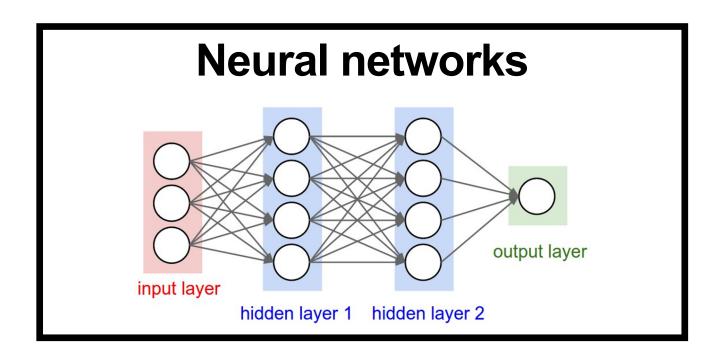


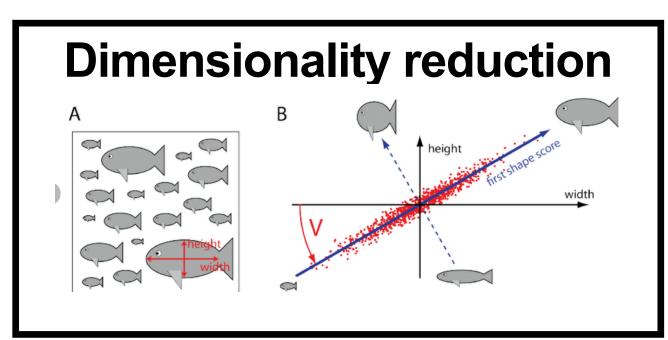
#### Outline

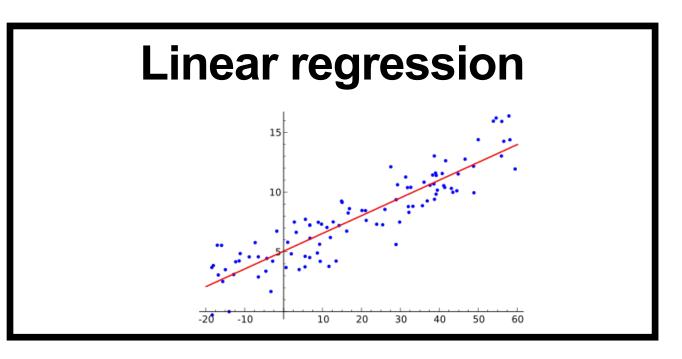
- Review of probability
  - Conditional distribution
  - Independence
  - Bayes rule
- Naive Bayes classifier
  - Finite-valued features
  - Continuous-valued features (Gaussian Naive Bayes classifier)
- Announcements:
  - Exercise hours: quiz 1

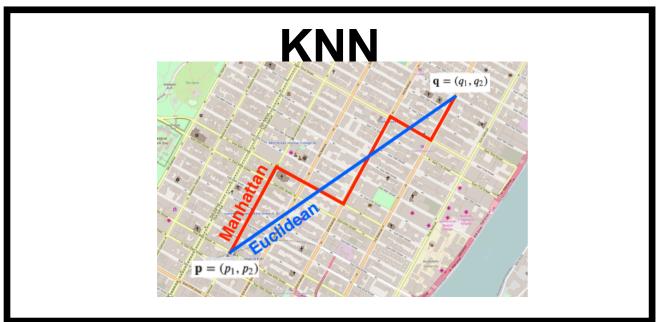


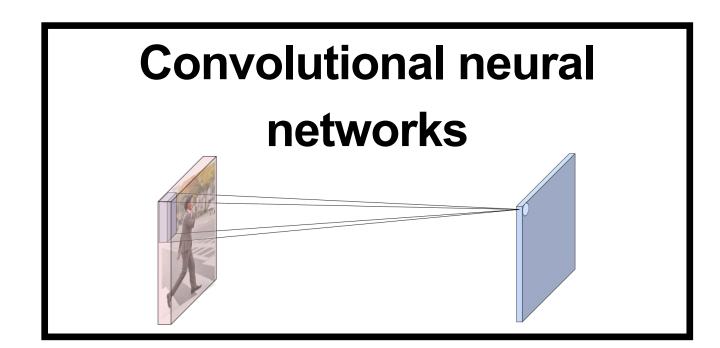


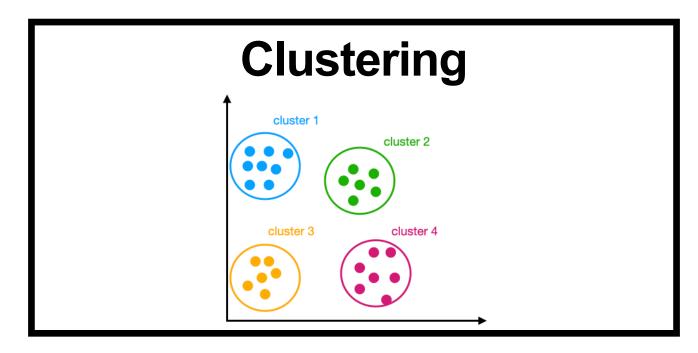


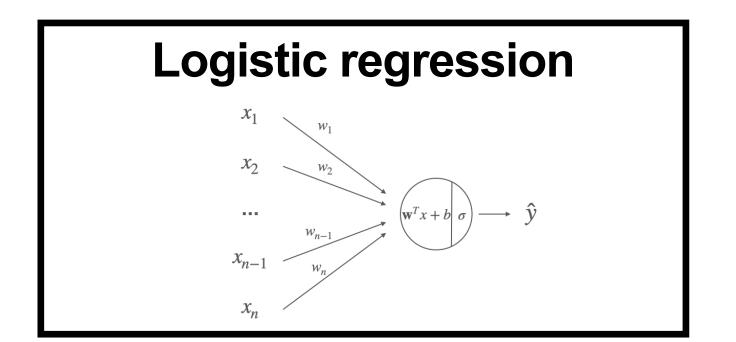


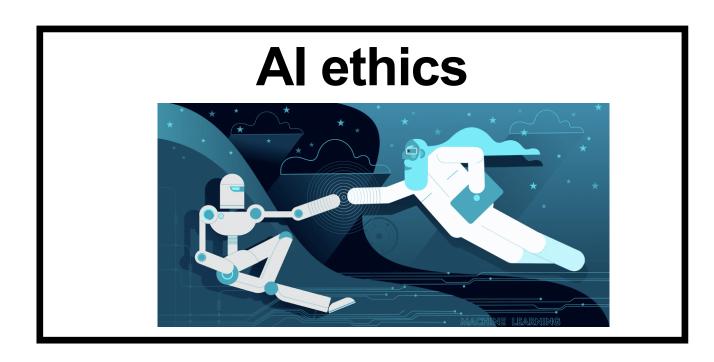


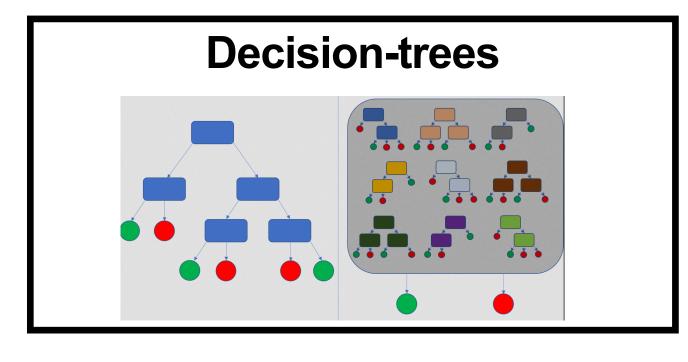


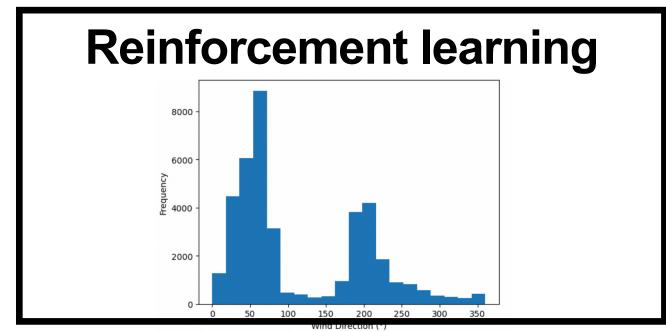














# Naive Bayes Classifier



#### Supervised learning - probabilistic classifier

Recall probabilistic interpretation of logistic regression

We will now discuss another approach to probabilistic classification



# Probability review

#### Probability, random variable

- probability of an event A, p, (A): how likely fire
  for event A to happen, ex: event A: die rolls a 6"
- Random variable X: a variable that takes different values with certain probabiliher. X: set of all values X

can take,  $\sum_{x \in X} P_v (X = x) = 1$ .  $x \in X_{all}$   $P_{vob} \cdot of X$  taking value x

X: the number the dir volus at  $P(X=x) = \frac{1}{6}$ 

dice example

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

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

	A
,	( - Pr(A)

"1"	"2"	"3"	( -
"4"	"5"	"6"	

"1"	"2"	"3"	$\sum_{x=1}^{6} P_{\nu} (X=x)$
"4"	"5"	"6"	= 7



## Probability review

#### Joint probability, independence



Marginalisation rule

$$\leq P(A, \chi = x) = P(A)$$

P(even) = /2

we call two random variables X, y indep if I values they can take P(X=x, Y=y)=P(X=x)P(Y=y)

P (	X=6)
-----	------

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

$$(1)^{2}$$
  $(2)^{2}$   $(3)^{2}$   $(4)^{2}$   $(5)^{2}$   $(6)^{2}$ 

# **Probability review**

• Conditional distribution  $P(A|B) = \frac{P(A|B)}{P(B)}$ 

#### conditional distribution

B: even

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

						1	/	
prob.	of	event	$\bigwedge$	gricen	that	event	Bhas	000000

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

P(6 | even) = 
$$\frac{1}{3}$$
 • Marginalization  $\leq P(x|B) = \frac{P(x,B)}{x \in X_{\alpha(i)}} = \frac{1}{2} \frac{\sum P(x,B)}{P(B)} = \frac{1}{2} \frac{\sum P(x,B)}{P(B)}$ 

Conditional independence

events A., Az are conditionally independent given event B:

• Product rule 
$$from * = P(A|B)P(B) = P(A|B) P(A)$$

#### **EPFL**

# Probability review Bayes rule

- Recall product rule P(A,B)=P(A|B)P(B) = P(B|A)P(A)

Bayes rule  $P(B|A) = \frac{P(A|B)P(B)}{P(A)}$ 

consider now 
$$P(B|A_1,A_2) = \frac{P(A_1,A_2|B)P(B)}{P(A_1,A_2)}$$

··· can be extended to {A; }! events.

Bayes rule for classification  $\{x', \eta'\}^{N}$   $\{y' \in \{1,2,...,K\}\}$ 

Probability to observe this  $\boldsymbol{x}$  knowing the label  $\boldsymbol{c}$ 



Probability to label with class c

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

Probability of the x

P(y) is called the prior probability

Describes the probability to encounter the data labeled y

# EPFL Bayes rule for classification

$$P(y=c|x) = \frac{P(x|y=c)P(y=c)}{P(x)}, \quad y = \arg\max_{c \in \{1,2,..,k\}} P(y=c|x)$$

$$e \times \operatorname{ample} y \in \{0,1\}$$

$$P(y=0|x) > P(y=1|x)$$

$$|x| = \arg\max_{c \in \{1,2,..,k\}} P(y=c|x)$$

$$\frac{P(x|y=0)P(y=0)}{P(x)} > \frac{P(x|y=1)P(y=1)}{P(x)} \iff P(x|y=0)P(y=0) > P(x|y=1)P(y=1)$$



#### Spam detection example

#### Bag of words to encode text to features

List all words encountered in the set of texts you have

Use a Boolean feature for the presence (1) or absence (0) of each word

oc; E {0,1} absence or presence of word;

class )

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

Suppose	J 001	emails	contein	Dear ", F	riend.	n. Money.
Data		Dear"	menel"	"Lunch"	"Money"	C 100 59
ema.					•	19n-5pan
email.	2					s pam



### Spam detection example

\ x', y'} = D

Use of Base's rule for spam detection

Need to calculate 
$$P(x|y="spam")P(y="spam")$$
 and  $P(x|y="no.spam")P(y="no.spam")$ .

$$P((1,0,0,1) | A = 1)$$
,  $P((1,0,0,1) | A = 0)$ 

#### **EPFL**

### Naive Bayes assumption

Conditional independence of different features given the label  $x \in \mathbb{R}^{d}$ 

# Naive Bayes classifier

 $x \in \{0,1\}^{c} = \{0,1\} \times \{0,1\} \times$ · · · × {0,1}

continued example . (i) 
$$P(x) = 1$$
  $P(y=1)$   $P(y=1)$  cl: # of features (here words)

to determine p(y=1|x) > p(y=0|x) we don't need to compate elenominate p(x) (same in both sider of the negucially)

Need to compare...

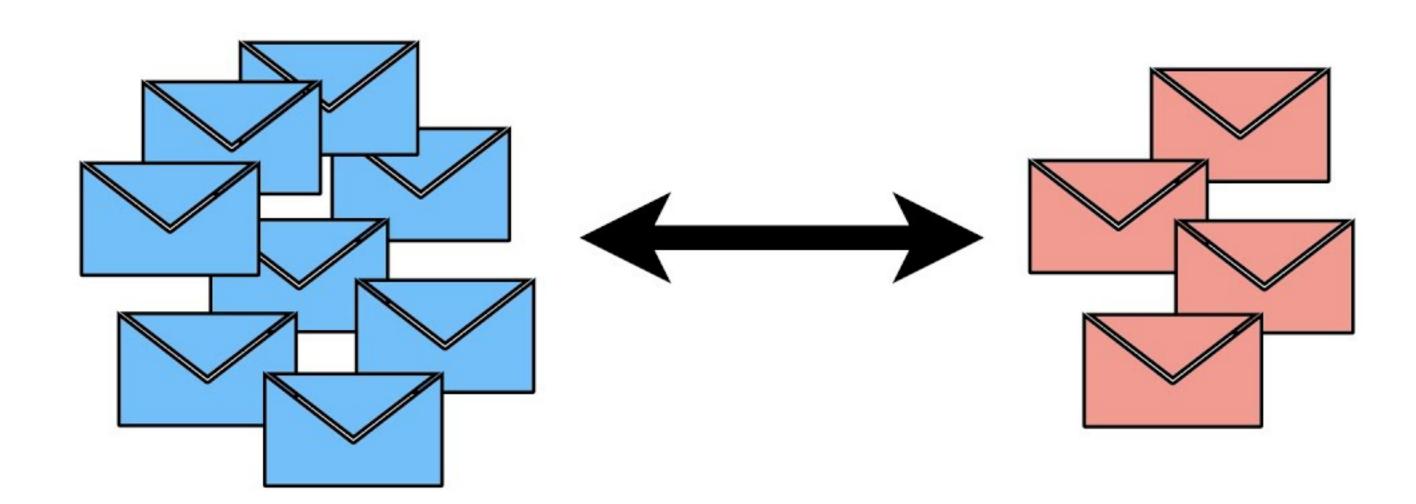
$$\frac{d}{dt} p(x_j | y=1) p(y=1) > \frac{d}{dt} p(x_j | y=0) p(y=0)$$

$$j=1$$



# Naive Bayes classifier for spam detection Example with worked out numbers

# Naive Bayes...



... Clearly Explained!!!

#### **EPFL**

# Activity

- Watch: <a href="https://www.youtube.com/watch?v=O2L2Uv9pdDA">https://www.youtube.com/watch?v=O2L2Uv9pdDA</a>
  - 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

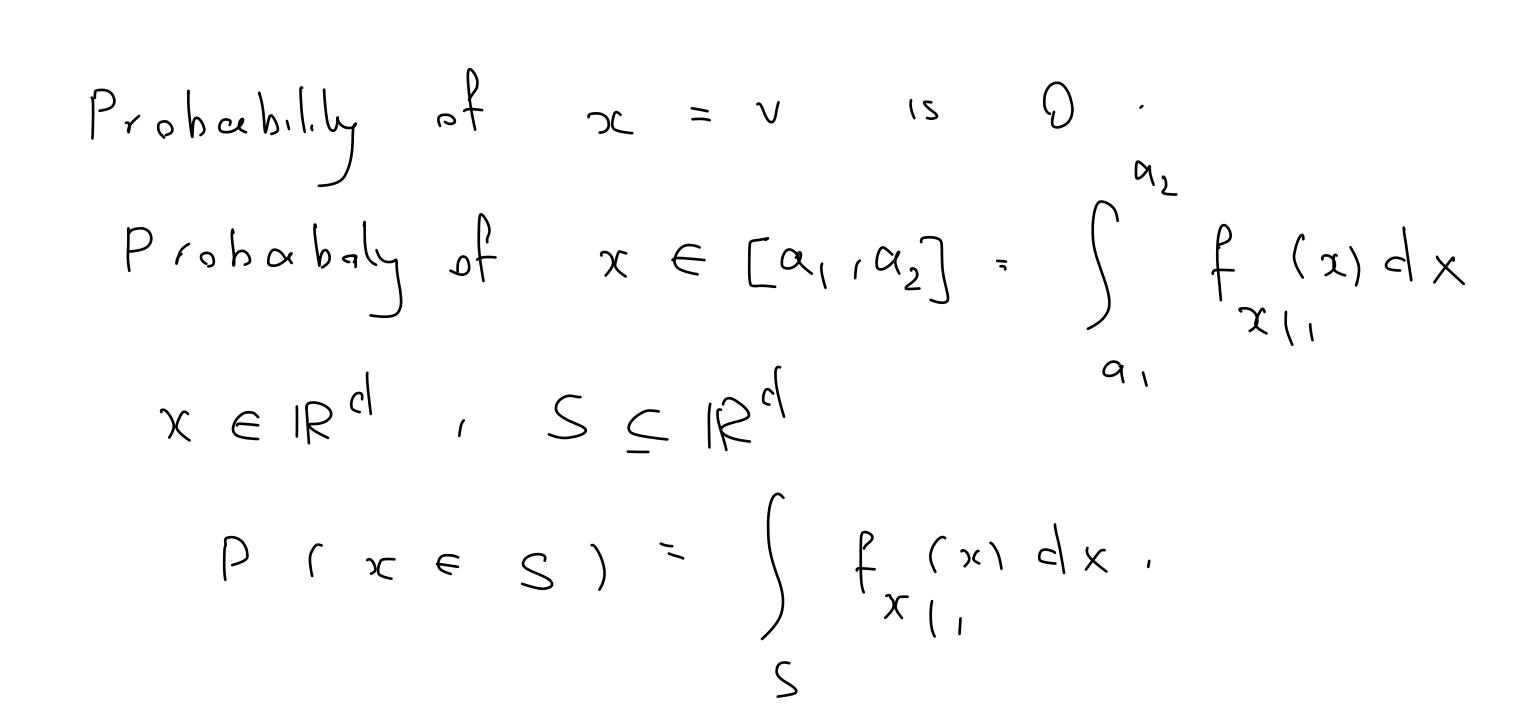


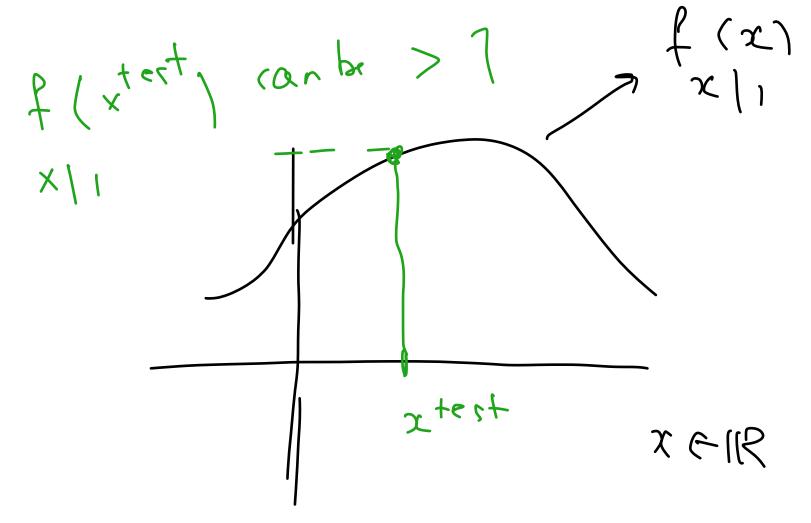
#### Bayes rule

#### Continuous-valued features

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

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







### Naive Bayes assumption

#### Continuous features

$$P(y=c \mid x) = \frac{f_{x\mid y=c}(x)P(y=c)}{f_{x}(x)} = \underbrace{\frac{\int_{z=c}^{x} f(x)}{\int_{z=c}^{x} f(x)} f(x)}_{C}$$

$$\frac{d}{\int_{j=1}^{\infty} f(x_j)} \int_{y=c}^{\infty} \int_{z=c}^{\infty} \int_{$$

Naive Bayer assumption: says given class c, the features are conclinenally independent.

$$x \in 18$$

$$x \in \mathbb{R}^d$$

$$f(x) : conclined density of x given clair c$$

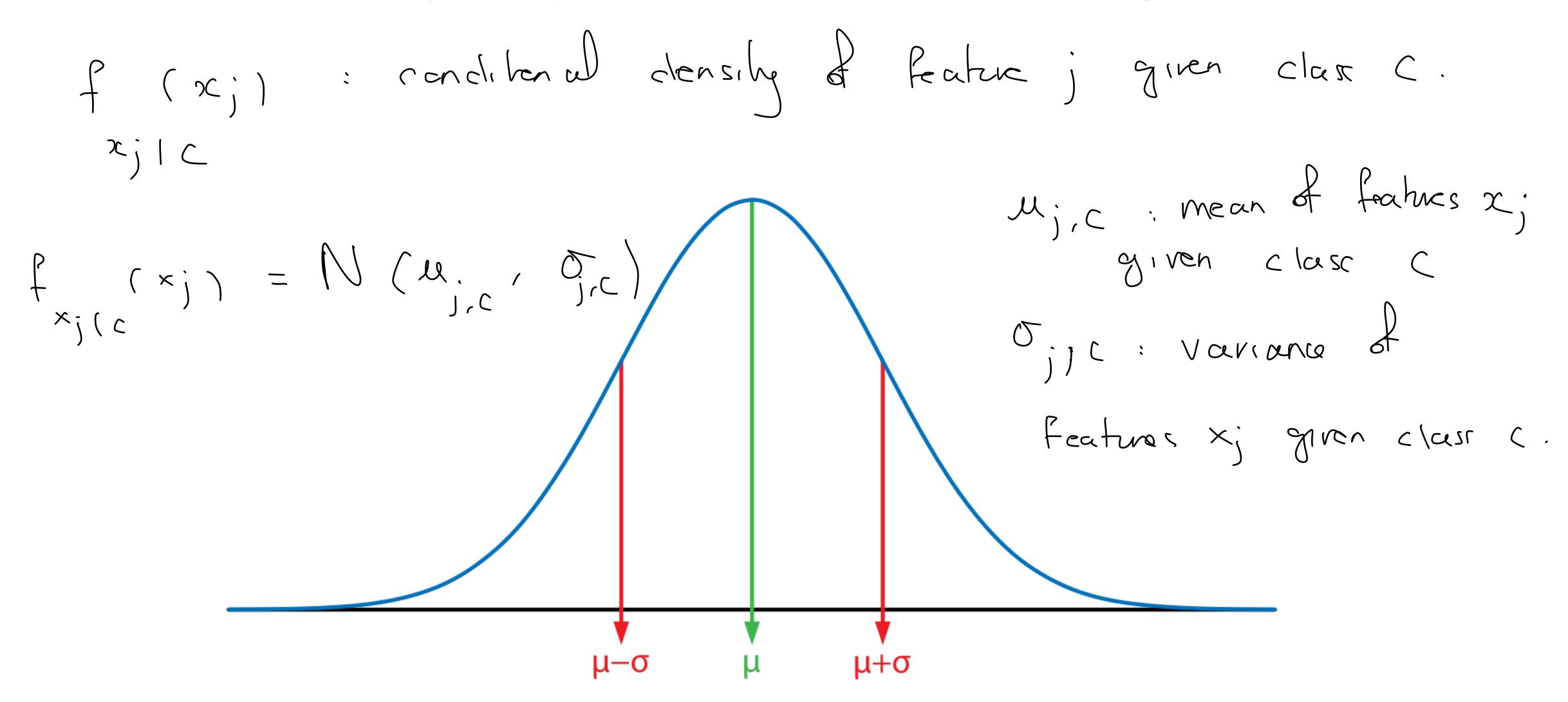
$$x|y=c$$

$$f(x) = \frac{d}{d}f(x)$$

$$= \frac{d}{d}f(x)$$

$$= \frac{d}{d}f(x)$$

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





$$\{x', y'\}$$
 $\{x' \in \mathbb{R}^d$ 

Estimating the conditional Gaussian distribution for each feature using sample data

Recall

$$u_{j,c} = \frac{1}{|I_c|} \sum_{i \in I_c} x_j$$
 $I_c : indices of class corresponds to class c.

 $V_c = \frac{\sum_{i \in I_c} (x_j' - \mu_{j,c})^2}{|I_c| - 1}$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$ 
 $V_c = \frac{1}{|I_c|} \sum_{i \in I_c} (x_i' - \mu_{j,c})^2$$ 

#### Classifier prediction

$$P(y=c \mid x) = \frac{f_{x\mid y=c}(x)P(y=c)}{f_{x}(x)} = \frac{\int_{x_{1}}^{d} f_{x_{1}}(x_{j}) P(y=c)}{f_{x_{1}}(x_{j})}$$
We used claba to approximate the density function.

$$f_{x_{j}\mid y=c}(x_{j}) = N(u_{j}, c \mid \sigma_{j,c})$$
8 to also cleterine  $P(y=c) = V(x_{j}, c \mid \sigma_{j,c})$ 

$$f_{x_{j}\mid y=c}(x_{j}) = N(u_{j}, c \mid \sigma_{j,c})$$
8 to also cleterine  $P(y=c) = V(x_{j}, c \mid \sigma_{j,c})$ 

$$f_{x_{j}\mid y=c}(x_{j}) = N(x_{j}) P(y=c)$$
8 to also cleterine  $P(y=c) = V(x_{j}) P(y=c)$ 
8 gives label  $V(x_{j}) = V(x_{j}) P(y=c)$ 
8 gives label  $V(x_{j}) = V(x_{j}) P(y=c)$ 

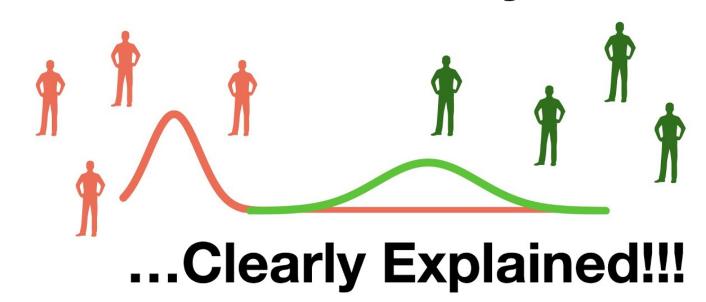


Naive Bayes assumption and final classifier

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

you can see a worked out example in the video below.

#### Guassian Naive Bayes....





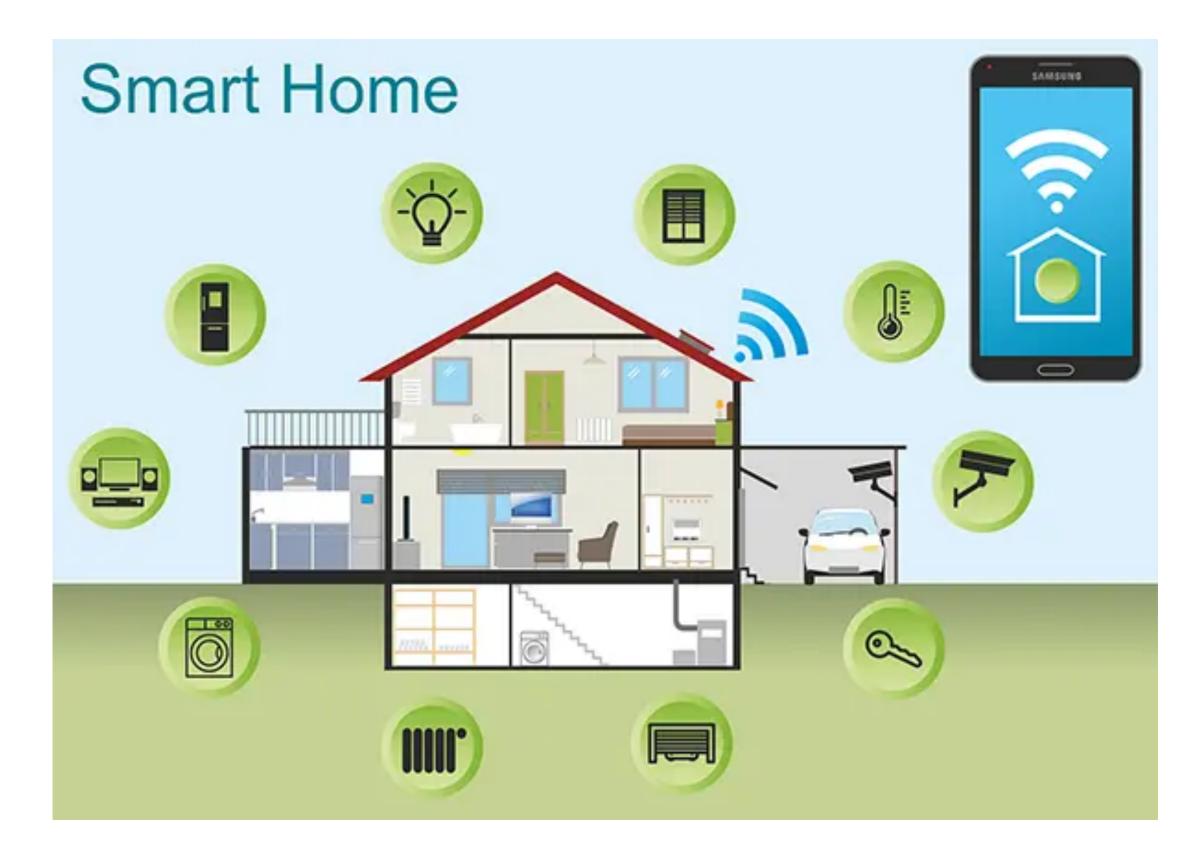
# Naive Bayes Summary

- 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 → Maire assumpted
  - In practice hard to verify the assumption but works well on certain problems



# Naive Bayes classifier Application in python exercises

- Dataset: Smart Home Device Efficiency
- Goals: Evaluate the efficiency of a smart home device



Quz an approach to AI & uses closed to ( Supervised leare) / xi, yi) learn predicter pattern (unsupervised leave) Super vised leavy yEIR regression classification y ∈ (1,2,..., K) { logistic regression

D= { x; y; } N

model fu, b (x)

yi ∈ { 1,2, ..., κ}