

# Lecture: Machine Learning for Data Science

Winter semester 2021/22

Lectures 5: Classification (Naïve Bayes classifiers)

Prof. Dr. Eirini Ntoutsi

# Recap

- Overfitting in general
- Overfitting in decision trees
- KNNs

# Happiness check

- Any feedback on the course (lectures/tutorials)?
- How happy are you with the material/pace?
- How happy are you with the hybrid format
- Any wishes (e.g., further readings, programming exercises)?
- Other issues you would like to discuss



# Outline

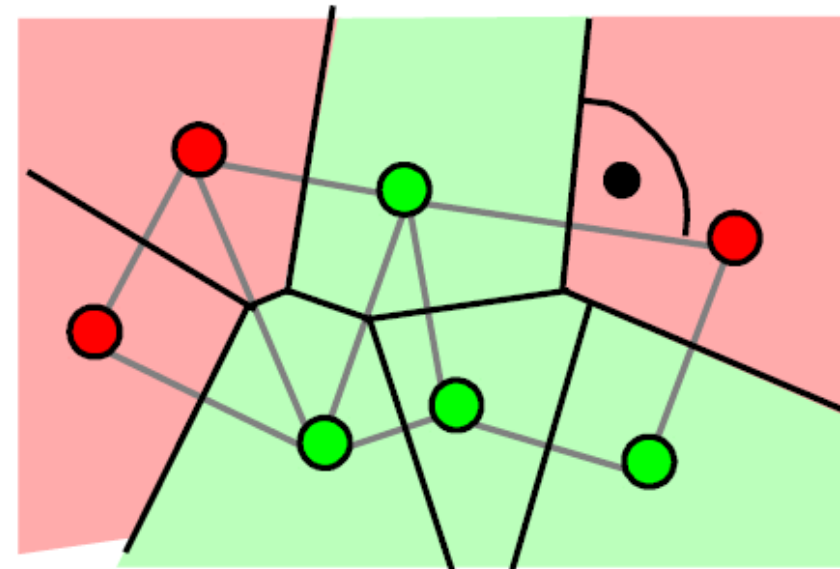
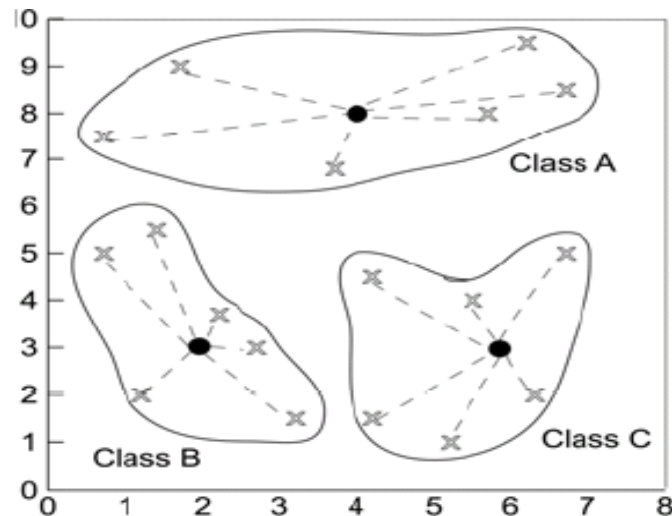
- A few more words on KNNs
- Generative vs Discriminative models
- Bayesian Classifiers
- Naïve Bayes classifiers
- Laplace correction
- Bayesian Belief Networks
- Things you should know from this lecture & reading material

# KNN classifiers: Inductive bias

- **Inductive bias**: the set of assumptions that, together with the training data, deductively justify the classifications assigned by the learner to future instances.
- **Inductive bias of KNN classifiers**: What is the policy by which a KNN classifier generalizes from observed training examples to classify unseen instances?
  - Similar instances have similar class labels.
  - All attributes/dimensions contribute equally.
  - ...

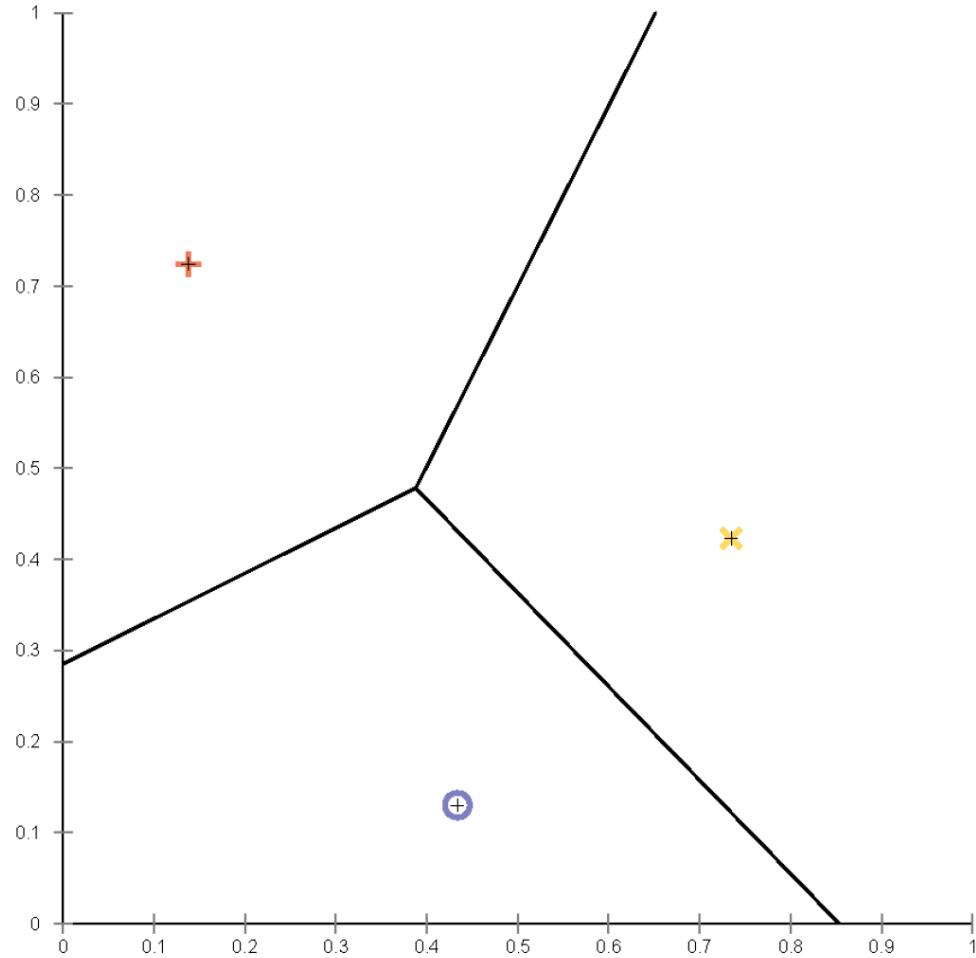
# KNN classifiers: Decision boundaries

- Nearest-neighbor classifiers can produce arbitrarily shaped decision boundaries
- We can visualize **class regions** by **Voronoi** cells
  - Each cell is a region, a set of points that are closer to the training example than to any other training example in the dataset



# KNN classifiers: Decision boundaries

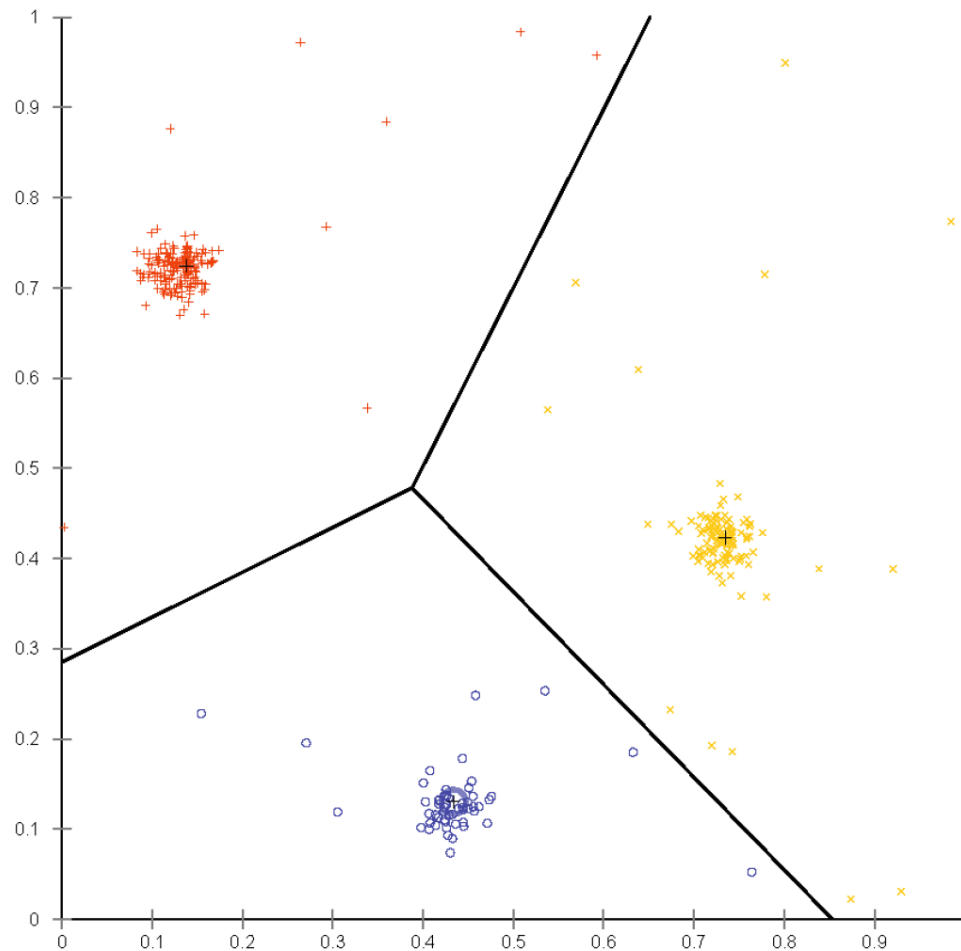
- Small training set → Simple decision boundaries



Slide by Arthur Zimek @SDU

# KNN classifiers: Decision boundaries

- Small training set → Simple decision boundaries

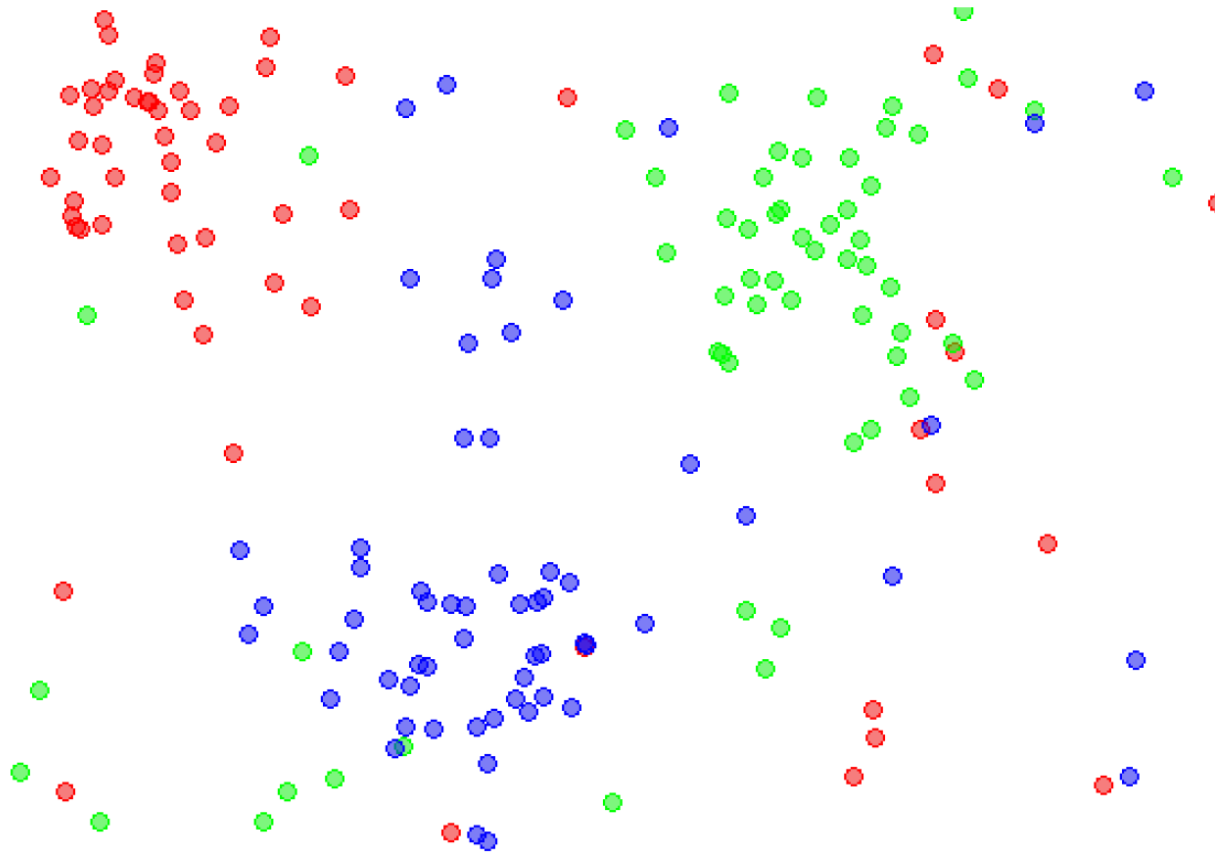


*Slide by Arthur Zimek @SDU*



# KNN classifiers: Decision boundaries

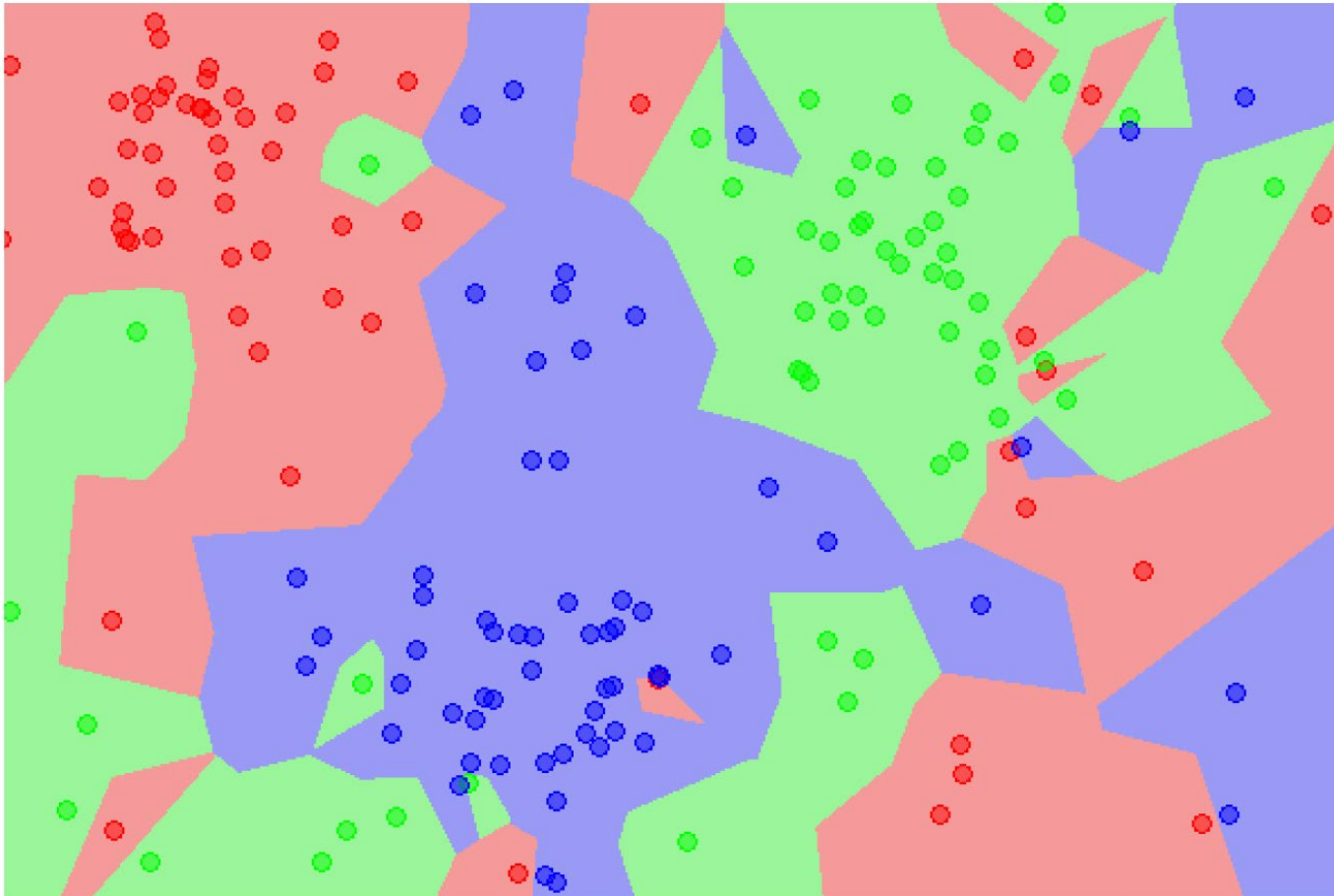
- Large training set → Potentially complex decision boundaries



*Slide by Arthur Zimek @SDU*

# KNN classifiers: Decision boundaries

- Large training set → Potentially complex decision boundaries ( $k=1$ )

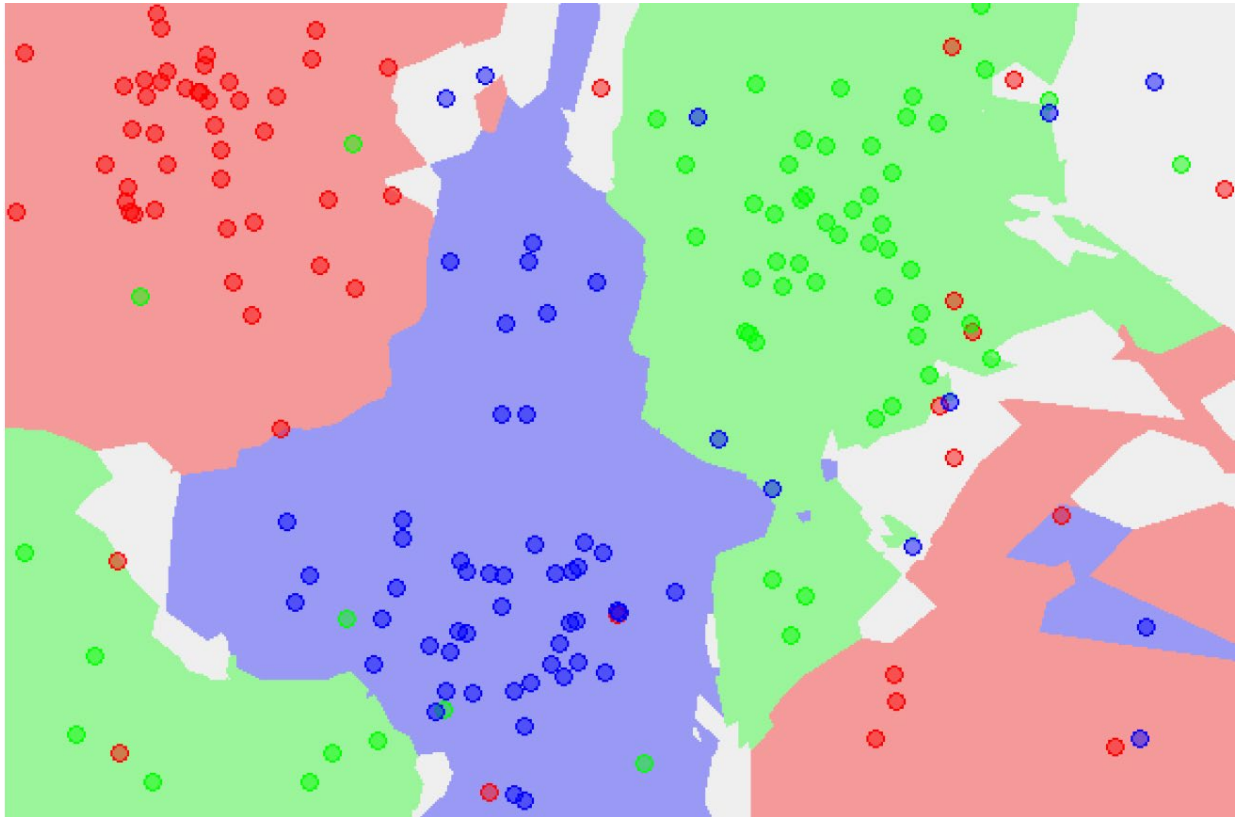


- Non-linear decision boundary
- Reflects the classes well

*Slide by Arthur Zimek @SDU*

# KNN classifiers: Decision boundaries

- Large training set → Potentially complex decision boundaries ( $k=5$ )



*Slide by Arthur Zimek @SDU*

# Outline

- A few more words on KNNs
- Generative vs Discriminative models
- Bayesian Classifiers
- Naïve Bayes classifier
- Laplace correction
- Bayesian Belief Networks
- Things you should know from this lecture & reading material

# Generative vs Discriminative models

- Thus far, we assumed that there is an underlying distribution  $P(X,Y)$  that generates our population/data but we don't have access to this distribution
  - Instead, we have access to training data  $D=\{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$  coming from this distribution
- **Discriminative models**
  - Try to find a mapping/function/hypothesis  $h(): X \rightarrow Y$  that **separates** the different classes
  - New instances are classified using  $h()$
- **Generative models**
  - Try to build a model for each individual class
    - It learns  $P(X|Y)$  and  $P(Y)$
  - New instances are tested against different models and the most likely model is the class (using Bayes' rule)  $\rightarrow P(Y|X)$

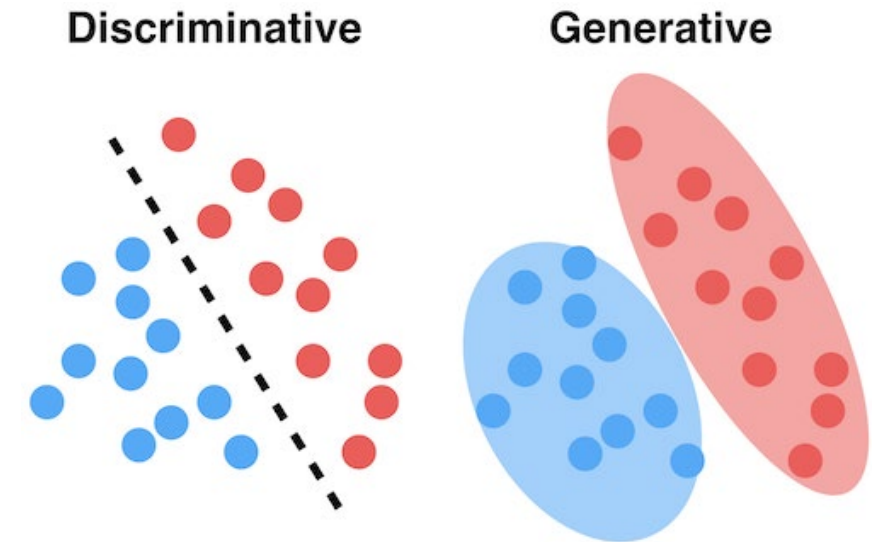


Image: [link](#)

# Generative vs Discriminative models

- The discriminative approach has the advantage of directly optimizing the predictive performance instead of learning the underlying distribution
  - *“When solving a given problem, try to avoid a more general problem as an intermediate step.” Vapnik*
- Usually it is harder to learn the underlying distribution than to learn an accurate predictor

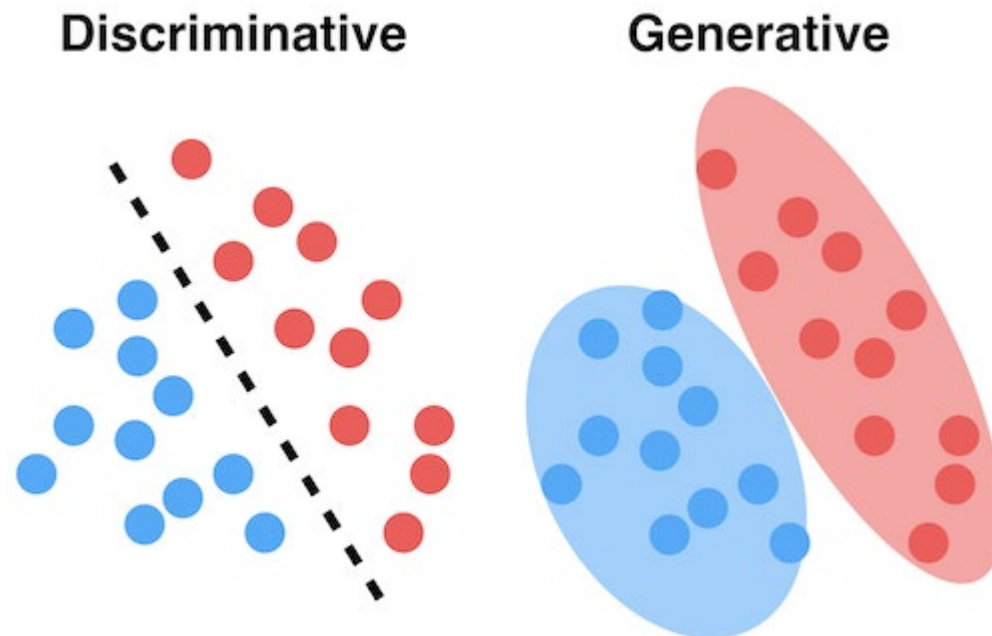


Image: [link](#)

# Outline

- Generative vs Discriminative models
- Bayesian Classifiers
- Naïve Bayes classifiers
- Laplace correction
- Bayesian Belief Networks
- Things you should know from this lecture & reading material

# Bayesian classifiers

- A **probabilistic framework** for solving classification problems
- Predict **class membership probabilities** for an instance
  - The class of an instance is the most likely class for the instance (**Maximum Likelihood classification**)
- Based on Bayes' rule
- Bayesian classifiers
  - **Naïve Bayes classifiers**
    - Assume class-conditional independence among attributes
  - **Bayesian Belief networks**
    - Graphical models
    - Model dependencies among attributes
- A popular method for e.g.: text classification, sentiment analysis



# Bayes' theorem

- The probability of an event  $C$  given an observation  $A$ :

$$P(C | A) = \frac{P(C)P(A | C)}{P(A)}$$

- $P(C)$ : prior
- $P(A|C)$ : likelihood
- $P(A)$ : evidence
- $P(C/A)$ : posterior

$$\textit{posterior} = \frac{\textit{prior} \times \textit{likelihood}}{\textit{evidence}}$$

# Bayes' theorem

- The probability of an event  $C$  given an observation  $A$ :

$$\begin{array}{ccc} & \nearrow & \\ \text{posterior} & P(C | A) = \frac{P(C)P(A | C)}{P(A)} & \nwarrow \text{likelihood} \\ & \nwarrow \text{prior} & \end{array}$$

- Example: Given that
  - A doctor knows that meningitis causes stiff neck 50% of the time
  - Prior probability of any patient having meningitis is  $P(M)=1/50,000$
  - Prior probability of any patient having stiff neck is  $P(S)=1/20$
- If a patient has stiff neck, what's the probability he/she has meningitis? **?**

$$P(M | S) = ?$$

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = ? \qquad \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

# Bayes classification

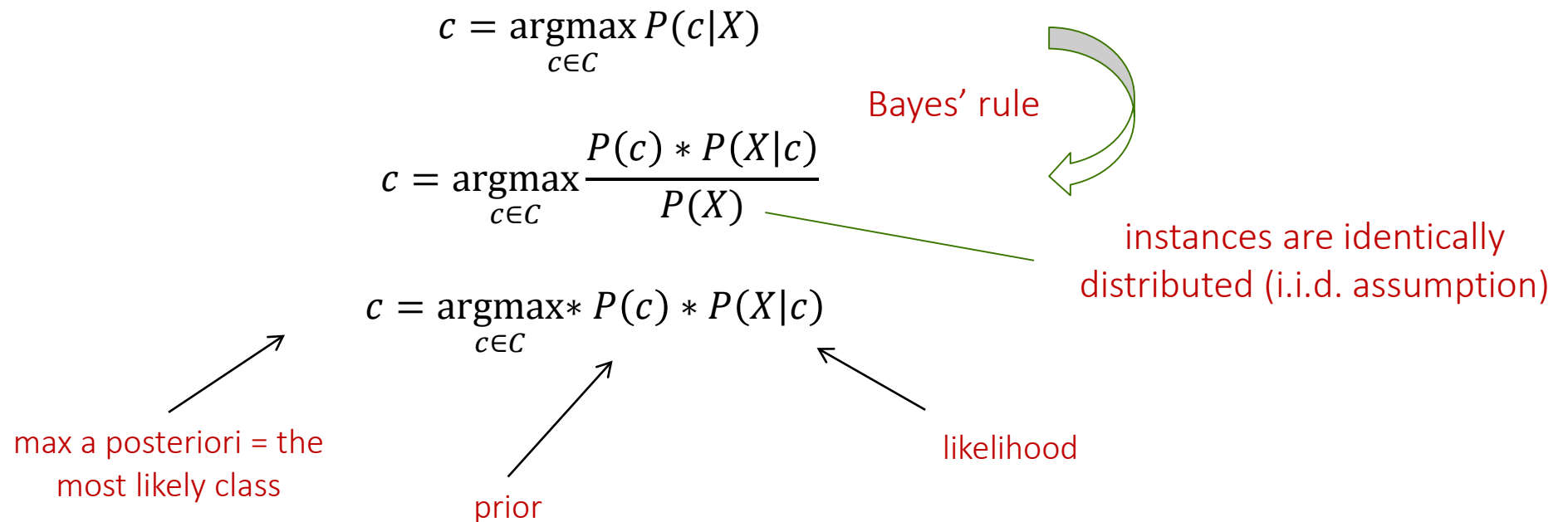
- Let  $C=\{c_1, c_2, \dots, c_k\}$  be the class attribute.
- Let  $X=(X_1, X_2, X_3, \dots, X_d)$  be a  $d$ -dimensional instance.
- Classification problem: What is the probability of a class  $c \in C$  given instance  $X$ ?



- The event  $C$  to be predicted is the class of the instance
  - The observation is the instance values  $X$
- The class of the instance is the class value with the highest probability:  $\operatorname{argmax}_{c \in C} P(c|X)$ 
  - $P(c_1/X)$ : posterior for  $c_1$
  - $P(c_2/X)$  : posterior for  $c_2$
  - ...
  - $P(c_k/X)$  : posterior for  $c_k$

# Bayes classification

- Consider each attribute and class label as random variables
- Given an instance  $X=(X_1, X_2, A_3, \dots, X_d)$ 
  - Goal is to predict class label  $c \in \mathcal{C}$
  - Specifically, we want to find the value  $c \in \mathcal{C}$  that maximizes  $P(c|X)$ , i.e.,  $\operatorname{argmax}_{c \in \mathcal{C}} P(c|X)$



# Bayes classification

- How can we estimate:

$$c = \operatorname{argmax}_{c \in C} P(c) * P(X|c)$$


- Short answer: from the data

- How to compute the class priors  $P(c)$ :

- Count the relative frequencies of the classes in the training set to estimate their priors

- Example: What is  $P(\text{Yes})$ ? 

- $P(\text{Yes}) = \frac{3}{10}$

- Example: What is  $P(\text{No})$ ? 

- $P(\text{No}) = \frac{7}{10}$

| <i>Tid</i> | Refund | Marital Status | Taxable Income | Evade |
|------------|--------|----------------|----------------|-------|
| 1          | Yes    | Single         | 125K           | No    |
| 2          | No     | Married        | 100K           | No    |
| 3          | No     | Single         | 70K            | No    |
| 4          | Yes    | Married        | 120K           | No    |
| 5          | No     | Divorced       | 95K            | Yes   |
| 6          | No     | Married        | 60K            | No    |
| 7          | Yes    | Divorced       | 220K           | No    |
| 8          | No     | Single         | 85K            | Yes   |
| 9          | No     | Married        | 75K            | No    |
| 10         | No     | Single         | 90K            | Yes   |

# Bayes classification

- How can we estimate:

$$c = \operatorname{argmax}_{c \in \mathcal{C}} P(c) * P(X|c)$$

- Short answer: from the data

- How to compute instance likelihood  $P(X/c)$ :

- What is the probability of an instance  $X$  given a class  $c$ ?

- $X=(X_1, X_2, \dots, X_d)$ , so,  $P(X/c)=P(X_1 \cap X_2 \cap \dots \cap X_d/c)$

- i.e., the probability of an instance given the class is equal to the probability of a set of features given the class

- So:  $c = \operatorname{argmax}_{c \in \mathcal{C}} P(c) * P(X_1 \cap X_2 \cap \dots \cap X_d|c)$

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |



## Simplification: For single-dimensional data ( $d=1$ )

- How to compute  $P(X|c) = P(X_1 \cap X_2 \cap \dots \cap X_d|c)$ ?
- Let us assume  $X$  is univariate (so only 1 dimension  $X_1$ , i.e.,  $d=1$ )

$$c = \operatorname{argmax}_{c \in C} P(c) * P(X_1|c)$$

- We can compute  $P(X_1/c)$  from the data
- Depending on the data type for  $X_1$  we distinguish between
  - Probability estimation for **categorical attributes** (e.g., hair color)
  - Probability estimation for **continuous attributes** (e.g., income)

# Probability estimation for categorical attributes

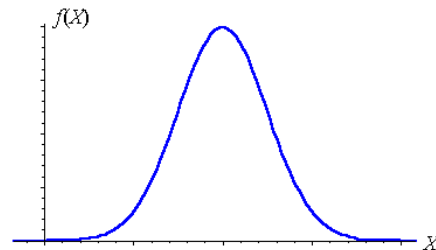
- Let  $X_1$  be a **categorical attribute**, what is  $P(X_1/c)$ ?
  - E.g., Marital status={Yes, No}, Color={green, blue, red}
- Based on the **relative frequency** of value  $X_1$  in class  $c$  in the training data
  - $P(X_1|c) = \frac{n_{1c}}{n_c}$
  - $n_c$ : #instances in class  $c$
  - $n_{1c}$ : #instances in class  $c$  having value  $X_1$
- Example: What is  $P(\text{Status}=\text{Married}|\text{No})$ ? 
  - $= \frac{4}{7}$
- Example: What is  $P(\text{Refund}=\text{No}|\text{Yes})$ ? 
  - $= \frac{3}{3}$

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |



# Probability estimation for continuous attributes

- Let  $X_1$  be a **continuous attribute**, what is  $P(X_1/c)$ ?
  - E.g., Income
- Idea 1: **Discretization**
  - Discretize the attribute → categorical attribute case
- Idea 2: **Probability density estimation**
  - Assume the attribute follows a known distribution
    - For example, Gaussian
  - Use data to estimate parameters of distribution
    - For example, mean and standard deviation
  - Once probability distribution is known, can be used to estimate the conditional probability  $P(X_1/c)$

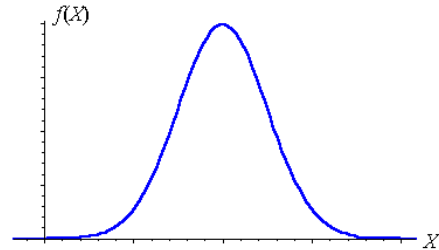


$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |

# Probability estimation for continuous attributes (Gaussian NBs)

- Example: Assume income follows the **Gaussian distribution**



- Use data to estimate the parameters of the distribution for each class, i.e., mean and standard deviation (See lecture 2)
  - For class No: mean income = 110
  - For class No: variance  $\sigma^2=2975$
- Once probability distribution is known, can be used to estimate the conditional probability  $P(X_1/c)$ 
  - E.g., what is the probability of income value 120 in class No?

$$P(\text{Income} = 120 \mid \text{No}) = \frac{1}{\sqrt{2\pi 2975}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

# Multidimensional data

- So far, we consider only one attribute for probability estimation
- For multi-dimensional data, we need to estimate the combined probabilities of specific attribute values:

$$P(X|c) = P(X_1 \cap X_2 \cap \dots \cap X_d|c)$$

- Example:

- $P(\text{Refund}=\text{yes} \cap \text{Marital status}=\text{single} | \text{No})$ 
  - $=1/10$
- $P(\text{Refund}=\text{no} \cap \text{Marital status}=\text{married} | \text{No})=$ 
  - $3/10$
- ...

- If we have  $d$  attributes each taking  $k$  values, we have  $k^d$  different combinations
- Typically, there are not enough training instances available to reliably estimate probabilities.

| <i>Tid</i> | Refund | Marital Status | Taxable Income | Evade |
|------------|--------|----------------|----------------|-------|
| 1          | Yes    | Single         | 125K           | No    |
| 2          | No     | Married        | 100K           | No    |
| 3          | No     | Single         | 70K            | No    |
| 4          | Yes    | Married        | 120K           | No    |
| 5          | No     | Divorced       | 95K            | Yes   |
| 6          | No     | Married        | 60K            | No    |
| 7          | Yes    | Divorced       | 220K           | No    |
| 8          | No     | Single         | 85K            | Yes   |
| 9          | No     | Married        | 75K            | No    |
| 10         | No     | Single         | 90K            | Yes   |

## Multidimensional data: Example

| ID | shape  | color  | class |
|----|--------|--------|-------|
| 1  | round  | orange | A     |
| 2  | round  | green  | A     |
| 3  | round  | yellow | A     |
| 4  | square | green  | A     |
| 5  | oval   | white  | B     |

$$\Pr(\text{shape}=\text{round} \cap \text{color}=\text{orange}|A) = \frac{1}{4}$$

$$\Pr(\text{shape}=\text{round} \cap \text{color}=\text{green}|A) = \frac{1}{4}$$

$$\Pr(\text{shape}=\text{round} \cap \text{color}=\text{yellow}|A) = \frac{1}{4}$$

$$\Pr(\text{shape}=\text{round} \cap \text{color}=\text{white}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{oval} \cap \text{color}=\text{orange}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{oval} \cap \text{color}=\text{green}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{oval} \cap \text{color}=\text{yellow}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{oval} \cap \text{color}=\text{white}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{square} \cap \text{color}=\text{orange}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{square} \cap \text{color}=\text{green}|A) = \frac{1}{4}$$

$$\Pr(\text{shape}=\text{square} \cap \text{color}=\text{yellow}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{square} \cap \text{color}=\text{white}|A) = \frac{0}{4}$$

$$\Pr(\text{shape}=\text{round} \cap \text{color}=\text{orange}|B) = \frac{0}{1}$$

- The probability estimates are unreliable, because the sample size is too small for each instance.

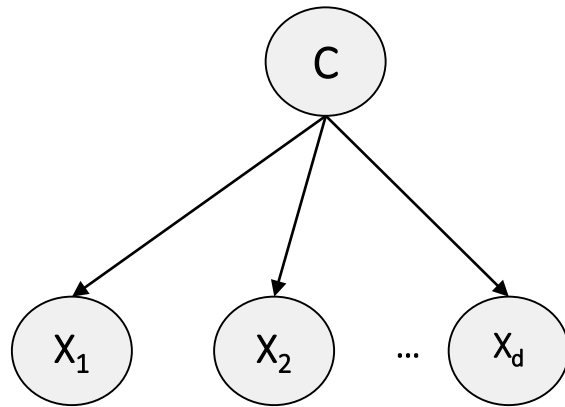
# Outline

- Generative vs Discriminative models
- Bayesian Classifiers
- Naïve Bayes classifier
- Laplace correction
- Bayesian Belief Networks
- Things you should know from this lecture & reading material

# Naïve Bayes classifier

- How to estimate instance likelihood  $P(X/c) = P(X_1 \cap X_2 \cap \dots \cap X_d / c)$ ?
- Assume independence among attributes  $X_i$  when class is given

$$P(X_1 \cap X_2 \cap \dots \cap X_d / c) = \prod P(X_i | c) = P(X_1 | c) P(X_2 | c) \dots P(X_d | c)$$



Strong class conditional independence assumption!!!

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |

- The class value is the hidden factor that explains all the dependencies between the attributes.
- In other words, once the class is observed an attribute does not give any information about other attributes



What class conditional independence means?  
E.g., in the context of sentiment analysis/ spam filtering?

# Naïve Bayes classifier

- Assume independence among attributes  $X_i$  when class is given:

$$P(X_1 \cap X_2 \cap \dots \cap X_d | c) = \prod (X_i | c) = P(X_1 | c) P(X_2 | c) \dots P(X_d | c)$$

Strong class conditional  
independence assumption!!!

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1   | Yes    | Single         | 125K           | No    |
| 2   | No     | Married        | 100K           | No    |
| 3   | No     | Single         | 70K            | No    |
| 4   | Yes    | Married        | 120K           | No    |
| 5   | No     | Divorced       | 95K            | Yes   |
| 6   | No     | Married        | 60K            | No    |
| 7   | Yes    | Divorced       | 220K           | No    |
| 8   | No     | Single         | 85K            | Yes   |
| 9   | No     | Married        | 75K            | No    |
| 10  | No     | Single         | 90K            | Yes   |

- Methodology
  - Using the independence assumption estimate  $P(X|c)$  for all classes  $c \in \mathcal{C}$  based on the training data
  - The instance  $X$  is finally classified into:

$$c = \operatorname{argmax}_{c \in \mathcal{C}} P(X|c) * P(c)$$

# Naive Bayes classifier: Example 1

## Training set

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis |
|-----|----------|-------------|----------|--------|------------|
| D1  | Sunny    | Hot         | High     | Weak   | No         |
| D2  | Sunny    | Hot         | High     | Strong | No         |
| D3  | Overcast | Hot         | High     | Weak   | Yes        |
| D4  | Rain     | Mild        | High     | Weak   | Yes        |
| D5  | Rain     | Cool        | Normal   | Weak   | Yes        |
| D6  | Rain     | Cool        | Normal   | Strong | No         |
| D7  | Overcast | Cool        | Normal   | Strong | Yes        |
| D8  | Sunny    | Mild        | High     | Weak   | No         |
| D9  | Sunny    | Cool        | Normal   | Weak   | Yes        |
| D10 | Rain     | Mild        | Normal   | Weak   | Yes        |
| D11 | Sunny    | Mild        | Normal   | Strong | Yes        |
| D12 | Overcast | Mild        | High     | Strong | Yes        |
| D13 | Overcast | Hot         | Normal   | Weak   | Yes        |
| D14 | Rain     | Mild        | High     | Strong | No         |

## Test instance X

| Outlook | Temperature | Humidity | Wind   | Play |
|---------|-------------|----------|--------|------|
| Sunny   | Cool        | High     | Strong | ?    |

$$P(\text{yes} | X) = \frac{P(X | \text{yes})P(\text{yes})}{P(X)} = \frac{P(O = \text{"sunny"} | \text{yes})P(T = \text{"cool"} | \text{yes})P(H = \text{"high"} | \text{yes})P(W = \text{"strong"} | \text{yes})P(\text{yes})}{P(X)}$$

$$P(O = \text{"sunny"} | \text{yes}) = \frac{2}{9} \quad P(T = \text{"cool"} | \text{yes}) = \frac{3}{9} \quad P(H = \text{"high"} | \text{yes}) = \frac{3}{9} \quad P(W = \text{"strong"} | \text{yes}) = \frac{3}{9}$$

$$P(\text{yes}) = \frac{9}{14}$$

$$P(\text{no} | X) = \frac{P(X | \text{no})P(\text{no})}{P(X)} = \frac{P(O = \text{"sunny"} | \text{no})P(T = \text{"cool"} | \text{no})P(H = \text{"high"} | \text{no})P(W = \text{"strong"} | \text{no})P(\text{no})}{P(X)}$$



# Naive Bayes classifier: Example 2

## Training set

| Name          | Give Birth | Can Fly | Live in Water | Have Legs | Class       |
|---------------|------------|---------|---------------|-----------|-------------|
| human         | yes        | no      | no            | yes       | mammals     |
| python        | no         | no      | no            | no        | non-mammals |
| salmon        | no         | no      | yes           | no        | non-mammals |
| whale         | yes        | no      | yes           | no        | mammals     |
| frog          | no         | no      | sometimes     | yes       | non-mammals |
| komodo        | no         | no      | no            | yes       | non-mammals |
| bat           | yes        | yes     | no            | yes       | mammals     |
| pigeon        | no         | yes     | no            | yes       | non-mammals |
| cat           | yes        | no      | no            | yes       | mammals     |
| leopard shark | yes        | no      | yes           | no        | non-mammals |
| turtle        | no         | no      | sometimes     | yes       | non-mammals |
| penguin       | no         | no      | sometimes     | yes       | non-mammals |
| porcupine     | yes        | no      | no            | yes       | mammals     |
| eel           | no         | no      | yes           | no        | non-mammals |
| salamander    | no         | no      | sometimes     | yes       | non-mammals |
| gila monster  | no         | no      | no            | yes       | non-mammals |
| platypus      | no         | no      | no            | yes       | mammals     |
| owl           | no         | yes     | no            | yes       | non-mammals |
| dolphin       | yes        | no      | yes           | no        | mammals     |
| eagle         | no         | yes     | no            | yes       | non-mammals |

## Test instance X

| Give Birth | Can Fly | Live in Water | Have Legs | Class |
|------------|---------|---------------|-----------|-------|
| yes        | no      | yes           | no        | ?     |

$$P(X | M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(X | N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

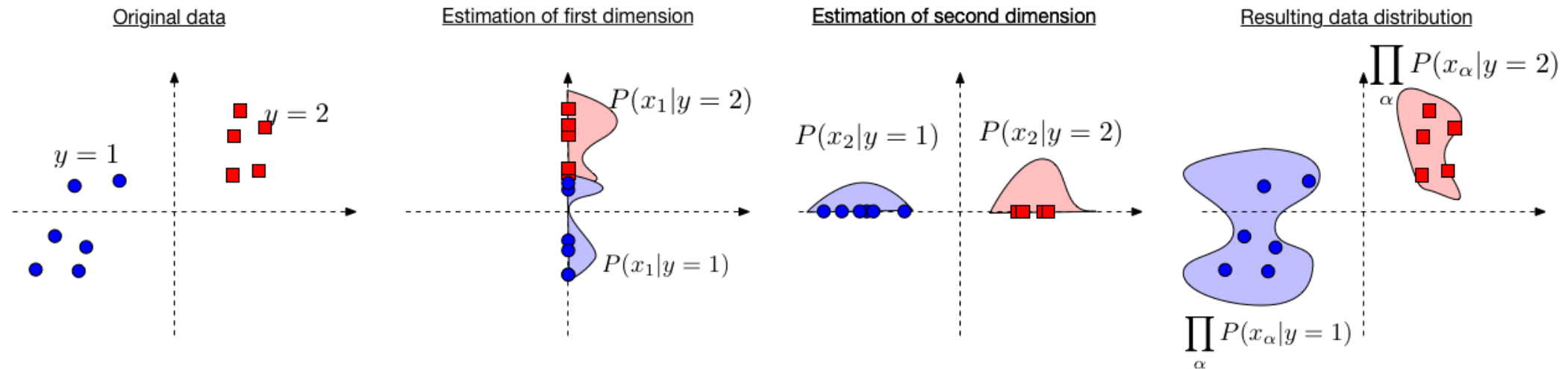
$$P(M | X) = P(X | M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(N | X) = P(X | N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

$P(M | X) > P(N | X) \rightarrow$  Mammals

# Naive Bayes classifier: Example 3

- We estimate  $P(X/c)$  independently in each dimension
- We multiple these probabilities under the conditional independence assumption



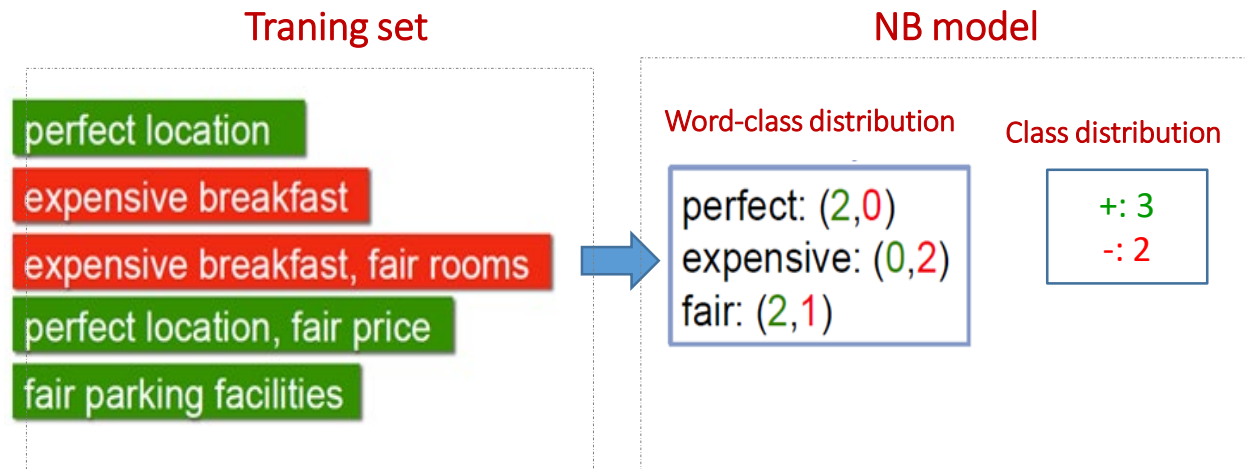
Source: [link](#)

# Outline

- Generative vs Discriminative models
- Bayesian Classifiers
- Naïve Bayes classifier
- Laplace correction
- Bayesian Belief Networks
- Things you should know from this lecture & reading material

# NB predictions

- Consider the following example



- Predict the class of the following instances:
  - X="Perfect breakfast"
  - X="Perfect weather"

# The zero frequency problem

- Naïve Bayesian prediction requires each conditional probability  $P(X_i/c)$  be non-zero. Otherwise, the predicted probability will be zero

$$c = \operatorname{argmax}_{c \in \mathcal{C}} P(X|c) * P(c)$$

- Solution: **Correction/Smoothing**:

$$\text{Original: } P(X_i | c) = \frac{n_{ic}}{n_c}$$

$$\text{Laplace: } P(X_i | c) = \frac{n_{ic} + 1}{n_c + k}$$

$k$ : number of classes

# The problem of 0-probabilities: example

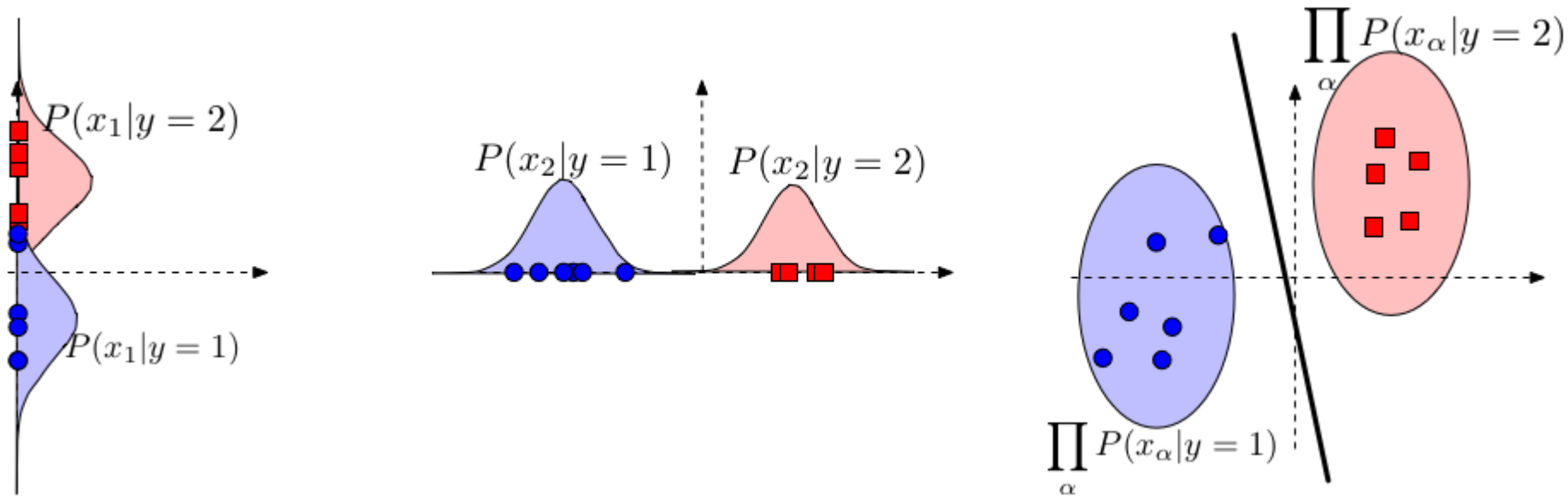
- Suppose a dataset with 1000 tuples:
  - income=low (0)
  - income= medium (990)
  - income = high (10)
- Use Laplacian correction (or Laplacian estimator): add 1 to each class value
  - $\text{Prob}(\text{income} = \text{low}) = 1/1003$
  - $\text{Prob}(\text{income} = \text{medium}) = 991/1003$
  - $\text{Prob}(\text{income} = \text{high}) = 11/1003$
- Result
  - The probabilities are never 0
  - The “corrected” prob. estimates are close to their “uncorrected” counterparts

# Inductive bias of NBs

- The assumption of independence can be seen as the bias inherent to the Naïve Bayes classifier.
- Relying on the bias, the classifier may have a tendency to be wrong (if the assumption does not hold).
- An unbiased probabilistic classifier is not practical due to a notorious lack of training examples.
  - In any practical scenario, it would overfit.
- As an example: For a binary classification problem, with 30 binary predictive attributes it would require more than 2 billion instances just to see each combination once (which is not sufficient for reliable probability estimates)
- So bias is necessary to make generalization feasible.

# Decision boundary

- Naive Bayes leads to a linear decision boundary in many common cases
- Illustrated here is the case of Gaussian NB where standard deviation is the same across all classes



Source: [Link](#)



# Naïve Bayes (NB) classifiers: overview

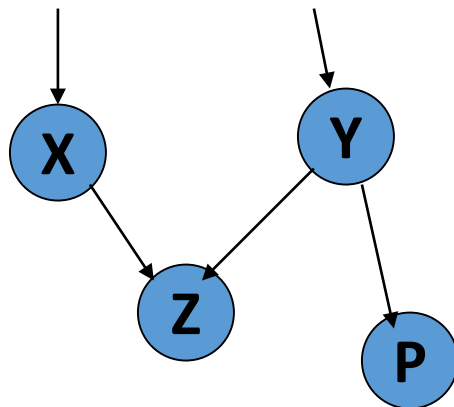
- (+) Easy to implement
- (+) It works surprisingly well in practice, although the independence assumption is too strong.
  - It does not require precise estimations of the probabilities
  - It is enough if the max probability belongs to the correct class
- (+) Robust to irrelevant attributes
- (+) Handles missing values by ignoring the value during probability estimate calculations
- (+) Robust to noise
- (+) Incremental
- (-) Strong independence assumption
- (-) Practically, there exist dependencies among variables
  - Such dependencies cannot be modeled by NB classifiers
  - Use other techniques such as Bayesian Belief Networks (BBN)

# Outline

- Generative vs Discriminative models
- Bayesian Classifiers
- Naïve Bayes classifier
- Laplace correction
- Bayesian Belief Networks
- Things you should know from this lecture & reading material

# Bayesian Belief Networks (BBN)

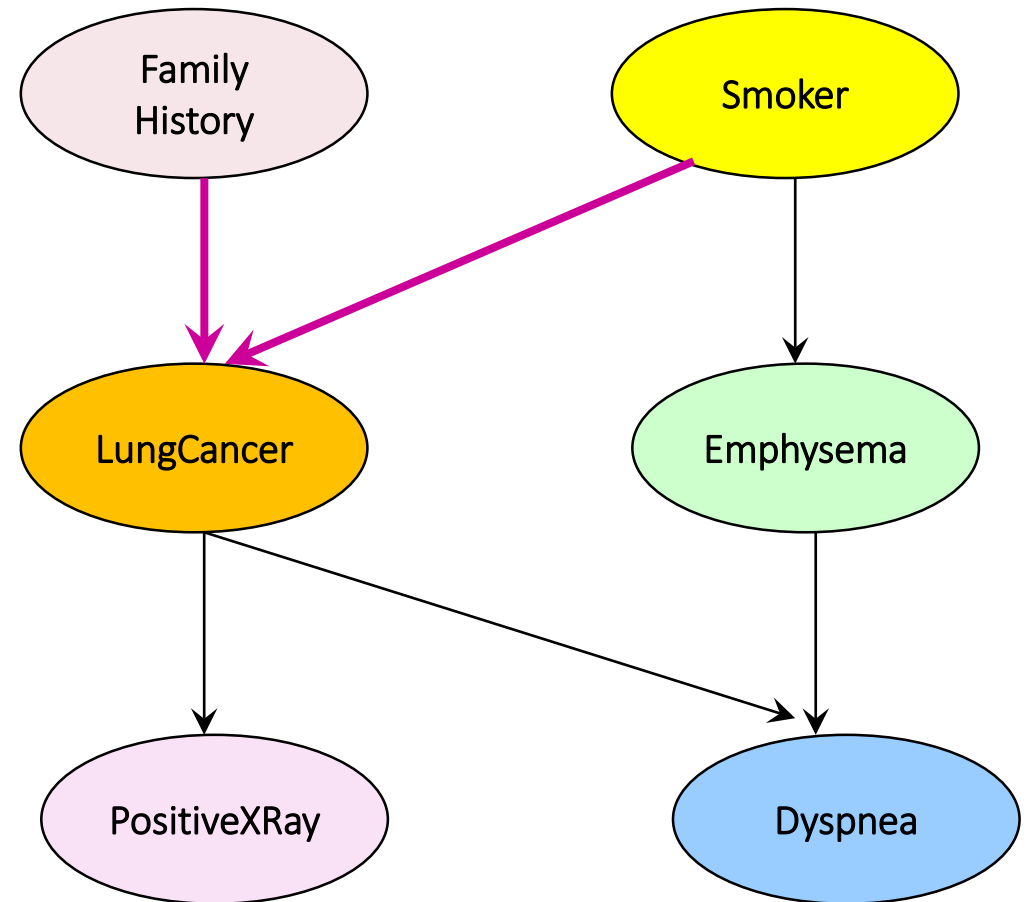
- Bayesian belief networks allow class conditional independence to be defined between subsets of variables instead of all variables (as in NB).
- They provide a graphical model of causal relationships, on which learning can be performed
- A belief network is defined by two components:
  - A **directed acyclic graph** of nodes encoding the **dependence relationships** among a set of variables.
  - A **set of conditional probability tables (CPT)** that associates each node to its immediate parent nodes.



- **Nodes:** random variables
- **Links:** dependency between variables
- X, Y are the parents of Z; Y is the parent of P.

# Conditional independence in BBN

- A node in a Bayesian network is **conditional independent** of its non-descendants, if its parents are known.
- In our example
  - having lung cancer is influenced by a person's family history and on whether or not the person is a smoker
  - PositiveXRay is independent of "family history" and "smoker" attributes once we know that the person has LungCancer.



# Bayesian Belief Networks

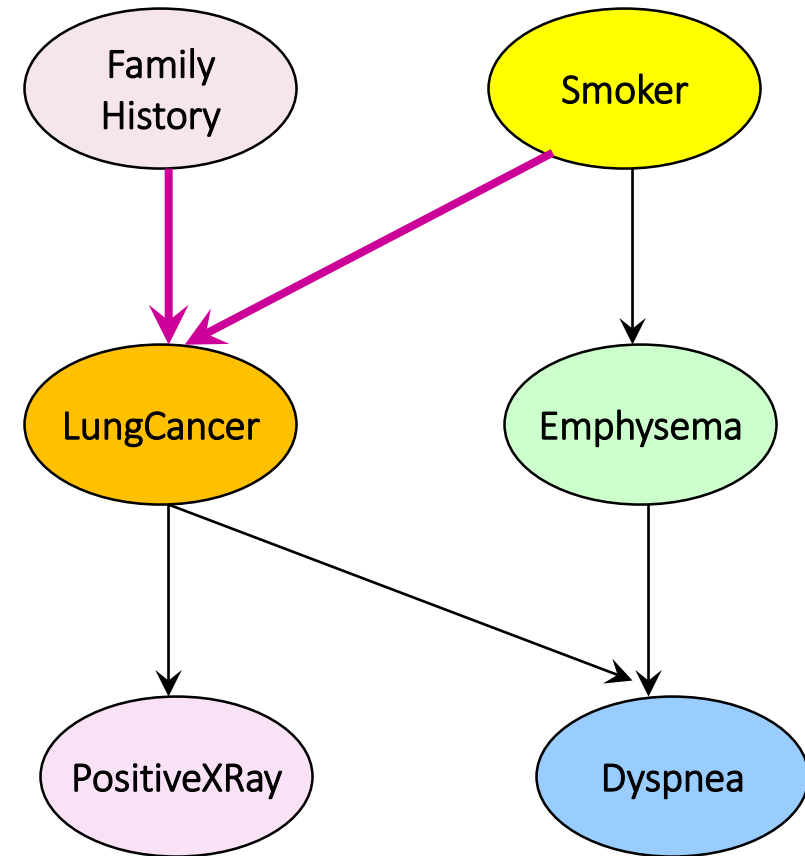
- Each variable  $X_i$  is associated with a conditional probability table (CPT)
- CPT of  $X_i$  specifies the conditional distribution  $P(X_i | Parents(X_i))$
- Let the conditional probability table (CPT) for variable LungCancer (LC):

|     | (FH, S) | (FH, ~S) | (~FH, S) | (~FH, ~S) |
|-----|---------|----------|----------|-----------|
| LC  | 0.8     | 0.5      | 0.7      | 0.1       |
| ~LC | 0.2     | 0.5      | 0.3      | 0.9       |

- Let a new instance  $X = (X_1, X_2, \dots, X_d)$ .
- The probability of  $X$  is given by:

$$P(X) = \prod_{i=1}^d P(X_i | Parents(X_i))$$

- Key challenge: How to get the CPTs



# Outline

- Bayesian Classifiers
  - Naïve Bayes classifier
  - Laplace correction
  - Bayesian Belief Networks
- Things you should know from this lecture & reading material

# Overview and Reading

- Overview

- Bayes rule
- Bayesian classifiers
- Naïve Bayes classifiers
- Laplace correction

- Reading

- Chapter 24: Generative Models, Understanding Machine Learning book by Shai Shalev-Schwartz and Shai Ben-David
- Chapter 5, Bayesian classifiers, Tan et al book

# Hands on experience



- Build a sentiment classifier for Twitter data (is a tweet positive or negative?)
  - Many datasets available
    - E.g., [Sentiment140 dataset](#)
- Classify newsgroup messages (20 newsgroups from politics to atheism)
  - Dataset: [20 newsgroup dataset](#)
- Spam filter (is tweet spam or not?)
  - Dataset: [HSPAM](#)
- If features do not represent categories but counts, we use multinomial distribution → Multinomial Naïve Bayes



# Acknowledgements

- The slides are based on
  - KDD I lecture at LMU Munich (Johannes Aßfalg, Christian Böhm, Karsten Borgwardt, Martin Ester, Eshref Januzaj, Karin Kailing, Peer Kröger, Eirini Ntoutsi, Jörg Sander, Matthias Schubert, Arthur Zimek, Andreas Züfle)
  - Introduction to Data Mining book slides at <http://www-users.cs.umn.edu/~kumar/dmbook/>
  - Pedro Domingos Machine Lecture course slides at the University of Washington
  - Machine Learning book by T. Mitchel slides at <http://www.cs.cmu.edu/~tom/mlbook-chapter-slides.html>
  - (DTs) J. Fürnkranz slides from TU Darmstadt (<https://www.ke.tu-darmstadt.de/lehre/archiv/ws0809/mlbm/>)
  - Thank you to all TAs contributing to their improvement, namely Vasileios Iosifidis, Damianos Melidis, Tai Le Quy, Han Tran.

Thank you

Questions/Feedback/Wishes?

# Acknowledgements

- The slides are based on
  - ❑ KDD I lecture at LMU Munich (Johannes Aßfalg, Christian Böhm, Karsten Borgwardt, Martin Ester, Eshref Januzaj, Karin Kailing, Peer Kröger, Eirini Ntoutsi, Jörg Sander, Matthias Schubert, Arthur Zimek, Andreas Züfle)
  - ❑ Introduction to Data Mining book slides at <http://www-users.cs.umn.edu/~kumar/dmbook/>
  - ❑ Pedro Domingos Machine Lecture course slides at the University of Washington
  - ❑ Machine Learning book by T. Mitchel slides at <http://www.cs.cmu.edu/~tom/mlbook-chapter-slides.html>
  - ❑ Arthur Zimek DMML lecture at SDU
  - ❑ Thank you to all TAs contributing to their improvement, namely Vasileios Iosifidis, Damianos Melidis, Tai Le Quy, Han Tran.