

AI 6102: Machine Learning Methodologies & Applications

L7: Bayesian Classifiers & K-NN Classifiers

Sinno Jialin Pan

Nanyang Technological University, Singapore

Homepage: <http://www.ntu.edu.sg/home/sinnopan>

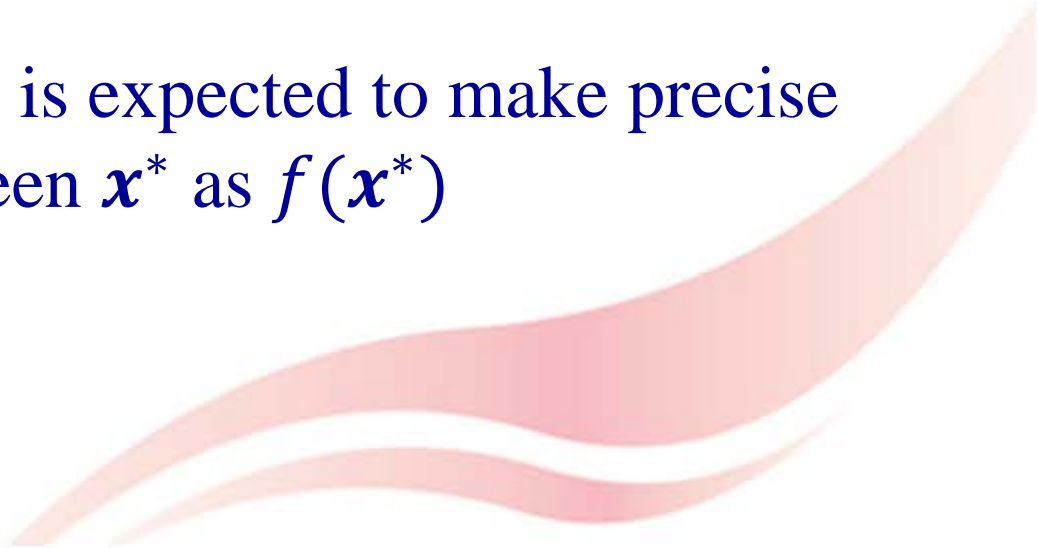
Outline

- Bayesian Classifiers
 - Naïve Bayes classifiers
- K Nearest neighbors (K -NN) classifiers
 - A lazy classifier



Supervised Learning: Recall

In mathematics

- Given: a set of $\{\mathbf{x}_i, y_i\}$ for $i = 1, \dots, N$, where $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{im}]$ is m -dimensional vector of numerical values, and y_i is a scalar
 - Aim to learn a mapping $f: \mathbf{x} \rightarrow y$ by requiring $f(\mathbf{x}_i) = y_i$
 - The learned mapping f is expected to make precise predictions on any unseen \mathbf{x}^* as $f(\mathbf{x}^*)$
- 

In Probability Point of View

- The mapping $f: \mathbf{x} \rightarrow y$ can be considered as a conditional probability $P(y|\mathbf{x})$
- Given a test data instance \mathbf{x}^*
$$y^* = c^* \text{ if } c^* = \arg \max_c P(y = c|\mathbf{x}^*), c \in \{0, \dots, C - 1\}$$
- In logistic regression, the conditional probabilities of different classes are assumed to be expressed as specific forms in terms of parameters \mathbf{w} , e.g., for binary classification (0 or 1)

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \quad P(y = 0|\mathbf{x}) = \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$


- Today, we introduce another way to estimate $P(y|\mathbf{x})$

Probability Review

- Let A be a random variable (a feature / a label in machine learning)
- Marginal probability $0 \leq P(A = a) \leq 1$

$$P(A = a)$$

refers to the probability that variable $A = a$

$$\sum_{a_i} P(A = a_i) = 1$$


Probability Review (cont.)

- Let A and B be a pair of random variables (features/labels in machine learning).
- Their joint probability

$$P(A = a, B = b)$$

refers to the probability that variable $A = a$, and at the same time variable $B = b$




Probability Review (cont.)

- Conditional probability

$$P(B = b|A = a)$$

refers to the probability that variable B will take on the value b , given that the variable A is observed to have the value a

$$\sum_{b_i} P(B = b_i|A = a) = 1$$


Sum Rule

- The connection between joint probability of A and B and marginal probability of A :

$$P(A = a) = \sum_{b_i} P(A = a, B = b_i) \quad \mathbf{OR} \quad P(A) = \sum_B P(A, B)$$

$$P(A = a) = \sum_{c_j} \sum_{b_i} P(A = a, B = b_i, C = c_j)$$

OR


$$P(A) = \sum_C \sum_B P(A, B, C)$$

Product Rule

- The connections between marginal, joint and conditional probabilities of A and B :

$$\begin{aligned}P(A = a, B = b) &= P(B = b|A = a) \times P(A = a) \\ &= P(A = a|B = b) \times P(B = b)\end{aligned}$$

OR

$$\begin{aligned}P(A, B) &= P(B|A) \times P(A) \\ &= P(A|B) \times P(B)\end{aligned}$$


Bayes Rule / Bayes Theorem

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

- The 1st and the 2nd equations are both based on product rule

$$P(A, B) = P(A|B) \times P(B) = P(B|A) \times P(A)$$

- Can be generalized to the case when **A** and **B** are a set of variables

$$P(A_1 \dots A_k | B_1 \dots B_p) = \frac{P(B_1 \dots B_p | A_1 \dots A_k) P(A_1 \dots A_k)}{P(B_1 \dots B_p)}$$

Bayesian Classifiers

- To estimate $P(y|\mathbf{x})$ from the training set $\{\mathbf{x}_i, y_i\} \ i = 1, \dots, N$, we can use the Bayes Rule

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

- Recall that to make a prediction on \mathbf{x}^*

$$y^* = c^* \text{ if } c^* = \arg \max_c P(y = c | \mathbf{x}^*), c \in \{0, \dots, C - 1\}$$

$P(\mathbf{x}^*)$ is a constant
w.r.t. different c

$$\begin{aligned} &= \arg \max_c \frac{P(\mathbf{x}^* | y = c)P(y = c)}{P(\mathbf{x}^*)} \\ &= \arg \max_c P(\mathbf{x}^* | y = c)P(y = c) \end{aligned}$$

An Example

- Suppose we aim to predict the class label (Repay = Yes or No) of the following data instance

Fixed Assets	Occupation	Income	Repay
Yes	Manager	125K	?

- That is we need to compute

V.S.

$$\begin{aligned}
 & P(\text{Rapy}=\text{Yes} \mid \text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K}) \\
 &= \frac{P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K} \mid \text{Rapy}=\text{Yes})P(\text{Rapy}=\text{Yes})}{\cancel{P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K})}} \\
 & P(\text{Rapy}=\text{No} \mid \text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K}) \\
 &= \frac{P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K} \mid \text{Rapy}=\text{No})P(\text{Rapy}=\text{No})}{\cancel{P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K})}}
 \end{aligned}$$

Bayesian Classifiers (cont.)

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

$$y^* = c^* \text{ if } c^* = \arg \max_c P(y = c | \mathbf{x}^*), c \in \{0, \dots, C - 1\}$$

$P(\mathbf{x}^*)$ is a constant
w.r.t. different c

$$= \arg \max_c \frac{P(\mathbf{x}^* | y = c)P(y = c)}{P(\mathbf{x}^*)}$$

$$= \arg \max_c \boxed{P(\mathbf{x}^* | y = c)P(y = c)}$$

Estimate these two types of
probabilities from training data

Bayesian Classifiers (cont.)

$$P(y|\mathbf{x}) \propto P(\mathbf{x}|y)P(y)$$

Class probabilities
from training data,
easy to estimate

In general, difficult to estimate
as all the possible combinations
need to be considered in training

- Consider the risk estimation task

$P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K} \mid \text{Rapy}=\text{Yes})$

$P(\text{Assets}=\text{No}, \text{Occ.}=\text{Manager}, \text{Income}=125\text{K} \mid \text{Rapy}=\text{Yes})$

$P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Engineer}, \text{Income}=125\text{K} \mid \text{Rapy}=\text{Yes})$

$P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Lawyer}, \text{Income}=125\text{K} \mid \text{Rapy}=\text{Yes})$

...

Bayesian Classifiers (cont.)

- In theory, the estimation of $P(\mathbf{x}|y)$ is computationally expensive
 - Need to consider all possible value combination of \mathbf{x} and y
 - How to make the estimation of $P(\mathbf{x}|y)$ computationally tractable?
- Two implementations of Bayesian classification methods
 - Naïve Bayes classifier
 - Based on a strong conditional independence assumption
 - Bayesian belief network
 - Based on a graph of dependence among variables

Not covered

Naïve Bayes Classifier

- Assume that the features are conditionally independent given the class label:

$$P(\mathbf{x}|y = c) = \prod_{i=1}^d P(x_i|y = c) \quad \text{where } \mathbf{x} = [x_1, x_2, \dots, x_d]$$

$$P(x_1, x_2, \dots, x_d|y = c) = \prod_{i=1}^d P(x_i|y = c)$$

Only need to different combinations of x_i and y , no need to consider all combinations of $[x_1, \dots, x_d]$ and y

For example:

$$P(\text{Assets=Yes, Occ.=Manager, Income=125K} \mid \text{Rapy=Yes})$$

$$= P(\text{Assets=Yes} \mid \text{Rapy=Yes})P(\text{Occ.=Manager} \mid \text{Rapy=Yes})P(\text{Income=125K} \mid \text{Rapy=Yes})$$

Independence

- Let A and B be two random variables
- A is said to be independent of B , if the following condition holds:

$$P(A, B) = \overset{P(A|B) = P(A)}{\underbrace{P(A|B)}} \times P(B) = \underbrace{P(A)} \times P(B)$$

- This can be generalized to the setting where \mathbf{A} and \mathbf{B} are two sets of random variables
- The variables in \mathbf{A} are said to be independent of \mathbf{B} , if the following condition holds:

$$P(\mathbf{A}, \mathbf{B}) = P(\mathbf{A}|\mathbf{B}) \times P(\mathbf{B}) = P(\mathbf{A}) \times P(\mathbf{B})$$

Conditional Independence

- Let **A**, **B**, and **C** be three sets of random variables
- The variables in **A** are said to be conditionally independent of **B**, given **C**, if the following condition holds:

$$P(\mathbf{A}|\mathbf{B}, \mathbf{C}) = P(\mathbf{A}|\mathbf{C})$$



$$P(\mathbf{x}|y = c) = \prod_{i=1}^d P(x_i|y = c)$$

Conditional Independence (cont.)

- The conditional independence between **A** and **B** given **C** can also be written as follows

$$P(\mathbf{A}, \mathbf{B}|\mathbf{C}) = \frac{P(\mathbf{A}, \mathbf{B}, \mathbf{C})}{P(\mathbf{C})} \quad \text{Product rule: } P(\mathbf{A}, \mathbf{B}|\mathbf{C})P(\mathbf{C}) = P(\mathbf{A}, \mathbf{B}, \mathbf{C})$$

$$= \frac{P(\mathbf{A}, \mathbf{B}, \mathbf{C})}{P(\mathbf{B}, \mathbf{C})} \times \frac{P(\mathbf{B}, \mathbf{C})}{P(\mathbf{C})} \quad \begin{array}{l} \text{Product rule:} \\ P(\mathbf{A}|\mathbf{B}, \mathbf{C})P(\mathbf{B}, \mathbf{C}) = P(\mathbf{A}, \mathbf{B}, \mathbf{C}) \\ P(\mathbf{B}|\mathbf{C})P(\mathbf{C}) = P(\mathbf{B}, \mathbf{C}) \end{array}$$

$$= P(\mathbf{A}|\mathbf{B}, \mathbf{C}) \times P(\mathbf{B}|\mathbf{C})$$

$$= P(\mathbf{A}|\mathbf{C}) \times P(\mathbf{B}|\mathbf{C})$$

Conditional independence:
 $P(\mathbf{A}|\mathbf{B}, \mathbf{C}) = P(\mathbf{A}|\mathbf{C})$

Naïve Bayes Classifier (cont.)

- The set of variables **A** and **B** are said to be independent given **C** if $P(\mathbf{A}, \mathbf{B} | \mathbf{C}) = P(\mathbf{A} | \mathbf{C}) \times P(\mathbf{B} | \mathbf{C})$
- Recall that naïve Bayes classifier assumes that the features are conditionally independent given the class label

$$\mathbf{A} = \{x_1, \dots, x_{d-1}\}, \mathbf{B} = \{x_d\}, \mathbf{C} = \{y = c\}$$

$$P(x_1, x_2, \dots, x_d | y = c) = P(x_1, \dots, x_{d-1} | y = c) P(x_d | y = c)$$



$$P(x_1, \dots, x_{d-1} | y = c) = P(x_1, \dots, x_{d-2} | y = c) P(x_{d-1} | y = c)$$



...

$$\begin{aligned} P(x_1, x_2, \dots, x_d | y = c) \\ = P(x_1 | y = c) P(x_2 | y = c) \dots P(x_d | y = c) = \prod_{i=1}^d P(x_i | y = c) \end{aligned}$$

Naïve Bayes Classifier (cont.)

- For any test data instance \mathbf{x}^*

$$\begin{aligned}c^* &= \arg \max_c P(y = c | \mathbf{x}^*) \\&= \arg \max_c \frac{P(\mathbf{x}^* | y = c) P(y = c)}{P(\mathbf{x}^*)} \\&= \arg \max_c P(\mathbf{x}^* | y = c) P(y = c) \\&= \arg \max_c P(y = c) \prod_{i=1}^d P(x_i^* | y = c)\end{aligned}$$

- In training, we need to estimate $P(y)$ for different classes, and for each class c and feature x_i , $P(x_i | y = c)$ for different possible values of x_i

Credit Risk Estimation

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

Testing

ID	Assets	Occupation	Income	Repay
11	No	Engineer	85K	?

Training

$$P(x_i | y = c)$$

$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No})$
 $P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No})$
 $P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes})$
 $P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes})$

 $P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No})$
 $P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No})$
 $P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No})$
 $P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes})$
 $P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes})$
 $P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes})$

 $P(\text{Income}=v \mid \text{Repay}=\text{Yes})$
 $P(\text{Income}=v \mid \text{Repay}=\text{No})$
 where $v \geq 0$

$$P(y = c)$$

$P(\text{Repay}=\text{Yes})$
 $P(\text{Repay}=\text{No})$

$$P(\text{No})P(\text{Assets}=\text{No} \mid \text{No})P(\text{Occu.}=\text{Engineer} \mid \text{No})P(\text{Income}=85\text{K} \mid \text{No})$$

V.S.

$$P(\text{Yes})P(\text{Assets}=\text{No} \mid \text{Yes})P(\text{Occu.}=\text{Engineer} \mid \text{Yes})P(\text{Income}=85\text{K} \mid \text{Yes})$$

Margin Probability of Class

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

Number of data instances of class c

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

Total number of training data instances

$$P(\text{Repay}=\text{Yes}) = \frac{7}{10}$$

$$P(\text{Repay}=\text{No}) = \frac{3}{10}$$

Conditional Probability on Discrete Features

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

$$P(\text{Assets} = \text{Yes} \mid \text{Repay} = \text{Yes})$$

$$= \frac{\#(\text{Assets} = \text{Yes} \wedge \text{Repay} = \text{Yes})}{\#(\text{Repay} = \text{Yes})} = \frac{3}{7}$$

$$P(\text{Occ.} = \text{Manager} \mid \text{Repay} = \text{No})$$

$$= \frac{\#(\text{Occ.} = \text{Manager} \wedge \text{Repay} = \text{No})}{\#(\text{Repay} = \text{No})} = \frac{2}{3}$$

Number of data instances of class c ,
whose values of the i -th feature are z

$$|(x_i = z) \wedge (y = c)|$$

$$P(x_i = z \mid y = c) = \frac{|(x_i = z) \wedge (y = c)|}{|y = c|}$$

Value of the i -th
feature equals to z

$$|y = c|$$

Number of data
instances of class c

Conditional Probability on Continuous Features

- Assume the values of a specific feature x_i given a specific class c follow a Gaussian distribution, i.e., $P(x_i|y = c)$ is a Gaussian distribution

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

- Use training data of class c to estimate parameters of the Gaussian distribution, i.e., mean μ and variance σ^2
 - Once the parameters are estimated, the Gaussian distribution is known, and we can use it to compute conditional probability
- Note: more methods will be introduced when introducing density estimation

Conditional Probability on Continuous Features (cont.)

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

{Income, Repay=Yes}, Gaussian distribution:

$$P(\text{Inc.}|\text{Yes}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\text{Inc.}-\mu)^2}{2\sigma^2}}$$

μ and σ^2 are the mean and variance of the income of the data instances whose labels are Yes (Repay=Yes)

$$\mu_x = \frac{1}{N} \sum_{k=1}^N x_k \quad \sigma_x = \sqrt{\frac{1}{N-1} \sum_{k=1}^N (x_i - \mu_x)^2}$$

$$\mu_{\{\text{inc.}, \text{Yes}\}} = 110$$

$$\sigma_{\{\text{Inc.}, \text{Yes}\}}^2 = 2975$$

$$\sigma_{\{\text{Inc.}, \text{Yes}\}} = 54.54$$

$$P(\text{Inc.}|\text{Yes}) = \frac{1}{\sqrt{2\pi} \times 54.54} e^{-\frac{(\text{Inc.}-110)^2}{2 \times 2975}}$$

Conditional Probability on Continuous Features (cont.)

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

{Income, Repay=No}, Gaussian distribution:

$$P(\text{Inc.}|\text{No}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\text{Inc.}-\mu)^2}{2\sigma^2}}$$

μ and σ^2 are the mean and variance of the income of the data instances whose labels are No (Repay=No)

$$\mu_{\{\text{inc.}, \text{No}\}} = 90$$

$$\sigma_{\{\text{Inc.}, \text{No}\}}^2 = 25$$

$$\sigma_{\{\text{Inc.}, \text{No}\}} = 5$$

$$P(\text{Inc.}|\text{No}) = \frac{1}{\sqrt{2\pi} \times 5} e^{-\frac{(\text{Inc.}-90)^2}{2 \times 25}}$$

Conditional Probability on Continuous Features (cont.)

$$P(\text{Inc.}|\text{No}) = \frac{1}{\sqrt{2\pi} \times 5} e^{-\frac{(\text{Inc.}-90)^2}{2 \times 25}}$$

$$P(\text{Inc.}|\text{Yes}) = \frac{1}{\sqrt{2\pi} \times 54.54} e^{-\frac{(\text{Inc.}-110)^2}{2 \times 2975}}$$

ID	Fixed Assets	Occupation	Income	Repay
11	No	Engineer	85k	?

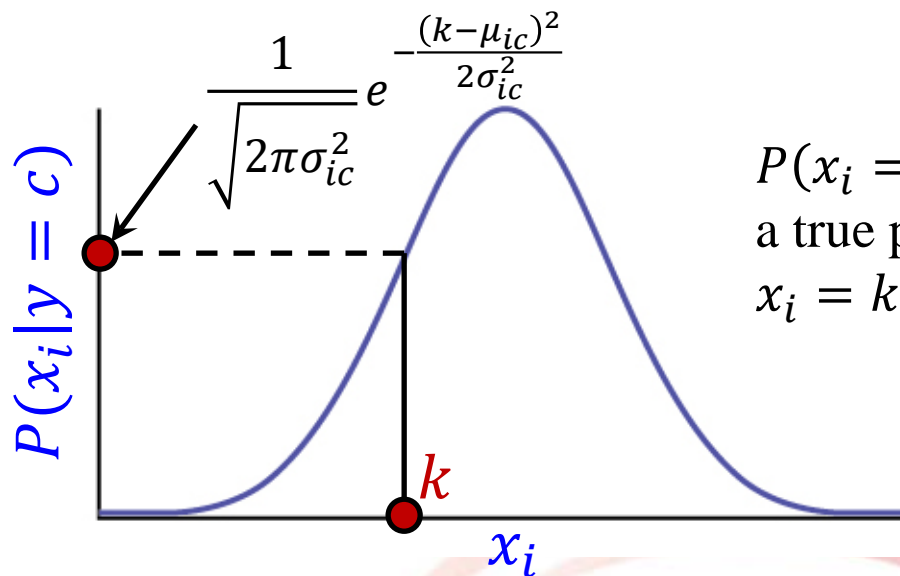
$$P(\text{Income}=85|\text{No}) = \frac{1}{\sqrt{2\pi} \times 5} e^{-\frac{(85-90)^2}{2 \times 25}} = 0.048$$

$$P(\text{Income}=85|\text{Yes}) = \frac{1}{\sqrt{2\pi} \times 54.54} e^{-\frac{(85-110)^2}{2 \times 2975}} = 0.007$$

Additional Notes

Probability density function $P(x_i|y = c) = \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(x_i - \mu_{ic})^2}{2\sigma_{ic}^2}}$

- The probability density function is continuous, the probability is defined as the area under the curve of the probability density function

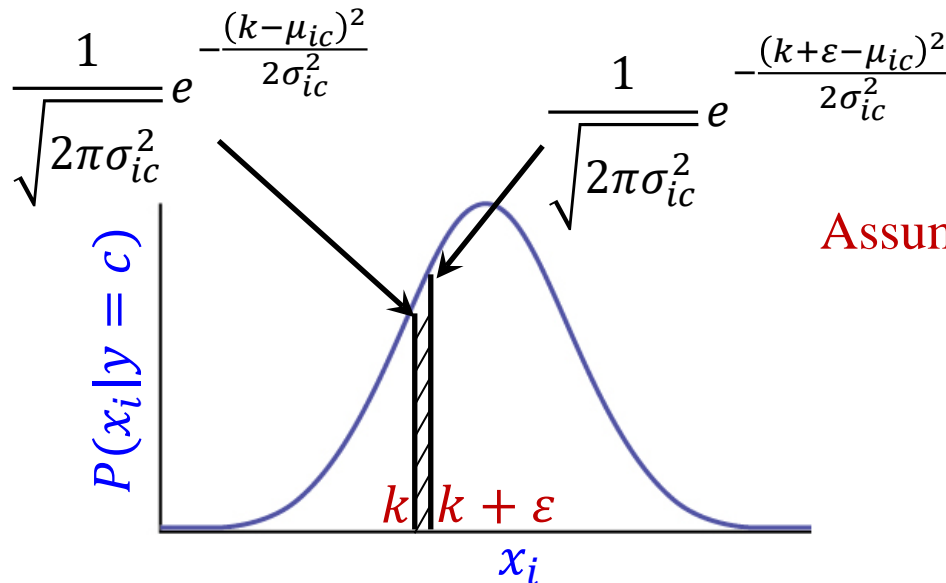


$P(x_i = k|y = c)$ is not a true probability that $x_i = k$ for class c

Additional Notes (cont.)

- Instead, we should compute

$$P(k \leq x_i \leq k + \underbrace{\varepsilon}_{\text{Small positive constant}} | y = c) = \int_k^{k+\varepsilon} \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(x_i - \mu_{ic})^2}{2\sigma_{ic}^2}} dx_i$$



$$\text{Assume } \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(k - \mu_{ic})^2}{2\sigma_{ic}^2}} \approx \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(k + \varepsilon - \mu_{ic})^2}{2\sigma_{ic}^2}}$$

$$\approx \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(k - \mu_{ic})^2}{2\sigma_{ic}^2}} \times \varepsilon$$

Additional Notes (cont.)

- Since ε appears as a constant multiplicative factor for each class, it cancels out when comparing posterior probabilities $P(y = c|\mathbf{x})$ for each class
- E.g., consider binary classification and instance is represented by a single feature of continuous values

$$P(y = 0|x = k) \quad \text{vs.} \quad P(y = 1|x = k)$$



$$P(x = k|y = 0)P(y = 0) \quad \text{vs.} \quad P(x = k|y = 1)P(y = 1)$$



$$\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(k-\mu_0)^2}{2\sigma_0^2}} \times \cancel{\varepsilon} \times P(y = 0) \quad \text{vs.} \quad \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(k-\mu_1)^2}{2\sigma_1^2}} \times \cancel{\varepsilon} \times P(y = 1)$$

Additional Notes (cont.)

- Therefore, we can still apply the following equation to approximate the probability of $x_i = k$ for class c

$$P(x_i = k|y = c) = \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(k-\mu_{ic})^2}{2\sigma_{ic}^2}}$$

Naïve Bayes Classifier: An Example

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

$$\arg \max_c P(y = c) \prod_{i=1}^d P(x_i^* | y = c)$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No}) = 0/3$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No}) = 3/3$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes}) = 3/7$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes}) = 4/7$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No}) = 2/3$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No}) = 0/3$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No}) = 1/3$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes}) = 2/7$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes}) = 4/7$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes}) = 1/7$$

$$P(\text{Income} \mid \text{Repay}=\text{Yes})$$

$$\mu_{\{\text{inc.}, \text{Yes}\}} = 110, \sigma_{\{\text{Inc.}, \text{Yes}\}}^2 = 2975$$

$$P(\text{Income} \mid \text{Repay}=\text{No})$$

$$\mu_{\{\text{inc.}, \text{No}\}} = 90, \sigma_{\{\text{Inc.}, \text{No}\}}^2 = 25$$

$$P(\text{Repay}=\text{Yes}) = 7/10$$

$$P(\text{Repay}=\text{No}) = 3/10$$

ID	Fixed Assets	Occupation	Income	Repay
11	No	Engineer	85k	?

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No}) = 0/3$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No}) = 3/3$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes}) = 3/7$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes}) = 4/7$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No}) = 2/3$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No}) = 0/3$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No}) = 1/3$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes}) = 2/7$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes}) = 4/7$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes}) = 1/7$$

$$P(\text{Income} \mid \text{Repay}=\text{Yes})$$

$$\mu_{\{\text{inc.}, \text{Yes}\}} = 110, \sigma_{\{\text{inc.}, \text{Yes}\}}^2 = 2975$$

$$P(\text{Income} \mid \text{Repay}=\text{No})$$

$$\mu_{\{\text{inc.}, \text{No}\}} = 90, \sigma_{\{\text{inc.}, \text{No}\}}^2 = 25$$

$$P(\text{Repay}=\text{Yes}) = 7/10$$

$$P(\text{Repay}=\text{No}) = 3/10$$

$$\begin{aligned}
 P(\mathbf{x}^*|\text{No}) &= P(\text{Assets}=\text{No} \mid \text{No}) \\
 &\quad \times P(\text{Occ.}=\text{Engineer} \mid \text{No}) \\
 &\quad \times P(\text{Income}=85 \mid \text{No}) \\
 &= 1 \times 0 \times 0.048 = 0
 \end{aligned}$$

one of the conditional probabilities is 0, the entire expression is 0

$$\begin{aligned}
 P(\mathbf{x}^*|\text{Yes}) &= P(\text{Assets}=\text{No} \mid \text{Yes}) \\
 &\quad \times P(\text{Occ.}=\text{Engineer} \mid \text{Yes}) \\
 &\quad \times P(\text{Income}=85 \mid \text{Yes}) \\
 &= 4/7 \times 4/7 \times 0.007 = 0.0023
 \end{aligned}$$

$$P(\mathbf{x}^*|\text{No}) \times P(\text{No}) = 0 \times 0.3 = 0$$

$$P(\mathbf{x}^*|\text{Yes}) \times P(\text{Yes}) = 0.0023 \times 0.7 = 0.0016$$

predict Repay=Yes

Laplace Estimate or Smoothing

Original: $P(x_i = z|y = c) = \frac{|(x_i = z) \wedge (y = c)|}{|y = c|}$

- Laplace: $P(x_i = z|y = c) = \frac{|(x_i = z) \wedge (y = c)| + 1}{|y = c| + n_i}$

Number of
possible
values of x_i

$P(\text{Engineer}|\text{No}) = \frac{\#(\text{Engineer} \wedge \text{No})}{\#(\text{No})} = \frac{0}{3}$

$P(\text{Engineer}|\text{No}) = \frac{\#(\text{Engineer} \wedge \text{No}) + 1}{\#(\text{No}) + 3} = \frac{1}{6}$

The same to $P(\text{Manager}|\text{No})$ and $P(\text{Lawyer}|\text{No})$

Extreme case - no training data:

$P(\text{Manager}|\text{No}) = P(\text{Engineer}|\text{No}) = P(\text{Lawyer}|\text{No}) = \frac{1}{3}$

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

More General Form

Laplace: $P(x_i = z|y = c) = \frac{|(x_i = z) \wedge (y = c)| + \boxed{\alpha}}{|y = c| + \alpha n_i}$

$\alpha > 0$ is a smoothing parameter

For example, $\alpha = 0.1$

$$P(\text{Engineer}|\text{No}) = \frac{\#(\text{Engineer} \wedge \text{No}) + 0.1}{\#(\text{No}) + 0.3} = \frac{1}{33}$$

For example, $\alpha = 10$

$$P(\text{Engineer}|\text{No}) = \frac{\#(\text{Engineer} \wedge \text{No}) + 10}{\#(\text{No}) + 30} = \frac{10}{33}$$

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	Engineer	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	No	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

Practice

Use Laplace smoothing with $\alpha = 1$ to re-estimate $P(\text{Assets}|\text{Repay})$ and $P(\text{Occ.}|\text{Repay})$

$$\begin{aligned}P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No}) &= 0/3 \\P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No}) &= 3/3 \\P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes}) &= 3/7 \\P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes}) &= 4/7\end{aligned}$$

$$\begin{aligned}P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No}) &= 2/3 \\P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No}) &= 0/3 \\P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No}) &= 1/3 \\P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes}) &= 2/7 \\P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes}) &= 4/7 \\P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes}) &= 1/7\end{aligned}$$

$$P(\text{Income} \mid \text{Repay}=\text{Yes})$$

$$\mu_{\{\text{inc.}, \text{Yes}\}} = 110, \sigma_{\{\text{Inc.}, \text{Yes}\}}^2 = 2975$$

$$P(\text{Income} \mid \text{Repay}=\text{No})$$

$$\mu_{\{\text{inc.}, \text{No}\}} = 90, \sigma_{\{\text{Inc.}, \text{No}\}}^2 = 25$$

$$P(\text{Repay}=\text{Yes}) = 7/10$$

$$P(\text{Repay}=\text{No}) = 3/10$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No}) = ?$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No}) = ?$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes}) = ?$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes}) = ?$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No}) = ?$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No}) = ?$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No}) = ?$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes}) = ?$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes}) = ?$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes}) = ?$$

$$P(\text{Income} \mid \text{Repay}=\text{Yes})$$

$$\mu_{\{\text{inc.}, \text{Yes}\}} = 110, \sigma_{\{\text{Inc.}, \text{Yes}\}}^2 = 2975$$

$$P(\text{Income} \mid \text{Repay}=\text{No})$$

$$\mu_{\{\text{inc.}, \text{No}\}} = 90, \sigma_{\{\text{Inc.}, \text{No}\}}^2 = 25$$

$$P(\text{Repay}=\text{Yes}) = 7/10$$

$$P(\text{Repay}=\text{No}) = 3/10$$

$$P(x_i = z | y = c) = \frac{|(x_i = z) \wedge (y = c)| + \alpha}{|y = c| + \alpha n_i}$$

$$\alpha = 1$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No}) = 0/3$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No}) = 3/3$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes}) = 3/7$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes}) = 4/7$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No}) = 2/3$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No}) = 0/3$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No}) = 1/3$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes}) = 2/7$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes}) = 4/7$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes}) = 1/7$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{No}) = \frac{0 + 1}{3 + 2} = \frac{1}{5}$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{No}) = \frac{3 + 1}{3 + 2} = \frac{4}{5}$$

$$P(\text{Assets}=\text{Yes} \mid \text{Repay}=\text{Yes}) = \frac{3 + 1}{7 + 2} = \frac{4}{9}$$

$$P(\text{Assets}=\text{No} \mid \text{Repay}=\text{Yes}) = \frac{4 + 1}{7 + 2} = \frac{5}{9}$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{No}) = \frac{2 + 1}{3 + 3} = \frac{1}{2}$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{No}) = \frac{0 + 1}{3 + 3} = \frac{1}{6}$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{No}) = \frac{1 + 1}{3 + 3} = \frac{1}{3}$$

$$P(\text{Occ.}=\text{Manager} \mid \text{Repay}=\text{Yes}) = \frac{2 + 1}{7 + 3} = \frac{3}{10}$$

$$P(\text{Occ.}=\text{Engineer} \mid \text{Repay}=\text{Yes}) = \frac{4 + 1}{7 + 3} = \frac{1}{2}$$

$$P(\text{Occ.}=\text{Lawyer} \mid \text{Repay}=\text{Yes}) = \frac{1 + 1}{7 + 3} = \frac{1}{5}$$

Naïve Bayes vs. Logistic Regression

- Both are probabilistic models for classification
- Use different ways to estimate $P(y|\mathbf{x})$

– Naïve Bayes:

Generative model

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

– Logistic Regression:

Discriminative model

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Binary classification

$$P(y = 0|\mathbf{x}) = \frac{\exp(-\mathbf{w}^T \mathbf{x})}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Deal with Missing Values

- In training, we only need to compute $P(x_i = z|y = c)$ for each feature independently
 - Ignore the missing value, e.g., when we compute $P(\text{Occ.} = z|\text{Repay} = \text{Yes})$ and $P(\text{Occ.} = z|\text{Repay} = \text{No})$, where $z \in \{\text{Manager}, \text{Engineer}, \text{Lawyer}\}$, we only consider the data instances without missing values of Occ.
 - No need to remove whole data instances or features from the training dataset

ID	Assets	Occupation	Income	Repay
1	Yes	Manager	125K	Yes
2	No	?	100K	Yes
3	No	Manager	70K	Yes
4	Yes	Engineer	120K	Yes
5	?	Lawyer	95K	No
6	No	Engineer	60K	Yes
7	Yes	Lawyer	220K	Yes
8	No	Manager	85K	No
9	No	Engineer	75K	Yes
10	No	Manager	90K	No

- In testing,

ID	Assets	Occupation	Income	Repay
11	?	Engineer	85k	?

$$\text{v.s.} \begin{cases} P(\text{No} \mid \text{Occ.}=\text{Engineer}, \text{Income}=85) \\ P(\text{Yes} \mid \text{Occ.}=\text{Engineer}, \text{Income}=85) \end{cases}$$

$$P(\text{No} \mid \text{Occ.}=\text{Engineer}, \text{Income}=85) \propto P(\text{Occ.}=\text{Engineer}, \text{Income}=85 \mid \text{No})P(\text{No})$$

$$= P(\text{Occ.}=\text{Engineer}, \text{Income}=85, \text{No})$$

By using the sum rule $\sum_B P(A, B) = P(A)$

$$= P(\text{Assets}=\text{No}, \text{Occ.}=\text{Eng.}, \text{Income}=85, \text{No}) + P(\text{Assets}=\text{Yes}, \text{Occ.}=\text{Eng.}, \text{Income}=85, \text{No})$$

$$= P(\text{Assets}=\text{No} \mid \text{No}) \times P(\text{Occ.}=\text{Engineer} \mid \text{No}) \times P(\text{Income}=85 \mid \text{No}) \times P(\text{No})$$

$$+ P(\text{Assets}=\text{Yes} \mid \text{No}) \times P(\text{Occ.}=\text{Engineer} \mid \text{No}) \times P(\text{Income}=85 \mid \text{No}) \times P(\text{No})$$

$$= (P(\text{Assets}=\text{No} \mid \text{No}) + P(\text{Assets}=\text{Yes} \mid \text{No})) \times P(\text{Occ.}=\text{Engineer} \mid \text{No}) \times P(\text{Income}=85 \mid \text{No}) \times P(\text{No})$$

$$= P(\text{Occ.}=\text{Engineer} \mid \text{No}) \times P(\text{Income}=85 \mid \text{No}) \times P(\text{No})$$

$$P(\text{Yes} \mid \text{Occ.}=\text{Engineer}, \text{Income}=85) \propto P(\text{Occ.}=\text{Engineer} \mid \text{Yes}) \times P(\text{Income}=85 \mid \text{Yes}) \times P(\text{Yes})$$

Summary

- Computationally efficient
- Computational efficiency is obtained based on a very strong assumption of conditional independence
 - The assumption may not hold in practice (most of the time)
 - That is why we call it “naïve”
 - It was widely used for text classification in the past



Implementation using scikit-learn

- API: `sklearn.naive_bayes`: Naive Bayes

https://scikit-learn.org/stable/modules/classes.html#module-sklearn.naive_bayes

`sklearn.naive_bayes`: Naive Bayes

The `sklearn.naive_bayes` module implements Naive Bayes algorithms. These are supervised learning methods based on applying Bayes' theorem with strong (naive) feature independence assumptions.

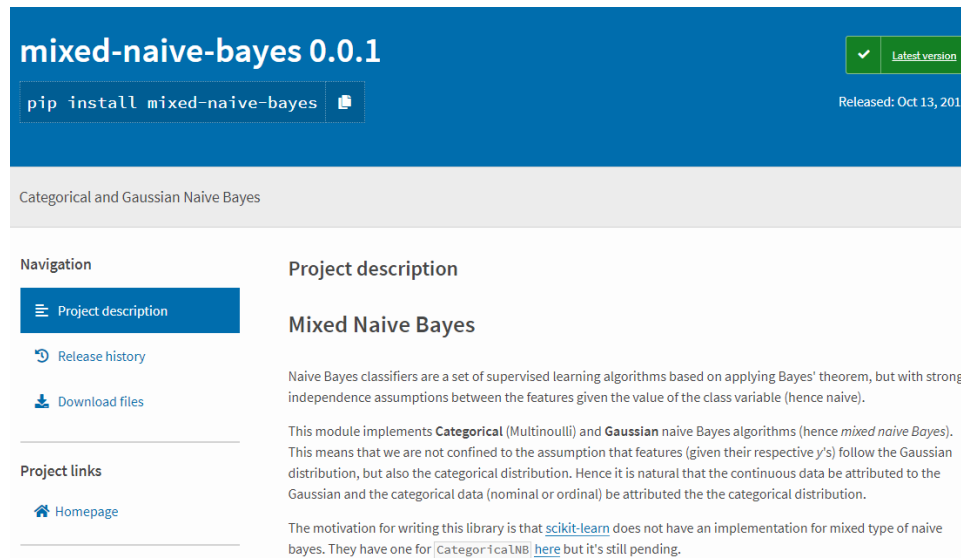
User guide: See the [Naive Bayes](#) section for further details.

<code>naive_bayes.BernoulliNB(* [, alpha, ...])</code>	Naive Bayes classifier for multivariate Bernoulli models.
<code>naive_bayes.CategoricalNB(* [, alpha, ...])</code>	Naive Bayes classifier for categorical features
<code>naive_bayes.ComplementNB(* [, alpha, ...])</code>	The Complement Naive Bayes classifier described in Rennie et al.
<code>naive_bayes.GaussianNB(* [, priors, ...])</code>	Gaussian Naive Bayes (GaussianNB)
<code>naive_bayes.MultinomialNB(* [, alpha, ...])</code>	Naive Bayes classifier for multinomial models

Documentation: https://scikit-learn.org/stable/modules/naive_bayes.html

Mixed Naïve Bayes Implementation

<https://pypi.org/project/mixed-naive-bayes/#installation>



The screenshot shows the PyPI page for the 'mixed-naive-bayes' package, version 0.0.1. The header is blue with the package name and version. A green button indicates it is the 'Latest version'. Below the header, there is a command to install the package: 'pip install mixed-naive-bayes'. The description states it is 'Categorical and Gaussian Naive Bayes'. The 'Project description' section explains that Naive Bayes classifiers are based on Bayes' theorem with strong independence assumptions. It also mentions that the module implements 'Categorical' (Multinoulli) and 'Gaussian' naive Bayes algorithms, hence 'mixed naive Bayes'. The motivation for writing this library is that 'scikit-learn' does not have an implementation for mixed type of naive bayes.

```
>>> from mixed_naive_bayes import MixedNB
```

```
>>> nbC = MixedNB(categorical_features=[0,1,3])
```

```
>>> nbC.fit(X, y)
```

```
>>> nbC.predict(X)
```

Specify which columns
are categorical features

Outline

- Bayesian Classifiers
 - Naïve Bayes classifiers
- K Nearest neighbors (K -NN) classifiers
 - A lazy classifier



Typical Learning Procedure

Inductive Learning

Labeled training data $\{\mathbf{x}_i, y_i\}, i = 1, \dots, N$

ID	Gender	Profession	Income	Saving	Repay
1	F	Engineer	60k	200k	Yes
2	M	Student	10k	20k	Yes
...
10	M	Student	8k	5k	No

Training phase

Test data \mathbf{x}^*

ID	Gender	Profession	Income	Saving
11	F	Lawyer	70k	100k



Test phase

Repay
Yes or No

Model
 $f: \mathbf{x} \rightarrow y$

Lazy Learning Procedure

Lazy Learning

Labeled training data $\{\mathbf{x}_i, y_i\}, i = 1, \dots, N$

ID	Gender	Profession	Income	Saving	Repay
1	F	Engineer	60k	200k	Yes
2	M	Student	10k	20k	Yes
...
10	M	Student	10k	20k	No

Training phase

- A model is not learned during “training phase”
- Instead, hashing table or indexing can be built

Test data \mathbf{x}^*

ID	Gender	Profession	Income	Saving
11	F	Lawyer	70k	100k

Test phase


- Retrieve similar training instances



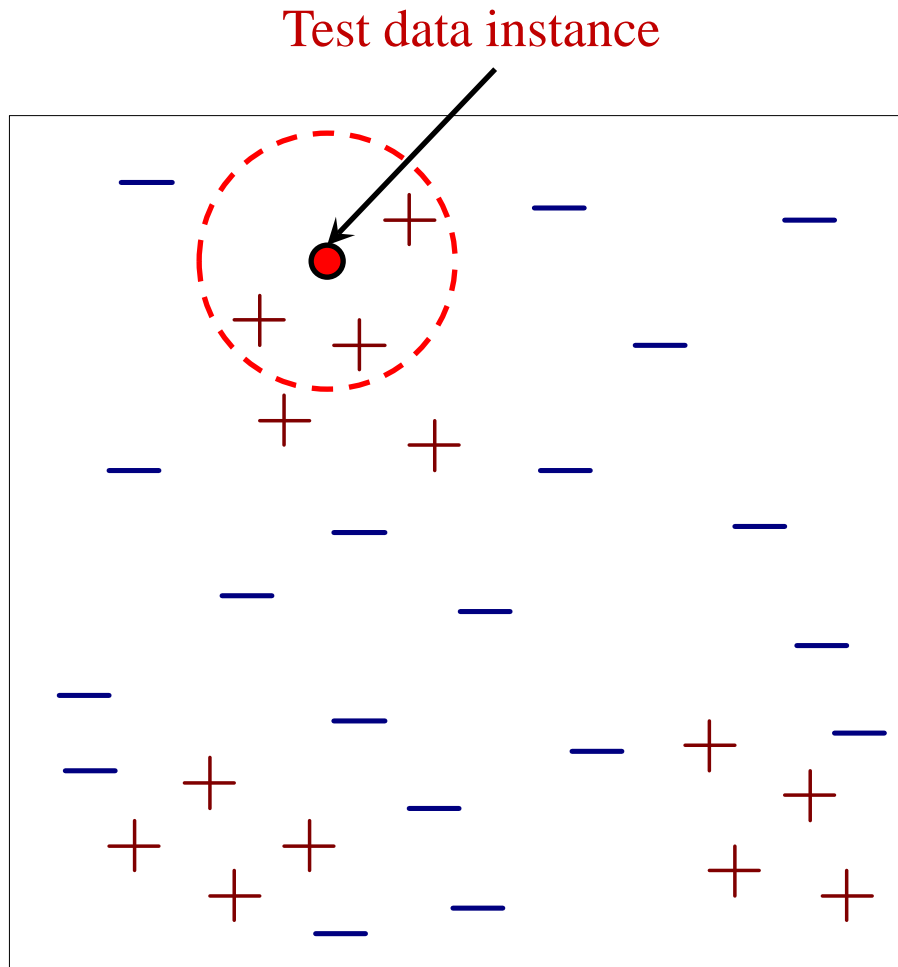
Based on the labels of the retrieved data instances

Repay
Yes or No

K -Nearest Neighbors Classifiers

- Algorithm:
 - For each test instance \mathbf{x}^* , retrieve K training data instances that are the most similar to \mathbf{x}^* (K nearest neighbors) from the training set
 - Based on the class labels of the K nearest neighbors to make a prediction on \mathbf{x}^*
- 

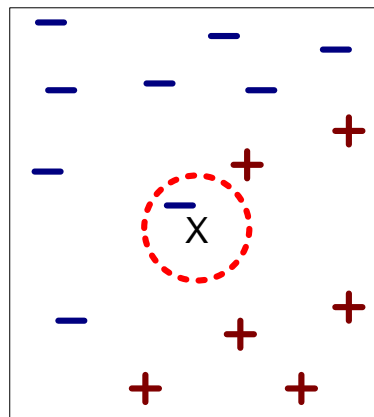
Illustration



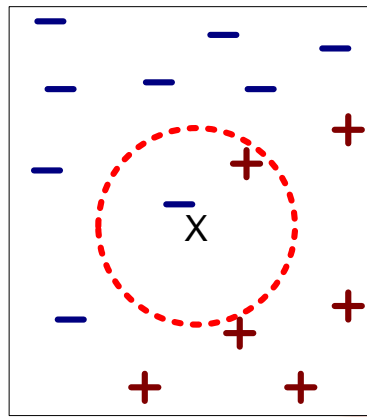
- Requirements:
 - A set of stored labeled training instances
 - Distance measure to compute distance between instances
 - The value of K , the number of nearest neighbors to retrieve
- To classify a test instance:
 - Compute distance to all the training instances
 - Identify K nearest neighbors
 - Use class labels of nearest neighbors to determine the class label of the test instance (e.g., using the majority class)

Distance & Nearest Neighbors

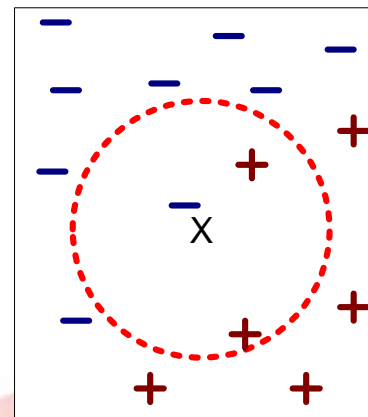
- Euclidean distance $d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{k=1}^m (x_{ik} - x_{jk})^2}$
 - The smaller is the value, the more similar are two data instances
 - K -nearest neighbors of an instance \mathbf{x} are data instances that have the K smallest distance to \mathbf{x}



(a) 1-nearest
neighbor



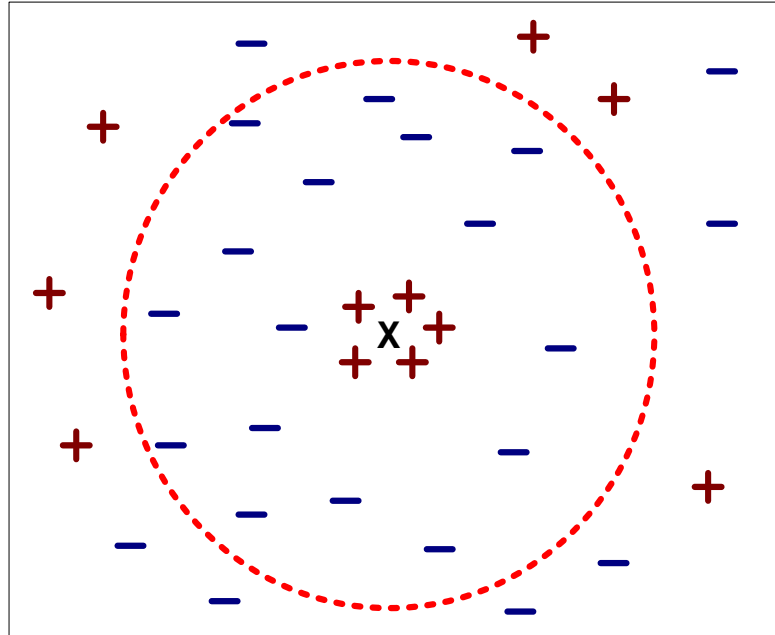
(b) 2-nearest
neighbor



(c) 3-nearest
neighbor

Value of K

- K is a hyper-parameter:
 - If K is too small, sensitive to noise points
 - If K is too large, neighborhood may include points from other classes



Determine Class Label

- Determine the class from nearest neighbor list
 - Take the majority vote of class labels among the K -nearest neighbors

- For majority voting:

$$y^* = \arg \max_c \sum_{(x_i, y_i) \in \mathcal{N}_{x^*}} I(c = y_i)$$

Indicator function that returns 1 if its input is true, otherwise 0

Nearest neighbors of the test instance \mathbf{x}^*

- Every neighbor has the same impact on the classification
- This makes the algorithm more sensitive to the choice of K

Weighted Voting

- Alternative scheme: distance-weight voting
 - Weight the influence of each nearest neighbor \mathbf{x}_i according to its distance to the test data

$$w_i = \frac{1}{d(\mathbf{x}^*, \mathbf{x}_i)^2}$$

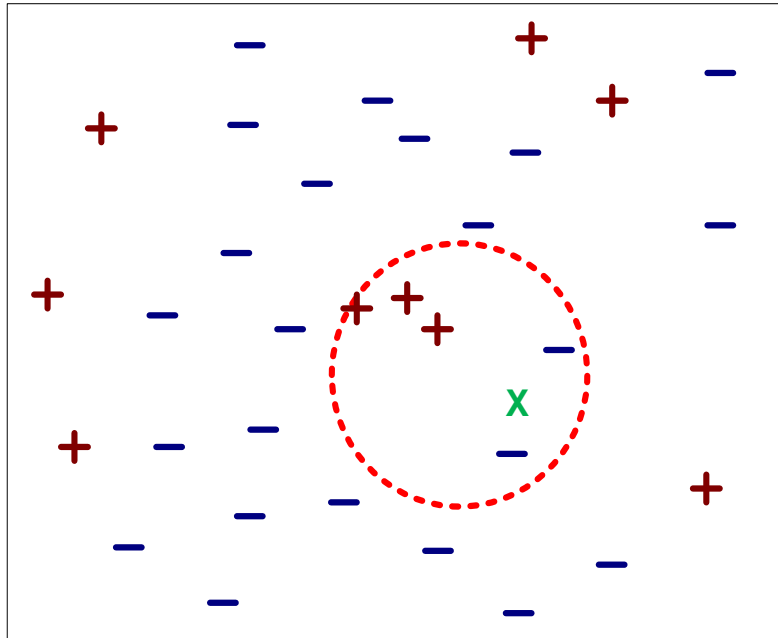
- That is

$$y^* = \arg \max_c \sum_{(\mathbf{x}_i, y_i) \in \mathcal{N}_{\mathbf{x}^*}} w_i I(c = y_i)$$

The larger is the distance to the test data,
the smaller influence of the corresponding
nearest neighbor to the vote

An Example

Consider a binary classification problem, and a 5-NN classifier



Instance ID	Class	Squared distance to test data
1	+	9
2	+	12.25
3	+	16
4	-	2.25
5	-	4

- Majority voting:

$$\textcircled{+ : 3} > - : 2$$

- Distance-weight voting:

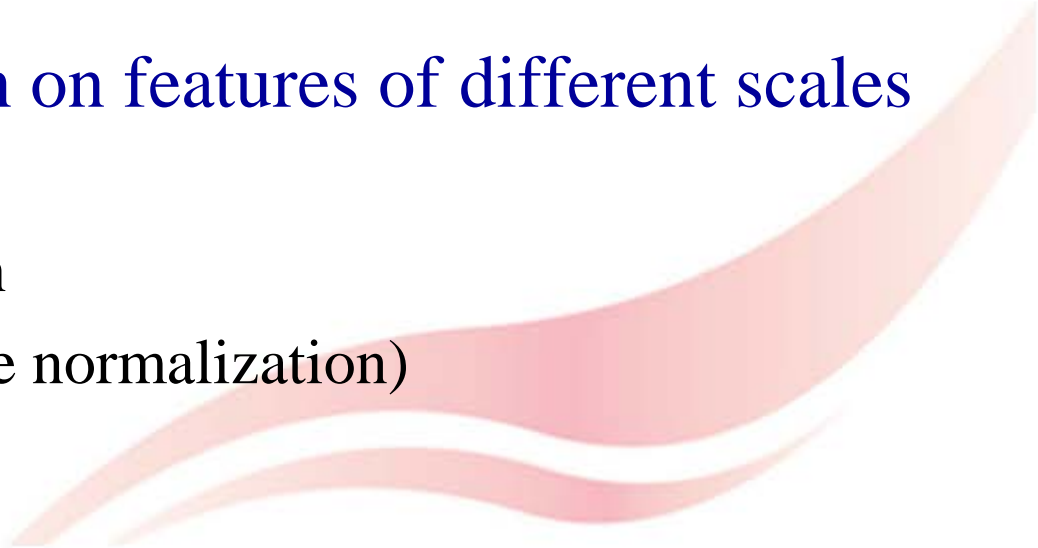
Distance-Weight votes for +:

$$\frac{1}{9} + \frac{1}{12.25} + \frac{1}{16} = 0.2552$$

Distance-Weight votes for -:

$$< \frac{1}{4} + \frac{1}{2.25} = 0.6944$$

Potential Issues

- Features may have very different scales, e.g.,
 - height of a person may vary from 1.5m to 1.8m
 - weight of a person may vary from 3kg to 200kg
 - income of a person may vary from \$10K to \$1M
 - Features need to be rescaled to prevent distance from being dominated by some features
 - Solution: normalization on features of different scales to the same scale
 - Min-Max normalization
 - Standardization (z-score normalization)
- 

Summary

- The *KNN* classifiers are a lazy learner
 - A classification model is not built explicitly
 - “Training” is very efficient
 - Classifying test instances is relatively expensive



Implementation using scikit-learn

- API: `sklearn.neighbors.KNeighborsClassifier`

<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>

sklearn.neighbors.KNeighborsClassifier

```
class sklearn.neighbors.KNeighborsClassifier(n_neighbors=5, *, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None, **kwargs)
```

[\[source\]](#)

Classifier implementing the k-nearest neighbors vote.

Read more in the [User Guide](#).

Parameters:	n_neighbors : int, default=5 Number of neighbors to use by default for <code>kneighbors</code> queries.
	weights : {'uniform', 'distance'} or callable, default='uniform' weight function used in prediction. Possible values: <ul style="list-style-type: none">• 'uniform' : uniform weights. All points in each neighborhood are weighted equally.• 'distance' : weight points by the inverse of their distance. in this case, closer neighbors of a query point will have a greater influence than neighbors which are further away.• [callable] : a user-defined function which accepts an array of distances, and returns an array of the same shape containing the weights.
	algorithm : {'auto', 'ball_tree', 'kd_tree', 'brute'}, default='auto' Algorithm used to compute the nearest neighbors: <ul style="list-style-type: none">• 'ball_tree' will use <code>BallTree</code>• 'kd_tree' will use <code>KDTree</code>• 'brute' will use a brute-force search.• 'auto' will attempt to decide the most appropriate algorithm based on the values passed to <code>fit</code> method.

Example

```
>>> from sklearn.neighbors import KNeighborsClassifier  
>>> import numpy as np
```

```
>>> n_samples, n_features = 10, 5  
>>> rng = np.random.RandomState(0)  
>>> y = rng.integers(2, n_samples)  
>>> X = rng.randn(n_samples, n_features)
```

```
>>> knnC = KNeighborsClassifier(n_neighbors=3)  
>>> knnC.fit(X, y)  
>>> pred = knnC.predict(X)
```

set number of neighbors

Build indices s.t. it is more efficient
when making predictions on test data

Thank you!

