

Supervised Learning: Classification with Naïve Bayes Classifier

Dr. Yuzana Win (Nagasaki University, Japan)

Lecturer

Department of Computer Engineering and
Information Technology

Lecture Objectives

- To introduce
 - Classification Techniques or Methods
 - What is Naïve Bayes Classifier?
 - How does the Naïve Bayes Classifier works?
 - Pros and Cons

Classification Techniques or Methods

- Naïve Bayes Classifier *
- Support Vector Machines *
- Neural Networks (Deep Learning)
- Decision Tree based Methods
- Rule-based Methods
- K-Nearest Neighbors (kNN) *

Classification Techniques or Methods

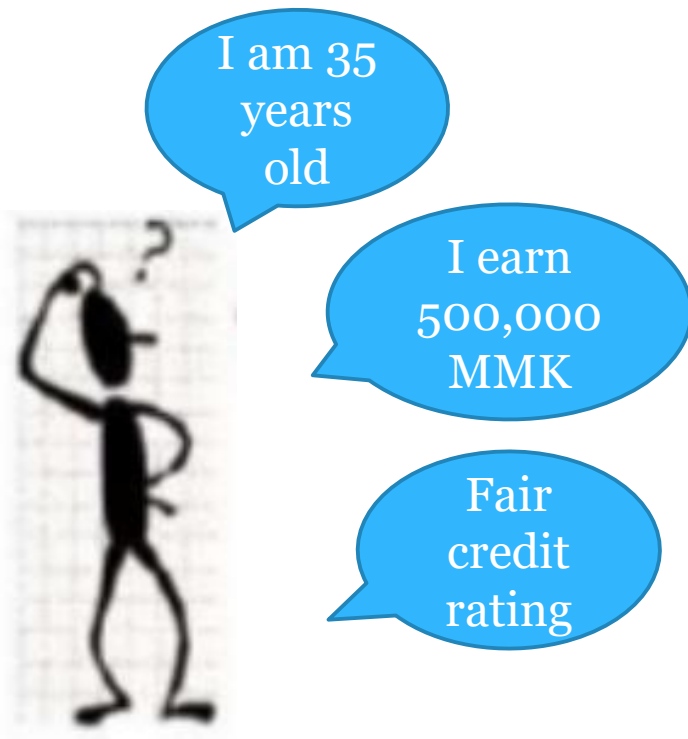
- Naïve Bayes Classifier *
- Support Vector Machines *
- Neural Networks (Deep Learning)
- Decision Tree based Methods
- Rule-based Methods
- K-Nearest Neighbors (KNN) *

Naïve Bayes Classifier

- **Supervised** Learning Method
- Statistical method for classification
- **Probabilistic model** under the Bayes theorem
- Solve problems involving **both categorical and continuous valued attributes**

How Naïve Bayes Classifier Works?

- X : 35 years old customer with an income of 500,000 MMK, and fair credit rating
- H: Hypothesis that the customer will buy a computer



Will he buy a computer?



Bayes Theorem

Probability that the customer is 35 years old, have fair credit rate and earns MMK 500000, given that he has bought computer (**Posterior Probability of A**)

Probability that the customer will buy a computer regardless of age, credit rating, income (**Prior Probability of C**)

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

Probability that the customer will buy a computer given that his age, credit rating and income (**Posterior Probability of C**)

Probability that a person from 35 years old, have fair credit rate and earns MMK 500000 (**Prior Probability of A**)

How Naïve Bayes Classifier Works?

- Given:
 - A doctor knows that **meningitis** causes stiff neck 50% of the time
 - Prior probability of any patient having meningitis is 1/50,000
 - Prior probability of any patient having stiff neck is 1/20

- If a patient has stiff neck, what's the probability he/she has meningitis?

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

Naïve Bayes Classifier

- A: Set of tuples
 - Each Tuple is an ‘n’ dimensional attribute vector
 - $X: (x_1, x_2, x_3, \dots, x_n)$
 - Where x_i is the value of attribute A_i
- Let there are ‘m’ Classes:
 - $C_1, C_2, C_3, \dots, C_m$
- Bayesian classifier predicts X belongs to Class C_i if
 - $P(C_i | X) > P(C_j | X)$ for $1 \leq j \leq m, j \neq i$
- Maximum Posterior Hypothesis
 - $$P(C_i | X) = \frac{P(X | C_i) P(C_i)}{P(X)}$$
 - Maximize $P(X | C_i) P(C_i)$ as $P(X)$ is constant

Naïve Bayes Classifier

- With many attributes, it is computationally expensive to Calculate $P(X|C_i)$
- Naïve Assumption of “class conditional independence”

$$\begin{aligned}
 P(X | C_i) &= P(x_1, x_2, \dots, x_n | C_i) \\
 &= P(x_1 | C_i) * P(x_2 | C_i) * \dots * P(x_n | C_i) \\
 &= \prod_{k=1}^n P(x_k | C_i)
 \end{aligned}$$

Naïve Bayes Classifier

- To calculate, $P(x_k|C_i)$
- $A_k \Rightarrow$ categorical:

$$P(x_k|C_i) = \frac{\text{the number of tuples of class } C_i \text{ in } D \text{ having the value } x_k \text{ for } A_k}{\text{the number of tuples of class } C_i \text{ in } D.}$$

- $A_k \Rightarrow$ continuous:

A continuous-valued attribute is typically assumed to have a Gaussian distribution with a mean μ and standard deviation σ

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

$$P(x_k|C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i}).$$

How to Solve “ZERO” Problem?

- If there is a Class, C_i and X has an attribute value x_k , such that none of the samples in C_i has that attribute value, how to solve this problem?
- In that case, we solve **the probability estimation equation** as below

$$\text{Original : } P(A_i | C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace : } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

$$\text{m - estimate : } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m}$$

c : number of classes

p : prior probability

m : parameter

Example of Naïve Bayes Classifier I

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

A: attributes

M: mammals

N: non-mammals

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

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

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

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

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

$P(A|M)P(M) > P(A|N)P(N)$

=> **Mammals**

Example of Naïve Bayes Classifier II

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

The Evidence relates all attributes without Exceptions.

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

Evidence E

Probability of class "yes"

$$\begin{aligned}
 \Pr[\text{yes} | E] &= \Pr[\text{Outlook} = \text{Sunny} | \text{yes}] \\
 &\quad \times \Pr[\text{Temperature} = \text{Cool} | \text{yes}] \\
 &\quad \times \Pr[\text{Humidity} = \text{High} | \text{yes}] \\
 &\quad \times \Pr[\text{Windy} = \text{True} | \text{yes}] \\
 &\quad \times \frac{\Pr[\text{yes}]}{\Pr[E]} \\
 &= \frac{\frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}}{\Pr[E]}
 \end{aligned}$$

Outlook			Temperature			Humidity			Windy			Play	
Yes	No		Yes	No		Yes	No		Yes	No	Yes	No	
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								

Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

For compute prediction for new day:

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

Likelihood of the two classes

$$\text{For "yes"} = 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$\text{For "no"} = 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0206$$

Conversion into a probability by normalization:

$$P(\text{"yes"}) = 0.0053 / (0.0053 + 0.0206) = \mathbf{0.205}$$

$$P(\text{"no"}) = 0.0206 / (0.0053 + 0.0206) = \mathbf{0.795}$$

Example of Naïve Bayes Classifier III

Training dataset

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
30...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

$P(\text{buys_computer} = \text{„yes“}) = 9/14$

$P(\text{buys_computer} = \text{„no“}) = 5/14$

Class:

C1:buys_computer= 'yes'

C2:buys_computer= 'no'

Data sample

X =(age<=30,
Income=medium,
Student=yes
Credit_rating=
Fair)

Example of Naïve Bayes Classifier III

- Compute $P(X|C_i)$ for each class

$$P(\text{age}=\text{"<30"} \mid \text{buys_computer}=\text{"yes"}) = 2/9=0.222$$

$$P(\text{age}=\text{"<30"} \mid \text{buys_computer}=\text{"no"}) = 3/5 = 0.6$$

$$P(\text{income}=\text{"medium"} \mid \text{buys_computer}=\text{"yes"}) = 4/9 = 0.444$$

$$P(\text{income}=\text{"medium"} \mid \text{buys_computer}=\text{"no"}) = 2/5 = 0.4$$

$$P(\text{student}=\text{"yes"} \mid \text{buys_computer}=\text{"yes"}) = 6/9 = 0.667$$

$$P(\text{student}=\text{"yes"} \mid \text{buys_computer}=\text{"no"}) = 1/5 = 0.2$$

$$P(\text{credit_rating}=\text{"fair"} \mid \text{buys_computer}=\text{"yes"}) = 6/9 = 0.667$$

$$P(\text{credit_rating}=\text{"fair"} \mid \text{buys_computer}=\text{"no"}) = 2/5 = 0.4$$

- $X = (\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair})$

$$\mathbf{P(X|C_i)} : P(X|\text{buys_computer}=\text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$$

$$P(X|\text{buys_computer}=\text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

$$\mathbf{P(X|C_i) * P(C_i)} : P(X|\text{buys_computer}=\text{"yes"}) * P(\text{buys_computer}=\text{"yes"}) = \mathbf{0.028}$$

$$P(X|\text{buys_computer}=\text{"no"}) * P(\text{buys_computer}=\text{"no"}) = \mathbf{0.007}$$

- X belongs to class "buys_computer=yes"

Application Area of Naïve Bayes Classifier

- Text Classification
- Spam Filtering
- Hybrid Recommender System
- Online Application
 - Cyberbullying Detection System
 - Simple Emotion Modeling

Text Classification with Naïve Bayes Classifier

- Naïve Bayes classification method based on Bayes rule
- Relies on very simple representation of document
 - **Bag of words**

The bag of words representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.

The bag of words representation

I **love** this movie! It's **sweet**, but with **satirical** humor. The dialogue is **great** and the adventure scenes are **fun...** It manages to be **whimsical** and **romantic** while **laughing** at the conventions of the fairy tale genre. I would **recommend** it to just about anyone. I've seen it **several** times, and I'm always **happy** to see it **again** whenever I have a friend who hasn't seen it yet.

The bag of words representation

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

Text Classification with Naïve Bayes Classifier

- **Naïve Bayes Text Classifier Algorithm** is used which is a probabilistic model.

$$c = \arg \max \Pr(C = c_j) \prod_{i=1}^{|A|} \Pr(A_i = a_i | C = c_j)$$

$$\Pr(A_i = a_i | C = c_j) = \frac{n_{ij} + \lambda}{n_j + \lambda m_i}$$

- Let A be the attributes of the training dataset and C be the class attribute
- n_{ij} be the number of examples that have both $A_i = a_i$ and $C = c_j$
- n_j be the total number of examples with $C=c_j$ in the training data set
- m_i is the number of values of attribute A_i
- λ is a multiplicative factor

Text Classification with Naïve Bayes Classifier

- Input:
 - a document d (Cyber-bullied terms document)
 - classes $C = \{c_1, c_2\} \Rightarrow$ bully or non-bully
 - A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$
- Output:
 - a learned classifier $\gamma: d \rightarrow c$

Doc	Text	Class
1.	['not', 'lazi', 'jobless', 'just', 'f**k', 'poor']	-
2.	['come', 'speed', 'good', 'movi']	+
3.	['dumbass', 'inde']	-
4.	['nah', 'neither', 'good']	+
5.	['agre', 'movi', 's**t', 'act']	-

Unique words 17

<not, lazi, jobless, just, f**k, poor, come, speed, good, movi, dumbass, inde, nah, neither, agre, s**t, act>

Doc	not	lazi	jobless	just	f**k	poor	come	speed	good	movi	dumbass	inde	nah	neither	agre	s**t	act	Class
1.	1	1	1	1	1	1												-
2.							1	1	1	1								+
3.											1	1						-
4.									1				1	1				+
5.										1					1	1	1	-

The probabilities are calculated using **naïve bayes** algorithm

Doc	not	lazi	jobless	just	f**k	poor	come	speed	good	movi	dumbass	inde	nah	neither	agre	s**t	act	Class
2.							1	1	1	1								+
4.									1				1	1				+

$$P(+)=2/5=0.4$$

$$P(\text{not}|+)=0+1/7+17=0.0416$$

$$P(\text{lazi}|+)=0+1/7+17=0.0416$$

$$P(\text{jobless}|+)=0+1/7+17=0.0416$$

$$P(\text{just}|+)=0+1/7+17=0.0416$$

$$P(\text{f**k}|+)=0+1/7+17=0.0416$$

$$P(\text{poor}|+)=0+1/7+17=0.0416$$

$$P(\text{come}|+)=1+1/7+17=0.0833$$

$$P(\text{speed}|+)=1+1/7+17=0.0833$$

$$P(\text{good}|+)=2+1/7+17=0.125$$

$$P(\text{movi}|+)=1+1/7+17=0.0833$$

$$P(\text{dumbass}|+)=0+1/7+17=0.0416$$

$$P(\text{inde}|+)=0+1/7+17=0.0416$$

$$P(\text{nah}|+)=1+1/7+17=0.0833$$

$$P(\text{neither}|+)=1+1/7+17=0.0833$$

$$P(\text{agre}|+)=0+1/7+17=0.0416$$

$$P(\text{s**t}|+)=0+1/7+17=0.0416$$

$$P(\text{act}|+)=0+1/7+17=0.0416$$

Doc	not	lazi	jobless	just	f**k	poor	come	speed	good	movi	dumbass	inde	nah	neither	agre	s**t	act	Class
1.	1	1	1	1	1	1												-
3.											1	1						-
5.										1					1	1	1	-

$$P(-) = 3/5 = 0.6$$

$$P(\text{not}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{lazi}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{jobless}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{just}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{f**k}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{poor}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{come}|-) = 0+1/12+17 = 0.0344$$

$$P(\text{speed}|-) = 0+1/12+17 = 0.0344$$

$$P(\text{good}|-) = 0+1/12+17 = 0.0344$$

$$P(\text{movi}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{dumbass}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{inde}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{nah}|-) = 0+1/12+17 = 0.0344$$

$$P(\text{neither}|-) = 0+1/12+17 = 0.0344$$

$$P(\text{agre}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{s**t}|-) = 1+1/12+17 = 0.0689$$

$$P(\text{act}|-) = 1+1/12+17 = 0.0689$$

Let's classify a new sentence according to:

$$V_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_{w \in \text{words}} P(w|v_j)$$

Where V stands for “class”

Example Sentence : “dumbass inde”

$$\begin{aligned} V = +; & P(+)\ P(\text{dumbass}|+)\ P(\text{inde}|+) \\ & = (0.4)\ (0.0416)\ (0.0416) \\ & = 6.92224 \times 10^{-4} \end{aligned}$$

$$\begin{aligned} V = -; & P(-)\ P(\text{dumbass}|-)\ P(\text{inde}|-) \\ & = (0.6)\ (0.0689)\ (0.0689) \\ & = 2.848326 \times 10^{-3} \end{aligned}$$

Greater

Therefore, the sentence is **negative**(contained **bullied** word).

Naïve Bayes Coding

```

train = [('not lazi jobless just fuck poor', 'neg'),
         ('come speed good movi', 'pos'),
         ('dumbass inde', 'neg'),
         ('nah neither good', 'pos'),
         ('agre movi shit act', 'neg'),
         ('yeah justin stop fuck forum troll', 'neg'),
         ('yeah like less blu ray version bought', 'pos'),
         ('yeah lot problem associ pretti much know elimin peopl think thatll never work', 'pos')]

from textblob.classifiers import NaiveBayesClassifier
cl = NaiveBayesClassifier(train)

a=cl.classify("come speed good movi")
print ("The input sentence has " + a + " sense")

b=cl.classify("jobless life give shit")
print ("The input sentence has " + b + " sense")

c=cl.classify("lazi peopl usually think good problem never work")
print ("The input sentence has " + c + " sense")

```

Output

```

===== RESTART: I:\My data\BE thesis\Step 2\Testing\Nb Testing\4.py =====
The input sentence has pos sense
The input sentence has neg sense
The input sentence has pos sense
>>> |

```

Advantages and Disadvantages of Naïve Bayes Classifier

- Advantages :
 - Easy to implement.
 - Good results obtained in most of the cases
- Disadvantages
 - Class conditional independence assumption, therefore loss of accuracy
 - Practically, dependencies exist among variables
 - E.g., hospital patients' profile: age, family history etc.,
 - Symptoms: fever, cough etc., Disease: lung cancer, diabetes etc.,

Conclusion

- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Widely used in text classification and spam filtering

Next Week Lecture

- Supervised Learning: Classification with Decision Tree