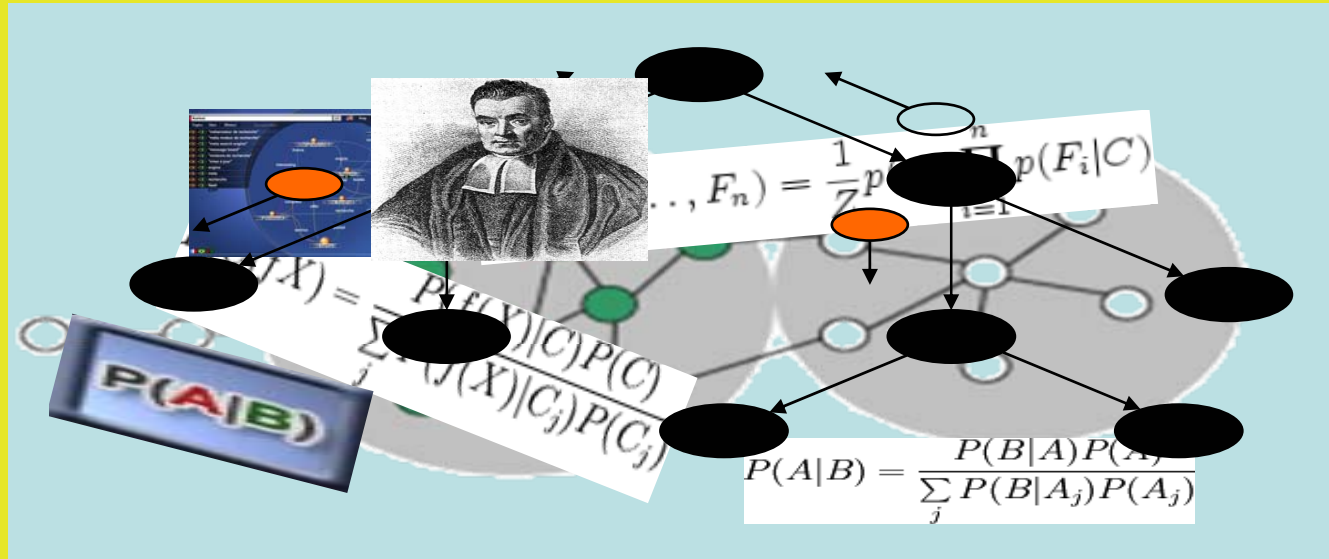


Statistical Data Mining



Application of Bayesian Statistics in Classification

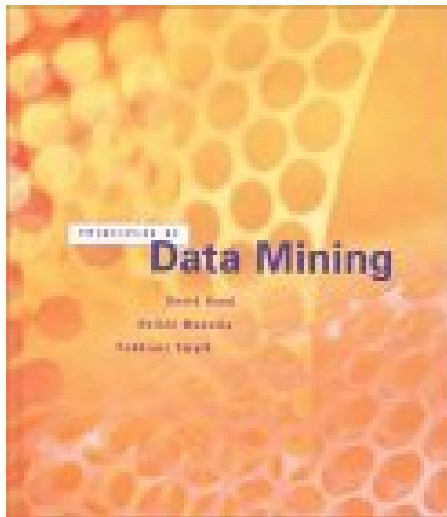
Naïve Bayes

Professor Dr. Gholamreza Nakhaeizadeh

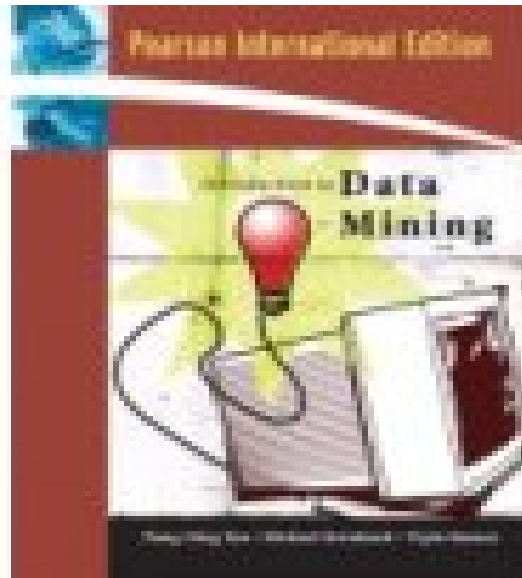
Content

- Literature used
- Introduction
- Bayes Theorem
- Application of Bayes-Theorem in classification
- Naïve Bayes

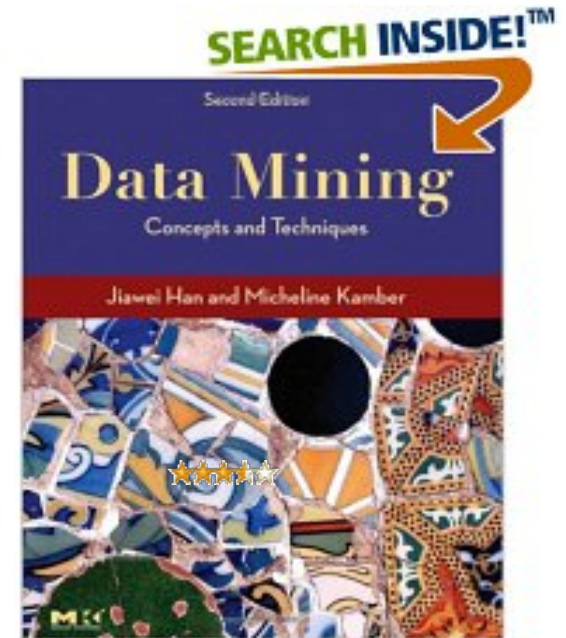
Literatur used



Principles of Data Mining
[David J. Hand](#), [Heikki Mannila](#),
[Padhraic Smyth](#)



Pang-Ning Tan,
Michael Steinbach,
Vipin Kumar



[Jiawei Han](#) and
[Micheline Kamber](#)

Naïve Bayes

Introduction

- In classification tasks, assigning the class value to an observation is often stochastic and not deterministic

Reasons:

- Noisy data
- Some of the relevant attributes are not considered that are stochastic

Example:

To handle such situations, one needs stochastic approaches



- Naïve Bayes
- Logistic Regression
-

income	car	gender	class
2000	yes	F	good
2000	yes	F	bad

Naïve Bayes

Bayes-Theorem

X and Y : Random variables

P (X,Y) : joint probability of X and Y

P(X|Y) : Conditional probability of X given Y

Then: $P(X,Y) = P(Y|X) \cdot P(X) = P(X|Y) \cdot P(Y)$

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)} \quad (1)$$

Relation (1) is known as Bayes-Theorem

Naïve Bayes

Application of Bayes-Theorem in classification

X : Attributes Vector, Y: Class Vector

X and Y : Random variables

P (Y|X) : Posterior probability, P(Y) : Prior probability

Bayesian Classification Task

- **Building the Classifier:**
Learning $P(Y|X)$ by using data on X and Y
- **Classification of new tuples:**
To each new tuple X' assign the class value that maximizes $p(Y'|X)$

Naïve Bayes

Application of Bayes-Theorem in classification

How can we compute $P(Y|X)$?

By using the Bayes-Theorem
$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

- $P(x)$ is independent of Y and can be ignored
- Computing of $P(Y)$ can be done easily by using the observations on Y
- To compute $P(X|Y)$ there are different alternative:

Naïve Bayes is one of them

Naïve Bayes

Naïve Bayes Classifier

- **Goal:** Estimating the class conditional probability
- **(Naïve) Assumption:** given the class label Y the attributes are conditionally independent

$$P(X|Y=y) = \prod_{i=1}^m P(X_i | Y=y), \quad X = (X_1, X_2, \dots, X_m) \quad (2)$$

Instead of computing the joint conditional probability of X , it is just necessary to compute the probability of each X_i given Y

Thus, from (1) and (2) we have

$$P(Y|X) = \frac{P(Y) \cdot \prod_{i=1}^m P(X_i | Y)}{P(X)} \quad (3)$$

Naïve Bayes Classification Rule:
Assign to the new Vector X' the class that maximizes the numerator of (3)

Naïve Bayes

Naïve Bayes Classifier

$$P(Y|X) = \frac{P(Y) \cdot \prod_{i=1}^m P(X_i | Y)}{P(X)}$$

$\alpha = 1/P(X) \longrightarrow \text{constant}$



$$P(Y|X) = \alpha \cdot P(Y) \cdot \prod_{i=1}^m P(X_i | Y) \quad (4)$$

Naïve Bayes

Example

Source: <http://www-users.itlabs.umn.edu/classes/Spring-2006/csci5523/index.php?page=lecture%20slides>

<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

Y = Evade

X= (Refund, Marital Status, Taxable Income)

New Record X':

**(Refund= no,
marital status=married,
taxable income = \$ 120 K)**

Based on the training data, we compute

$P(\text{yes} \mid X')$ and $P(\text{no} \mid X')$

The new record is classified as “yes” if

$P(\text{yes} \mid X') > P(\text{no} \mid X')$

Otherwise it is classified as “no”

Naïve Bayes

Estimation conditional probabilities $P(X|Y)$

A. Attribute is nominal

1. Choose a value of Y
2. Determine the values of the nominal attribute X that corresponds to this selected value of Y
3. Determine the fraction of these values

Example:

1. We choose $Y = \text{Evade} = \text{No}$

2. The values of the attribute "Refund" that corresponds to $\text{Evade} = \text{no}$ are:

Refund	Evade
Yes	No
No	No
No	No
Yes	No
No	No
Yes	No
No	No

3. Determining of the fractions

$$P(X = \text{Yes} \mid Y = \text{No}) = 3/7$$

$$P(X = \text{No} \mid Y = \text{No}) = 4/7$$

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

Naïve Bayes

Estimation conditional probabilities $P(X|Y)$

B. Attribute is continuous-valued

Alternative 1:
Discretization of the continuous-valued Attribute. The rest of the procedure is similar to the case A

Alternative 2:
Assume a certain conditional distribution for the continuous-valued attribute(e.g. normal distribution)

$$P(X_i | Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

The distributions parameters can be estimated by using the observations on X and Y

Naïve Bayes

Estimation conditional probabilities P (X|Y)

B. Attribute is continuous-valued

Example:

Taxable Income	Evade
125	No
100	No
70	No
120	No
60	No
220	No
75	No

$$\bar{X} = (125 + 100 + 70 + \dots + 75) / 7 = 110$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$S^2 = 17850 / 6 = 2975$$

$$S = \sqrt{2975} = 54.54$$

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

$$P(X_i | Y = \text{No}) = \frac{1}{\sqrt{2\pi} \cdot 54.54} e^{-\frac{(X_i - 110)^2}{2 \cdot 2975}}$$

$$X_i = 120 \longrightarrow P(X_i | Y = \text{No}) = 0.0072$$

Naïve Bayes

Example:

Determine the class of a new Record X:

(Refund= no, marital status=married, taxable income = \$ 120 K)

$$P(Y = \text{No}) = 7/10 \quad P(Y = \text{Yes}) = 3/10$$

$$P(X | Y = \text{No}) =$$

$$P(\text{Refund} = \text{No} | Y = \text{No}) \cdot P(\text{status} = \text{married} | Y = \text{No}) \cdot P(\text{T. Income} = 120 | Y = \text{No})$$

Refund	Evade
Yes	No
No	No
No	No
Yes	No
No	No
Yes	No
No	No

$$\rightarrow P(\text{Refund} = \text{No} | Y = \text{No}) = 4/7$$

$$P(\text{T. Income} = 120 | Y = \text{No}) = 0.0072$$

(see the last slide)

M. status	Evade
single	No
married	No
single	No
married	No
married	No
divorced	No
married	No

$$\rightarrow P(\text{Status} = \text{married} | Y = \text{No}) = 4/7$$

$$P(X | Y = \text{No}) = 4/7 \cdot 4/7 \cdot 0.0072 = 0.0024$$

$$\text{From (4) we have: } P(\text{No} | X) = \alpha \quad P(Y = \text{No}) \cdot P(X | Y = \text{No}) = 7/10 \cdot 0.0024 \quad \alpha = 0.0016 \quad \alpha$$

$$\text{Using the same method } \rightarrow P(\text{Yes} | X) = 0 \rightarrow P(\text{No} | X) > P(\text{Yes} | X)$$

The class value of the new record is computed as No

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

Naïve Bayes

Strength and Weakness

- Robust to noise and irrelevant attributes
- Independence assumption may not hold for some attributes