

Naïve Bayes Classifier

Outline

- Background
- Probability Basics
- Probabilistic Classification
- Naïve Bayes
- Example: Play Tennis
- Relevant Issues
- Conclusions

Background

- There are three methods to establish a classifier
 - a) Model a classification rule directly*

Examples: k-NN, decision trees, perceptron, SVM
 - b) Model the probability of class memberships given input data*

Example: multi-layered perceptron with the cross-entropy cost
 - c) Make a probabilistic model of data within each class*

Examples: naive Bayes, model based classifiers
- *a)* and *b)* are examples of **discriminative** classification
- *c)* is an example of **generative** classification
- *b)* and *c)* are both examples of **probabilistic** classification

Probability Basics

- Prior, conditional and joint probability
 - Prior probability: $P(X)$
 - Conditional probability: $P(X_1 | X_2), P(X_2 | X_1)$
 - Joint probability: $\mathbf{X} = (X_1, X_2), P(\mathbf{X}) = P(X_1, X_2)$
 - Relationship: $P(X_1, X_2) = P(X_2 | X_1)P(X_1) = P(X_1 | X_2)P(X_2)$
 - Independence: $P(X_2 | X_1) = P(X_2), P(X_1 | X_2) = P(X_1), P(X_1, X_2) = P(X_1)P(X_2)$
- Bayesian Rule

$$P(C | \mathbf{X}) = \frac{P(\mathbf{X} | C)P(C)}{P(\mathbf{X})} \quad \text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

Probabilistic Classification

- Establishing a probabilistic model for classification

- Discriminative model

$$P(C | \mathbf{X}) \quad C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_n)$$

- Generative model

$$P(\mathbf{X} | C) \quad C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_n)$$

- MAP classification rule

- **MAP: M**aximum **A** Posterior

- Assign x to c^* if $P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x}) \quad c \neq c^*, c = c_1, \dots, c_L$

- Generative classification with the MAP rule

- Apply Bayesian rule to convert:
$$P(C | \mathbf{X}) = \frac{P(\mathbf{X} | C)P(C)}{P(\mathbf{X})} \propto P(\mathbf{X} | C)P(C)$$

Naïve Bayes

- Bayes classification

$$P(C | \mathbf{X}) \propto P(\mathbf{X} | C)P(C) = P(X_1, \dots, X_n | C)P(C)$$

Difficulty: learning the joint probability $P(X_1, \dots, X_n | C)$

- Naïve Bayes classification

- Making the assumption that **all input attributes are independent**

$$\begin{aligned} P(X_1, X_2, \dots, X_n | C) &= \underline{P(X_1 | X_2, \dots, X_n; C)} P(X_2, \dots, X_n | C) \\ &= \underline{P(X_1 | C)} \underline{P(X_2, \dots, X_n | C)} \\ &= \underline{P(X_1 | C) P(X_2 | C) \dots P(X_n | C)} \end{aligned}$$

- MAP classification rule

$$[P(x_1 | c^*) \dots P(x_n | c^*)]P(c^*) > [P(x_1 | c) \dots P(x_n | c)]P(c), \quad c \neq c^*, c = c_1, \dots, c_L$$

Naïve Bayes

- Naïve Bayes Algorithm (for discrete input attributes)

- Learning Phase: Given a training set S ,

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in S ;

For every attribute value a_{jk} of each attribute x_j ($j = 1, \dots, n; k = 1, \dots, N_j$)

$\hat{P}(X_j = a_{jk} | C = c_i) \leftarrow$ estimate $P(X_j = a_{jk} | C = c_i)$ with examples in S ;

Output: conditional probability tables; for x_j , $N_j \times L$ elements

- Test Phase: Given an unknown instance $\mathbf{X}' = (a'_1, \dots, a'_n)$

Look up tables to assign the label c^* to \mathbf{X}' if

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_n | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_n | c)] \hat{P}(c), \quad c \neq c^*, c = c_1, \dots, c_L$$

Example

- Example: Play Tennis

PlayTennis: training examples

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

Example

- Learning Phase

Outlook	Play=Yes	Play=No
<i>Sunny</i>	2/9	3/5
<i>Overcast</i>	4/9	0/5
<i>Rain</i>	3/9	2/5

Temperature	Play=Yes	Play=No
<i>Hot</i>	2/9	2/5
<i>Mild</i>	4/9	2/5
<i>Cool</i>	3/9	1/5

Humidity	Play=Yes	Play=No
<i>High</i>	3/9	4/5
<i>Normal</i>	6/9	1/5

Wind	Play=Yes	Play=No
<i>Strong</i>	3/9	3/5
<i>Weak</i>	6/9	2/5

$$P(\text{Play=Yes}) = 9/14 \quad P(\text{Play=No}) = 5/14$$

Example

- Test Phase

- Given a new instance,

$\mathbf{x}' = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$

- Look up tables

$$P(\text{Outlook}=\text{Sunny} \mid \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temperature}=\text{Cool} \mid \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Humidity}=\text{High} \mid \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Wind}=\text{Strong} \mid \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

$$P(\text{Outlook}=\text{Sunny} \mid \text{Play}=\text{No}) = 3/5$$

$$P(\text{Temperature}=\text{Cool} \mid \text{Play}=\text{No}) = 1/5$$

$$P(\text{Humidity}=\text{High} \mid \text{Play}=\text{No}) = 4/5$$

$$P(\text{Wind}=\text{Strong} \mid \text{Play}=\text{No}) = 3/5$$

$$P(\text{Play}=\text{No}) = 5/14$$

- MAP rule

$$P(\text{Yes} \mid \mathbf{x}'): [P(\text{Sunny} \mid \text{Yes})P(\text{Cool} \mid \text{Yes})P(\text{High} \mid \text{Yes})P(\text{Strong} \mid \text{Yes})]P(\text{Play}=\text{Yes}) = 0.0053$$

$$P(\text{No} \mid \mathbf{x}'): [P(\text{Sunny} \mid \text{No})P(\text{Cool} \mid \text{No})P(\text{High} \mid \text{No})P(\text{Strong} \mid \text{No})]P(\text{Play}=\text{No}) = 0.0206$$

Given the fact $P(\text{Yes} \mid \mathbf{x}') < P(\text{No} \mid \mathbf{x}')$, we label \mathbf{x}' to be “No”.

Relevant Issues

- Violation of Independence Assumption
 - For many real world tasks, $P(X_1, \dots, X_n | C) \neq P(X_1 | C) \dots P(X_n | C)$
 - Nevertheless, naïve Bayes works surprisingly well anyway!
- Zero conditional probability Problem
 - If no example contains the attribute value $X_j = a_{jk}$, $\hat{P}(X_j = a_{jk} | C = c_i) = 0$
 - In this circumstance, $\hat{P}(x_1 | c_i) \dots \hat{P}(a_{jk} | c_i) \dots \hat{P}(x_n | c_i) = 0$ during test
 - For a remedy, conditional probabilities estimated with

$$\hat{P}(X_j = a_{jk} | C = c_i) = \frac{n_c + mp}{n + m}$$

n_c : number of training examples for which $X_j = a_{jk}$ and $C = c_i$

n : number of training examples for which $C = c_i$

p : prior estimate (usually, $p = 1/t$ for t possible values of X_j)

m : weight to prior (number of "virtual" examples, $m \geq 1$)

Relevant Issues

- Continuous-valued Input Attributes
 - Numberless values for an attribute
 - Conditional probability modeled with the normal distribution

$$\hat{P}(X_j | C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$

σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

- Learning Phase: for $\mathbf{X} = (X_1, \dots, X_n)$, $C = c_1, \dots, c_L$
Output: $n \times L$ normal distributions and $P(C = c_i)$ $i = 1, \dots, L$
- Test Phase: for $\mathbf{X}' = (X'_1, \dots, X'_n)$
 - Calculate conditional probabilities with all the normal distributions
 - Apply the MAP rule to make a decision

Conclusions

- Naïve Bayes based on the independence assumption
 - Training is very easy and fast; just requiring considering each attribute in each class separately
 - Test is straightforward; just looking up tables or calculating conditional probabilities with normal distributions
- A popular generative model
 - Performance competitive to most of state-of-the-art classifiers even in presence of violating independence assumption
 - Many successful applications, e.g., spam mail filtering
 - Apart from classification, naïve Bayes can do more...