

Introduction to Big Data and Machine Learning Classification

Dr. Mihail

September 19, 2019

Goal

- Goal of classification: take an input vector x and assign it to one of K discrete classes \mathcal{C}_k where $k = 1, \dots, K$
- The input space is therefore divided into *decision regions* whose boundaries are called “*decision boundaries*” or “*decision surfaces*”
- Here, we will consider linear models where the decision boundaries are linear functions of the input vector “ x ” and hence are defined by $D - 1$ -dimensional hyperplanes within the D -dimensional input space
- Data sets that can be separated exactly by linear decision surfaces are said to be “*linearly separable*”

- For probabilistic models, the most convenient, in the case of two-class problems is the binary representation, in which there is a single target variable $t = \{0, 1\}$
- For $K > 2$ classes, it is convenient to use 1 – of – K coding scheme, in which t is a vector of length K such that if the class is \mathcal{C}_j , then all elements of t_k are zero except t_j .
- For instance if we have 5 classes, the a patter from class 2 would be given by the target vector $t = (0, 1, 0, 0, 0)^T$

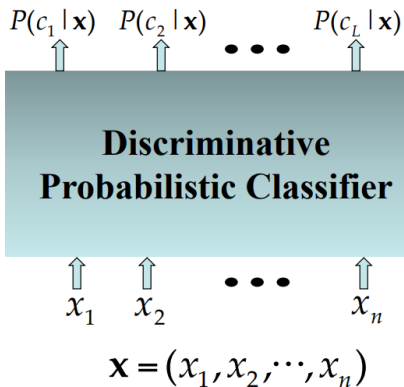
Using Bayes Theorem

- Model posterior class conditional probability: $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$
- **Notice** the denominator is not a function of C
- Prior class distribution: $p(C_k)$
- Class conditional density: $p(x|C_k)$

Discriminative model

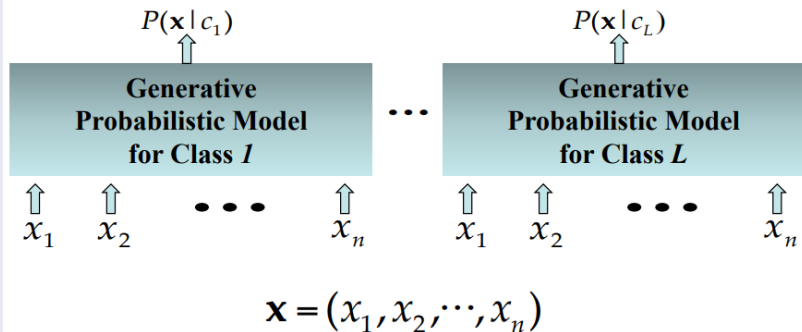
- $P(c|x)$
- To train a discriminative classifier, all training examples of different classes must be jointly used to build up a single discriminative classifier
- Output K probabilities for K class labels in probabilistic classifiers, while a single label is produced by non-probabilistic classifier

Discriminative classifier



- $P(x|c)$, $c = c_1, \dots, c_K$, $x = (x_1, \dots, x_n)$
- K probabilistic models have to be trained **independently**
- Each is trained on only the examples of the same label
- Output K probabilities for a given input with K models
- “Generative” means that model can produce data via distribution sampling

Generative classifier



Maximum a-posteriori (MAP)

- For an input x , find the largest one from K probabilities output by a discriminative probabilistic classifier $P(c_1|x), \dots, P(c_K|x)$
- Assign x to label c^* if $P(c^*|x)$ is the largest
- Generative classification with the MAP rule:

$$P(c_i|x) = \frac{P(x|c_i)P(c_i)}{P(x)} \propto P(x|c_i)P(c_i) \quad (1)$$

Bayes classification

$$P(c|x) \propto P(x|c)P(c) = P(x_1, \dots, x_n|c)P(c) \quad (2)$$

for $c = c_1, \dots, c_K$

Bayes classification

$$P(c|x) \propto P(x|c)P(c) = P(x_1, \dots, x_n|c)P(c) \quad (2)$$

for $c = c_1, \dots, c_K$

Problem

The joint probability $P(x_1, \dots, x_n|c)$ is not feasible to learn.

Bayes classification

$$P(c|x) \propto P(x|c)P(c) = P(x_1, \dots, x_n|c)P(c) \quad (2)$$

for $c = c_1, \dots, c_K$

Problem

The joint probability $P(x_1, \dots, x_n|c)$ is not feasible to learn.

Solution

Assume all input features are class conditionally independent!

Bayes model

$$\begin{aligned}P(x_1, x_2, \dots, x_n | c) &= P(x_1 | x_2, \dots, x_n, c) P(x_2, \dots, x_n | c) \\ &= P(x_1 | c) P(x_2, \dots, x_n | c) \\ &= P(x_1 | c) P(x_2 | c) \dots P(x_n | c)\end{aligned}\tag{3}$$

Discrete valued features

Learning phase: Given a training set S of F features and K classes,

- For each target value of c_i ($c_i = c_1, \dots, c_K$):
 - $\hat{P}(c_i) \leftarrow$ estimate $P(c_i)$ with examples in S
 - For every feature value x_{jk} of each feature x_j ($j = 1, \dots, F; k = 1, \dots, N$):
 - $\hat{P}(x_j = x_{jk} | c_i) \leftarrow$ estimate $P(x_{jk} | c_i)$ with samples in S

Output: $F \times K$ conditional probabilistic (generative) models.

Test phase: Given an unknown instance $x' = (a'_1, \dots, a'_n)$ assign label c^* to x' if

$$[\hat{P}(a'_1 | c^*) \dots \hat{P}(a'_n | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c_i) \dots \hat{P}(a'_n | c_i)] \hat{P}(c_i) \quad (4)$$

for $c_i \neq c^*, c_i = c_1, \dots, c_K$

Example

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

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

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

- Given a new instance, predict its label:

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

- Look up tables:

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

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

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

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

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

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

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

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

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

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

- Make decision with the MAP rule:

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

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

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