

Naïve Bayes Classification

Introduction

- The Bayes Classifier requires probability structure of the problem to be known.
- Density estimation (using non-parametric or parametric methods) is one way to handle the problem.
- There are several problems

Problems with density estimation

- Large datasets are needed.
- Numeric valued features are required.
- In practice these two may not be satisfied.

How to overcome the problem

- One has to work with the given data set.
- So, Probability estimations needs to be done using the given data only.
- Often marginal probabilities can be better estimated than the joint probabilities.
- Also, marginal probabilities are easy to compute.

Play-tennis data

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

- $P(\langle \text{sunny}, \text{cool}, \text{high}, \text{false} \rangle | N) = 0$
- But, $P(\text{sunny} | N) = 3/5$, $P(\text{cool} | N) = 1/5$,
 $P(\text{high} | N) = 4/5$, $P(\text{false} | N) = 2/5$.

- $P(\langle \text{sunny, cool, high, false} \rangle | N) = 0$
- This may be because of the smaller dataset.
- If we increase the dataset size, this may become a positive number.
- This problem is often referred to as “the curse of dimensionality”.

Assumption

- Make the assumption that for a given class, features are independent of each other.
- In practice, this assumption holds very often.
- Then $P(\langle \text{sunny}, \text{cool}, \text{high}, \text{false} \rangle | N) = P(\text{sunny} | N) \cdot P(\text{cool} | N) \cdot P(\text{high} | N) \cdot P(\text{false} | N) = 3/5 \cdot 1/5 \cdot 4/5 \cdot 2/5 = 24/625$.

Naïve Bayesian Classification

- Naïve assumption: for a given class, features are independent of each other

$$P(\langle x_1, \dots, x_k \rangle | C) = P(x_1 | C) \cdot \dots \cdot P(x_k | C)$$

- $P(x_i | C)$ is estimated as the relative freq of samples having value x_i as i -th attribute in class C
- It often makes the problem a feasible and easy one to solve.

Play-tennis example: estimating $P(x_i | C)$

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

$$P(p) = 9/14$$

$$P(n) = 5/14$$

outlook	
$P(\text{sunny} p) = 2/9$	$P(\text{sunny} n) = 3/5$
$P(\text{overcast} p) = 4/9$	$P(\text{overcast} n) = 0$
$P(\text{rain} p) = 3/9$	$P(\text{rain} n) = 2/5$
temperature	
$P(\text{hot} p) = 2/9$	$P(\text{hot} n) = 2/5$
$P(\text{mild} p) = 4/9$	$P(\text{mild} n) = 2/5$
$P(\text{cool} p) = 3/9$	$P(\text{cool} n) = 1/5$
humidity	
$P(\text{high} p) = 3/9$	$P(\text{high} n) = 4/5$
$P(\text{normal} p) = 6/9$	$P(\text{normal} n) = 2/5$
windy	
$P(\text{true} p) = 3/9$	$P(\text{true} n) = 3/5$
$P(\text{false} p) = 6/9$	$P(\text{false} n) = 2/5$

Play-tennis example: classifying X

- An unseen sample $X = \langle \text{rain, hot, high, false} \rangle$
- $P(X|p) \cdot P(p) =$
 $P(\text{rain}|p) \cdot P(\text{hot}|p) \cdot P(\text{high}|p) \cdot P(\text{false}|p) \cdot P(p) = 3/9 \cdot 2/9 \cdot 3/9 \cdot 6/9 \cdot 9/14 =$
 0.010582
- $P(X|n) \cdot P(n) =$
 $P(\text{rain}|n) \cdot P(\text{hot}|n) \cdot P(\text{high}|n) \cdot P(\text{false}|n) \cdot P(n) = 2/5 \cdot 2/5 \cdot 4/5 \cdot 2/5 \cdot 5/14 =$
 0.018286
- Sample X is classified in class n (don't play)

With Continuous features

- In order to use the Naïve Bayes classifier, the features has to be discretized appropriately (otherwise what happens?)

With Continuous features

- In order to use the Naïve Bayes classifier, the features has to be discretized appropriately (otherwise what happens?)
- Height = 4.234 will not occur anywhere in that column; but 4.213, 4.285 may be occurring. If you discretize (eg., rounding) then frequency ratio's are meaningful.
- Clustering of feature values of a feature may be done to achieve a better discretization.