

Machine Learning

Naive Bayes Classifier

(NB)

$$P(y|x_1, x_2, \dots, x_n) = \frac{P(x_1, x_2, \dots, x_n|y)P(y)}{P(x_1, x_2, \dots, x_n)}$$

Dr. Ali Valinejad

valinejad.ir

valinejad@umz.ac.ir

University of Mazandaran

Probabilistic Classification

Data Set:

$(\mathbf{x}^{(i)})^T$	$y^{(i)}$	
$(\mathbf{x}^{(1)})^T$	$y^{(1)}$	
$(\mathbf{x}^{(2)})^T$	$y^{(2)}$	$\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}]^T$
\vdots	\vdots	
$(\mathbf{x}^{(m)})^T$	$y^{(m)}$	
		$y^{(i)} \in \{c_1, c_2, \dots, c_K\}$

Problem statement:

Given features (x_1, x_2, \dots, x_n)

Predict a label $y = c_1, c_2, \dots, c_K$

A good strategy is to predict:

$$\underset{y}{\operatorname{argmax}} P(y|x_1, x_2, \dots, x_n)$$

Car Theft Example

Data Set:

Example No.	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

Features:

Color, Type, Origin

$(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic})$

$(x_1: \text{Color}, x_2: \text{Type}, x_3: \text{Origin})$

Label:

Stolen (can be either Yes or No). $(y = \text{Yes or No})$

Car Theft Example

Problem statement:

Given features x_1, x_2, \dots, x_n

Predict a label $y = \text{Yes}, \text{No}$

if $x_1 = \text{Red}$, $x_2 = \text{SUV}$, $x_3 = \text{Domestic}$, **predict a label for Stolen** ($y = \text{Yes or No}$)

$$\underset{y}{\operatorname{argmax}} P(y | x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic})$$

Probabilistic Classification

Establishing a probabilistic model for classification

❖ Discriminative model

$$P(y|X) \quad y = c_1, c_2, \dots, c_K, X = (x_1, x_2, \dots, x_n)$$

❖ Generative model

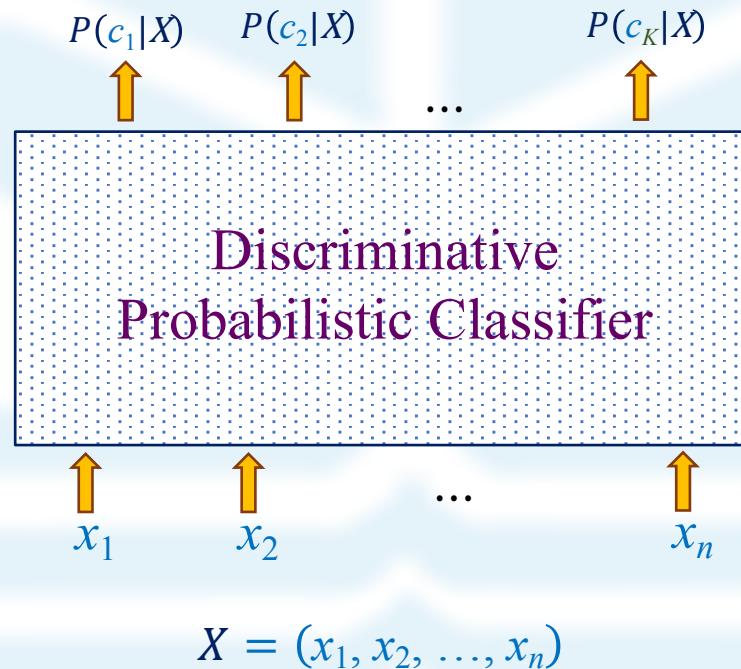
$$P(X|y) \quad y = c_1, c_2, \dots, c_K, X = (x_1, x_2, \dots, x_n)$$

Probabilistic Classification

Establishing a probabilistic model for classification

❖ Discriminative model

$$P(y|X) \quad y = c_1, c_2, \dots, c_K, X = (x_1, x_2, \dots, x_n)$$

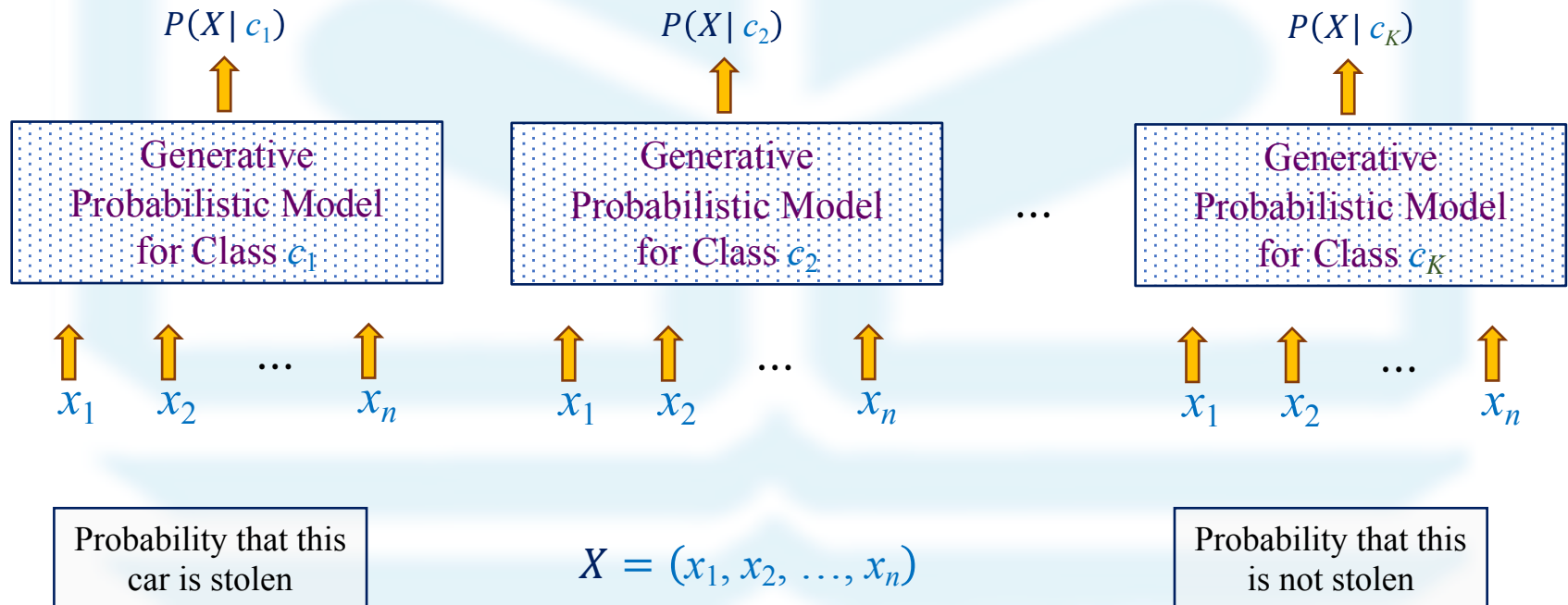


Probabilistic Classification

Establishing a probabilistic model for classification

❖ Generative model

$$P(X|y) \quad y = c_1, c_2, \dots, c_K, X = (x_1, x_2, \dots, x_n)$$



The Bayes Classifier

Establishing a probabilistic model for classification

❖ Generative model

$$P(X|y) \quad y = c_1, c_2, \dots, c_K, X = (x_1, x_2, \dots, x_n)$$

$$P(y|X) = \frac{P(X|y)P(y)}{P(X)}$$

Normalization Constant

Likelihood Prior

$$P(y|x_1, x_2, \dots, x_n) = \frac{P(x_1, x_2, \dots, x_n|y)P(y)}{P(x_1, x_2, \dots, x_n)}$$

Why did this help?

We might be able to specify how features are “generated” by the class label.

The Bayes Classifier

$$\begin{aligned} P(y = c | x_1, x_2, \dots, x_n) &= \frac{P(x_1, x_2, \dots, x_n | y = c) P(y = c)}{P(x_1, x_2, \dots, x_n)} \\ &= \frac{P(x_1, x_2, \dots, x_n | y = c) P(y = c)}{P(x_1, x_2, \dots, x_n | y = \text{Yes}) P(y = \text{Yes}) + P(x_1, x_2, \dots, x_n | y = \text{No}) P(y = \text{No})} \end{aligned}$$

$$\begin{aligned} &P(y = \text{Yes} | x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic}) \\ &= \frac{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{Yes}) P(y = \text{Yes})}{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic})} \\ &= \frac{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{Yes}) P(y = \text{Yes})}{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{Yes}) P(y = \text{Yes}) + P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{No}) P(y = \text{No})} \end{aligned}$$

$$\begin{aligned} &P(y = \text{No} | x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic}) \\ &= \frac{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{No}) P(y = \text{No})}{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic})} \\ &= \frac{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{No}) P(y = \text{No})}{P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{No}) P(y = \text{No}) + P(x_1 = \text{Red}, x_2 = \text{SUV}, x_3 = \text{Domestic} | y = \text{Yes}) P(y = \text{Yes})} \end{aligned}$$

To classify, we'll simply compute these two probabilities and predict based on which one is greater

The Bayes Classifier

$$P(y|X) = \frac{P(X|y)P(y)}{P(X)}$$

Likelihood Prior

Normalization Constant

$X = (x_1, x_2, \dots, x_n)$

For the Bayes classifier, we need to “**learn**” two functions,

- the **likelihood**, $P(X|y)$
- the **prior**, $P(y)$

Model Parameters

- 1- How many parameters are required to specify the **prior**?
- 2- How many parameters are required to specify the **likelihood**?

The problem with explicitly modeling $P(X|y)$ is that there are usually too many parameters:

We'll run out of space

We'll run out of time

And we'll need tons of training data (which is usually not available)

MAP classification rule

MAP classification rule

- **MAP: M**aximum **A** Posterior
- Assign $X = (x_1, x_2, \dots, x_n)$ to c^* if

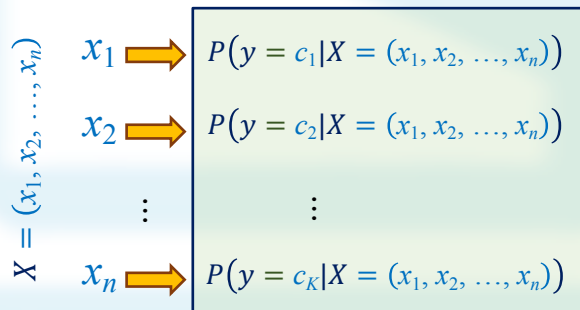
$$P(y = c^* | X = (x_1, x_2, \dots, x_n)) > P(y = c | X = (x_1, x_2, \dots, x_n)) \\ c \neq c^*, c = c_1, c_2, \dots, c_K$$

Algorithm. *MAP classification rule*

Step1: Compute $P(y = c_i | X = (x_1, x_2, \dots, x_n))$ for $i=1, 2, \dots, K$

Step2: Compute $c^* = \operatorname{argmax}_{c=c_1, c_2, \dots, c_K} P(y = c | X = (x_1, x_2, \dots, x_n))$

Step3: Assign $X = (x_1, x_2, \dots, x_n)$ to c^* as it's predicted label

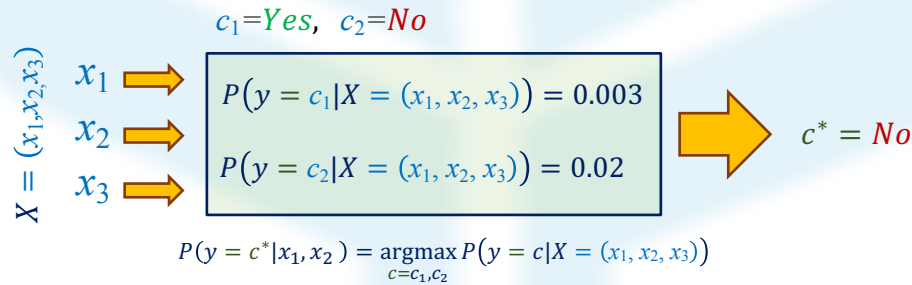


$$P(y = c^* | x_1, x_2, \dots, x_n) = \operatorname{argmax}_{c=c_1, c_2, \dots, c_K} P(y = c | X = (x_1, x_2, \dots, x_n))$$

MAP classification rule

Example:

Assume $X = (x_1=\text{Red}, x_2=\text{SUV}, x_3=\text{Domestic})$,
if $P(y = \text{Yes} | x_1=\text{Red}, x_2=\text{SUV}, x_3=\text{Domestic}) = 0.003$
and
 $P(y = \text{No} | x_1=\text{Red}, x_2=\text{SUV}, x_3=\text{Domestic}) = 0.02$
then
 $c^* = \text{No}$



Method of *Generative classification* with the **MAP** rule

1. Apply Bayesian rule to convert them into posterior probabilities

$$P(y = c_i | x_1, x_2, \dots, x_n) = \frac{P(x_1, x_2, \dots, x_n | y = c_i)P(y = c_i)}{P(x_1, x_2, \dots, x_n)} \propto P(x_1, x_2, \dots, x_n | y = c_i)P(y = c_i)$$

for $i=1, 2, \dots, K$

2. Then apply the MAP rule:

$$P(y = c^* | x_1, x_2, \dots, x_n) > P(y = c | x_1, x_2, \dots, x_n) \quad c \neq c^*, \quad c = c_1, c_2, \dots, c_K$$

Algorithm. *Generative classification* with the **MAP** rule

Step1: Compute $P(y = c_i)$ for $i=1, 2, \dots, K$

Step2: Compute $P(x_1, x_2, \dots, x_n | y = c_i)$ for $i=1, 2, \dots, K$

Step3: Apply Bayesian rule to compute $P(y = c_i | X = (x_1, x_2, \dots, x_n))$ for $i=1, 2, \dots, K$

Step4: Compute $c^* = \underset{c=c_1, c_2, \dots, c_K}{\operatorname{argmax}} P(y = c | X = (x_1, x_2, \dots, x_n))$

Step5: Assign $X = (x_1, x_2, \dots, x_n)$ to c^* as it's predicted label

Bayes classification

$$P(y|X) \propto P(X|y)P(y) = P(x_1, x_2, \dots, x_n|y)P(y)$$

Difficulty

learning the joint probability

The Naïve Bayes Model

- The *Naïve Bayes Assumption*:

Assume that all features are *independent* **given the class label y**

$$P(x_1, x_2, \dots, x_n | y) = \prod_{i=1}^n P(x_i | y)$$

Naïve Bayes Algorithm (for *discrete valued Features*)

Input: Data set $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}$, for $i=1,2,\dots, m$, such that $\mathbf{x}^{(i)}, y^{(i)} \in \{c_1, c_2, \dots, c_K\}$ and each feature $x_j^{(i)}$ can be take it's discrete value from set $\{x_{j1}^{(i)}, x_{j2}^{(i)}, \dots, x_{jn_j}^{(i)}\}$, for $j = 1, 2, \dots, n$.

Output:

- Given an unknown instance $X^{new} = (x_1^{new}, x_2^{new}, \dots, x_n^{new})$, predict it's label.

Training Phase:

- Estimate $P(y = c_i)$ with examples in D for $i=1,2,\dots, K$:

$$\hat{P}(y = c_i) \approx P(y = c_i)$$

- For every feature value x_{jk} of each feature x_j , estimate $P(x_j = x_{jk} | y = c_i)$ for $k=1,2,\dots, n_j$ with examples in D:

$$\hat{P}(x_j = x_{jk} | y = c_i) \approx P(x_j = x_{jk} | y = c_i)$$

Output of training phase: Conditional probability table for each feature x_j , for $j = 1, 2, \dots, n$.

Testing Phase:

Given an unknown instance $X^{new} = (x_1^{new}, x_2^{new}, \dots, x_n^{new})$, predict it's label:

- Look up tables to assign the label c^* to X^{new} :

$$c^* = \underset{c=c_1, c_2, \dots, c_K}{\operatorname{argmax}} \quad P(y = c) \prod_{j=1}^n \hat{P}(x_j^{new} | y = c)$$

Training is easy, just
create probability tables.

Classification is easy, just
multiply probabilities

Naïve Bayes Algorithm

(for *discrete valued Features*)

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

$$P(\text{Outlook} = \text{Sunny}, \text{Play} = \text{Yes}) = \frac{\#(\text{Outlook}=\text{Sunny}, \text{Play}=\text{Yes})}{\# \text{ total samples}} = \frac{2}{14}$$

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

$$P(\text{Outlook} = \text{Sunny}, \text{Play} = \text{No}) = \frac{\#(\text{Outlook}=\text{Sunny}, \text{Play}=\text{No})}{\# \text{ total samples}} = \frac{3}{14}$$

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

Naïve Bayes Algorithm

(for *discrete valued Features*)

Example: Play Tennis

Cont.

Training Phase

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

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

Temperature

$$P(\text{Hot}|\text{Yes}) = \frac{P(\text{Hot}, \text{Yes})}{P(\text{Yes})} = \frac{\frac{2}{14}}{\frac{9}{14}} = \frac{2}{9}$$

$$P(\text{Cool}|\text{Yes}) = \frac{P(\text{Cool}, \text{Yes})}{P(\text{Yes})} = \frac{\frac{3}{14}}{\frac{9}{14}} = \frac{3}{9}$$

$$P(\text{Hot}|\text{No}) = \frac{P(\text{Hot}, \text{No})}{P(\text{No})} = \frac{\frac{2}{14}}{\frac{5}{14}} = \frac{2}{5}$$

$$P(\text{Cool}|\text{No}) = \frac{P(\text{Cool}, \text{No})}{P(\text{No})} = \frac{\frac{1}{14}}{\frac{5}{14}} = \frac{1}{5}$$

$$P(\text{Mild}|\text{Yes}) = \frac{P(\text{Mild}, \text{Yes})}{P(\text{Yes})} = \frac{\frac{4}{14}}{\frac{9}{14}} = \frac{4}{9}$$

$$P(\text{Mild}|\text{No}) = \frac{P(\text{Mild}, \text{No})}{P(\text{No})} = \frac{\frac{2}{14}}{\frac{5}{14}} = \frac{2}{5}$$

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

Naïve Bayes Algorithm

(for *discrete valued Features*)

Example: Play Tennis

Cont.

Training Phase

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

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

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

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

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

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

Naïve Bayes Algorithm (for *discrete valued Features*)

Example: Play Tennis

Cont.

Testing Phase

Given a new instance $X^{new} = (\text{Outlook} = \text{Sunny}, \text{Temperature} = \text{Cool}, \text{Humidity} = \text{High}, \text{Wind} = \text{Strong})$, what is it's label?

- Look up tables achieved in the learning phrase:

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

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

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

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

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

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

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

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

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

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

- Decision making with the MAP rule:

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

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

Given the fact $P(\text{Yes} \mid X^{new}) < P(\text{No} \mid X^{new})$, we label X^{new} to be *No*.

Naïve Bayes Algorithm

(Zero conditional probability)

If no example contains the feature value

m-estimate of probability

- ❖ Note that we estimated conditional probabilities $P(A|B)$ by $\frac{n_c}{n}$ where n_c is the number of times A and B happened in same time and n is the number of times B happened in the training data
- ❖ *This can cause trouble if $n_c=0$.*
- ❖ To avoid this, we fix the numbers p and m beforehand:
 - A nonzero prior estimate p for $P(A|B)$, and
 - A number m that says how confident we are of our prior estimate p , as measured in number of samples.
- ❖ Then instead of using $\frac{n_c}{n}$ for the estimate, use $\frac{n_c+pm}{n+m}$.
- ❖ Just think of this as adding a bunch of samples to start the whole process
- ❖ If we don't have any knowledge of p , assume the attribute is uniformly distributed over all possible values

$$\hat{P}(x_{jk}|y = c_t) \approx \frac{n_c + pm}{n + m}$$

n_c : number of training examples for which $x_j = x_{jk}$ and $y = c_t$.

n : number of training examples for $y = c_t$.

p : prior estimate (usually, $p = \frac{1}{n_j}$, $(x_j \in \{x_{j1}^{(i)}, x_{j2}^{(i)}, \dots, x_{jn_j}^{(i)}\})$).

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

The parameter m is also known as equivalent sample size. It prevents the probabilities from being 0.

Naïve Bayes Algorithm

(for *Continuous valued Features*)

Input: Data set $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}$, for $i=1,2,\dots, m$, such that $\mathbf{x}^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{c_1, c_2, \dots, c_K\}$.

Output:

- Given an unknown instance $X^{new} \in \mathbb{R}^n$, predict it's label.

Training Phase:

- Estimate $P(y = c_i)$ with examples in D for $i=1,2,\dots, K$:
- Numberless values taken by a continuous-valued feature
- For each feature x_j estimate conditional probability $P(x_j|y = c_t)$ using *normal distribution*

$$\hat{P}(x_j|y = c_t) \approx P(x_j|y = c_t) = \frac{1}{\sqrt{2\pi} \sigma_{jt}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_{jt}^2}}$$

μ_j : mean(average) of feature values x_j of examples for which $y = c_t$

σ_{jt} : standard deviation of feature values x_j of examples for which $y = c_t$

Output of training phase: Conditional probability table for each feature x_j , for $j = 1, 2, \dots, n$.

Testing Phase:

Given an unknown instance $X^{new} = (x_1^{new}, x_2^{new}, \dots, x_n^{new}) \in \mathbb{R}^n$, predict it's label:

- Instead of looking-up tables, calculate conditional probabilities with all the normal distributions achieved in the learning phase
- Apply the MAP rule to assign the label c^* to X^{new} :

$$c^* = \operatorname{argmax}_{c=c_1, c_2, \dots, c_K} P(y = c) \prod_{j=1}^n \hat{P}(x_j^{new} | y = c)$$

Naïve Bayes Algorithm (for *Continuous valued Features*)

Cont.

Example

- **Temperature** is naturally of continuous value.

Yes: 25.2, 19.3, 18.5, 21.7, 20.1, 24.3, 22.8, 23.1, 19.8

No: 27.3, 30.1, 17.4, 29.5, 15.1

Training Phase:

Estimate mean and variance for each class

$$\mu = \frac{1}{N} \sum_{n=1}^N x_n \quad \Rightarrow \quad \mu_{Yes} = 21.64, \quad \mu_{No} = 23.88$$

$$\sigma^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu)^2 \quad \Rightarrow \quad \sigma_{Yes} = 2.35, \quad \sigma_{No} = 7.09$$

Output: two Gaussian models for $P(temp|c)$

$$\hat{P}(x_j|Yes) \approx \frac{1}{2.35\sqrt{2\pi}} e^{-\frac{(x_j-21.64)^2}{2 \times 2.35^2}} = \frac{1}{2.35\sqrt{2\pi}} e^{-\frac{(x_j-21.64)^2}{11.09}}$$

$$\hat{P}(x_j|No) \approx \frac{1}{7.09\sqrt{2\pi}} e^{-\frac{(x_j-23.88)^2}{2 \times 7.09^2}} = \frac{1}{7.09\sqrt{2\pi}} e^{-\frac{(x_j-23.88)^2}{50.25}}$$

References

- 1- <https://cse.sc.edu/~rose/587/PPT/NaiveBayes.ppt>
- 2- <http://euclid.nmu.edu/~mkowalc/cs495/notes/NaiveBayesSlides/lesson004.pdf>