

UVA CS 6316 – Fall 2015 Graduate: Machine Learning

Lecture 12: Naïve Bayes Classifier

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

10/12/15

1

Where are we ? →

Five major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
- Graphical models

10/7/14

Where are we ? →

Three major sections for classification

- We can divide the large variety of classification approaches into **roughly three major types**

1. Discriminative

- directly estimate a decision rule/boundary
- e.g., support vector machine, decision tree



2. Generative:

- build a generative statistical model
- e.g., **naïve bayes classifier**, Bayesian networks

3. Instance based classifiers

- Use observation directly (no models)
- e.g. K nearest neighbors

10/7/14

Last Lecture Recap: Probability Review

- The big picture : data \leftrightarrow probabilistic model
- Sample space, Events and Event spaces
- Random variables
- Joint probability, Marginal probability, conditional probability,
- Chain rule, Bayes Rule, Law of total probability, etc.
- Structural properties
 - Independence, conditional independence

10/12/15

4

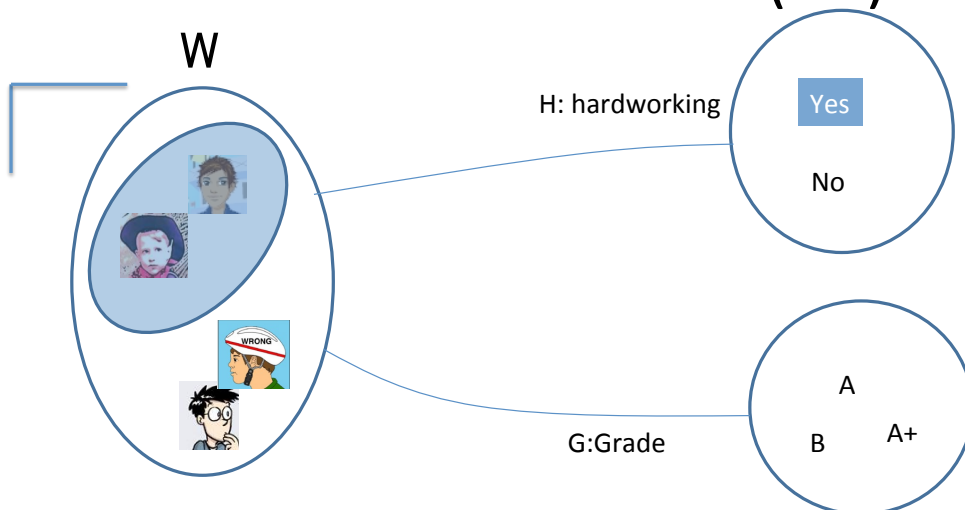
Sample space and Events

- **W: Sample Space,**
 - result of an experiment / set of all outcomes
 - If you toss a coin twice $W = \{HH, HT, TH, TT\}$
- **Event:** a subset of W
 - First toss is head = $\{HH, HT\}$
- **S: event space, a set of events:**
 - Contains the empty event and W

10/12/15

5

Random Variables (RV)



$$P(H = \text{Yes}) = P(\{\text{all students who is working hard on the course}\})$$

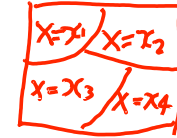
- “functions” that maps from W to an attribute space T .

10/12/15

6

Probability of Discrete RV

- Probability mass function (pmf): $P(X = x_i)$
- Easy facts about pmf
 - $\sum_i P(X = x_i) = 1$
 - $P(X = x_i \cap X = x_j) = 0$ if $i \neq j$
 - $P(X = x_i \cup X = x_j) = P(X = x_i) + P(X = x_j)$ if $i \neq j$
 - $P(X = x_1 \cup X = x_2 \cup \dots \cup X = x_k) = 1$



10/7/14

e.g. Coin Flips cont.

- You flip a coin
 - Head with probability p
 - Binary random variable
 - **Bernoulli trial** with success probability p
- You flip k coins
 - How many heads would you expect
 - Number of heads X : discrete random variable
 - **Binomial distribution** with parameters k and p

$$X \sim \text{Bin}(k, p) \quad P(X = i) = \binom{k}{i} p^i (1-p)^{k-i}$$

10/7/14

Joint prob: e.g., Coin Flips by Two Persons

- Your friend and you both flip coins
 - Head with probability 0.5
 - You flip 50 times; your friend flip 100 times
 - How many heads will both of you get
- Given two discrete RVs X and Y , their **joint distribution** is the distribution of X and Y together
 - E.g. $P(\text{You get 21 heads AND you friend get 70 heads})$

$$P((X=21) \wedge (Y=70))$$

10/12/15

9

Conditional Probability

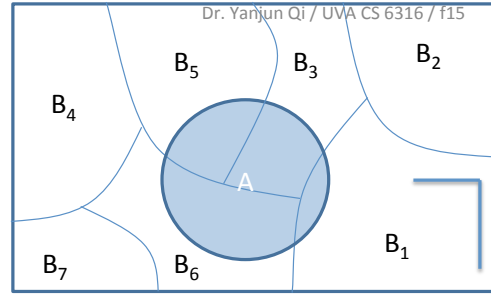
- $P(X = x | Y = y)$ is the probability of $X = x$, given the occurrence of $Y = y$
 - E.g. you get 0 heads, given that your friend gets 61 heads

$$P(\underbrace{X = x} | \underbrace{Y = y}) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$

10/12/15

10

Marginalization



Marginal Probability

Law of Total Probability

$$\begin{aligned}
 P(X = x_i) &= \sum_j P(X = x_i \cap Y = y_j) \\
 &= \sum_j P(X = x_i | Y = y_j) P(Y = y_j)
 \end{aligned}$$

Conditional Probability
Marginal Probability

10/12/15

11

Bayes Rule

- X and Y are discrete RVs...

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

$$P(X = x_i | Y = y_j) = \frac{P(Y = y_j | X = x_i) P(X = x_i)}{\sum_k P(Y = y_j | X = x_k) P(X = x_k)}$$

10/12/15

12

Bayes Rule cont.

- You can condition on more variables

$$P(x|y,z) = \frac{P(x|z)P(y|x,z)}{P(y|z)}$$

$$P(x|y|z) = \frac{P(x,y|z)}{P(y|z)}$$

Conditional Probability Example

What is the probability that the 2nd ball drawn from the set $\{r,r,r,b\}$ will be red?

Using marginalization, $P(B_2 = r) = P(B_2 = r \wedge B_1 = r) + P(B_2 = r \wedge B_1 = b)$

$$= P(B_1 = r)P(B_2 = r|B_1 = r) + P(B_1 = b)P(B_2 = r|B_1 = b)$$

$$\begin{bmatrix} P(B_2=r) \\ P(B_2=b) \end{bmatrix} = \begin{bmatrix} P(B_2=r|B_1=r)P(B_1=r) + P(B_2=r|B_1=b)P(B_1=b) \\ P(B_2=b|B_1=r)P(B_1=r) + P(B_2=b|B_1=b)P(B_1=b) \end{bmatrix}$$

→ Matrix Notation

$$\begin{bmatrix} P(B_2=r|B_1=r), & P(B_2=r|B_1=b) \\ P(B_2=b|B_1=r), & P(B_2=b|B_1=b) \end{bmatrix} \begin{bmatrix} P(B_1=r) \\ P(B_1=b) \end{bmatrix}$$

$P(B_2)$ $P(B_2|B_1)$ $P(B_1)$

⇒ matrix notation form $\Pi_2 = G^T \Pi_1$

For short, we write this using vectors and a **stochastic matrix**:

Today : Naïve Bayes Classifier

- ✓ Probability review
 - Structural properties, i.e., Independence, conditional independence
- ✓ Naïve Bayes Classifier
 - Spam email classification

Independent RVs

- Intuition: X and Y are independent means that $X = x$ **neither** makes it **more or less** probable that $Y = y$

- Definition: X and Y are independent *iff*

$$P(X = x \cap Y = y) = P(X = x)P(Y = y)$$

$$P(X = x | Y = y) = P(X = x)$$

10/7/14

More on Independence

- $P(X = x \cap Y = y) = P(X = x)P(Y = y)$

$$P(X = x | Y = y) = P(X = x) \quad P(Y = y | X = x) = P(Y = y)$$

- **E.g.** no matter how many heads you get, your friend will not be affected, and vice versa

10/7/14

More on Independence

- X is independent of Y means that knowing Y does not change our belief about X.
 - $P(X|Y=y) = P(X)$
 - $P(X=x, Y=y) = P(X=x) P(Y=y)$

- The above should hold for all x_i, y_j
- It is symmetric and written as $X \perp Y$

$X \perp Y$

10/7/14

Conditionally Independent RVs

- Intuition: X and Y are conditionally independent given Z means that once Z is **known**, the value of X does not add any **additional** information about Y
- Definition: X and Y are conditionally independent given Z **iff**

$\left\{ \begin{array}{l} X: \text{sale} \\ \text{ice-cream} \\ Y: \text{rate of Dr} \\ Z: \text{Weather} \end{array} \right.$

$$P(X=x \cap Y=y | Z=z) = P(X=x | Z=z) P(Y=y | Z=z)$$

If holding for all x_i, y_j, z_k

$X \perp Y | Z$

10/7/14

More on Conditional Independence

$$P(X = x \cap Y = y | Z = z) = P(X = x | Z = z) P(Y = y | Z = z)$$

$$\frac{P(X=x, Y=y | Z=z)}{P(Y=y | Z=z)}$$

$$P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$

$$P(Y = y | X = x, Z = z) = P(Y = y | Z = z)$$

10/7/14

Today : Naïve Bayes Classifier

- ✓ Probability review
 - Structural properties, i.e., Independence, conditional independence
- ➔ ✓ Naïve Bayes Classifier
 - Spam email classification

10/7/14

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

Output as Discrete
Class Label
 C_1, C_2, \dots, C_L

$$P(C | \mathbf{X})$$

- **Data**/points/instances/examples/samples/records: [rows]
- **Features**/attributes/dimensions/independent variables/covariates, predictors/regressors: [columns, except the last]
- **Target**/outcome/response/label/dependent variable: special column to be predicted [last column]

10/7/14

Bayes classifiers

- Treat each feature attribute and the class label as random variables.
- Given a sample \mathbf{x} with attributes (x_1, x_2, \dots, x_p) :
 - Goal is to predict its class C .
 - Specifically, we want to find the value of C_i that maximizes $p(C_i | x_1, x_2, \dots, x_p)$.
- Can we estimate $p(C_i | \mathbf{x}) = p(C_i | x_1, x_2, \dots, x_p)$ directly from data?

10/7/14

Bayes classifiers

→ MAP classification rule

- Establishing a probabilistic model for classification

→ MAP classification rule

- **MAP: Maximum A Posterior**
- Assign x to c^* if

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

$$\left. \begin{array}{l} P(C=c_1 | \mathbf{x}) \\ P(C=c_2 | \mathbf{x}) \\ P(C=c_3 | \mathbf{x}) \end{array} \right\} \max \Rightarrow c_i$$

10/7/14

Adapt from Prof. Ke Chen NB slides

Bayes Classification Rule – (1)

- Establishing a probabilistic model for classification
- **(1) Discriminative model**

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

$$P(c_1 | \mathbf{x}) \quad P(c_2 | \mathbf{x}) \quad \dots \quad P(c_L | \mathbf{x})$$

**Discriminative
Probabilistic Classifier**

$$x_1 \quad x_2 \quad \dots \quad x_n$$

$$\mathbf{X} = (x_1, x_2, \dots, x_n)$$

*logistic
regression*

10/7/14

Adapt from Prof. Ke Chen NB slides

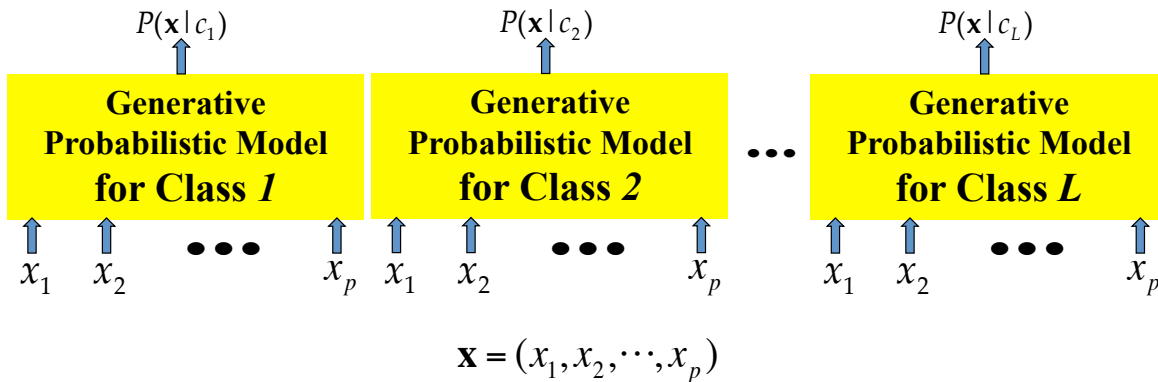
Bayes Classification Rule – (2)

- Establishing a probabilistic model for classification (cont.)

- (2) **Generative model**

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

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



10/7/14

Adapt from Prof. Ke Chen NB slides

Review : Bayes' Rule – (2)

$$P(C, X) = P(C | X)P(X) = P(X | C)P(C)$$

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

Posterior $\text{max} \Rightarrow C^*$

Prior

$P(C_1), P(C_2), \dots, P(C_L)$

$P(C_1|x), P(C_2|x), \dots, P(C_L|x)$

$$P(C | \mathbf{X}) = \frac{P(\mathbf{X} | C)P(C)}{P(\mathbf{X})}$$

10/12/15

28

Review: Bayes Rule – (2)

- Prior, conditional and marginal probability
 - Prior probability: $P(C)$ $P(C_1), P(C_2), \dots, P(C_L)$
 - Likelihood (through a generative model): $P(\mathbf{X} | C)$
 - Evidence (marginal prob. of sample): $P(\mathbf{X})$
 - Posterior probability: $P(C | \mathbf{X})$ $P(C_1|x), P(C_2|x), \dots, P(C_L|x)$
- Bayes Rule

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

10/7/14

Adapt from Prof. Ke Chen NB slides

(2) Generative classification with the MAP rule

- MAP classification rule
 - **MAP: Maximum A Posterior**
 - Assign x to c^* if
- Generative classification with the MAP rule
 - Apply Bayes rule to convert them into posterior probabilities

$$P(C = c_i | \mathbf{X} = \mathbf{x}) = \frac{P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i)}{P(\mathbf{X} = \mathbf{x})}$$

$\propto P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i)$
 for $i = 1, 2, \dots, L$

$P(C = c_i)$
 $P(\mathbf{X} = \mathbf{x} | C = c_i)$

- Then apply the MAP rule

10/7/14

Adapt from Prof. Ke Chen NB slides

Naïve Bayes Classifier

- Bayes classification

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

$P(X_1, \dots, X_p | C)$

Difficulty: learning the joint probability

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

given class ✓

10/7/14

Naïve Bayes Classifier

- Naïve Bayes classification

$$P(C) P(\mathbf{x} | C)$$

$$P(X_1, \dots, X_p | C)$$

- Assumption that **all input attributes are conditionally independent!**

$$\begin{aligned}
 P(X_1, X_2, \dots, X_p | C) &= P(X_1 | X_2, \dots, X_p, C) P(X_2, \dots, X_p | C) \quad \leftarrow \text{Chain Rule} \\
 &= P(X_1 | C) P(X_2, \dots, X_p | C) \quad \leftarrow \text{con.I.H.R} \\
 &= P(X_1 | C) P(X_2 | C) \cdots P(X_p | C)
 \end{aligned}$$

- MAP classification rule: for $\mathbf{x} = (x_1, x_2, \dots, x_n)$

$$[P(x_1 | c^*) \cdots P(x_p | c^*)] P(c^*) > [P(x_1 | c) \cdots P(x_p | c)] P(c),$$

10/7/14

$$c \neq c^*, c = c_1, \dots, c_L$$

Naïve Bayes Classifier

- Naïve Bayes classification $P(X_1, \dots, X_p | C)$
 - Assumption that **all input attributes are conditionally independent!**
 - MAP classification rule: for a testing sample

$$[P(x_1 | c^*) \cdots P(x_p | c^*)]P(c^*) > [P(x_1 | c) \cdots P(x_p | c)]P(c),$$

$$c \neq c^*, c = c_1, \dots, c_L$$

$$P(X_j | c_i) \left\{ \begin{array}{l} \text{Bern} \\ \text{Bio} \\ \text{Multi} \\ \text{Gaussian} \end{array} \right.$$

10/7/14

Naïve Bayes (for discrete input attributes) - training

- 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 ; → L

For every attribute value x_{jk} of each attribute X_j ($j = 1, \dots, p; k = 1, \dots, K_j$)

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

Output: conditional probability tables; for $X_j, K_j \times L$ elements

$\left\{ \begin{array}{l} K_1, K_2, \dots, K_p \\ (X_1, X_2, \dots, X_p) \end{array} \right\}$

$\left\{ \begin{array}{l} K_1 \times L + \\ K_2 \times L + \\ \dots + \\ K_p \times L \end{array} \right.$

10/7/14

Naïve Bayes (for discrete input attributes) - testing

- Naïve Bayes Algorithm (for discrete input attributes)

- **Test Phase:** Given an unknown instance $\mathbf{X}' = (a'_1, \dots, a'_p)$
Look up tables to assign the label c^* to \mathbf{X}' if

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > \underbrace{[\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)},$$

$c \neq c^*, c = c_1, \dots, c_L$

$$P(\mathbf{X}' | c_i) P(c_i)$$

$$= P(a'_1 | c_i) P(a'_2 | c_i) \dots P(a'_p | c_i) P(c_i)$$

$i = 1, 2, \dots, L$

10/7/14

$P(X_1, X_2, X_3, X_4 | \text{Yes})$
 $P(X_1, X_2, X_3, X_4 | \text{No})$

Example

• Example: Play Tennis $\rightarrow 36 \times 2$

3 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

$$P(c = c_i)$$

$$\{ \text{Yes} \}$$

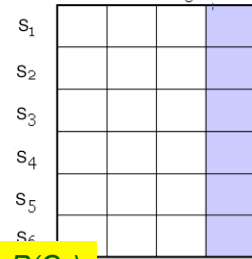
$$\{ \text{No} \}$$

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

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

$$3 \times 3 \times 2 \times 2 = 36$$

10/7/14



(a) Generative Bayes Classifier

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

- Learning Phase

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

$P(C_1), P(C_2), \dots, P(C_L)$

$P(X_1, X_2, \dots, X_p | C_1), P(X_1, X_2, \dots, X_p | C_2)$

$p=4 \Rightarrow L \times 2^p$

Outlook (3 values)	Temperature (3 values)	Humidity (2 values)	Wind (2 values)	Play=Yes	Play=No
sunny	hot	high	weak	0/9	1/5
sunny	hot	high	strong	.../9	.../5
sunny	hot	normal	weak	.../9	.../5
sunny	hot	normal	strong	.../9	.../5
...
...
...
...

$3 \times 3 \times 2 \times 2$ [conjunctions of attributes] * 2 [two classes] = 72 parameters

(b) Naïve Bayes Classifier $P(X_1, X_2, X_3, X_4 | C_i)$

Estimate $P(X_j = x_{jk} | C = c_i)$ with examples in S;

- Learning Phase

$P(X_2 | C_1), P(X_2 | C_2)$

X_1 3	Outlook	Play=Yes	Play=No	Temperature	Play=Yes	Play=No	
	Sunny	2/9	3/5		Hot	2/9	2/5
	Overcast	4/9	0/5		Mild	4/9	2/5
	Rain	3/9	2/5	Cool	3/9	1/5	

2	Humidity	Play=Yes	Play=No	Wind	Play=Yes	Play=No	
	High	3/9	4/5		Strong	3/9	3/5
	Normal	6/9	1/5		Weak	6/9	2/5

$3+3+2+2$ [naïve assumption] * 2 [two classes] = 20 parameters

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

$P(C_1), P(C_2), \dots, P(C_L)$

$P(C_i)$

P_L

(b) Naïve Bayes Classifier

$$[\hat{P}(a'_1|c^*) \cdots \hat{P}(a'_p|c^*)] \hat{P}(c^*) > [\hat{P}(a'_1|c) \cdots \hat{P}(a'_p|c)] \hat{P}(c)$$

- Test Phase

- Given a new instance,

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

$$\begin{aligned} &\rightarrow P(c_1) P(\text{Sunny} | c_1) P(\text{Cool} | c_1) P(\text{High} | c_1) P(\text{Strong} | c_1) \\ &= \frac{9}{14} \times \frac{2}{9} \cdots \cdots = \\ &\rightarrow P(c_2) P(\text{Su} | c_2) P(\text{Co} | c_2) P(\text{hi} | c_2) P(\text{St} | c_2) \\ &= \frac{5}{14} \times \frac{3}{5} \times \cdots \cdots = \end{aligned}$$

(b) Naïve Bayes Classifier

$$[\hat{P}(a'_1|c^*) \cdots \hat{P}(a'_p|c^*)] \hat{P}(c^*) > [\hat{P}(a'_1|c) \cdots \hat{P}(a'_p|c)] \hat{P}(c)$$

- 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

$P(\text{Outlook}=\text{Sunny} \text{Play}=\text{Yes}) = 2/9$	$P(\text{Outlook}=\text{Sunny} \text{Play}=\text{No}) = 3/5$
$P(\text{Temperature}=\text{Cool} \text{Play}=\text{Yes}) = 3/9$	$P(\text{Temperature}=\text{Cool} \text{Play}=\text{No}) = 1/5$
$P(\text{Humidity}=\text{High} \text{Play}=\text{Yes}) = 3/9$	$P(\text{Humidity}=\text{High} \text{Play}=\text{No}) = 4/5$
$P(\text{Wind}=\text{Strong} \text{Play}=\text{Yes}) = 3/9$	$P(\text{Wind}=\text{Strong} \text{Play}=\text{No}) = 3/5$
$P(\text{Play}=\text{Yes}) = 9/14$	$P(\text{Play}=\text{No}) = 5/14$

- MAP rule

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

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



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

Naïve Bayes Assumption

- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_p | c_j)$
 - $O(|X|^p \cdot |C|)$ parameters
 - Could only be estimated if a very, very large number of training examples was available.

If no naïve assumption

Naïve Bayes Conditional Independence Assumption:

$$|C| \cdot P \cdot |X|$$

- Assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities $P(x_i | c_j)$.

10/12/15

Adapt From Manning' textCat tutorial⁴¹

References

- Prof. Andrew Moore' s review tutorial
- Prof. Ke Chen NB slides
- Prof. Carlos Guestrin recitation slides

10/7/14