



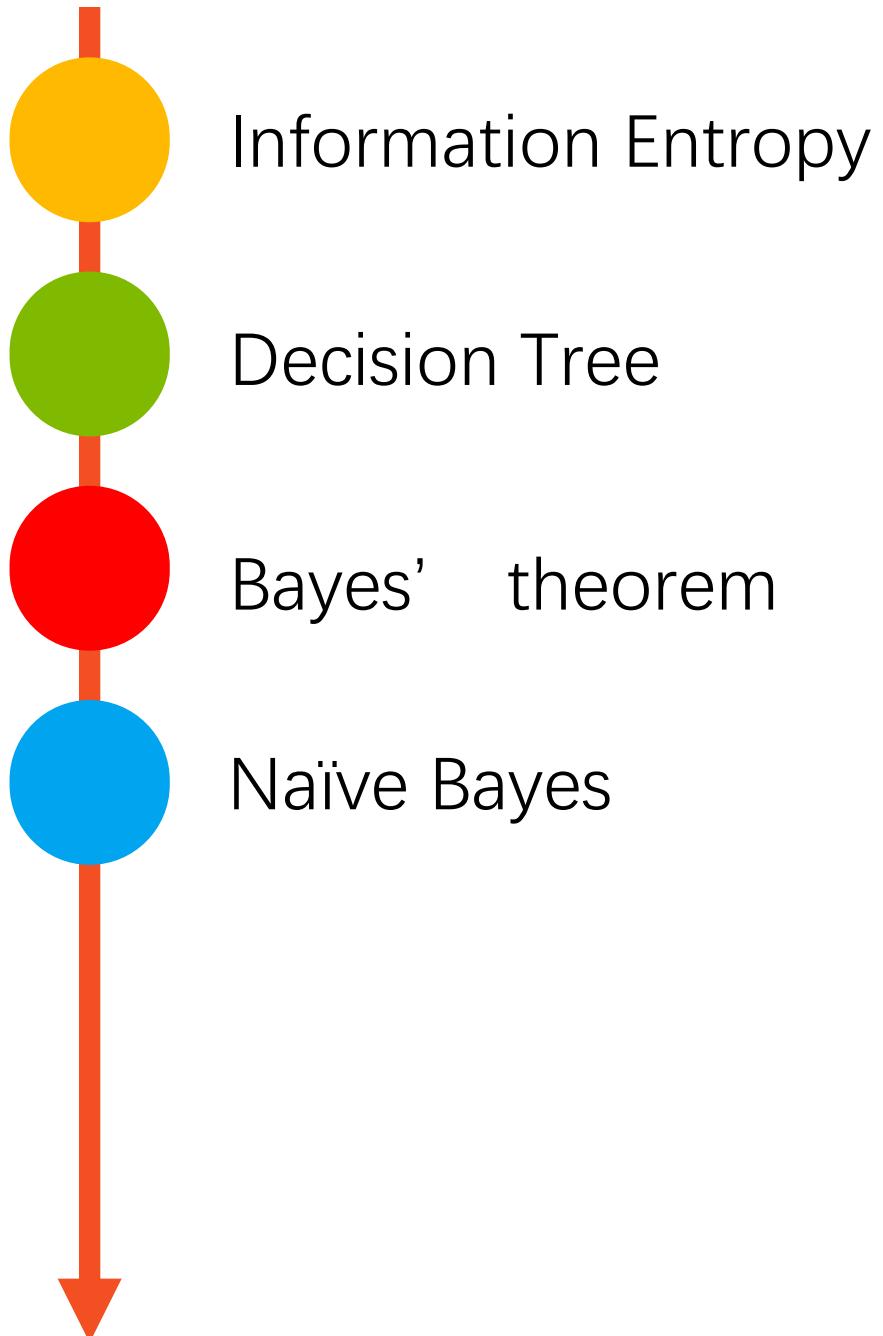
2018 TJMSC Tech. Courses

---

# Machine Learning Algorithm

Yiran Zhuo  
Tongji Microsoft Student Club

Nov 4, 2018  
Room 516, Ji Shi Building  
SSE, Tongji Univ



# Information Entropy

---

Information contained by an event can be described as

$$f(P) = \log \frac{1}{p} = -\log p$$

Thus, when the event has a lower-probability value (i.e., when a low-probability event occurs), the event carries more "information" ("surprisal") than when the event has a higher-probability value

# Information Entropy

---

The information carried by a set of events can be described as:

$$H(X) = - \sum_{x \in X} P(x) \log P(x))$$

# Decision trees

---

Problem: decide whether to wait for a table at a restaurant, based on the following attributes:

1. Alternate: is there an alternative restaurant nearby?
2. Bar: is there a comfortable bar area to wait in?
3. Fri/Sat: is today Friday or Saturday?
4. Hungry: are we hungry?
5. Patrons: number of people in the restaurant (None, Some, Full)
6. Price: price range (\$, \$\$, \$\$\$)
7. Raining: is it raining outside?
8. Reservation: have we made a reservation?
9. Type: kind of restaurant (French, Italian, Thai (泰国菜), Burger)
10. WaitEstimate: estimated waiting time (0-10, 10-30, 30-60, >60)

# Attribute-based representations

Examples described by **attribute values** (Boolean, discrete, continuous)

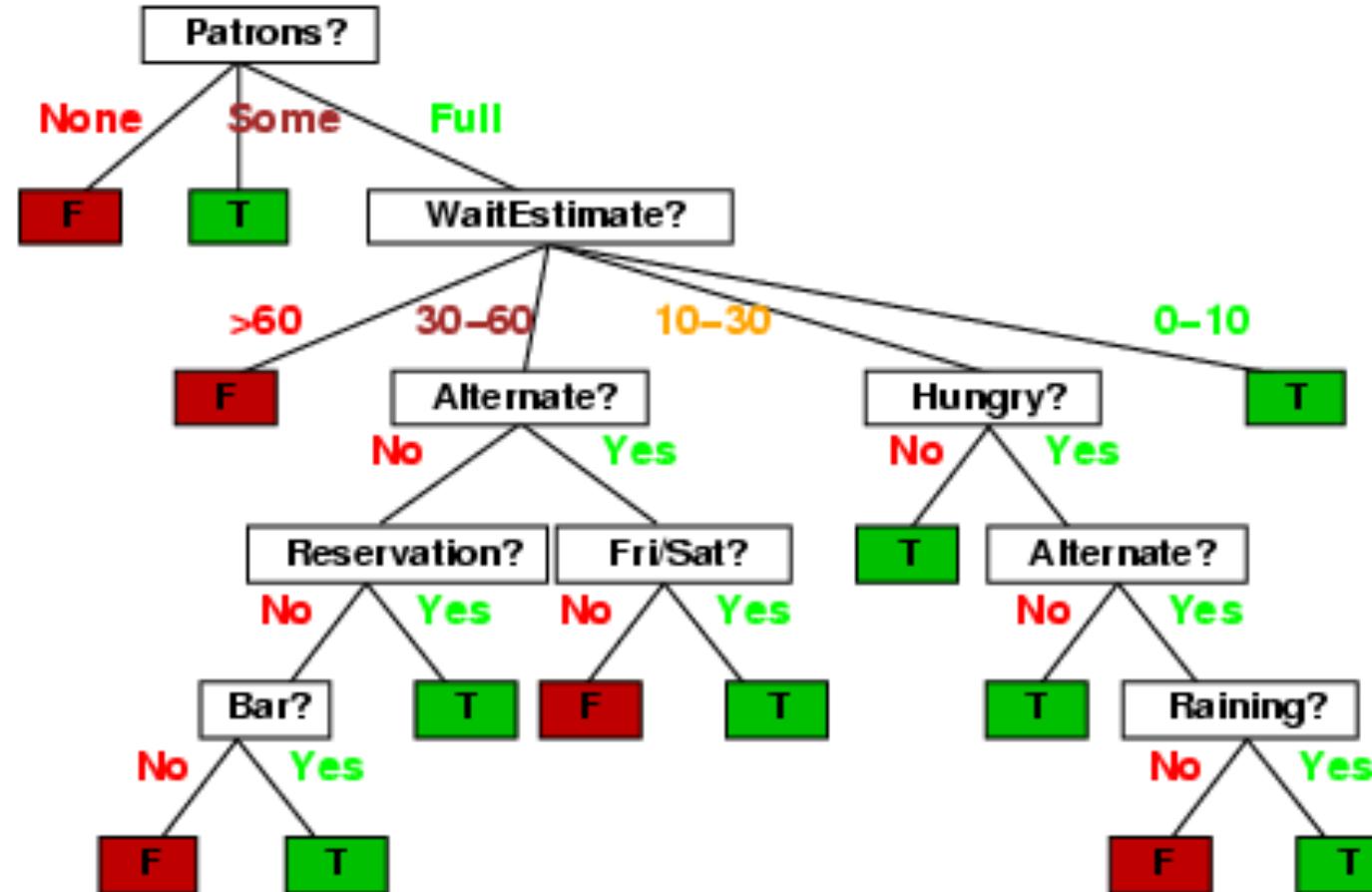
E.g., situations where I will/won't wait for a table:

Example	Attributes											Target <i>Wait</i>
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>		
$X_1$	T	F	F	T	Some	\$\$\$	F	T	French	0–10	T	
$X_2$	T	F	F	T	Full	\$	F	F	Thai	30–60	F	
$X_3$	F	T	F	F	Some	\$	F	F	Burger	0–10	T	
$X_4$	T	F	T	T	Full	\$	F	F	Thai	10–30	T	
$X_5$	T	F	T	F	Full	\$\$\$	F	T	French	>60	F	
$X_6$	F	T	F	T	Some	\$\$	T	T	Italian	0–10	T	
$X_7$	F	T	F	F	None	\$	T	F	Burger	0–10	F	
$X_8$	F	F	F	T	Some	\$\$	T	T	Thai	0–10	T	
$X_9$	F	T	T	F	Full	\$	T	F	Burger	>60	F	
$X_{10}$	T	T	T	T	Full	\$\$\$	F	T	Italian	10–30	F	
$X_{11}$	F	F	F	F	None	\$	F	F	Thai	0–10	F	
$X_{12}$	T	T	T	T	Full	\$	F	F	Burger	30–60	T	

Classification of examples is **positive** (T) or **negative** (F)

# Decision trees

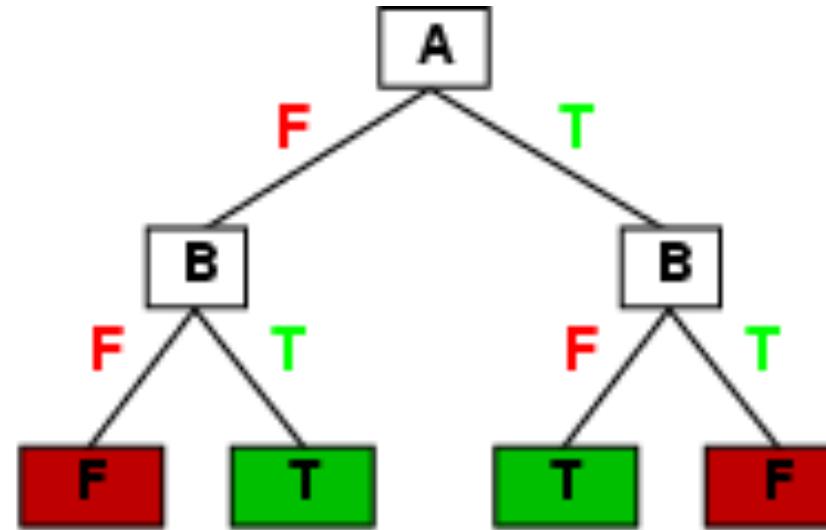
- One possible representation for hypotheses
- E.g., here is the “true” tree for deciding whether to wait:



# Expressiveness

- Decision trees can **express any function** of the input attributes.
- E.g., for Boolean functions, truth table row → path to leaf:

A	B	$A \text{ xor } B$
F	F	F
F	T	T
T	F	T
T	T	F



- Trivially, there is a consistent decision tree for any training set with one path to leaf for each example (unless  $f$  nondeterministic in  $x$ ) but it probably won't generalize to new examples
- Prefer to find more **compact** decision trees

# Decision tree learning

- Aim: find a small tree consistent with the training examples
- Idea: (recursively) choose "most significant" attribute as root of (sub)tree

检测数据集中的每个子项是否属于同一分类:

If so return 类标签

Else

    寻找划分数据集的最好特征

    划分数据集

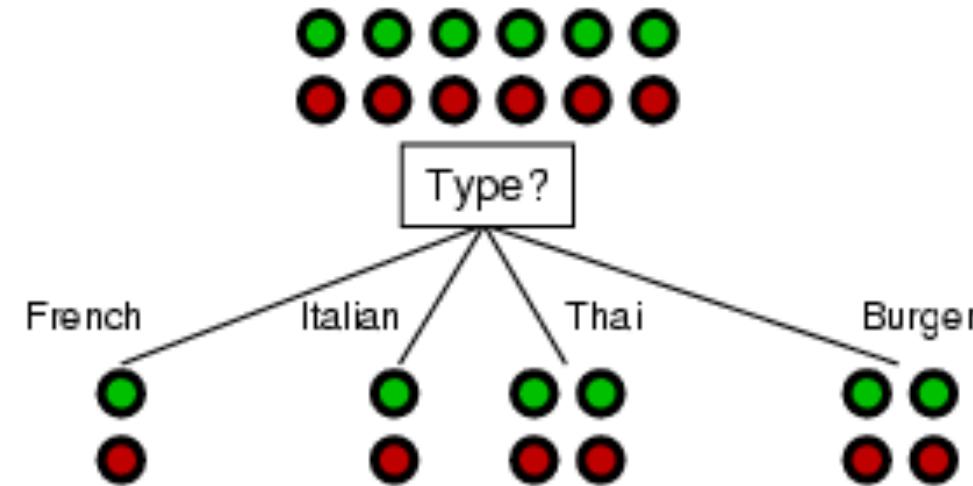
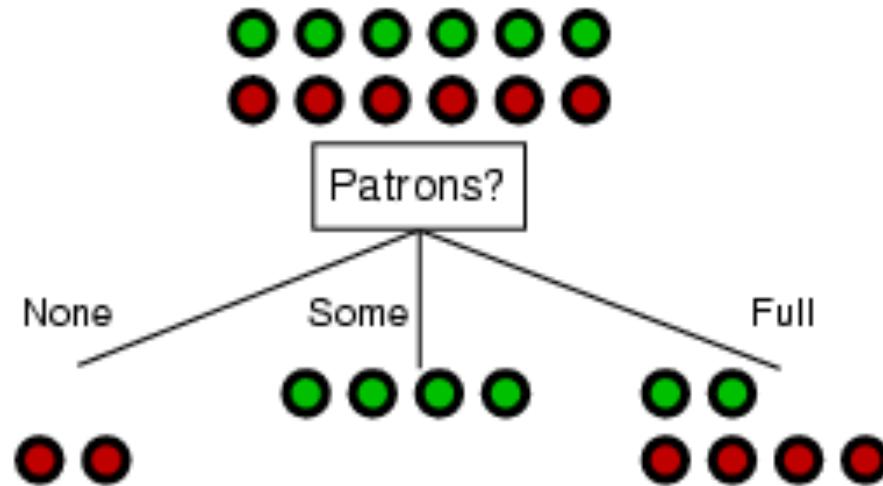
    创建分支节点

        for 每个划分的子集调用函数 createBranch 并增加返回结果到分支节点中

    return 分支节点

# Choosing an attribute

- Idea: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



- Patrons?* is a better choice

# Using information theory

---

- To implement Choose-Attribute in the DTL algorithm
- Information Content (Entropy):

$$I(P(v_1), \dots, P(v_n)) = \sum_{i=1} -P(v_i) \log_2 P(v_i)$$

- For a training set containing  $p$  positive examples and  $n$  negative examples:

$$I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

# Information gain

- A chosen attribute  $A$  divides the training set  $E$  into subsets  $E_1, \dots, E_v$  according to their values for  $A$ , where  $A$  has  $v$  distinct values.

$$\text{remainder}(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right)$$

- Information Gain (IG) or reduction in entropy from the attribute test:

$$IG(A) = I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) - \text{remainder}(A)$$

- Choose the attribute with the largest IG

# Information gain

For the training set,  $p = n = 6$ ,  $I(6/12, 6/12) = 1$  bit

Consider the attributes *Patrons* and *Type* (and others too)

Example	Attributes											Target <i>Wait</i>
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>		
$X_1$	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T	
$X_2$	T	F	F	T	Full	\$	F	F	Thai	30-60	F	
$X_3$	F	T	F	F	Some	\$	F	F	Burger	0-10	T	
$X_4$	T	F	T	T	Full	\$	F	F	Thai	10-30	T	
$X_5$	T	F	T	F	Full	\$\$\$	F	T	French	>60	F	
$X_6$	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T	
$X_7$	F	T	F	F	None	\$	T	F	Burger	0-10	F	
$X_8$	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T	
$X_9$	F	T	T	F	Full	\$	T	F	Burger	>60	F	
$X_{10}$	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F	
$X_{11}$	F	F	F	F	None	\$	F	F	Thai	0-10	F	
$X_{12}$	T	T	T	T	Full	\$	F	F	Burger	30-60	T	

# Information gain

Pat:

4 Some : 4 T 0 F

6 Full : 2 T 4 F

2 None : 0 T 2 F

Type:

2 French : 1 T 1 F

2 Italian : 1 T 1 F

4 Thai : 2 T 2 F

4 Burger: 2 T 2 F

.....

*Patrons* has the highest IG of all attributes and so is chosen by the DTL algorithm as the root

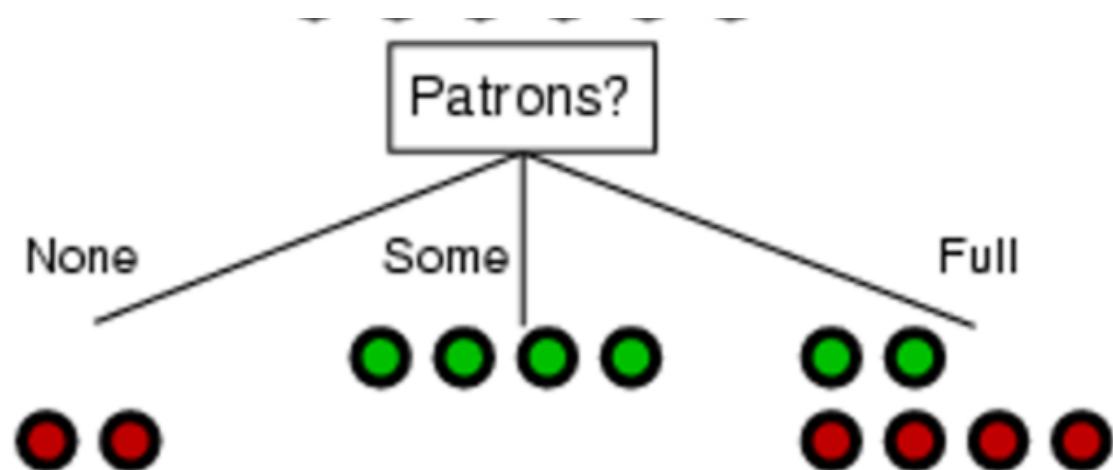
Example	Attributes											Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	Wait	
$X_1$	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T	
$X_2$	T	F	F	T	Full	\$	F	F	Thai	30-60	F	
$X_3$	F	T	F	F	Some	\$	F	F	Burger	0-10	T	
$X_4$	T	F	T	T	Full	\$	F	F	Thai	10-30	T	
$X_5$	T	F	T	F	Full	\$\$\$	F	T	French	>60	F	
$X_6$	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T	
$X_7$	F	T	F	F	None	\$	T	F	Burger	0-10	F	
$X_8$	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T	
$X_9$	F	T	T	F	Full	\$	T	F	Burger	>60	F	
$X_{10}$	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F	
$X_{11}$	F	F	F	F	None	\$	F	F	Thai	0-10	F	
$X_{12}$	T	T	T	T	Full	\$	F	F	Burger	30-60	T	

$$IG(Patrons) = 1 - \left[ \frac{2}{12} I(0,1) + \frac{4}{12} I(1,0) + \frac{6}{12} I\left(\frac{2}{6}, \frac{4}{6}\right) \right] = .0541 \text{ bits}$$

$$IG(Type) = 1 - \left[ \frac{2}{12} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{12} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{4}{12} I\left(\frac{2}{4}, \frac{2}{4}\right) + \frac{4}{12} I\left(\frac{2}{4}, \frac{2}{4}\right) \right] = 0 \text{ bits}$$

# Information gain

After choosing *Patrons* as the root, we still need to choose another node  
*at the situation 'Full'*



# Information gain

Hun:

4 T : 2 T 2 F

2 F : 0 T 2 F

Type:

1 French : 0 T 1 F

1 Italian : 0 T 1 F

2 Thai : 1 T 1 F

2 Burger: 1 T 1 F

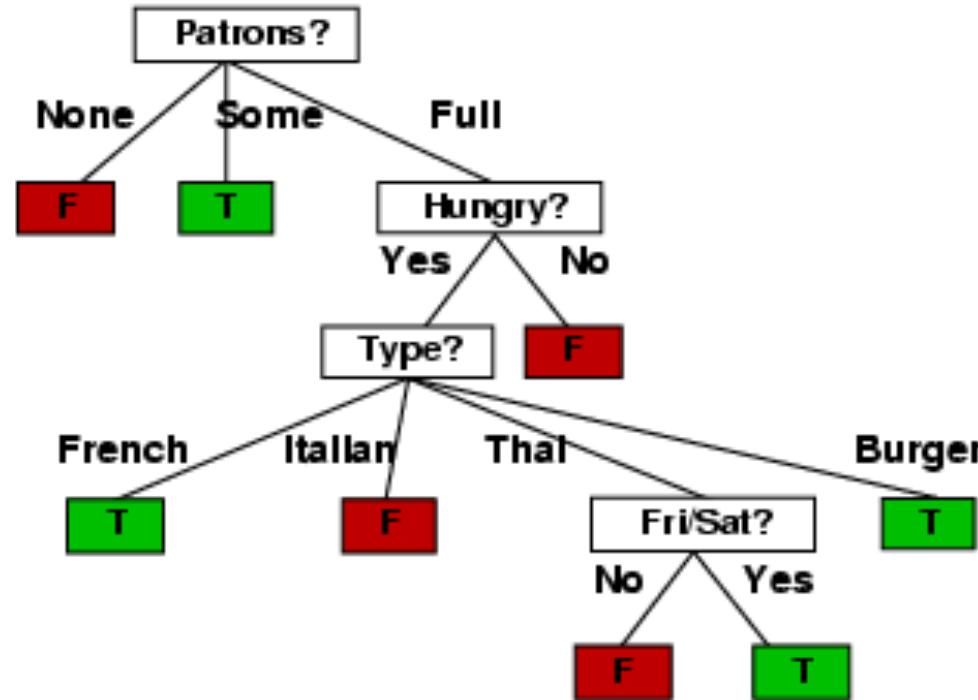
.....

*Hungrey has the highest IG of all attributes and so is chosen by the DTL algorithm as the root*

Example	Attributes											Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	Wait	
$X_2$	T	F	F	T	Full	\$	F	F	Thai	30–60	F	
$X_4$	T	F	T	T	Full	\$	F	F	Thai	10–30	T	
$X_5$	T	F	T	F	Full	\$\$\$	F	T	French	>60	F	
$X_9$	F	T	T	F	Full	\$	T	F	Burger	>60	F	
$X_{10}$	T	T	T	T	Full	\$\$\$	F	T	Italian	10–30	F	
$X_{12}$	T	T	T	T	Full	\$	F	F	Burger	30–60	T	

# Example contd.

- Decision tree learned from the 12 examples:



- Substantially simpler than “true” tree---a more complex hypothesis isn’t justified by small amount of data

# Bayes' theorem

---

Bayes' theorem describes the probability of an event, based on prior knowledge of conditions that might be related to the event

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

# Naïve Bayes

---

## Preliminaries and Notations

$\omega_i \in \{ \omega_1, \omega_2, \dots, \omega_c \}$  : a state of nature

$x$  : feature vector

$p(\omega_i)$  : prior probability

$p(x|\omega_i)$  : class-conditional density

$p(\omega_i|x)$  : posterior probability

# Naïve Bayes

---

## Fish Sorting Example

-Define  $\omega$  as the type of fish we observe (state of nature, class) where

$\omega = \omega_1$  for sea bass;

$\omega = \omega_2$  for salmon.

$-p(\omega_1)$  is the a priori probability that the next fish is a sea bass.

$-p(\omega_2)$  is the a priori probability that the next fish is a salmon.

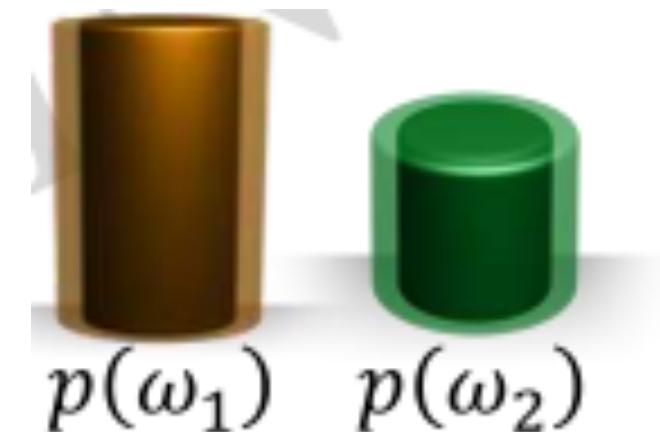
# Naïve Bayes

---

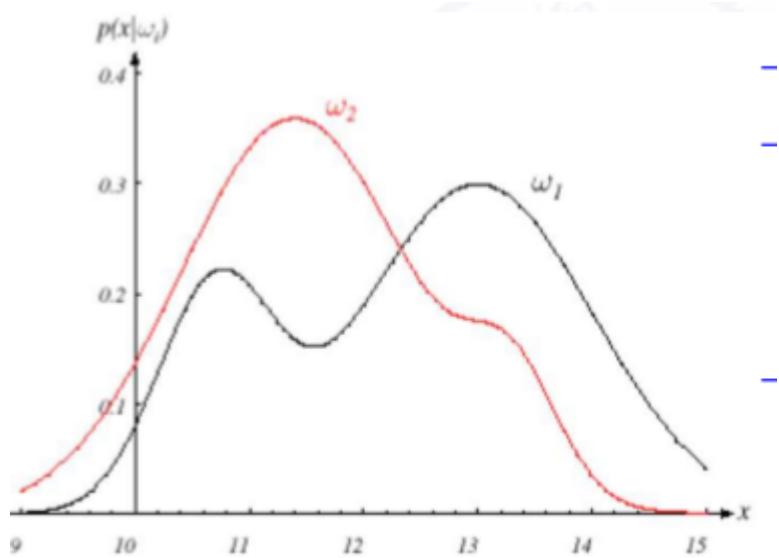
If we make decision with only the prior information

$$\text{Decision} \begin{cases} \omega_1 & p(\omega_1) > p(\omega_2) \\ \omega_2 & p(\omega_1) \leq p(\omega_2) \end{cases}$$

$$p(\text{error}) = \min\{ p(\omega_1), p(\omega_2) \}$$



# Naïve Bayes



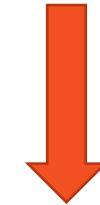
$x$  : continuous random variable. In this example,  $x$  is the lightness of fish.

- $p(x|\omega_i)$  : the class-conditional probability density (probability of  $x$  given that the state of nature is  $\omega_i$  for  $i = 1, 2$ )
- $p(x|\omega_1)$  and  $p(x|\omega_2)$  : the difference in lightness between populations of sea bass and salmon.

# Naïve Bayes

How can we make a decision after observing the value of  $x$ ?

$$\text{Decision} \begin{cases} \omega_1 & p(\omega_1|x) > p(\omega_2|x) \\ \omega_2 & p(\omega_1|x) \leq p(\omega_2|x) \end{cases}$$

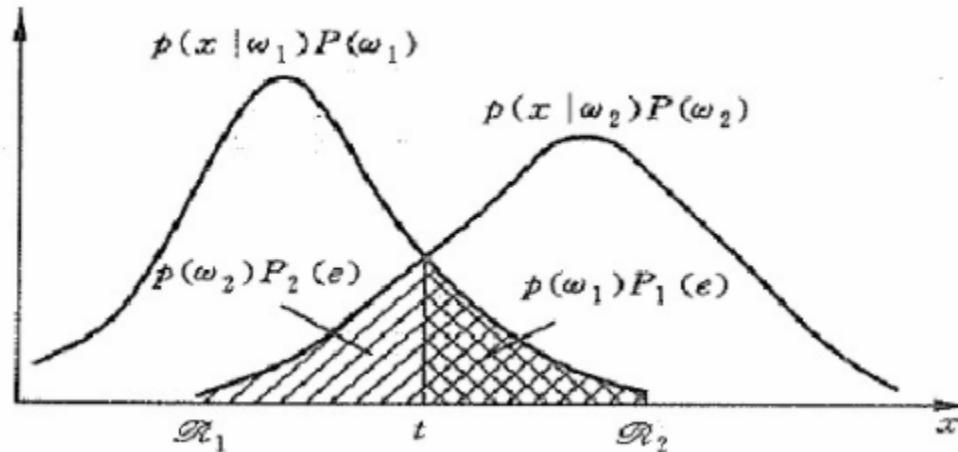


$$p(\omega_i|x) = \frac{p(x|\omega_i)p(\omega_i)}{p(x)}$$

$$\text{Decision} \begin{cases} \omega_1 & p(x|\omega_1)p(\omega_1) > p(x|\omega_2)p(\omega_2) \\ \omega_2 & p(x|\omega_1)p(\omega_1) \leq p(x|\omega_2)p(\omega_2) \end{cases}$$

# Naïve Bayes

$$p(\text{error}) = \begin{cases} p(\omega_1) & \text{decision} = \omega_2 \\ p(\omega_2) & \text{decision} = \omega_2 \end{cases} = \int p(\text{error}) dx$$



# Naïve Bayes

---

## Exercise1:

Two types of cells : normal( $\omega_1$ ) and abnormal( $\omega_2$ ), the prior probabilities are known :

Normal( $\omega_1$ ) :  $p(\omega_1) = 0.9$

Abnormal( $\omega_2$ ) :  $p(\omega_2) = 0.1$

For the given cell  $x$  , we know its class-conditional probabilities:

$p(x|\omega_1) = 0.2$  ,  $p(x|\omega_2) = 0.4$

This cell  $x$  is normal or abnormal?

# Naïve Bayes

---

## Exercise1:

$$p(\omega_1) = 0.9$$

$$p(\omega_2) = 0.1$$

$$p(x|\omega_1) = 0.2, p(x|\omega_2) = 0.4$$

$$\text{Decision} \begin{cases} \omega_1 & p(x|\omega_1)p(\omega_1) > p(x|\omega_2)p(\omega_2) \\ \omega_2 & p(x|\omega_1)p(\omega_1) \leq p(x|\omega_2)p(\omega_2) \end{cases}$$

$$p(x|\omega_1)p(\omega_1) = 0.9 \times 0.2 = 0.18$$

$$p(x|\omega_2)p(\omega_2) = 0.1 \times 0.4 = 0.04$$

Normal

# Naïve Bayes

Given the title of an article, decide whether it is about sports or not

the training set is:

Text	Category
A great game	Sports
The election was over	Not Sports
Very clean match	Sports
A clean but forgettable game	Sports
It was a close election	Not sports

Is the <A very close game> about sports or not?

# Naïve Bayes

$P(\text{a very close game} | \text{Sports})$

$$= P(a | \text{Sports}) \times P(\text{very} | \text{Sports}) \times P(\text{close} | \text{Sports}) \times P(\text{game} | \text{Sports})$$

$P(\text{a very close game} | \text{Not Sports})$

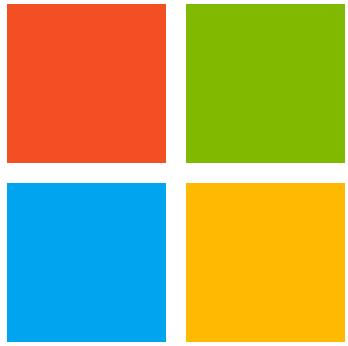
$$= P(a | \text{NS}) \times P(\text{very} | \text{NS}) \times P(\text{close} | \text{NS}) \times P(\text{game} | \text{NS})$$

Word	$P(\text{word}   \text{Sports})$	$P(\text{word}   \text{Not Sports})$
a	$\frac{2 + 1}{11 + 14}$	$\frac{1 + 1}{9 + 14}$
very	$\frac{1 + 1}{11 + 14}$	$\frac{0 + 1}{9 + 14}$
close	$\frac{0 + 1}{11 + 14}$	$\frac{1 + 1}{9 + 14}$
game	$\frac{2 + 1}{11 + 14}$	$\frac{0 + 1}{9 + 14}$

# Why Naive?

---

We assume that every word in a sentence has nothing to do with other words. This means that we do not look at the whole sentence, but the single word.



Microsoft



微软学生俱乐部

扫一扫二维码，加入该群。



关注TJMSC微信公众号