

بە نام پەگانە معۇد بىخىنلە مەربان

# مبانی یادگیری ماشین

## Machine Learning Foundations

گروه هوش مصنوعی، دانشکده مهندسی کامپیوتر، دانشگاه اصفهان

ترم اول سال تحصیلی ۱۴۰۲

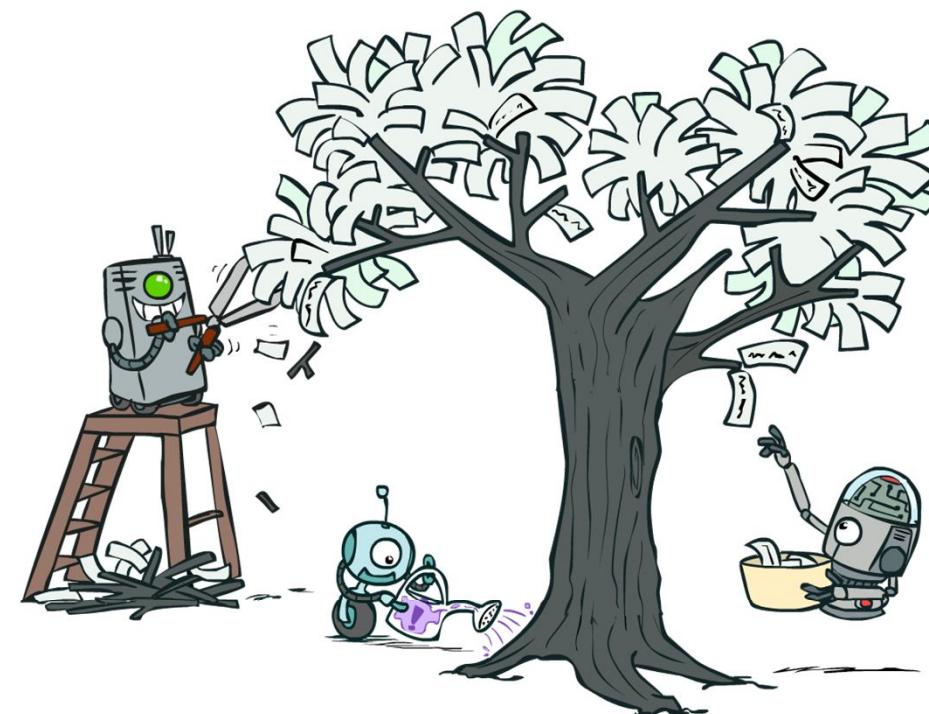
ارائه دهنده: پیمان ادبی

درخت تصمیم

Decision Tree

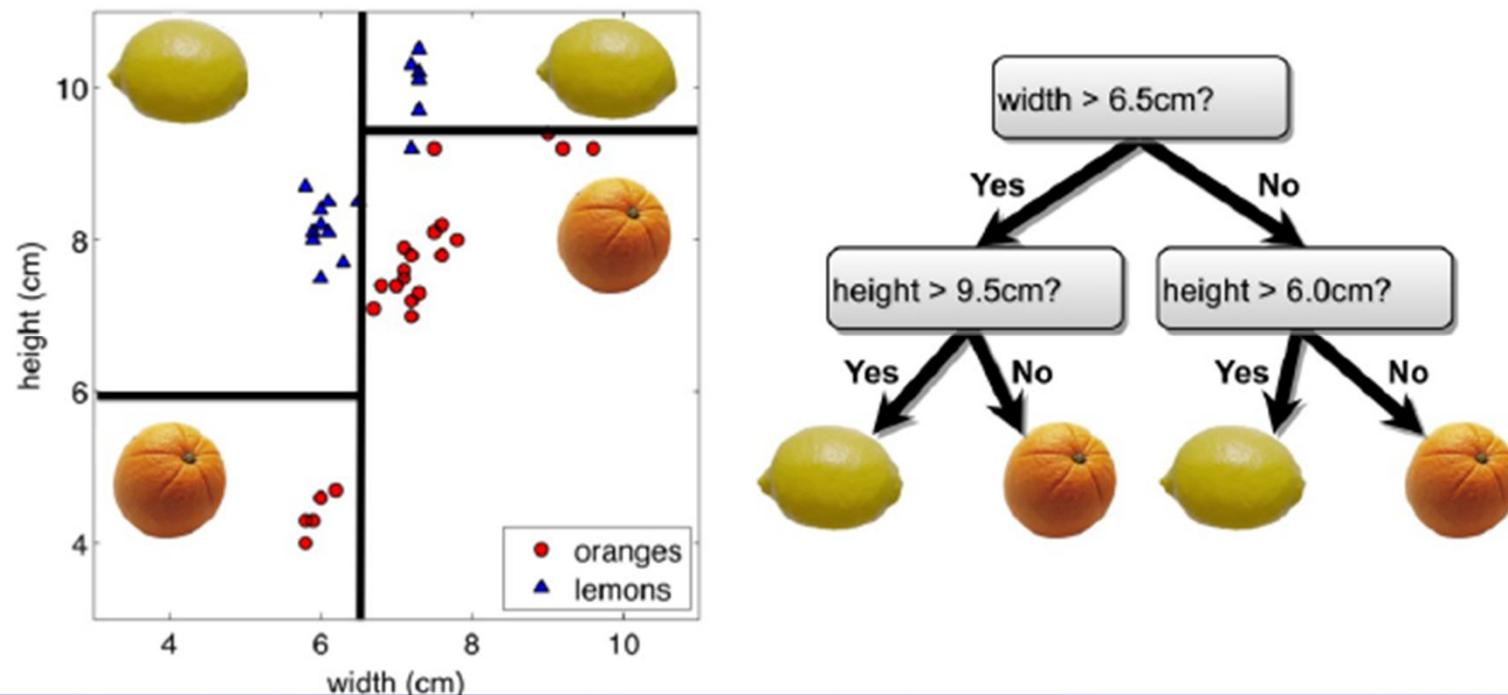
# درخت های تصمیم

- Decision Trees
  - ▶ entropy
  - ▶ information gain



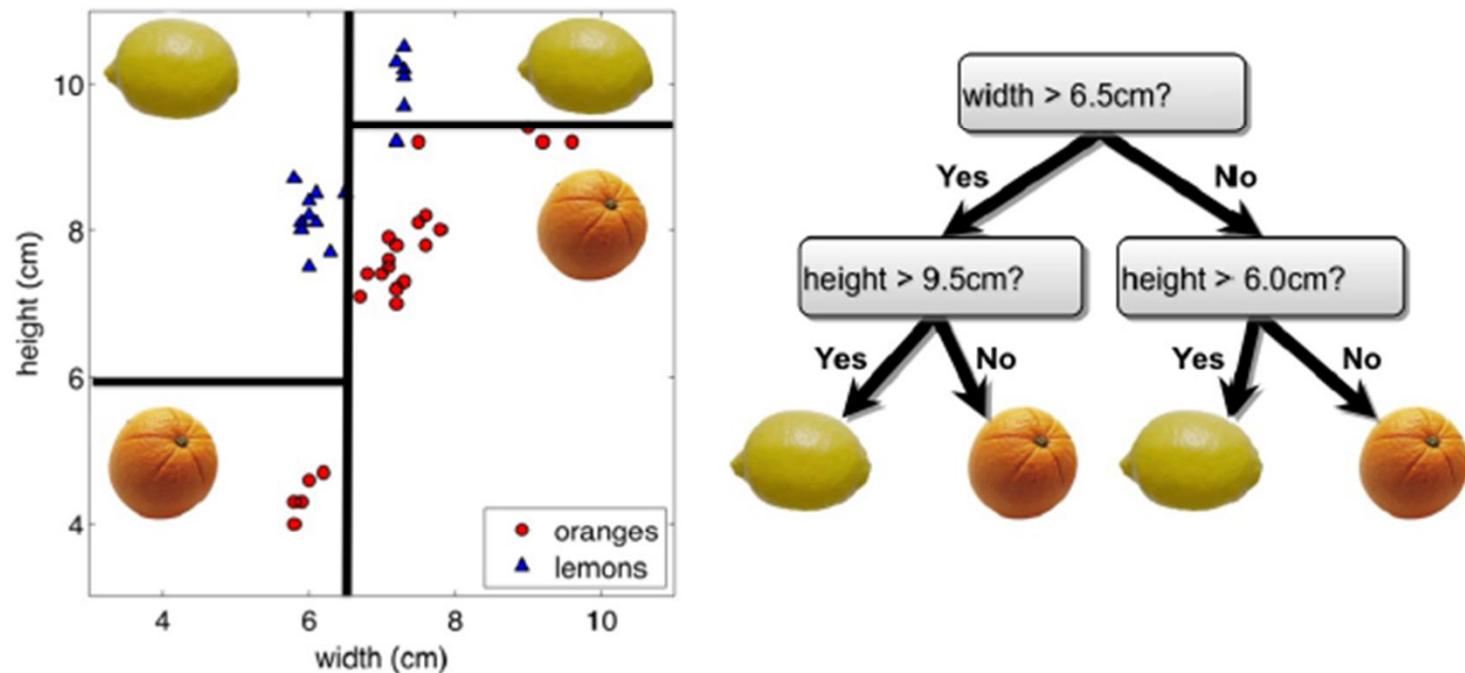
# ایده دیگری برای دسته بندی

- We tried linear classification (eg, logistic regression), and nearest neighbors.  
Any other idea?
- Pick an attribute, do a simple test
- Conditioned on a choice, pick another attribute, do another test
- In the leaves, assign a class with majority vote
- Do other branches as well

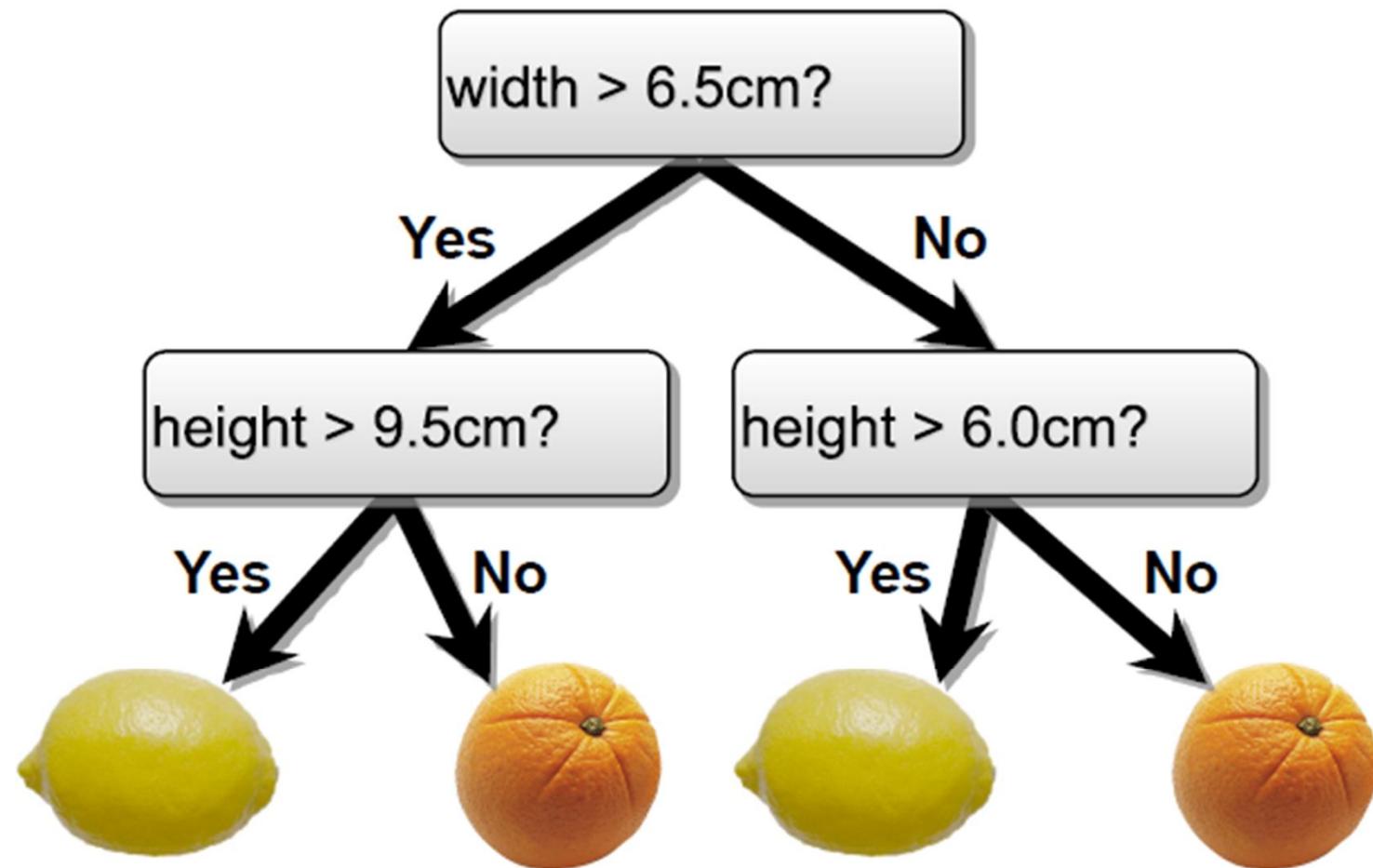


# ایده دیگری برای دسته بندی

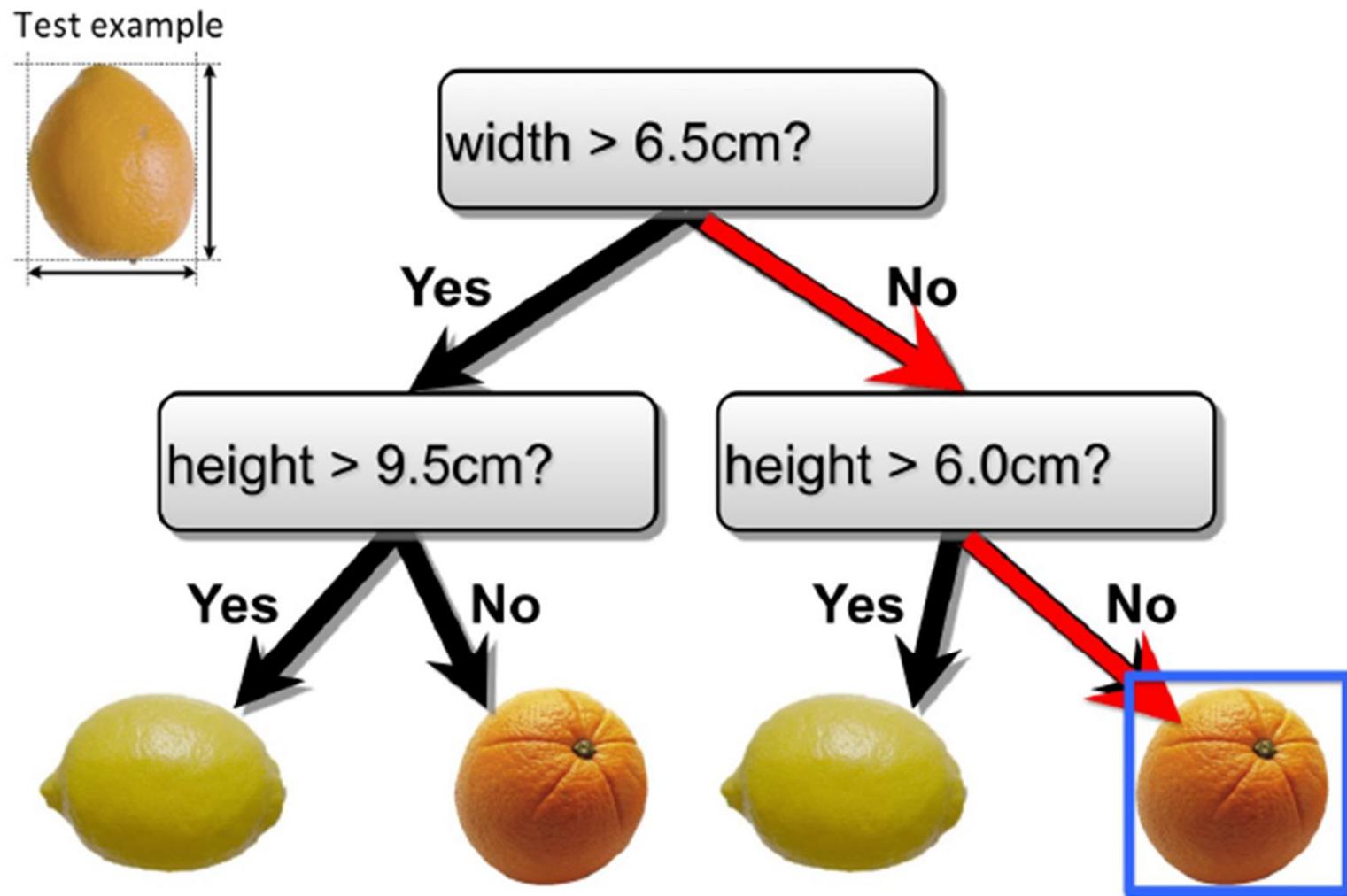
- Gives axes aligned decision boundaries



# درخت تصمیم - مثال



# درخت تصمیم - مثال



# درخت تصمیم - مثال با ویژگی های گستته

- What if the attributes are discrete?

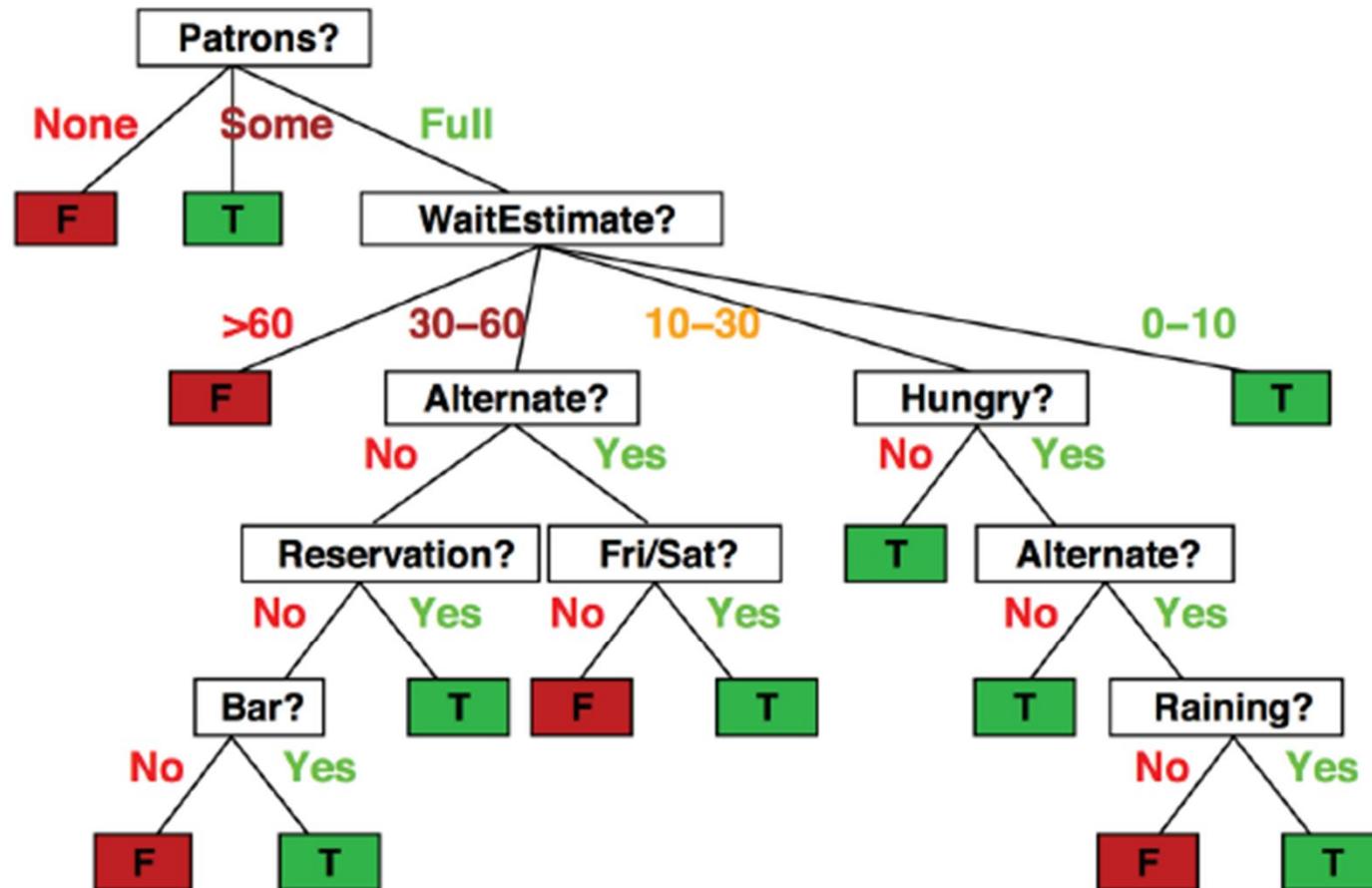
Example	Input Attributes										Goal WillWait
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
$x_1$	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0-10	$y_1 = \text{Yes}$
$x_2$	Yes	No	No	Yes	Full	\$	No	No	Thai	30-60	$y_2 = \text{No}$
$x_3$	No	Yes	No	No	Some	\$	No	No	Burger	0-10	$y_3 = \text{Yes}$
$x_4$	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10-30	$y_4 = \text{Yes}$
$x_5$	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = \text{No}$
$x_6$	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0-10	$y_6 = \text{Yes}$
$x_7$	No	Yes	No	No	None	\$	Yes	No	Burger	0-10	$y_7 = \text{No}$
$x_8$	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0-10	$y_8 = \text{Yes}$
$x_9$	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$
$x_{10}$	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10-30	$y_{10} = \text{No}$
$x_{11}$	No	No	No	No	None	\$	No	No	Thai	0-10	$y_{11} = \text{No}$
$x_{12}$	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30-60	$y_{12} = \text{Yes}$

Attributes:

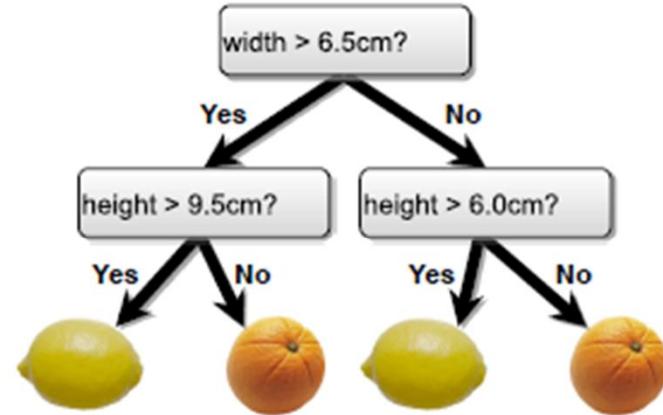
1. Alternate: whether there is a suitable alternative restaurant nearby.
2. Bar: whether the restaurant has a comfortable bar area to wait in.
3. Fri/Sat: true on Fridays and Saturdays.
4. Hungry: whether we are hungry.
5. Patrons: how many people are in the restaurant (values are None, Some, and Full).
6. Price: the restaurant's price range (\$, \$\$, \$\$\$).
7. Raining: whether it is raining outside.
8. Reservation: whether we made a reservation.
9. Type: the kind of restaurant (French, Italian, Thai or Burger).
10. WaitEstimate: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60).

# درخت تصمیم - مثال با ویژگی های گستته

- The tree to decide whether to wait (T) or not (F)



# درخت تصمیم - ساختار



- Internal nodes **test attributes**
- Branching is determined by **attribute value**
- Leaf nodes are **outputs** (class assignments)

# درخت تصمیم - الگوریتم

- Choose an attribute on which to descend at each level.
- Condition on earlier (higher) choices.
- Generally, restrict only one dimension at a time.
- Declare an output value when you get to the bottom
- In the orange/lemon example, we only split each dimension once, but that is not required.

# درخت تصمیم - دسته بندی و رگرسیون

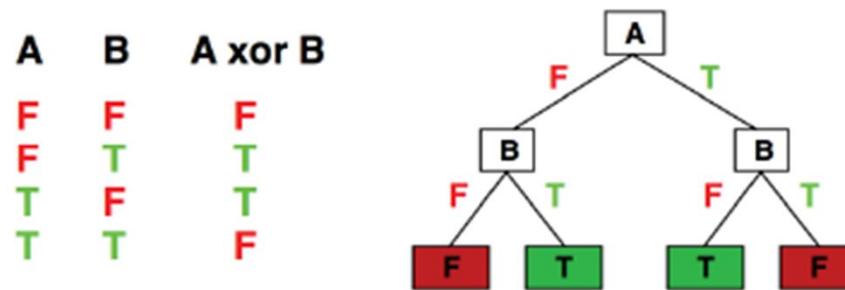
- Each path from root to a leaf defines a region  $R_m$  of input space
- Let  $\{(x^{(m_1)}, t^{(m_1)}), \dots, (x^{(m_k)}, t^{(m_k)})\}$  be the training examples that fall into  $R_m$
- Classification tree:
  - ▶ discrete output
  - ▶ leaf value  $y^m$  typically set to the most common value in  $\{t^{(m_1)}, \dots, t^{(m_k)}\}$
- Regression tree:
  - ▶ continuous output
  - ▶ leaf value  $y^m$  typically set to the mean value in  $\{t^{(m_1)}, \dots, t^{(m_k)}\}$

Note: We will only talk about classification

[Slide credit: S. Russell]

# درخت تصمیم - قدرت بیان

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



- Continuous-input, continuous-output case:
  - ▶ Can approximate any function arbitrarily closely
  - Trivially, there is a consistent decision tree for any training set w/ one path to leaf for each example (unless  $f$  nondeterministic in  $x$ ) but it probably won't generalize to new examples

Need some kind of regularization to ensure more **compact** decision trees

[Slide credit: S. Russell]

# درخت تصمیم - یادگیری

- How do we construct a useful decision tree?

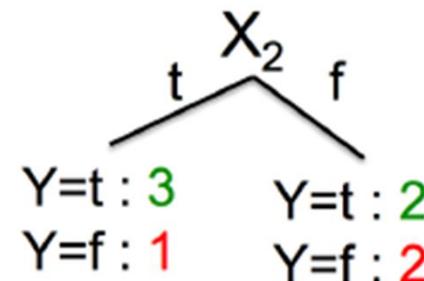
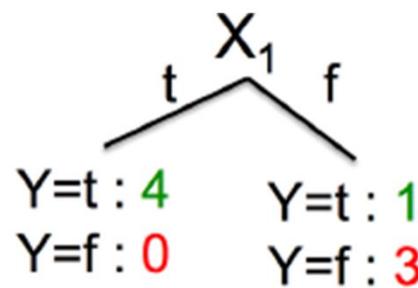
Learning the simplest (smallest) decision tree is an NP complete problem [if you are interested, check: Hyafil & Rivest'76]

- Resort to a greedy heuristic:
  - ▶ Start from an empty decision tree
  - ▶ Split on next best attribute
  - ▶ Recurse
- What is **best** attribute?
- We use [information theory](#) to guide us

[Slide credit: D. Sonntag]

# انتخاب یک ویژگی خوب

- Which attribute is better to split on,  $X_1$  or  $X_2$ ?



$X_1$	$X_2$	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F
F	T	F
F	F	F

Idea: Use counts at leaves to define probability distributions, so we can measure uncertainty

# انتخاب یک ویژگی خوب

- Which attribute is better to split on,  $X_1$  or  $X_2$ ?
  - ▶ Deterministic: good (all are true or false; just one class in the leaf)
  - ▶ Uniform distribution: bad (all classes in leaf equally probable)
  - ▶ What about distributions in between?

Note: Let's take a slight detour and remember concepts from information theory

[Slide credit: D. Sonntag]

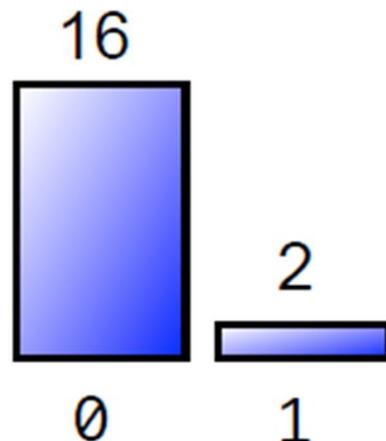
# مثال: پرتاب دو سکه

Sequence 1:

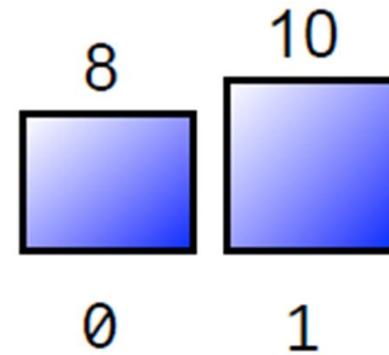
0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 1 0 0 ... ?

Sequence 2:

0 1 0 1 0 1 1 1 0 1 0 0 1 1 0 1 0 1 0 1 ... ?



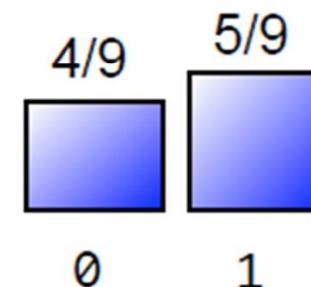
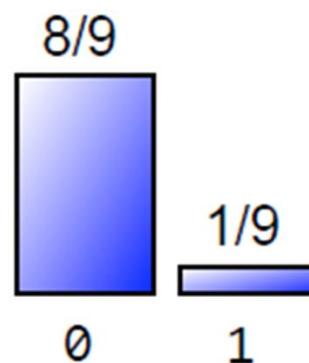
versus



# مثال: پرتاب دو سکه - کمی سازی عدم قطعیت

Entropy  $H$ :

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$



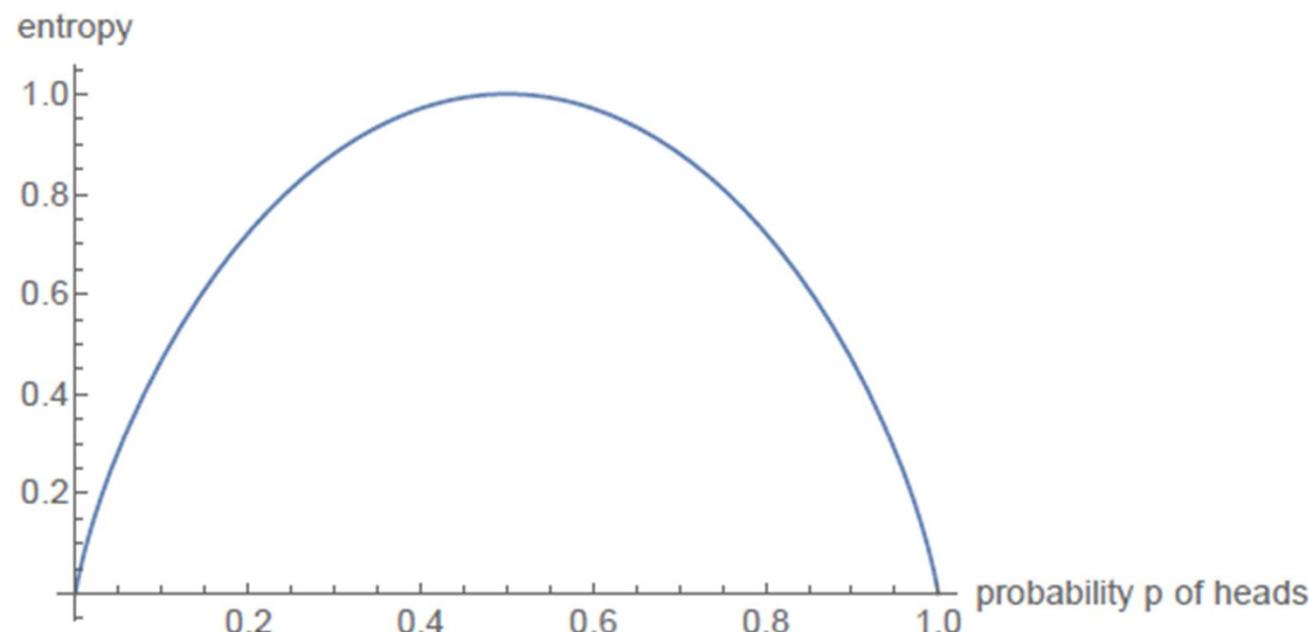
$$-\frac{8}{9} \log_2 \frac{8}{9} - \frac{1}{9} \log_2 \frac{1}{9} \approx \frac{1}{2}$$

$$-\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} \approx 0.99$$

- How surprised are we by a new value in the sequence?
- How much information does it convey?

# مثال: پرتاب دو سکه - کمی سازی عدم قطعیت

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$



# انتروپی

- “High Entropy”:
  - ▶ Variable has a uniform like distribution
  - ▶ Flat histogram
  - ▶ Values sampled from it are less predictable
- “Low Entropy”
  - ▶ Distribution of variable has many peaks and valleys
  - ▶ Histogram has many lows and highs
  - ▶ Values sampled from it are more predictable

[Slide credit: Vibhav Gogate]

# انتروپی توأم

- Example:  $X = \{\text{Raining, Not raining}\}$ ,  $Y = \{\text{Cloudy, Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

$$\begin{aligned} H(X, Y) &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y) \\ &= -\frac{24}{100} \log_2 \frac{24}{100} - \frac{1}{100} \log_2 \frac{1}{100} - \frac{25}{100} \log_2 \frac{25}{100} - \frac{50}{100} \log_2 \frac{50}{100} \\ &\approx 1.56 \text{ bits} \end{aligned}$$

# انتروپی شرطی مشخص

- Example:  $X = \{\text{Raining, Not raining}\}$ ,  $Y = \{\text{Cloudy, Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- What is the entropy of cloudiness  $Y$ , given that it is raining?

$$\begin{aligned} H(Y|X=x) &= -\sum_{y \in Y} p(y|x) \log_2 p(y|x) \\ &= -\frac{24}{25} \log_2 \frac{24}{25} - \frac{1}{25} \log_2 \frac{1}{25} \\ &\approx 0.24 \text{ bits} \end{aligned}$$

- We used:  $p(y|x) = \frac{p(x,y)}{p(x)}$ , and  $p(x) = \sum_y p(x,y)$  (sum in a row)

# انتروپی شرطی

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- The expected conditional entropy:

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} p(x)H(Y|X=x) \\ &= -\sum_{x \in X} \sum_{y \in Y} p(x,y) \log_2 p(y|x) \end{aligned}$$

# انتروپی شرطی

- Example:  $X = \{\text{Raining, Not raining}\}$ ,  $Y = \{\text{Cloudy, Not cloudy}\}$

	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- What is the entropy of cloudiness, given the knowledge of whether or not it is raining?

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} p(x)H(Y|X=x) \\ &= \frac{1}{4}H(\text{cloudy}| \text{is raining}) + \frac{3}{4}H(\text{cloudy}| \text{not raining}) \\ &\approx 0.75 \text{ bits} \end{aligned}$$

# انتروپی شرطی

- Some useful properties:
  - ▶  $H$  is always non-negative
  - ▶ Chain rule:  $H(X, Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$
  - ▶ If  $X$  and  $Y$  independent, then  $X$  doesn't tell us anything about  $Y$ :  
 $H(Y|X) = H(Y)$
  - ▶ But  $Y$  tells us everything about  $Y$ :  $H(Y|Y) = 0$
  - ▶ By knowing  $X$ , we can only decrease uncertainty about  $Y$ :  
 $H(Y|X) \leq H(Y)$

# بھرہ اطلاعاتی

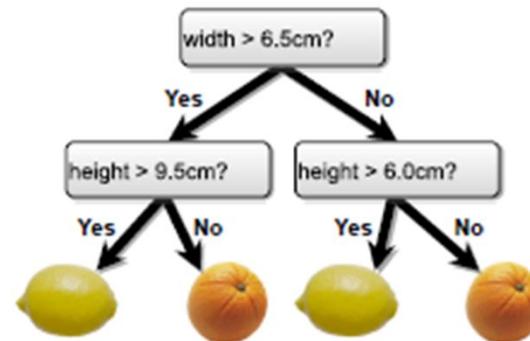
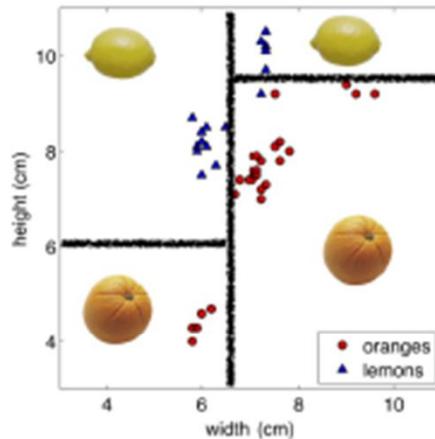
	Cloudy	Not Cloudy
Raining	24/100	1/100
Not Raining	25/100	50/100

- How much information about cloudiness do we get by discovering whether it is raining?

$$\begin{aligned} IG(Y|X) &= H(Y) - H(Y|X) \\ &\approx 0.25 \text{ bits} \end{aligned}$$

- Also called **information gain** in  $Y$  due to  $X$
- If  $X$  is completely uninformative about  $Y$ :  $IG(Y|X) = 0$
- If  $X$  is completely informative about  $Y$ :  $IG(Y|X) = H(Y)$
- How can we use this to construct our decision tree?

# ساخت درخت های تصمیم



- I made the fruit data partitioning just by eyeballing it.
- We can use the **information gain** to automate the process.
- At each level, one must choose:
  - Which variable to split.
  - Possibly where to split it.
- Choose them based on how much information we would gain from the decision! (choose attribute that gives the highest gain)

# ساخت درخت های تصمیم - الگوریتم

- Simple, greedy, recursive approach, builds up tree node-by-node
  1. pick an attribute to split at a non-terminal node
  2. split examples into groups based on attribute value
  3. for each group:
    - ▶ if no examples – return majority from parent
    - ▶ else if all examples in same class – return class
    - ▶ else loop to step 1

# درخت تصمیم - مثال قبلی

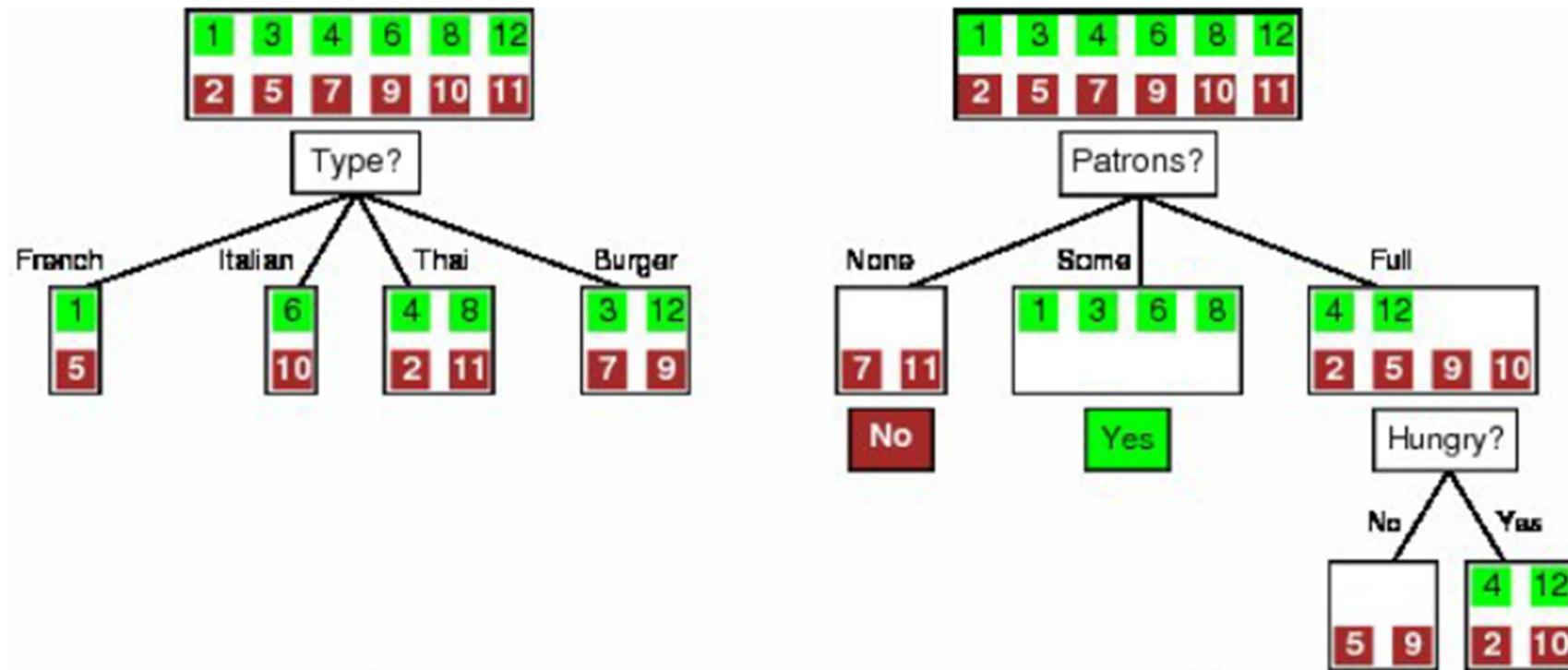
Example	Input Attributes										Goal <i>WillWait</i>
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
x <sub>1</sub>	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0-10	$y_1 = \text{Yes}$
x <sub>2</sub>	Yes	No	No	Yes	Full	\$	No	No	Thai	30-60	$y_2 = \text{No}$
x <sub>3</sub>	No	Yes	No	No	Some	\$	No	No	Burger	0-10	$y_3 = \text{Yes}$
x <sub>4</sub>	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10-30	$y_4 = \text{Yes}$
x <sub>5</sub>	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = \text{No}$
x <sub>6</sub>	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0-10	$y_6 = \text{Yes}$
x <sub>7</sub>	No	Yes	No	No	None	\$	Yes	No	Burger	0-10	$y_7 = \text{No}$
x <sub>8</sub>	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0-10	$y_8 = \text{Yes}$
x <sub>9</sub>	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$
x <sub>10</sub>	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10-30	$y_{10} = \text{No}$
x <sub>11</sub>	No	No	No	No	None	\$	No	No	Thai	0-10	$y_{11} = \text{No}$
x <sub>12</sub>	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30-60	$y_{12} = \text{Yes}$

1. Alternate: whether there is a suitable alternative restaurant nearby.
2. Bar: whether the restaurant has a comfortable bar area to wait in.
3. Fri/Sat: true on Fridays and Saturdays.
4. Hungry: whether we are hungry.
5. Patrons: how many people are in the restaurant (values are None, Some, and Full).
6. Price: the restaurant's price range (\$, \$\$, \$\$\$).
7. Raining: whether it is raining outside.
8. Reservation: whether we made a reservation.
9. Type: the kind of restaurant (French, Italian, Thai or Burger).
10. WaitEstimate: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60).

Attributes:

[from: Russell & Norvig]

# انتخاب ویژگی ها

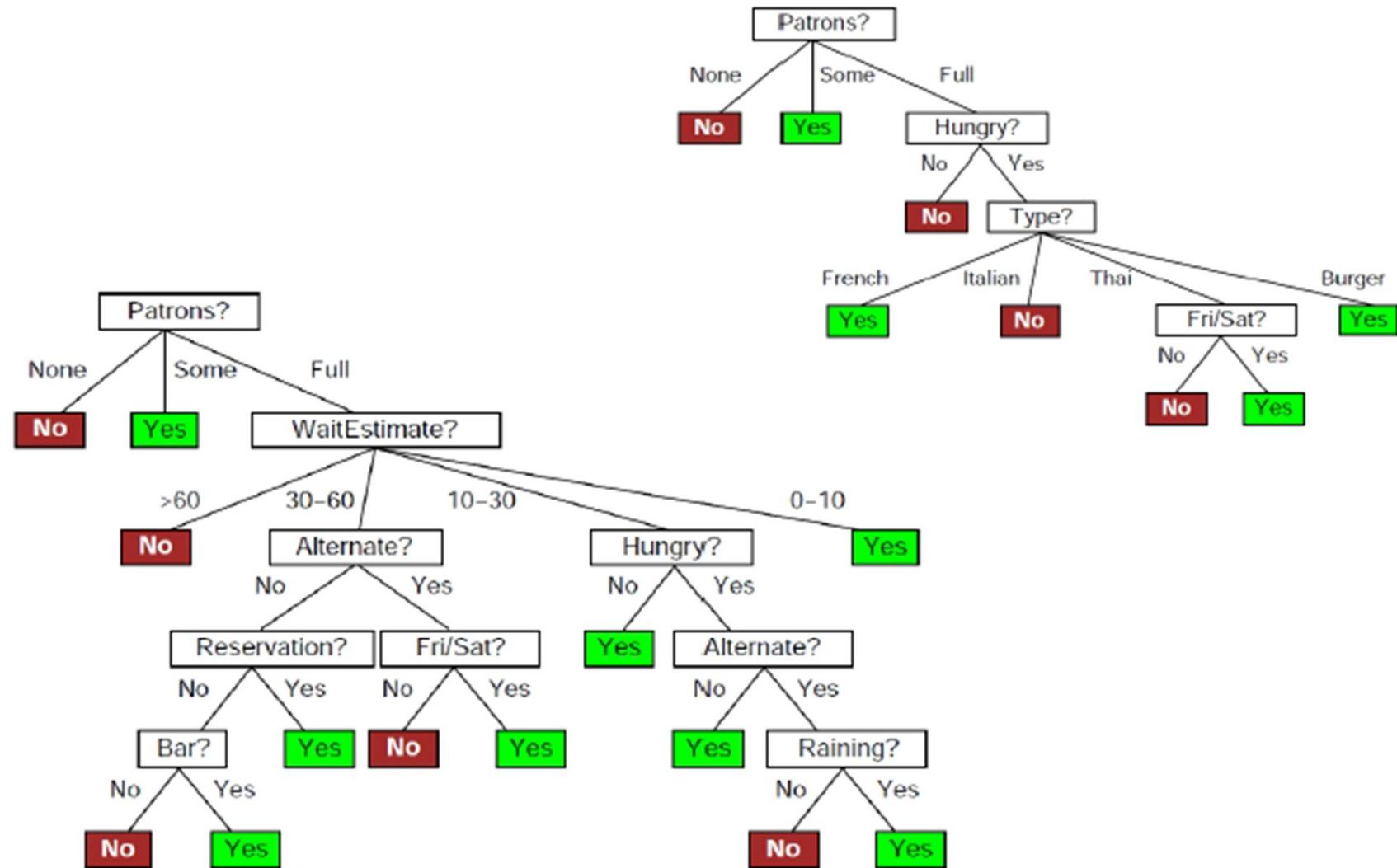


$$IG(Y) = H(Y) - H(Y|X)$$

$$IG(type) = 1 - \left[ \frac{2}{12}H(Y|Fr.) + \frac{2}{12}H(Y|It.) + \frac{4}{12}H(Y|Thai) + \frac{4}{12}H(Y|Bur.) \right] = 0$$

$$IG(Patrons) = 1 - \left[ \frac{2}{12}H(0, 1) + \frac{4}{12}H(1, 0) + \frac{6}{12}H\left(\frac{2}{6}, \frac{4}{6}\right) \right] \approx 0.541$$

# کدام درخت بهتر است؟



✓382

# نکاتی پیرامون درخت تصمیم

## قاعده تیغ اکام (Occam's Razor)

- یافتن ساده‌ترین فرضیه که با مشاهدات مطابقت داشته باشد.

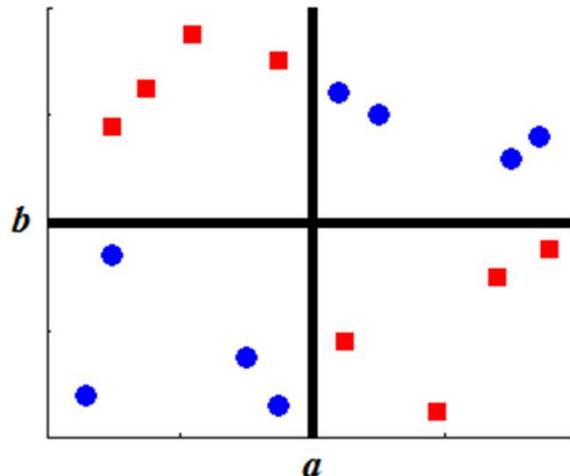
نتیجه

▪ پیشگیری از بیش برآذش

▪ کاهش محاسبات

## مشکلات

- تعداد نمونه‌ها در سطوح پایین تر بصورت نمایی کم می‌شود



- بزرگ بودن درخت میتواند موجب بیش برآذش شود

▪ الگوریتمهای حریصانه الزاماً بهینه سراسری

را به دست نمی‌آورند ←

## تنظیم

- جلوگیری از زیاد شدن عمق درخت

▪ مثلًاً با هرس کردن براساس داده اعتبار سنجی

در حین ساخت درخت

- قابل استفاده برای رگرسیون البته با فرمولاسیون متفاوت

# نکاتی پیرامون درخت تصمیم

## ■ مزایا

- مناسب برای بیان هر تابع بولین، بخصوص اگر تابع به تعداد ویژگی اندکی وابستگی جدی داشته باشد

## ■ معایب

- برای توابع بولین اکثریت (majority) و توازن (parity) و ویژگی های پیوسته خیلی مناسب نیست

## ■ برخی کاربردهای موفق

- شناسایی بی درنگ موقعیت بخش های بدن از تک تصویر عمق (XBox ↓)
- شبیه ساز پرواز (نود هزار نمونه از عملکرد خلبانان خبره ← درخت پرواز خودکار)

- چالش رتبه دهی بازیابی اطلاعات یا هو جنگل های تصادفی

... □

