# COMP 472: Artificial Intelligence
## Machine Learning *part #2*
## Decision Trees *video #4*

- Russell & Norvig: Sections 19.3

# Today

# Guess Who?

# Decision Trees

*Ross Quinlan* [handwritten annotation]
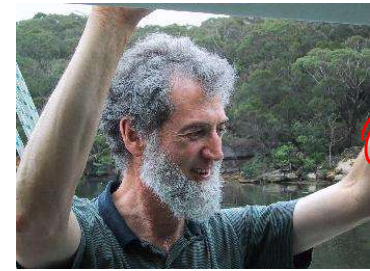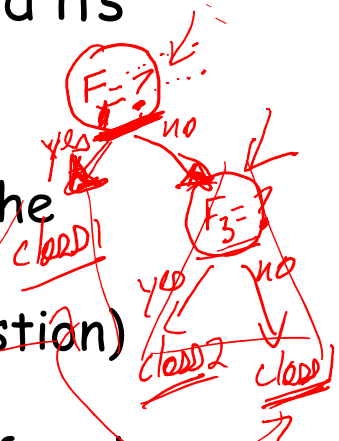
- Simplest, but most successful form of learning algorithm
- Very well-know algorithm is ID3 (Quinlan, 1987) and its successor C4.5 *J-48* [handwritten]

1. Rank features based on how good they are to indicate the result
2. Put the most discriminating feature as a node (as a question) of a tree *training* [handwritten]
3. Split the examples so that those with different values for the chosen feature are in a different set
4. Repeat the same process with the nest most discriminating feature
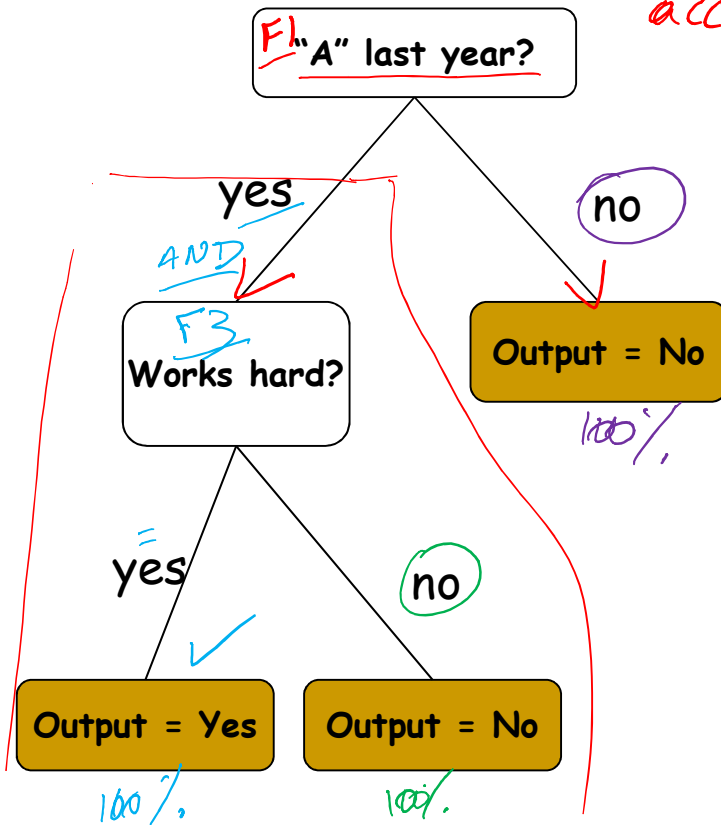
http://www.rulequest.com/Personal/

# Example 1

Info on last year's students to determine if a student will get an 'A' this year

| | Features (X) | | | | Output f(X) class |
|---|---|---|---|---|---|
| Student | 'A' last year? F1 | Black hair? F2 | Works hard? F3 | Drinks? F4 | 'A' this year? |
| X1: Richard | Yes | Yes | No | Yes | No |
| X2: Alan | Yes | Yes | Yes | No | Yes |
| X3: Alison | No | No | Yes | No | No |
| X4: Jeff | No | Yes | No | Yes | No |
| X5: Gail | Yes | No | Yes | Yes | Yes |
| X6: Simon | No | Yes | Yes | Yes | No |

# Example 1

A random decision tree that fits the dataset ~~with 100% accuracy~~



F1 "A" last year?

yes     AND     no

F3 Works hard?

Output = No   100%

yes   no

Output = Yes   100%

Output = No   100%

*same data set as previous slide*

| | Features | | | | Output f(X) |
|---|---|---|---|---|---|
| Student | 'A' last year? F1 | Black hair? | Works hard? F3 | Drinks ? | 'A' this year? |
| Richard | Yes | Yes | No | Yes | No . |
| Alan | Yes | Yes | Yes | No | Yes . |
| Alison | No | No | Yes | No | No |
| Jeff | No | Yes | No | Yes | No |
| Gail | Yes | No | Yes | Yes | Yes |
| Simon | No | Yes | Yes | Yes | No |

- fits the training set with 100% accuracy
- but the features are chosen at random so there might be a shorter DT
  — // better

6

# Example 2: The Restaurant

- Goal: learn whether one should wait for a table
- Attributes

1. Alternate: another suitable restaurant nearby
2. Bar: comfortable bar for waiting
3. Fri/Sat: true on Fridays and Saturdays
4. Hungry: whether one is hungry
5. Patrons: how many people are present (none, some, full)
6. Price: price range ($, $$, $$$)
7. Raining: raining outside
8. Reservation: reservation made
9. Type: kind of restaurant (French, Italian, Thai, Burger)
10. WaitEstimate: estimated wait by host (0-10 mins, 10-30, 30-60, >60)

*10 features*

# Example 2: The Restaurant
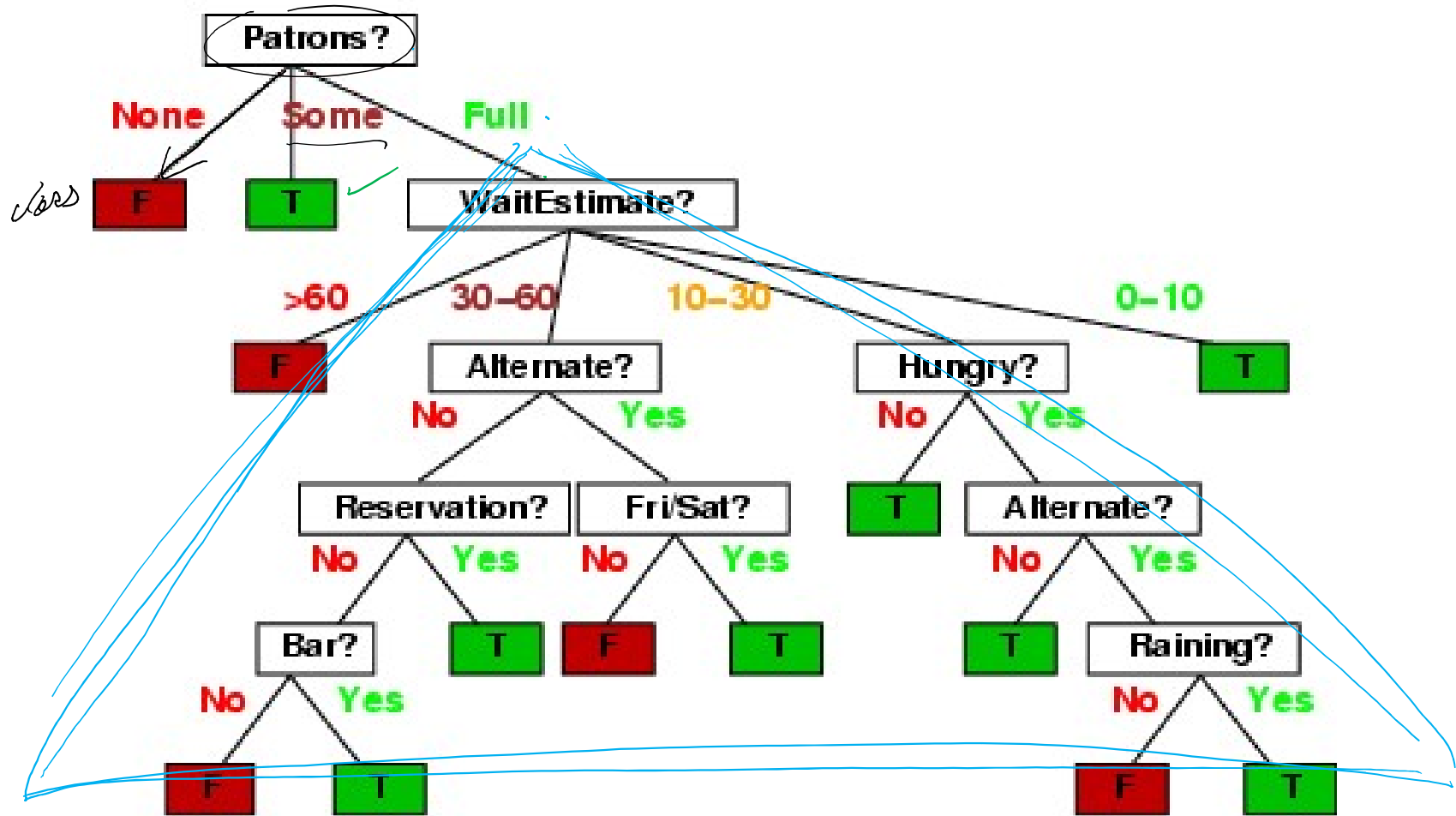
■ Training data:

*10 features.*

*f(x) class*

*12 instances*

| 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 |

source: Norvig (2003)

# A First Decision Tree
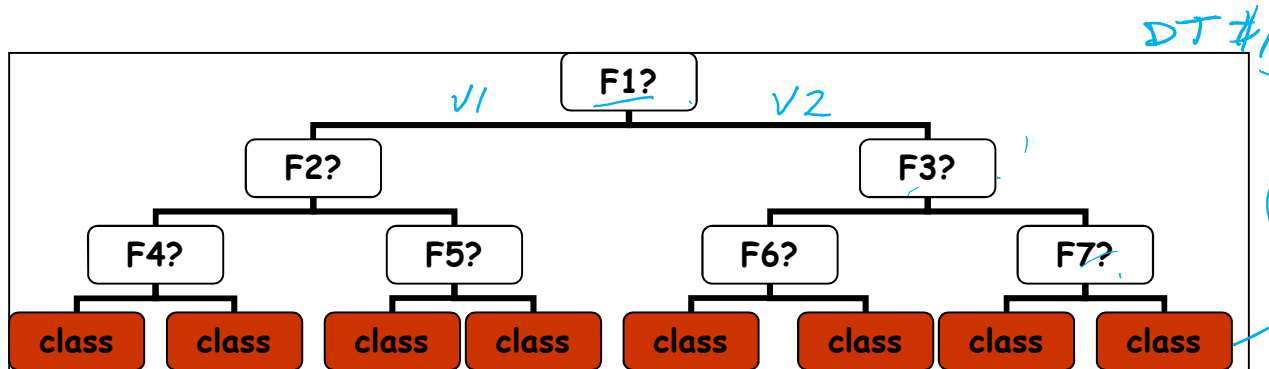


- But is it the best decision tree we can build?

# Ockham's Razor

*It is vain to do more than can be done with less… Entities should not be multiplied beyond necessity. [Ockham, 1324]*
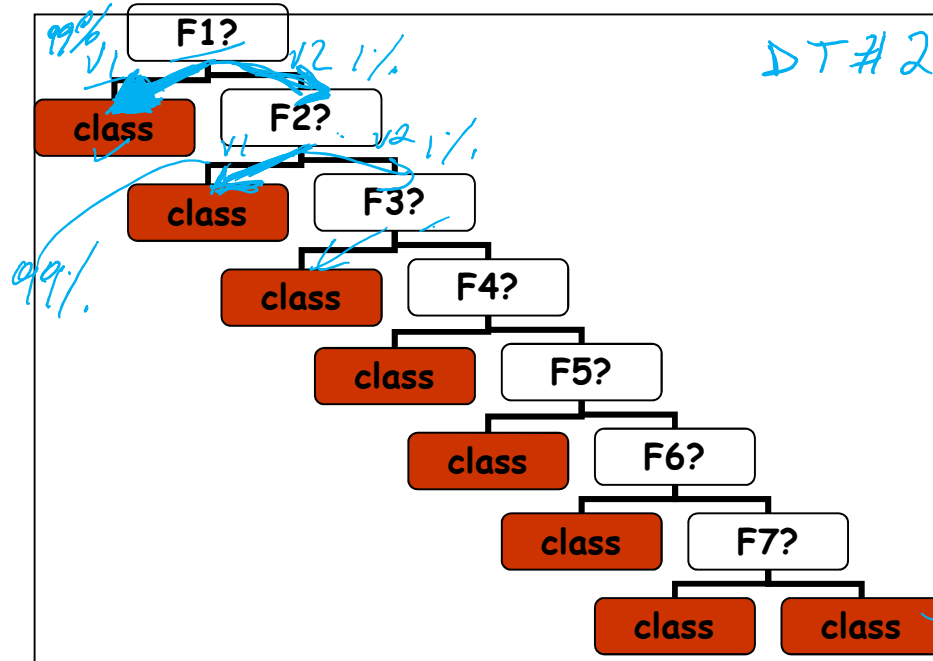
- In other words… always favor the simplest answer that correctly fits the training data

- i.e. the smallest tree on average

- This type of assumption is called inductive bias
  - inductive bias = making a choice beyond what the training instances contain
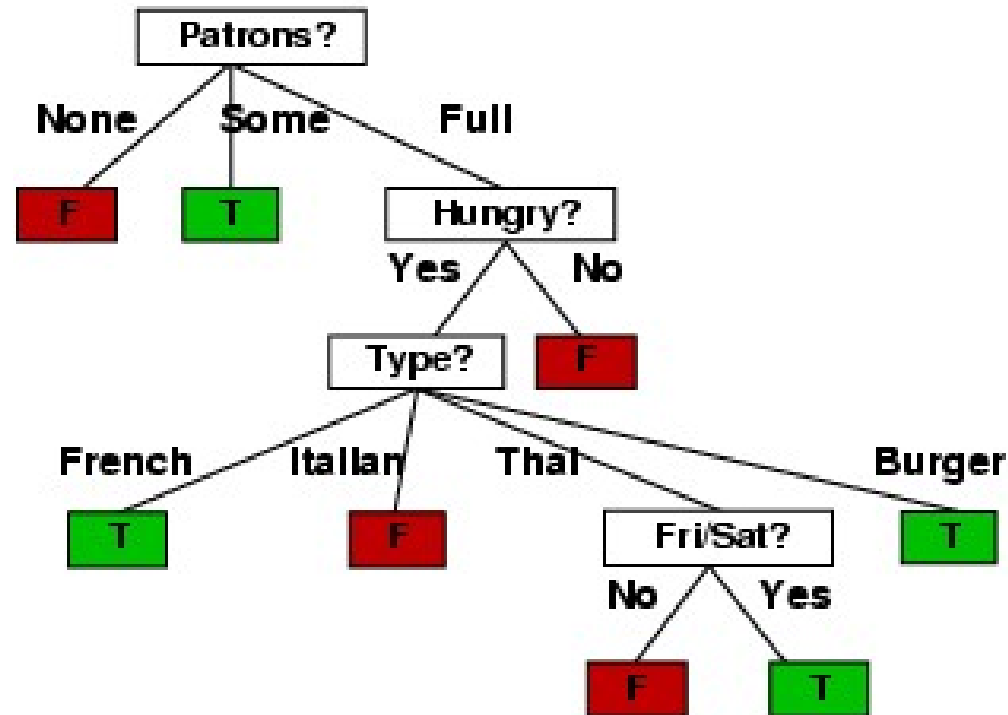
# Which Tree is Best?



DT #1

always asks 3 questions

DT #2

may be the shortest tree on average depending on the probability of each branch

F1?
V1    V2

F2?          F3?

F4?    F5?    F6?    F7?

class class class class class class class class

99%
V1    V2 1%

class
F2?
V1    V2 1%

class
F3?

class
F4?

class
F5?

class
F6?

class
F7?

class class

99%

# A Better Decision Tree

- 4 tests instead of 9
- 11 branches instead of 21



source: Norvig (2003)

# Choosing the Next Feature

- The key problem is choosing which feature to split a given set of examples

- Different measures have been proposed
- ID3 uses Maximum Information-Gain
  - i.e. we choose the feature that has the largest information gain
  - we expect this feature to result in the smallest tree on average
  - based on information theory

# Essential Information Theory

- Developed by *Claude* Shannon in the 40s

- Shannon developed the notion of entropy (aka information content) of a random variable (RV)

- Entropy measures how "predictable" a RV is
  - If you already have a good idea about the answer (e.g. 90/10 split)
    - → low entropy   *// sure thing*
  - If you have no idea about the answer (e.g. 50/50 split)
    - → high entropy   *// Total chaos*
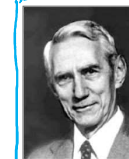


**Dartmouth Conference: The Founding Fathers of AI**

John McCarthy    Marvin Minsky    Claude Shannon    Ray Solomonoff

Alan Newell    Herbert Simon    Arthur Samuel

And three others…
Oliver Selfridge
  (Pandemonium theory)
Nathaniel Rochester
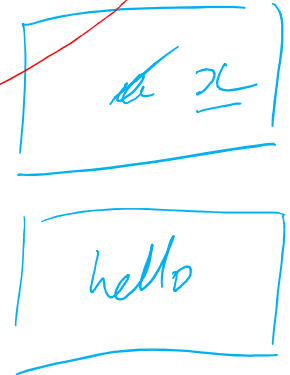  (IBM, designed 701)
Trenchard More
  (Natural Deduction)

# Entropy

- Let X be a discrete RV with *i* possible outcomes $x_i$
- Entropy (or information content) of X

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2 p(x_i)$$

- measures:
  - the *amount of information* in a RV
  - *average uncertainty* of a RV
  - *average length of a message* needed to transmit an outcome $x_i$ of that RV when encoded optimally over a binary channel
- measured in bits

base 2 always

hello

# Entropy of a Coin Toss

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

Entropy (or information content)

$$H(\text{fair coin toss}) = -\sum_{x_i \in X} p(x_i) \log_2 p(x_i) = H\left(\frac{1}{2}, \frac{1}{2}\right)$$

prob $x_1 =$ "head"   prob $x_2 =$ "tail"

$$= -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right) = 1 \text{ bit}$$

head       tail

entropy of a fair coin toss (the RV) with 2 possible outcomes, each with a probability of 1/2

a RV with only 2 outcomes $x_1$ and $x_2$
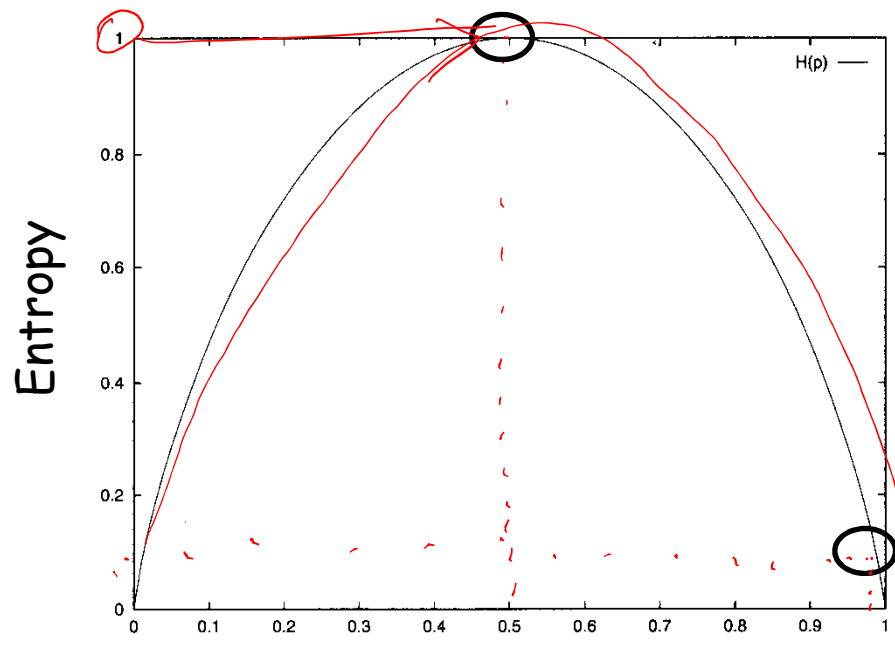will have $1 \geq H(X) \geq 0$            // sure thing

// total chaos

# Example: The Coin Toss

- Fair coin: $H(X) = -\sum_{i=1}^{n} p(x_i)\log_2 p(x_i) = -\left(\dfrac{1}{2}\log_2\dfrac{1}{2} + \dfrac{1}{2}\log_2\dfrac{1}{2}\right) = 1\,\text{bit}$

  *head*   *tail*

- Rigged coin: $H(X) = -\sum_{i=1}^{n} p(x_i)\log_2 p(x_i) = -\left(\dfrac{99}{100}\log_2\dfrac{99}{100} + \dfrac{1}{100}\log_2\dfrac{1}{100}\right) = 0.08\,\text{bits}$

  *head*



fair coin -> high entropy

rigged coin -> low entropy

17

# So what?

*Annotation (handwritten): for each attribute/feature I will compute how much knowing its values will reduce the entropy of the class*

- ## Training data:

*Annotation (handwritten, red): to features*

*Annotation (handwritten, red): RV*

| 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 |

*Annotation (handwritten, red): $x_1$  $x_2$*

*Annotation (handwritten, red): entropy (Target Wait) = 1*

source: Norvig (2003)

18

# Information Gain

- **information gain**
  - measure the entropy reduction of a RV, once a piece of information is known
  - used to measure the "discriminating power" of an attribute A given a data set S

  - Let Values(A) = the set of values that attribute A can take
  - Let $S_v$ = the set of examples in the data set which have value v for attribute A (for each value v from Values(A) )

// see slide 15

information gain (or entropy reduction)

$$gain(S, A) = H(S) - H(S|A)$$

$$= H(S) - \sum_{v \in values(A)} \frac{|S_v|}{|S|} \times H(S_v)$$

# Some Intuition

$$H(output) = 1 \quad // \text{ total chaos}$$

3 features

| Size | Color | Shape | Output |
|------|-------|-------|--------|
| Big | Red | Circle | + |
| Small | Red | Circle | + |
| Small | Red | Square | - |
| Big | Blue | Circle | - |

- *Size* is the least discriminating attribute (i.e. smallest information gain)

- *Shape* and *color* are the most discriminating attributes (i.e. highest information gain)
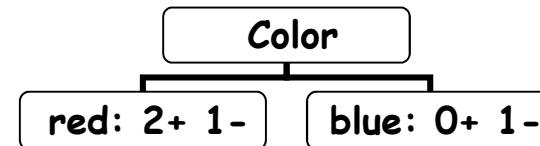
shape
color } more discriminating then size

# A Small Example (1)

*let's try F2 = color*

F2

| Size | Color | Shape | Output |
|------|-------|-------|--------|
| Big | Red | Circle | + |
| Small | Red | Circle | + |
| Small | Red | Square | – |
| Big | Blue | Circle | – |

Values(Color) = {red, blue}

Color

red: 2+ 1–    blue: 0+ 1–

$$gain(S, Color) = H(S) - \sum_{v \in values(Color)} \frac{|S_v|}{|S|} \times H(S_v)$$

*entropy of Output*

$$H(S) = -\left(\frac{2}{4}\log_2\frac{2}{4} + \frac{2}{4}\log_2\frac{2}{4}\right) = 1$$  // *total chaos, not knowing anything*

for each v of Values(Color)

• $H(S|Color = red) = H\left(\frac{2}{3}, \frac{1}{3}\right) = -\left(\frac{2}{3}\log_2\frac{2}{3} + \frac{1}{3}\log_2\frac{1}{3}\right) = 0.918$

• $H(S|Color = blue) = H(1,0) = -\left(\frac{1}{1}\log_2\frac{1}{1}\right) = 0$  // *sure thing*

*knowing the color*

$$H(S|Color) = \frac{3}{4}(0.918) + \frac{1}{4}(0) = 0.6885$$

Note: by definition,
- $Log_2 0 = -\infty$
- $0 log_2 0$ is 0

$$gain(Color) = H(S) - H(S|Color) = 1 - 0.6885 = 0.3115$$

*1 time out of 4 we have a blue*

*3 times out of 4 we have a red*

21

# A Small Example (2)

F3

| Size | Color | Shape | Output |
|------|-------|-------|--------|
| Big | Red | Circle | + |
| Small | Red | Circle | + |
| Small | Red | Square | - |
| Big | Blue | Circle | - |

Shape

circle: 2+ 1-    square: 0+ 1-

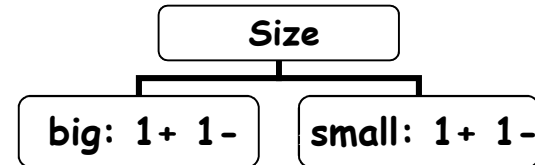$$H(S) = -\left(\frac{2}{4}\log_2\frac{2}{4} + \frac{2}{4}\log_2\frac{2}{4}\right) = 1$$

$$H(S\,|\,Shape) = \frac{3}{4}(0.918) + \frac{1}{4}(0) = 0.6885$$

$$gain(Shape) = H(S) - H(S\,|\,Shape) = 1 - 0.6885 = 0.3115$$

# A Small Example (3)

F1

| Size | Color | Shape | Output |
|------|-------|-------|--------|
| Big | Red | Circle | + |
| Small | Red | Circle | + |
| Small | Red | Square | - |
| Big | Blue | Circle | - |

Size

big: 1+ 1-    small: 1+ 1-

$$H(S) = -\left(\frac{2}{4}\log_2\frac{2}{4} + \frac{2}{4}\log_2\frac{2}{4}\right) = 1$$

$$H(S|Size) = \frac{2}{4}(1) + \frac{2}{4}(1)$$
$$gain(Size) = H(S) - H(S|Size) = 1-1=0$$

$H(Output)$ without knowing any feature

$H(Output/Size)$ knowing the size

# A Small Example (4)

| Size | Color | Shape | Output |
|------|-------|-------|--------|
| Big | Red | Circle | + |
| Small | Red | Circle | + |
| Small | Red | Square | - |
| Big | Blue | Circle | - |

$$gain(Shape) = 0.3115$$
$$gain(Color) = 0.3115$$
$$gain(Size) = 0$$

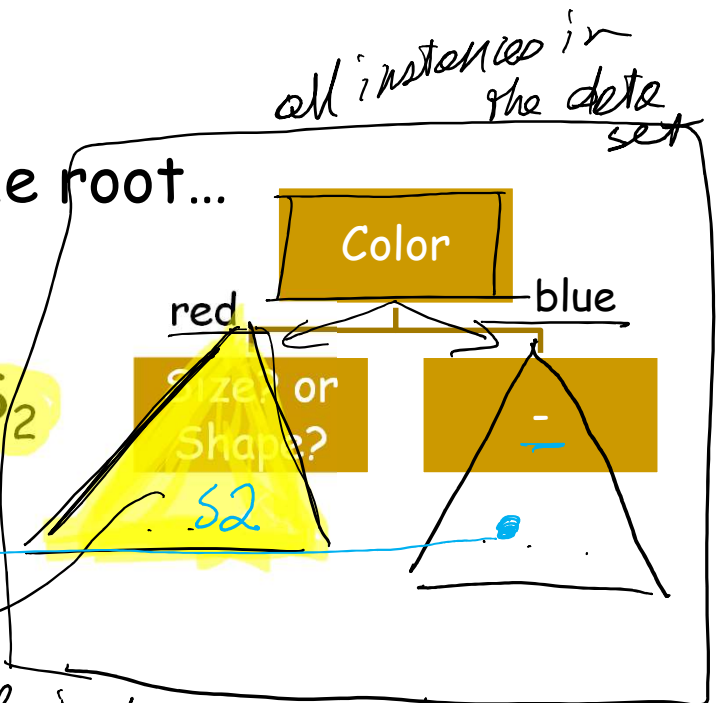- So first separate according to either _color_ or _shape_ (root of the tree)

# A Small Example (4)

■ Let's assume we pick *Color* for the root...

| Size | Color | Shape | Output |
|------|-------|-------|--------|
| Big | Red | Circle | + |
| Small | Red | Circle | + |
| Small | Red | Square | - |
| Big | Blue | Circle | - |

*[handwritten annotations: "S2" bracket next to first three rows, "S₂"]*

*[handwritten diagram: Color node branching to red → "Size or Shape?" and blue → "-", with "all instances in the data set", "only use the instances with color = red"]*

$$H(S_2) = -\left( \frac{2}{3}\log_2 \frac{2}{3} + \frac{1}{3}\log_2 \frac{1}{3} \right)$$

*[handwritten: "only use the instances with color = blue"]*

for each v of Values(Size)

$$H(S_2 \mid Size = big) = H\left( \frac{1}{1}, \frac{0}{1} \right) = 0$$

$$H(S_2 \mid Size = small) = H\left( \frac{1}{2}, \frac{1}{2} \right) = 1$$

$$H(S_2 \mid Size) = \frac{1}{3}(0) + \frac{2}{3}(1)$$

$$gain(Size) = H(S_2) - H(S_2 \mid Size)$$

for each v of Values(Shape)

$$H(S_2 \mid Shape = circle) = H\left( \frac{2}{2}, \frac{0}{2} \right) = 0$$

$$H(S_2 \mid Shape = square) = H\left( \frac{0}{1}, \frac{1}{1} \right) = 0$$

$$H(S_2 \mid Shape)$$

$$gain(Shape) = H(S_2) - H(S_2 \mid Shape)$$

# Back to the Restaurant

pick max info gain

- Training data:

| 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 |

source: Norvig (2003)

# The Restaurant Example

$$gain(alt) = ... \quad gain(bar) = ... \quad gain(fri) = ... \quad gain(hun) = ...$$

$$gain(pat) = 1 - \left( \frac{2}{12} \times H\left(\frac{0}{2}, \frac{2}{2}\right) + \frac{4}{12} \times H\left(\frac{0}{4}, \frac{4}{4}\right) + \frac{6}{12} \times H\left(\frac{2}{6}, \frac{4}{6}\right) \right)$$
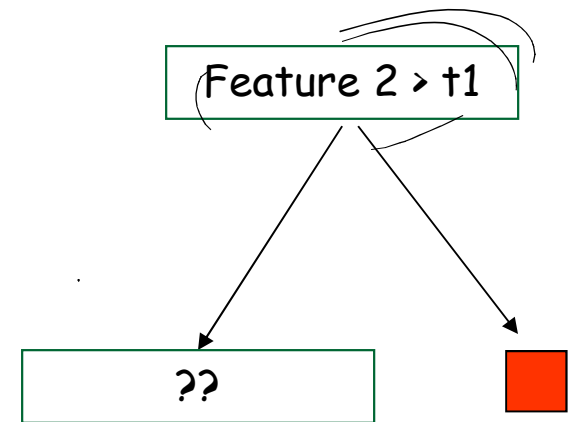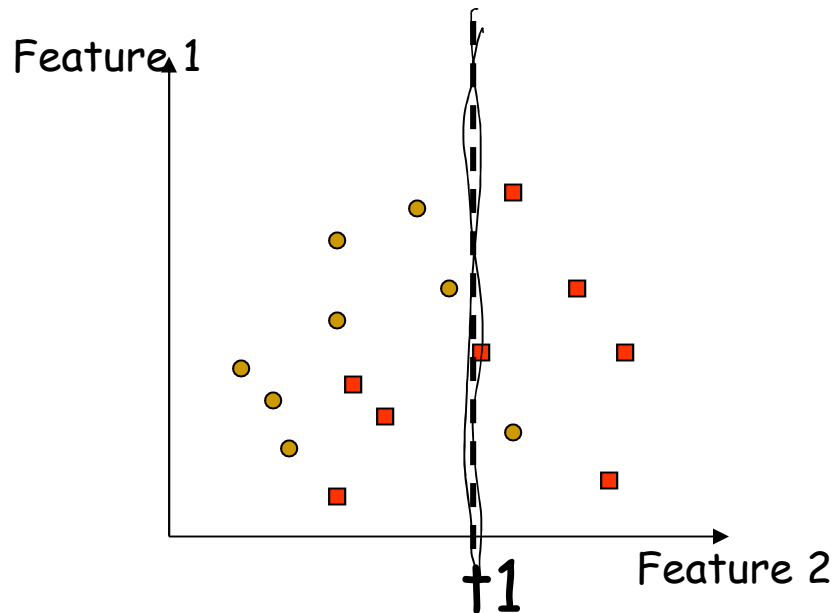
$$= 1 - \left( \frac{2}{12} \times -\left( \frac{0}{2}\log_2\frac{0}{2} + \frac{2}{2}\log_2\frac{2}{2} \right) + \frac{4}{12} \times -\left( \frac{0}{4}\log_2\frac{0}{4} + \frac{4}{4}\log_2\frac{4}{4} \right) + ... \right) \approx 0.541 \text{bits}$$

$$gain(price) = ... \quad gain(rain) = ... \quad gain(res) = ...$$

$$gain(type) = 1 - \left( \frac{2}{12} \times H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{12} \times H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{4}{12} \times H\left(\frac{2}{4}, \frac{2}{4}\right) + \frac{4}{12} \times H\left(\frac{2}{4}, \frac{2}{4}\right) \right) = 0 \text{ bits}$$

$$gain(est) = ...$$

*features*

- Attribute pat (Patron) has the highest gain, so root of the tree should be attribute *Patrons*
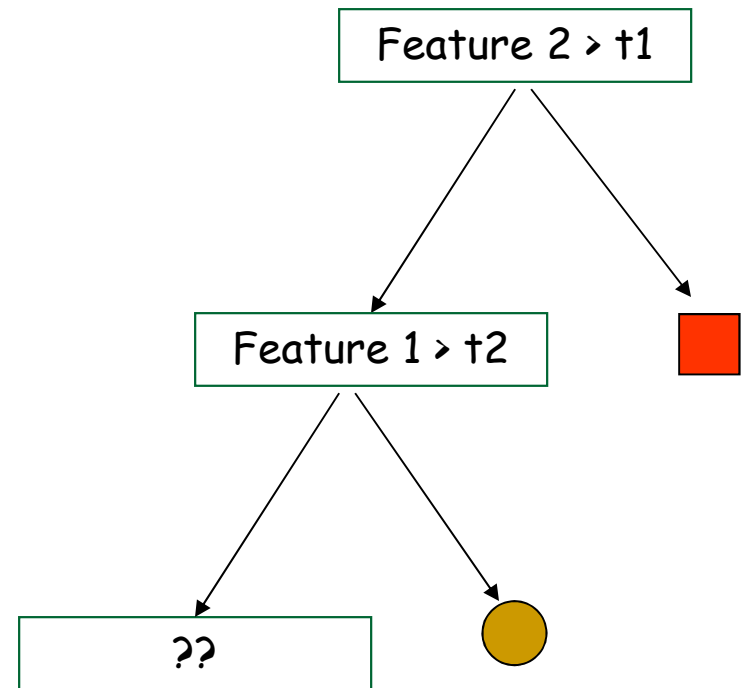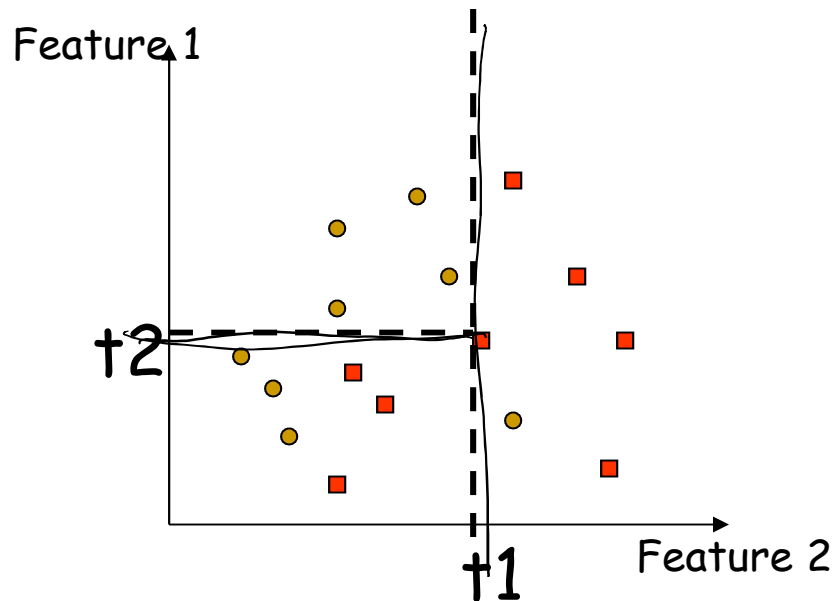- do recursively for subtrees
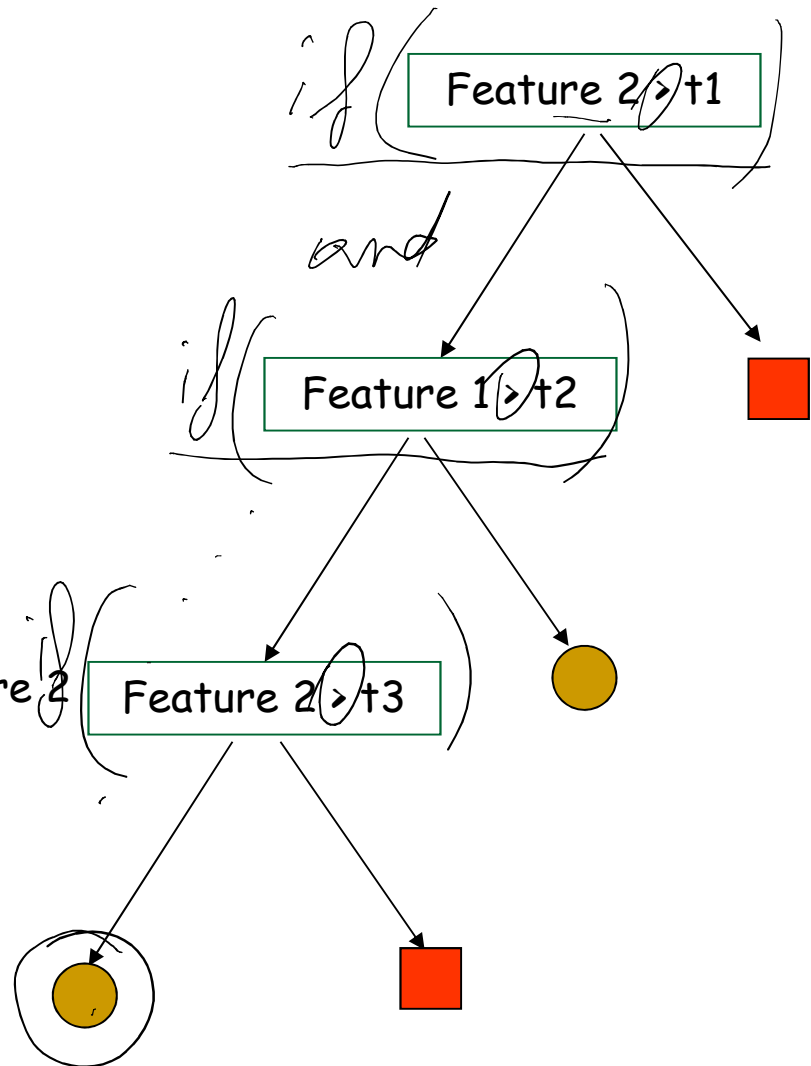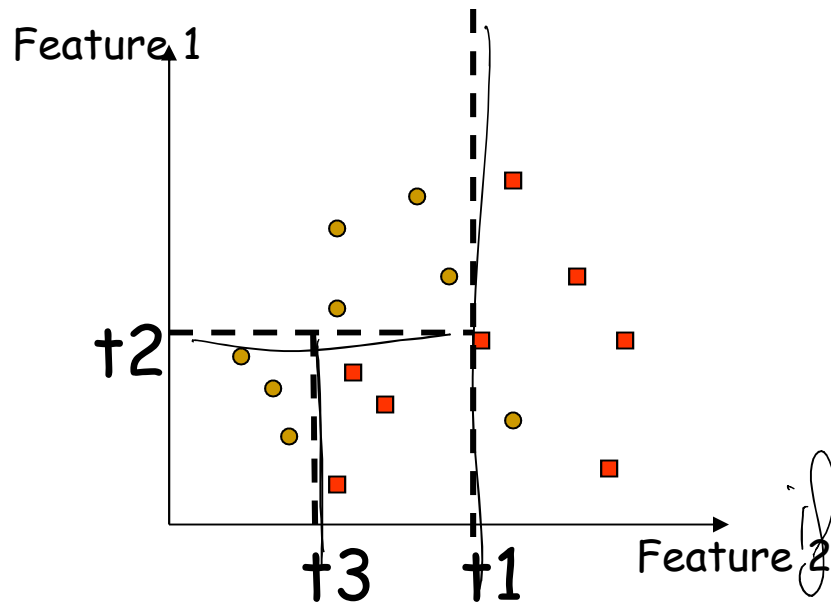
# Decision Boundaries of Decision Trees



Feature 1

class 2   class 1

Feature 2

# Decision Boundaries of Decision Trees



Feature 1

Feature 2

t1

Feature 2 > t1

??

# Decision Boundaries of Decision Trees

Feature 1

t2

t1

Feature 2

Feature 2 > t1

Feature 1 > t2

??

# Decision Boundaries of Decision Trees



Feature 1

t2

t3    t1    Feature 2

Feature 2 > t1

*if* ... *and*

Feature 1 > t2

*if* ...

Feature 2 > t3

*if* ...

# Applications of Decision Trees

- One of the most widely used learning methods in practice
  - Fast
  - Simple
  - Traceable (<-- very important!)

# Today

1. Introduction to ML ✓
2. Naïve Bayes Classification ✓
   a. Application to Spam Filtering ✓
3. Decision Trees ✓ *video #4*
4. ( Evaluation
5. Unsupervised Learning )
6. Neural Networks
   a. Perceptrons
   b. Multi Layered Neural Networks

# Up Next

Supervised Learning