

CS540 Introduction to Artificial Intelligence

Lecture 6

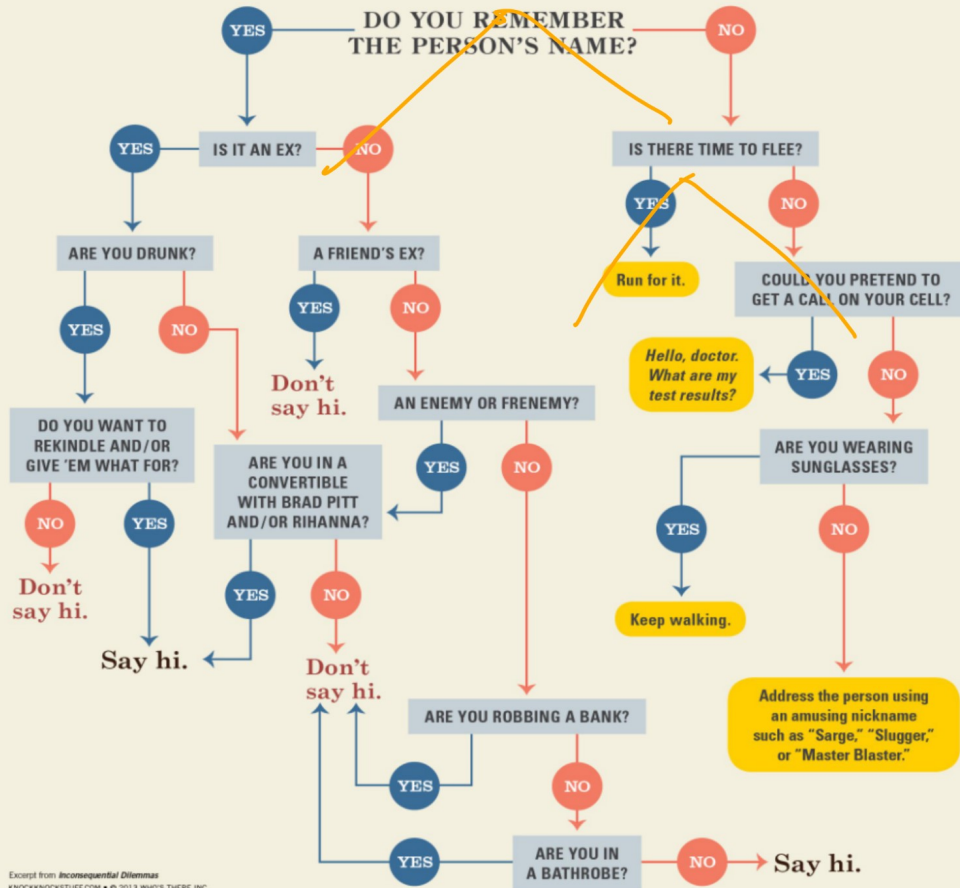
Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

June 10, 2020

I just saw someone I know.

DO I SAY HI?



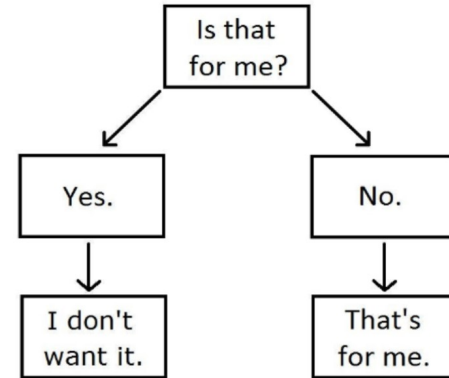
Excerpt from *Inconsequential Dilemmas*
KNOCKKNOCKSTUFF.COM • © 2013 WHO'S THERE INC.

n Forrest

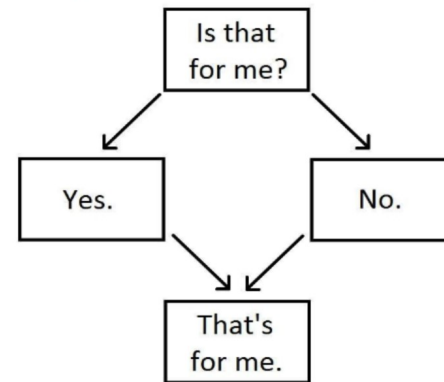
Nearest Neighbor
○○○○○○○

ion T...
vation

My Cat's Decision-Making Tree.

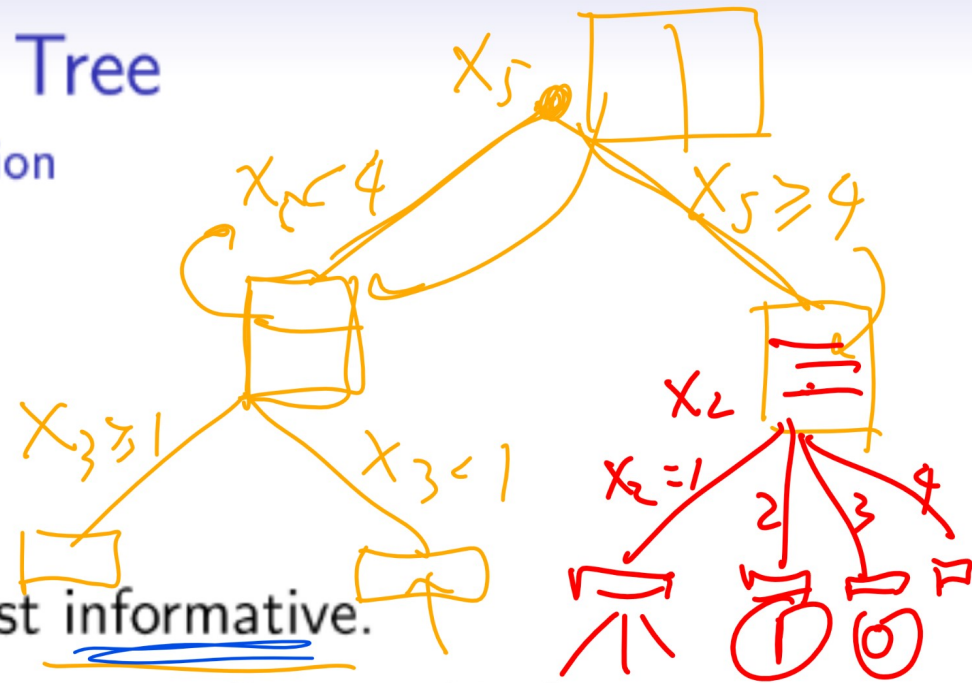
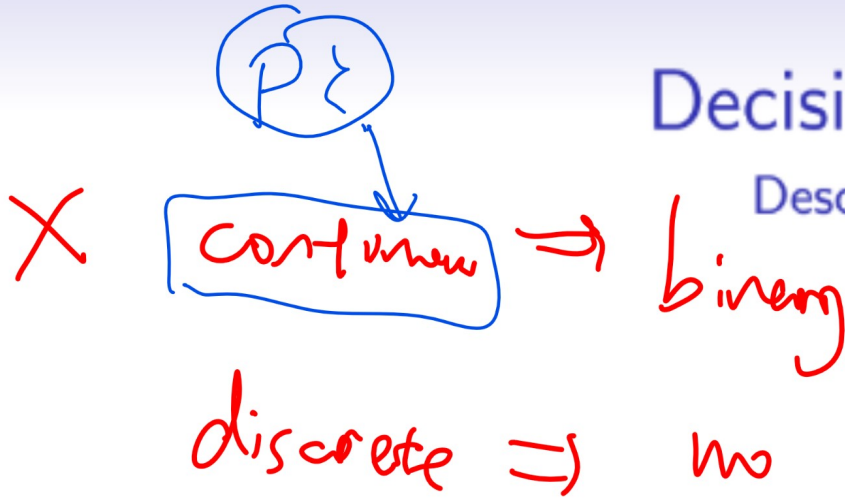


My Cat's Decision-Making Tree.



Decision Tree

Description



- Find the feature that is the most informative.
- Split the training set into subsets according to this feature.
- Repeat on the subsets until all the labels in the subset are the same.

Binary Entropy

Definition

- Entropy is the measure of uncertainty.
- The value of something uncertain is more informative than the value of something certain.
- For binary labels, $y_i \in \{0, 1\}$, suppose p_0 fraction of labels are 0 and $1 - p_0 = p_1$ fraction of the training set labels are 1, the entropy is:

$$H(Y) = p_0 \log_2 \left(\frac{1}{p_0} \right) + p_1 \log_2 \left(\frac{1}{p_1} \right)$$
$$= -p_0 \log_2 (p_0) - p_1 \log_2 (p_1)$$



Entropy

Definition

- If there are K classes and p_y fraction of the training set labels are in class y , with $y \in \{1, 2, \dots, K\}$, the entropy is:

$$\begin{aligned} H(Y) &= \sum_{y=1}^K p_y \log_2 \left(\frac{1}{p_y} \right) \\ &= - \sum_{y=1}^K p_y \log_2 (p_y) \end{aligned}$$

Conditional Entropy

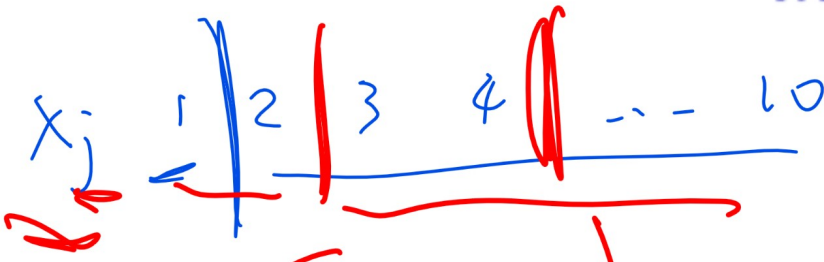
Definition

- Conditional entropy is the entropy of the conditional distribution. Let K_X be the possible values of a feature X and K_Y be the possible labels Y . Define p_x as the fraction of the instances that is x , and $p_{y|x}$ as the fraction of the labels that are y among the ones with instance x .

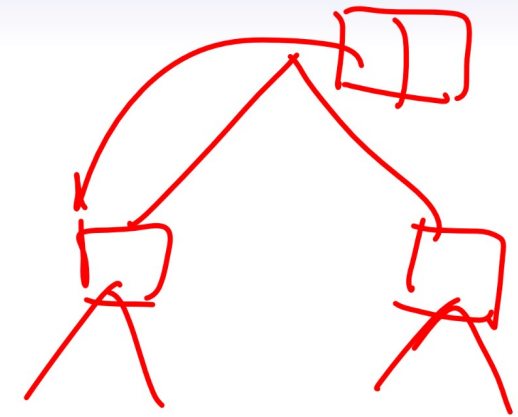
$$H(Y|X=x) = - \sum_{y=1}^{K_Y} p_{y|x} \log_2(p_{y|x})$$
$$H(Y|X) = \sum_{x=1}^{K_X} p_x H(Y|X=x)$$

Information Gain

Definition



for $ch_{j,c}$



- The information gain is defined as the difference between the entropy and the conditional entropy.

$$I(Y|X) = H(Y) - H(Y|X).$$

Annotations: "given" points to $H(Y)$; "value of X " points to $H(Y|X)$.

- The larger than information gain, the larger the reduction in uncertainty, and the better predictor the feature is.

Splitting Discrete Features

Definition

- The most informative feature is the one with the largest information gain.

$$\arg \max_j I(Y|X_j)$$

- Splitting means dividing the training set into K_{X_j} subsets.

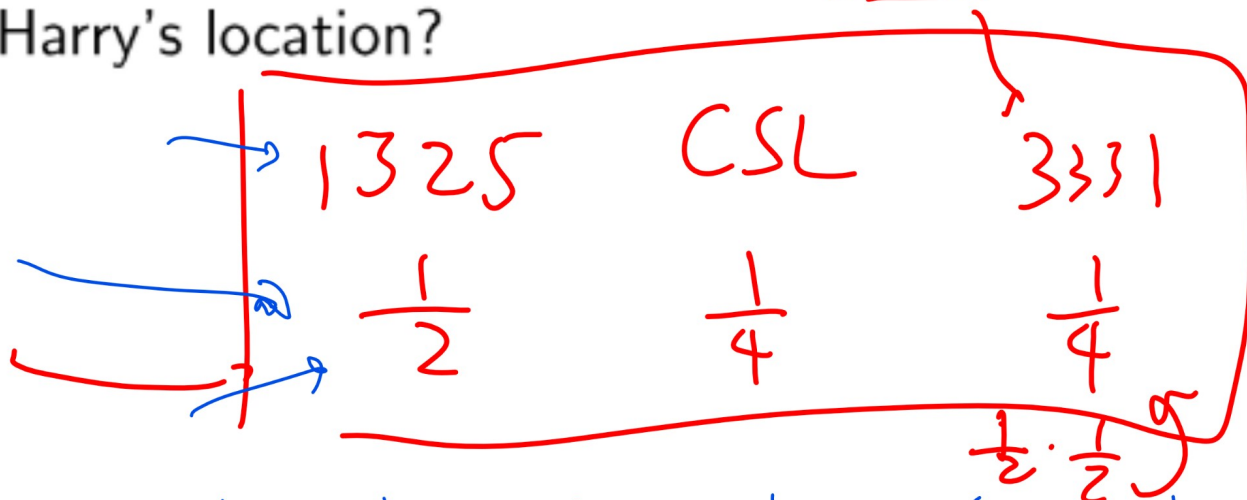
$$\{(x_i, y_i) : x_{ij} = 1\}, \{(x_i, y_i) : x_{ij} = 2\}, \dots, \{(x_i, y_i) : x_{ij} = K_{X_j}\}$$

Entropy

Quiz

- Fall 2010 Final Q10
- Running from You-Know-Who, Harry enters the CS building on the 1st floor. He flips a fair coin: if it is heads he hides in room 1325; otherwise, he climbs to the 2nd floor. In that case he flips the coin again: if it is heads he hides in CSL; otherwise, he climbs to the 3rd floor and hides in 3331. What is the entropy of Harry's location?

- A: 0.75
- B: 1
- C: 1.5
- D: 1.75
- E: None of the above.



$$\frac{1}{2} + \frac{1}{2} + \frac{1}{2} = 1.5$$

$$\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{4} \log_2 \frac{1}{4}$$

Entropy Math

Quiz

Entropy 2

Quiz

~~Q5~~ last.
Q6

- A bag contains a red ball, a green ball, a blue ball, and a black ball. Randomly draw a ball from the bag with equal probability. What is the entropy of the outcome?
- A: 1
- B: $\log_2(3)$
- C: 1.5
- **D: 2**
- E: 4

$$p_i \log_2 \frac{1}{p_i} + (1-p_i) \log_2 \frac{1}{1-p_i}$$

Handwritten notes:
- $p_i \log_2 \frac{1}{p_i}$
- $(1-p_i) \log_2 \frac{1}{1-p_i}$
- $p_i \log_2 p_i$
- $(1-p_i) \log_2 (1-p_i)$
- $p_i \log_2 p_i = 0$
- $(1-p_i) \log_2 (1-p_i) = 0$

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

Pruning Diagram

Disucssion

Baigging and Boosting Diagram

Discussion

K Nearest Neighbor

Description

- Given a new instance, find the K instances in the training set that are the closest.
- Predict the label of the new instance by the majority of the labels of the K instances.

Distance Function

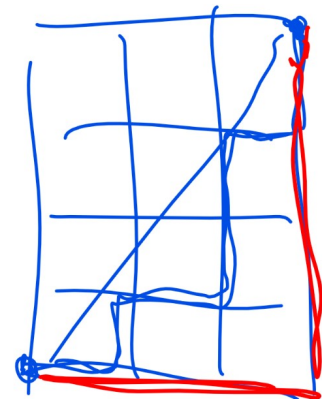
Definition

- Many distance functions can be used in place of the Euclidean distance.

$$\rho(x, x') = \|x - x'\|_2 = \sqrt{\sum_{j=1}^m (x_j - x'_j)^2}$$

- An example is Manhattan distance.

$$\rho(x, x') = \sum_{j=1}^m |x_j - x'_j|$$



1 Nearest Neighbor

Quiz

Q1

- Spring 2018 Midterm Q7
- Find the 1 Nearest Neighbor label for $\begin{bmatrix} 3 \\ 6 \end{bmatrix}$ using Manhattan distance.

x_1	1	1	3	5	2
x_2	1	7	3	4	5
y	0	1	1	0	0

- A: 0
- B: 1

dis: 7 3 3 4 2

3 Nearest Neighbor

Quiz

- Find the 3 Nearest Neighbor label for $\begin{bmatrix} 3 \\ 3 \end{bmatrix}$ using Manhattan distance.

new point

training set

x_1	1	1	3	5	2
x_2	1	7	3	4	5
y	0	1	1	0	0

dist = 4 6 | 0 3 3

- A: 0
- B: 1

Cross Validation

Quiz

- accuracy* *k-fold*
- Given the following training data. What is the 2 fold cross validation accuracy if 1 nearest neighbor classifier with Manhattan distance is used? The first fold is the first five data points.

x	1	1	2	2	3	3	4	4	5	5
y	1	2	3	3	2	2	3	3	2	1

fold 2 ↓ *fold 1*

↑ *2* *2* *2* *2* *2* *2* *2* *2* *2* *2*

2/5 *2/5*

label accuracy

ZCV accuracy = 4/10 = 40%

Cross Validation Diagram

Quiz

Cross Validation 2

Quiz

Q2

- Given the following training data. What is the 10 fold cross validation accuracy if 1 nearest neighbor classifier with Manhattan distance is used?

x	1	1	2	2	3	3	4	4	5	5
y	1	2	3	3	2	2	3	3	2	1

- A: 20 percent, B: 40, C: 60, D: 80, E: 100

2 1 3 3 2 2 3 3 1 2

6/10 correct.