

CS540 Introduction to Artificial Intelligence

Lecture 6

Young Wu

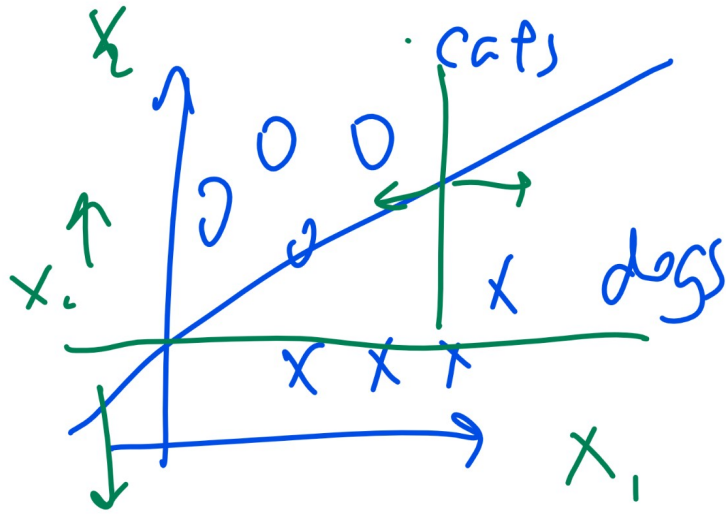
Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

June 6, 2022

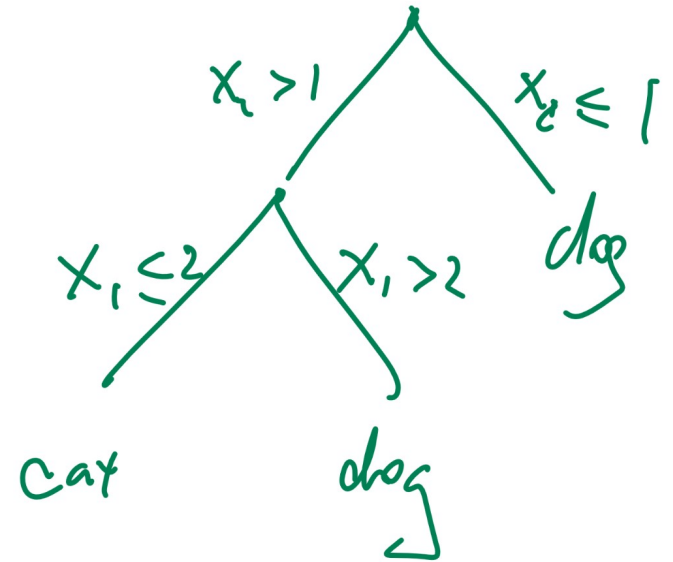
Axes Aligned Decision Boundary

Motivation

$$f(x^T w + b)$$



Decision 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

50% H
 50% T

Flipping

coin is not informative

100% H
 0% T

- 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:

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

Entropy

Quiz

- 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?

$$-\sum p_i \log_2 p_i$$

1325	CSL	3331
$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$

$$-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{4} \log_2 \frac{1}{4} = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = 1.5$$

Entropy 2

Quiz

- 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?

Q10

- A : 1
- B : $\log_2(3)$
- C : 1.5
- **D : 2**
- E : 4

$$- \sum P_i \log_2 P_i$$

~~$$\frac{1}{4} \log_2 \frac{1}{4}$$~~

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

back @ 7:55

⊕

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 are x , and $p_{y|x}$ as the fraction of the labels that are y among the ones with instance x .

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

$$\underline{H(Y|X)} = \sum_{x=1}^{K_X} p_x \underline{H(Y|X = x)}$$

$p_{y=1|x}$

$p_{y=2|x}$

Aside: Cross Entropy

Definition

$$y \log a + (1-y) \log(1-a)$$

- Cross entropy measures the difference between two distributions.

$$H(Y, X) = - \sum_{z=1}^K p_{Y=z} \log_2(p_{X=z})$$

- It is used in logistic regression to measure the difference between actual label Y_i and the predicted label A_i for instance i , and at the same time, to make the cost convex.

$$H(Y_i, A_i) = -y_i \log(a_i) - (1 - y_i) \log(1 - a_i)$$

Information Gain

Definition

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

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

informativeness of Y

info of Y given X.

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

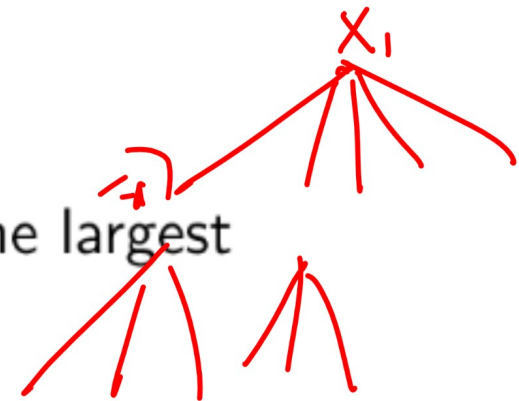
Splitting Discrete Features

Definition

$X_1 = 1, 2, 3, \dots, 5$

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

$$\operatorname{argmax}_j I(Y|X_j)$$



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

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

Splitting Continuous Variables Diagram

Definition

1, 3, 4, 5, 7



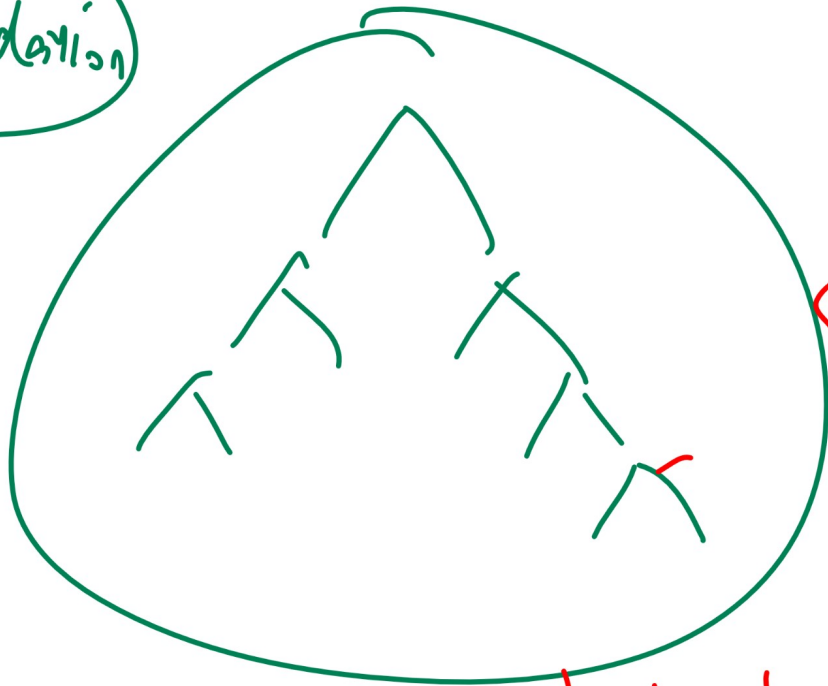
$\left. \begin{array}{l} \leq 2 \quad \geq 2 \\ \leq 3.75 \quad \geq 3.75 \\ \leq 5.75 \quad \geq 5.75 \end{array} \right\}$

find the one with largest IG
into gain

Pruning Diagram

Discussion

Validation



regularization
DT overfitting



which has higher validation accuracy?

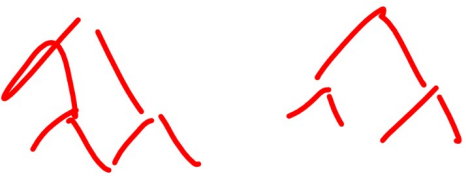
Bagging Diagram

Discussion

Bootstrap
Aggregating



⇒ NN



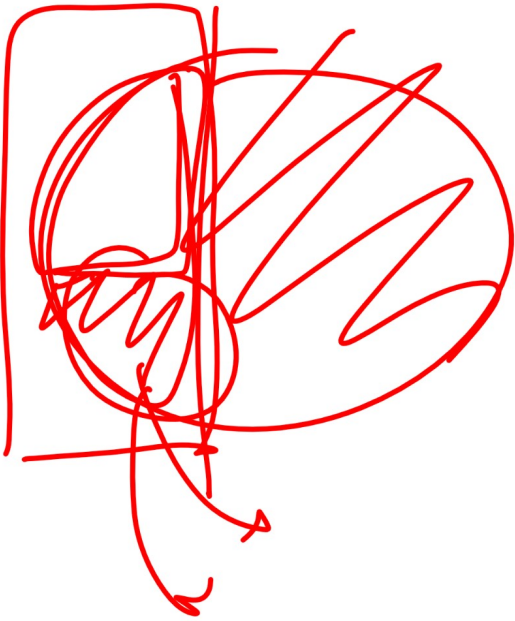
DT

random forest

{ Subset of data
subset of feature

Boosting Diagram

Discussion



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.

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

Leave One Out Cross Validation

Discussion

- If $K = n$, each time exactly one training instance is left out as the validation set. This special case is called Leave One Out Cross Validation (LOOCV).

Cross Validation

Quiz

- 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

Cross Validation 2

Quiz

- Given the following training data. What is the 10 fold cross-validation (LOOCV) 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

