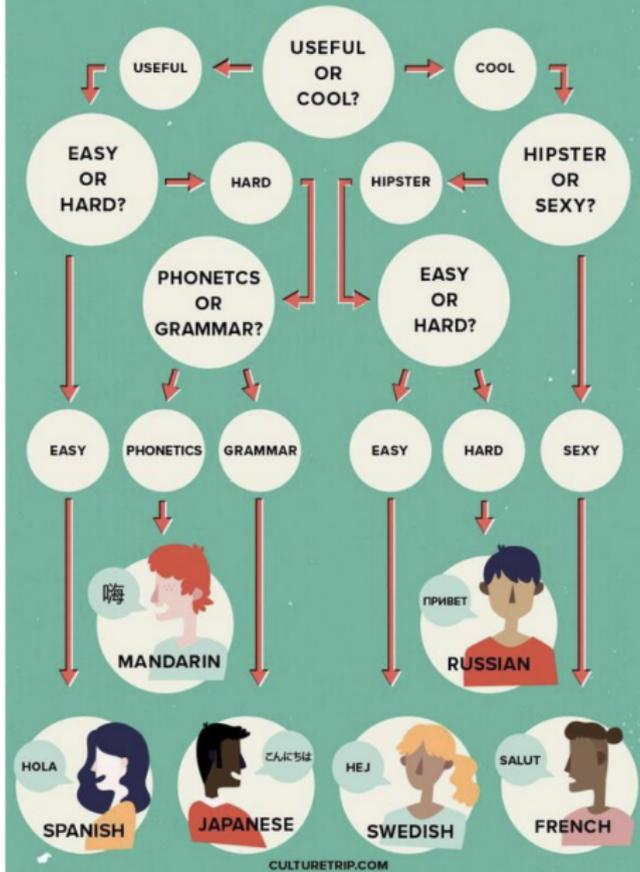


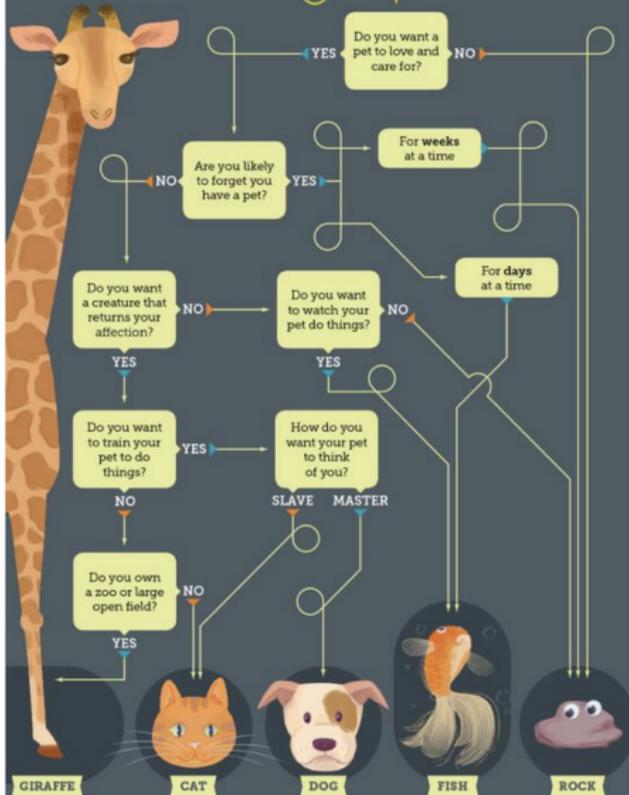
# Decision Trees

Sujay Sanghavi

# WHICH LANGUAGE SHOULD YOU LEARN IN 2018?



# What Kind of Pet is Right for You?



# Tree Methods

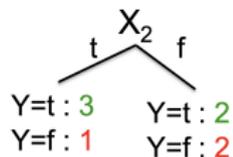
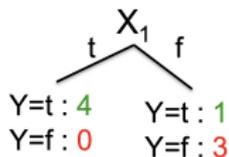
**Idea:** To get from  $x$  (the features) to  $y$  (the label / target to be predicted), do a **series of if-else questions** until you reach a decision (for classification) or prediction (for regression)

# Tree Methods

**Idea:** To get from  $x$  (the features) to  $y$  (the label / target to be predicted), do a **series of if-else questions** until you reach a decision (for classification) or prediction (for regression)

**Step 1:** Choosing which feature to split on.

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| T     | T     | T   |
| T     | F     | T   |
| T     | T     | T   |
| T     | F     | T   |
| F     | T     | T   |
| F     | F     | F   |
| F     | T     | F   |
| F     | F     | F   |

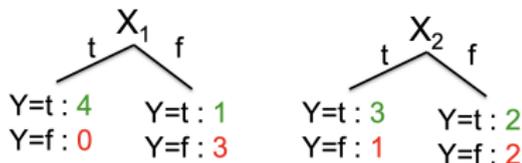


# Tree Methods

**Idea:** To get from  $x$  (the features) to  $y$  (the label / target to be predicted), do a **series of if-else questions** until you reach a decision (for classification) or prediction (for regression)

**Step 1:** Choosing which feature to split on.

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| T     | T     | T   |
| T     | F     | T   |
| T     | T     | T   |
| T     | F     | T   |
| F     | T     | T   |
| F     | F     | F   |
| F     | T     | F   |
| F     | F     | F   |



**Choose the feature that leaves us the most certain after splitting on it.**

# Measuring Uncertainty

Which of the following two distributions is more uncertain ?

|                |                |                |                |
|----------------|----------------|----------------|----------------|
| $P(Y=A) = 1/2$ | $P(Y=B) = 1/4$ | $P(Y=C) = 1/8$ | $P(Y=D) = 1/8$ |
|----------------|----------------|----------------|----------------|

|                |                |                |                |
|----------------|----------------|----------------|----------------|
| $P(Y=A) = 1/4$ | $P(Y=B) = 1/4$ | $P(Y=C) = 1/4$ | $P(Y=D) = 1/4$ |
|----------------|----------------|----------------|----------------|

# Tree Methods

**Basic idea:** First partition the training data samples. Then, to find prediction  $\hat{y}(x)$  for a query  $x$ :

- Find the partition that it goes into.
- Then

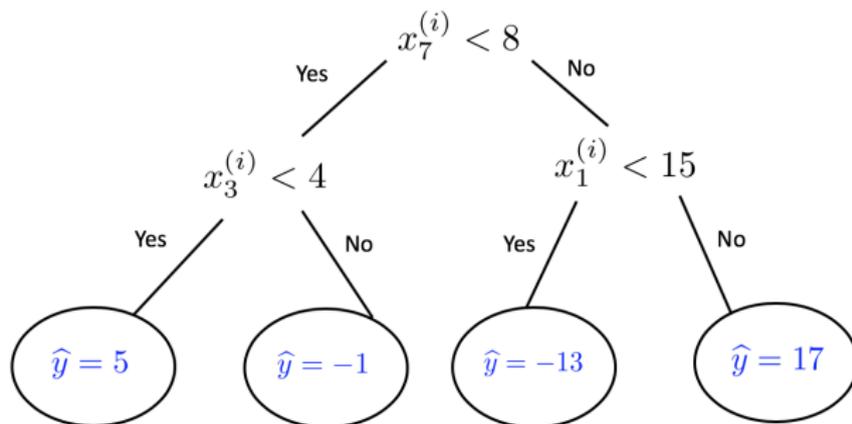
$$\begin{aligned}\hat{y}(x) &= \text{average of all training sample } y\text{'s in that partition} \\ &\quad (\text{for regression}) \\ &= \text{majority of all training sample } y\text{'s in that partition} \\ &\quad (\text{for classification})\end{aligned}$$

**Tree models** involve recursively partitioning the samples by

- Choose a feature  $x_j$  and a threshold  $\tau$ .
- For each sample  $x^{(i)}$ , put it in the “left bucket” if  $x_j^{(i)} < \tau$  and “right bucket” otherwise

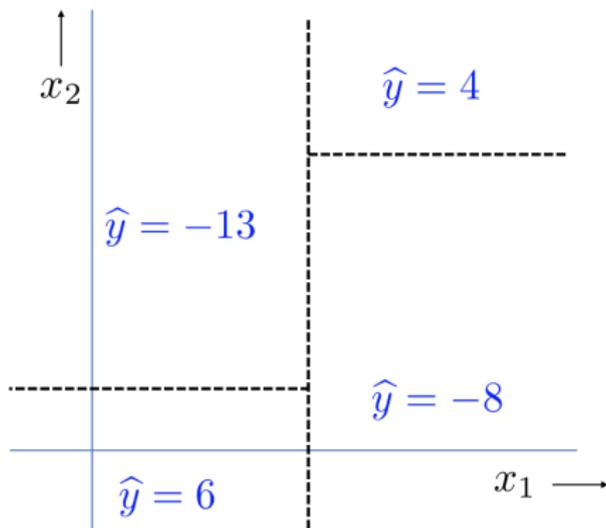
# Decision / Regression Trees

Decision rule represented by a tree



# Decision / Regression Trees

Decision regions can be visualized as a series of “stumps”



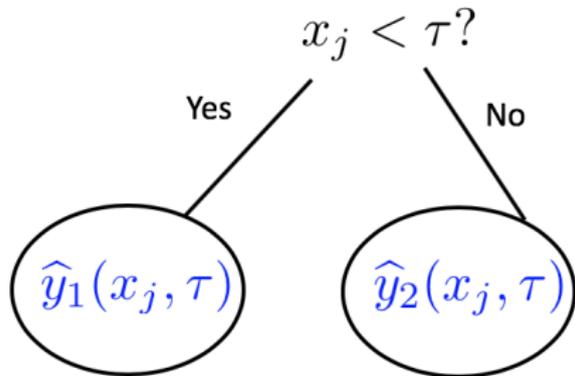
# Decision / Regression Trees

**Q: How to choose the feature and thresholds for each split ?**

**A1: Greedy:** at each step, choose the feature and threshold that gives the lowest final squared error.

If feature  $x_j$  and threshold  $\tau$ , final squared error:

$$\sum_{i \in S_1(x_j, \tau)} \left( y^{(i)} - \hat{y}_1(x_j, \tau) \right)^2 + \sum_{i \in S_2(x_j, \tau)} \left( y^{(i)} - \hat{y}_2(x_j, \tau) \right)^2$$



However, for a large number of features, this is very cumbersome.

Also: bigger depth == lower training error, but more likely to overfit ...