**Q** Estimate in-game win probability for American Football

↳ as a function of game-state

**\* Why?**
— Betting
— player valuation
— Strategic decision making:
make the decision which maximizes win probability

**\* Mathematical Models:**
—dynamic programming state-space models popular 20+ years ago

**\* Statistical Models:**
— approach: in "similar" situations across football history, what proportion of the time did the team with possession win the game?

— play-by-play NFL data is easily accessible today machine learning models can be fit quickly today Hence today, everyone does this
— we will focus on statistical models

\* <u>outcome var.</u> whether the team with possession on this play won or lost the game

<u>variables</u>
yardline
down
distance
score differential
game seconds remaining
off. and def. team quality → use point spread
Receive 2<sup>nd</sup> half kickoff
timeouts

\* Do these variables have a linear or additive relationship?

No → Nonlinear and Interacting variables

→ | Machine Learning |

Learn an arbitrarily complex relationship b/t the variables from the data.

\* Going to foug an Strategic decision making

Q  What fourth down decision should I make
   in {Punt, Go, FG} given the game-state?

\* The entire fourth down decision making
  process can be defined in terms of
  the win probability if you have a
  first down with ten yards to go.
  Ignoring down & yards to go will simplify our
  modeling process.

## Fourth down decision model

$$V_1(x) = \text{Value (win probability) of a } 1^{st} \text{ and } 10$$
$$\text{at game-state } x$$
$$(\text{or, for simplicity, yardline } y)$$

\* Value of punting:

$$V_{punt} = -\sum_{y'} V_1(\text{yardline } y') \cdot P(\text{yardline after punting is } y')$$

$$= -\mathbb{E}_{punt}\left[V_1(\text{yardline } y')\right]$$

$$\approx -V_1(\text{yardline } \mathbb{E}_{punt}[y']) \longrightarrow \text{linear regression function of yardline (spline)}$$

**\* Value of kicking a FG:**

$$V_{FG} = \mathbb{P}\binom{make}{FG} \cdot V\binom{make}{FG} + \left(1 - \mathbb{P}\binom{make}{FG}\right) \cdot \left(-V_1\binom{yardline}{100-y}\right)$$

<span style="color:blue">logistic Regression function of kicker quality and yardline</span>

$$= -\mathbb{E}_{kickoff}\left[V_1(yardline)\right]$$

$$\approx -V_1(yardline\ 75)$$

**\* Value of going for it:**

4th down and $z$ yards to go at yardline $y$

$$V_{GO} = \sum_{\Delta \geq z} V_1(yardline\ y-\Delta) \cdot \mathbb{P}(gain\ \Delta\ yards)$$

$$- \sum_{\Delta < z} V_1(yardline\ 100-(y-\Delta)) \cdot \mathbb{P}(gain\ \Delta\ yards)$$

$$\approx \mathbb{P}(convert) \cdot V_1(yardline\ y-z) + (1-\mathbb{P}(convert)) \cdot V_1\binom{yardline}{100-y}$$

<span style="color:blue">logistic Regression function of team qualities, yardline, and yards to go</span>

\* So, to obtain $V_{GO}, V_{FG}, V_{punt}$ we need

$V_1(x)$ = value (W.P.) of a $1^{st}$ and $10$ at game-state $x$

and $\underbrace{\mathbb{E}_{punt}[y'], \mathbb{P}\binom{make}{FG}, \mathbb{P}(convert)}$

<span style="color:blue">HW:</span> estimate these 3 functions of game-state

## <u>Task</u> Estimate $V_1(x) = WP(x \mid 1^{st}$ down and $10)$.
using Machine Learning.

## <u>Game-state X</u> yardline, score differential, timeouts
game seconds Remaining, Point spread, Receive $2^{nd}$ half kickoff

## <u>Model Setup</u> $i$ = index of $i^{th}$ play in dataset of NFL plays
$y_i = 1$ if team with possession on play $i$ wins, else $0$
$x_i$ = game-state vector of play $i$
$y_i = Logistic(f(x_i)) + \varepsilon_i$
Obtain the best possible $\hat{f}(x_i)$ ( best predictive performance )
( (logloss) )

\* We will focus on <u>tree machine learning</u> models,
which are easy to use in R, fairly quick to fit,
and are widely used for NFL win probabilities.
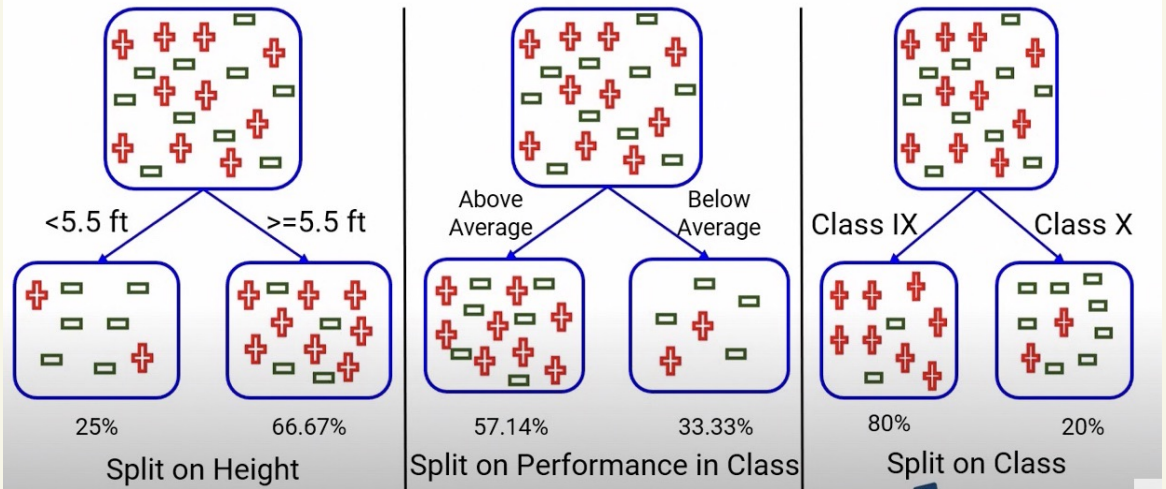
# Decision Trees
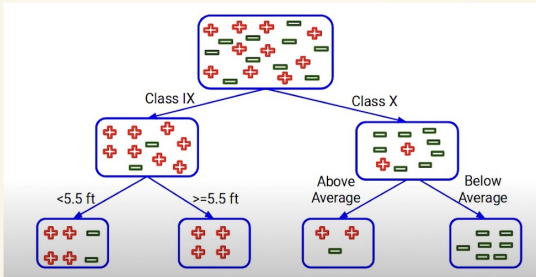
Ex 20 students, 10 play cricket
  Variables: Height, grades, class
  Fit a model to predict whether a student plays cricket.

Idea split on an X variable to classify the y variable



|  | Split on Height | Split on Performance in Class | Split on Class |
|---|---|---|---|

Which split is best? The Rightmost one.
Want to separate the classes as best as possible.



Can have multiple
splits in a decision tree!

✱ decision trees allow variables to interact
✱ but how to fit such a tree from data?

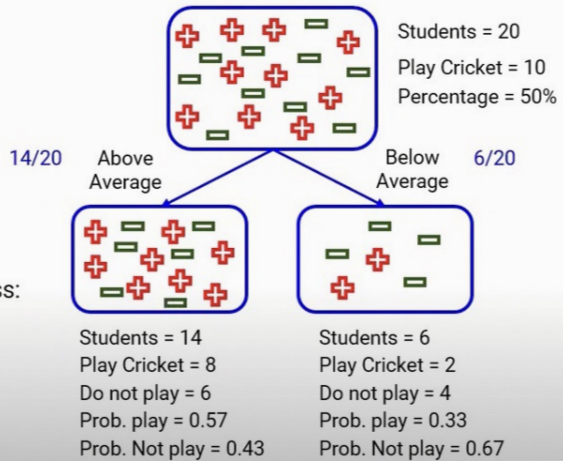# ✳ How to select the best split point?

→ select the split which results in the most homogeneous
                                                        sub-nodes

## 1. Grini Impurity



**Split on Performance in Class**

- Gini Impurity: sub-node Above Average:
  1 - [(0.57)*(0.57) + (0.43)*(0.43)] = 0.49

- Gini Impurity: sub-node Below Average:
  1 - [(0.33)*(0.33) + (0.67)*(0.67)] = 0.44

- Weighted Gini Impurity: Performance in Class:
  (14/20)*0.49 + (6/20)*0.44 = 0.475

Students = 20
Play Cricket = 10
Percentage = 50%

14/20    Above          Below      6/20
         Average        Average

Students = 14          Students = 6
Play Cricket = 8       Play Cricket = 2
Do not play = 6        Do not play = 4
Prob. play = 0.57      Prob. play = 0.33
Prob. Not play = 0.43  Prob. Not play = 0.67

Gini = sum of square probabilities for
         each outcome category, within a node

$$= P_+^2 + P_-^2$$

Gini Impurity of a split

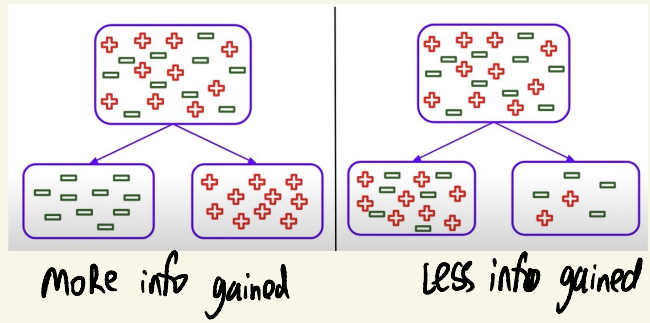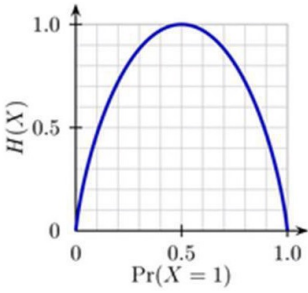   = weighted Gini Impurity of both sub-nodes of that split

   $$= W_L (1 - Gini_L) + W_R (1 - Gini_R)$$

| Split | Weighted Gini Impurity |
|---|---|
| Performance in Class | 0.475 |
| Class | 0.32 |

→ split on Class first!

# 2. Information Gain



More info gained        Less info gained

Information Gain = 1 - Entropy

Entropy = $-P\log_2 P - (1-P)\log_2(1-P)$

where
$P = \mathbb{P}(+) = \mathbb{P}(y=1)$



% Play = 0.50
% Not play = 0.50
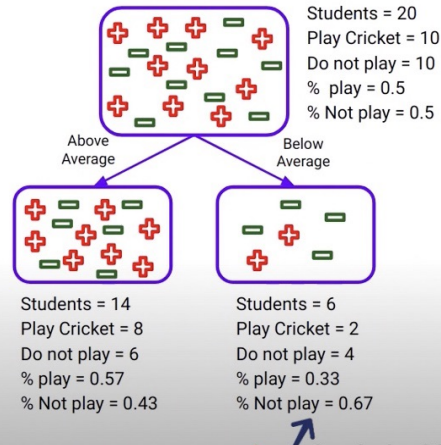
Entropy = - (0.5) * log2(0.5) - (0.5) * log2(0.5)

= 1

% Play = 0
% Not play = 1

Entropy = - (0) * log2(0) - (1) * log2(1)

= 0

Entropy of a split = weighted average of entropy of each Resulting sub-node

### Split on Performance in Class

- Entropy for Parent node:
  $-(0.5)*\log_2(0.5) -(0.5)*\log_2(0.5) = 1$

- Entropy for sub-node Above Average:
  $-(0.57)*\log_2(0.57) -(0.43)*\log_2(0.43) = 0.98$

- Entropy for sub-node Below Average:
  $-(0.33)*\log_2(0.33) -(0.67)*\log_2(0.67) = 0.91$

- Weighted Entropy: Performance in Class:
  $(14/20)*0.98 + (6/20)*0.91 = 0.959$

Students = 20
Play Cricket = 10
Do not play = 10
% play = 0.5
% Not play = 0.5

Above Average        Below Average

Students = 14
Play Cricket = 8
Do not play = 6
% play = 0.57
% Not play = 0.43

Students = 6
Play Cricket = 2
Do not play = 4
% play = 0.33
% Not play = 0.67

| Split | Entropy | Information Gain |
|---|---|---|
| Performance in Class | 0.959 | 0.041 |
| Class | 0.722 | 0.278 |

→ split on Class first!

**\* Our outcome variable** (cricket vs. no cricket) **is categorical,** (win vs. loss)
**so Gini Impurity and Information gain work well.**
**If outcome were continuous (e.g. height or money),**
**need another way to measure splits:** **3. Reduction in Variance.**
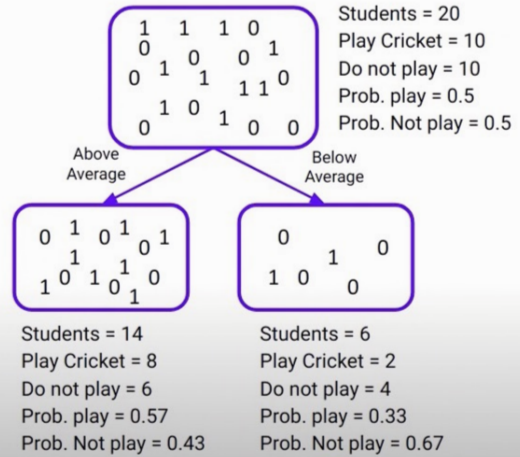
Variance = $\Sigma [(X - \mu)^2] / n$

| 2 | 6 | 7 |
|---|---|---|
| 4 | 7 | 9 |

Variance ~ 6

| 1 | 1 | 1 |
|---|---|---|
| 1 | 1 | 1 |

Variance = 0
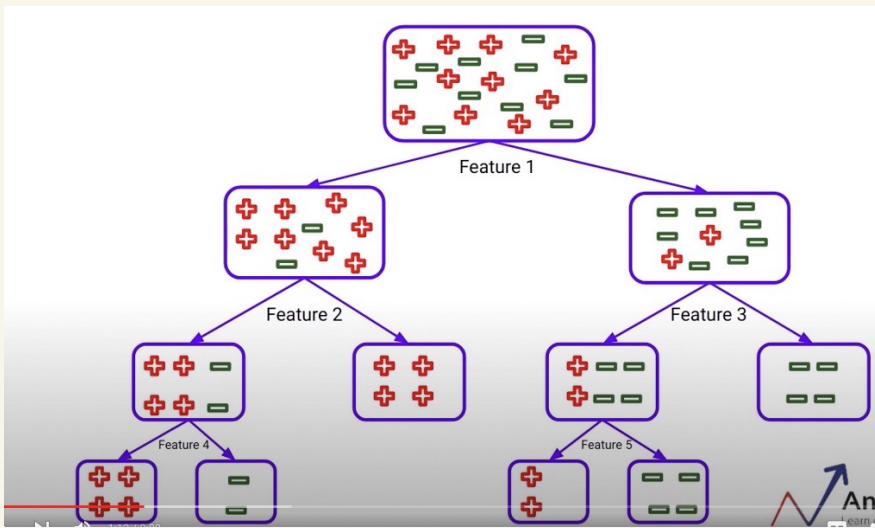
**Split with lower variance is selected.**

- Above Average node:
  - Mean = (8*1 + 6*0) / 14 = 0.57
  - Variance =
    $[8*(1-0.57)^2 + 6*(0-0.57)^2] / 14 = 0.245$

- Below Average node:
  - Mean = (2*1 + 4*0) / 6 = 0.33
  - Variance =
    $[2*(1-0.33)^2 + 4*(0-0.33)^2] / 6 = 0.222$

- Variance: Performance in Class:
  (14/20)*0.245 + (6/20)*0.222 = 0.238

Students = 20
Play Cricket = 10
Do not play = 10
Prob. play = 0.5
Prob. Not play = 0.5

| 1 | 1 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | | 1 | 1 | 1 | 0 |
| | 1 | 0 | 1 | | 0 |
| 0 | | | 0 | 0 |

Above Average → Below Average

| 0 | 1 | 0 | 1 | 0 | 1 |
| | 1 | | 1 | |
| 1 | 0 | 1 | 0 | | 0 |
| | | | 1 | |

Students = 14
Play Cricket = 8
Do not play = 6
Prob. play = 0.57
Prob. Not play = 0.43

| 0 | | | | 0 |
| | | 1 | |
| 1 | 0 | | 0 | |

Students = 6
Play Cricket = 2
Do not play = 4
Prob. play = 0.33
Prob. Not play = 0.67

| Split | Variance |
|---|---|
| Performance in Class | 0.238 |
| Class | 0.16 |

**→ split on Class first!**

**\* Now, we know how to select the best split point!**

Feature 1
Feature 2
Feature 3
Feature 4
Feature 5

Ana

* Could *iteratively make splits* until all nodes are "pure" but this would overfit (memorize noise in the training data)

* Tuning parameters (to control overfitting & underfitting)

    max depth of tree
    min. samples for node split
    min. samples for a terminal node
    max number of terminal nodes

* Straightforward to fit a decision tree in R

* HW: fit a decision tree NFL Win Prob. model. What goes wrong?

* Decision trees are unstable & prone to overfitting.

* How to Reduce overfitting? Ask Leo Breiman...