

Decision Trees

Blackbox Machine Learning

Q Estimate in-game win probability for American Football

↳ as a function of game-state

* Why?

- Betting
- Player valuation
- Strategic decision making

$$\left\{ \begin{array}{l} P(Td) = 0.2 \\ P(Tg) = 0.2 \\ P(Turnover) = 0.1 \\ P(Punt) = 0.5 \end{array} \right.$$

make the decision which maximizes our win probability.

* Mathematical Models:

— state-space models & dynamic programming

- States: $\left\{ \begin{array}{l} \# \text{ Possession} \\ \text{Score of Team A} \\ \text{Score of Team B} \end{array} \right\}$
- transition probabilities

$$P(2, 7, 0 \mid 3, 0, 0) = P(Td) = 0.2$$

* Statistical Models:

- learn entirely from historical data
- across "similar" situations over the entire (recent) history of football, what proportion of times did the team with possession win the game?

— play by play NFL data is easily accessible today (since 2017)

NFLFastR

— easy to fit your favorite ML model today

- thy today everyone uses Statistical models for WP

* We want to fit a ML WP model from historical data.
↳ as a function of game-state

variables

outcome variable — whether the team with possession won or lost the game
1/0

game state vars {
yardline
down
distance
score differential
game seconds remaining
off. and def. team quality
Receive 2nd half Kickoff
Timeouts

$$\beta_{11} \cdot ydl + \beta_{12} \cdot ydl^2 + \beta_{13} \cdot ydl^3 + \dots$$

→ point spread

categorical

$$\beta_{21} \text{ down } 1 + \beta_{22} \text{ down } 2 + \beta_{23} \text{ down } 3 + \beta_{2n} \text{ down } n 4$$

* Last week: Logistic Regression

$$P(\text{Win} = 1 | \text{game-state } x) = \frac{1}{1 + \exp(-(\beta_1 \cdot ydl + \beta_2 \cdot \text{down} + \beta_3 \cdot \text{dist} + \beta_4 \cdot \text{time rem.} + \beta_5 \cdot \text{game sel. rem.} + \dots))}$$

What's wrong with this?

* These variables are all interacting but regression models are fundamentally additive (non-interacting)

ex Scorediff and time remaining are def. interacting

$$\beta_3 \cdot \text{sign} \left(\frac{\text{score diff}}{\text{time rem}} \right) \cdot e^{-(\text{score diff}) \cdot (\text{time rem})}$$

$$\beta_1 \cdot \text{score diff} + \beta_2 \cdot \text{time rem} + \beta_3 \cdot (\text{score diff}) \cdot (\text{time rem})$$

* YOU can model interactions in additive regression settings but this is extremely difficult with 2 numeric/continuous vars and it is easy with indicator / binary vars

linear: $\beta \cdot X$

Stame $\beta_1 \cdot D_1 + \beta_2 \cdot D_2 + \beta_3 \cdot D_1 \cdot D_2$

$D_1 = \text{indicator}$

$$P(\text{out}) = \frac{1}{1 + \exp(-(\quad))}$$

All of these variables are interacting and nonlinear.

Logistic regression is a ~~terrible~~ bad idea.

→ Machine Learning!

Learn an arbitrarily complex relationship between variables, given enough data.

* Focus on using WP for strategic decision making.

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

* The entire 4th down decision process can be defined in terms of the Win Probability if you have a 1st down and 10 yards to go

def $V_1(x) =$ value (win probability)
of a 1st and 10
at game-state x

* value of kicking a FG on 4th down:

$$V_{FG}(x) = P(\text{make FG}) \cdot V(\text{make FG}) + P(\text{miss FG}) \cdot V(\text{miss FG})$$

$$1 - V_1 \left(\begin{array}{l} x \\ \text{except} \\ 75 \text{ yardline} \\ \text{score diff} + 3 \end{array} \right)$$

Logistic Regression

outcome 1/0 made/miss FG

yardline, Kicker quality,

weather, time rem.

HW

$$1 - V_1 \left(\begin{array}{l} x \\ \text{except} \\ \text{flip} \\ \text{ydl} \end{array} \right)$$

$$1 - P(\text{make FG})$$

$$P(\text{make FG} \mid ydl, Kq)$$

$Kq =$ Kicks made over expected

$$= \sum_{\text{his prev kicks } i} \left[\begin{cases} 1 & \text{if kick } i \text{ made} \\ 0 & \text{if not} \end{cases} - \underbrace{P(\text{kick } i \text{ made} \mid ydl)} \right]$$

$$P(\text{make FG} \mid ydl)$$

punter qual = punt yardlines added over expected

$$= \sum_{\text{his prev punts } i} \left[\begin{pmatrix} ydl & y'_i \\ \text{after} & \text{punt } i \end{pmatrix} - \underbrace{E \left(\begin{pmatrix} ydl & y'_i \\ \text{after} & \text{punt } i \end{pmatrix} \mid ydl \right)} \right]$$

ignore punter quality

* Value of Punting on 4th down:

$$V_{\text{punt}}(x) = \sum_{y'} P(\text{ydl after punting is } y') \cdot \left(1 - V_1 \left(\begin{matrix} x \\ \text{ydl } y' \end{matrix} \right) \right)$$

$$= 1 - \mathbb{E}_{\text{punt } y'} V_1 \left(\begin{matrix} x \\ \text{ydl } y' \end{matrix} \right)$$

$$\approx 1 - V_1 \left(\begin{matrix} x \\ \text{ydl, } \underbrace{\mathbb{E}_{\text{punt}}(y')} \end{matrix} \right)$$

- bin by ydl and ^{expected} take avg ^{next} yardline after punting
- linear regression ^{next} on ydl

$\mathbb{E} f(Y) \neq f(\mathbb{E} Y)$ in general on ydl + punter quality

but in this specific case = is close enough because f , $y \mapsto V_1(y)$ is nearly Linear and so $\mathbb{E} f(Y) = f(\mathbb{E} Y)$

↓
HW

* Value of going for it on 4th down:

4th down and z yards to go
 at yardline y
 game-state x

$$V_{go}(x) \approx$$

$$P(\text{convert}) \cdot V_1 \left(\begin{matrix} x \\ \text{except} \\ \text{new ydl} \end{matrix} \right) + (1 - P(\text{convert})) \cdot (1 - V_1 \left(\begin{matrix} x \\ \text{except} \\ \text{opp. ydl} \end{matrix} \right))$$

\downarrow $y - z$
 ydl. $y = y$ yards from td

\downarrow $100 - y$

$P(\text{gain} \geq z \text{ yards})$

off. def, team quality
 dist
 yardline
 outcome var: 1/0, 1 = conversion

HW

Logistic Regression

* V_{go} , V_{fg} , V_{punt} on 4th down
as a function of game-state x
can be written in terms of

$V_1(x) = WP$ if a team has a
1st and 10 at x
Machine Learning
and $P(\text{make}_{fg}|x)$, $P(\text{convert}|x)$, $E_{punt}[y'|x]$
Regression (HW)

Task Estimate $V_1(x) = WP(x | 1^{\text{st}} \text{ and } 10)$
Using Machine Learning.

Game-state x yardline, score diff, timeouts,
game sec. rem., point spread, Receive 2nd half kicks

Dataset

i = index of 1st and 10 play i in dataset of plays

x_i = game-state vector of play i

y_i = 1 if team with possession on play i wins the game else 0

Want

$$P(y_i=1|x_i) = E[y_i|x_i] = \hat{y}_i = \hat{f}(x_i)$$

\hat{f} could be an arbitrarily complex nonlinear interacting function



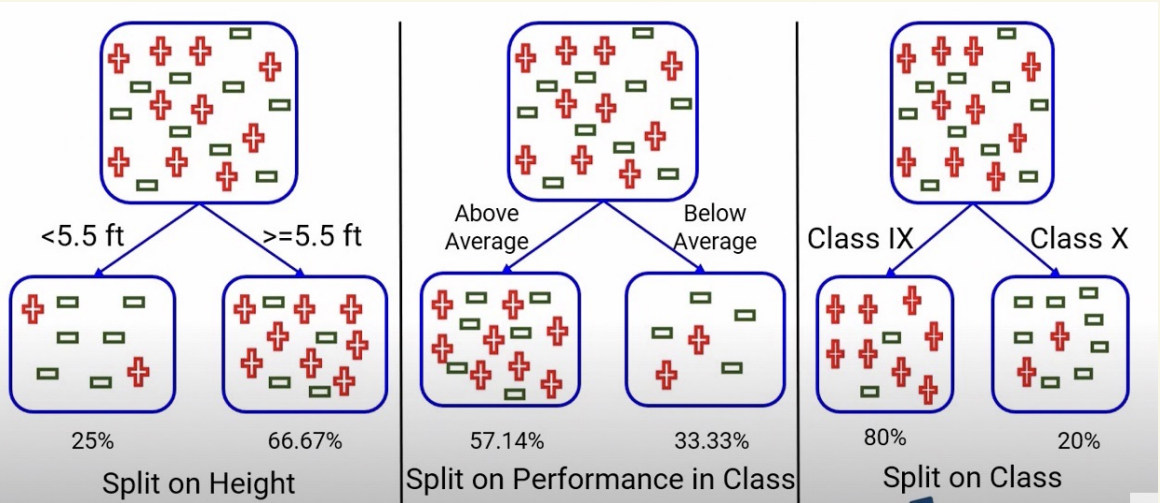
Tree Machine Learning

- decision trees / CART
classification & regression trees
- Random forests
- XGBoost (extreme gradient boosting)

Decision Trees

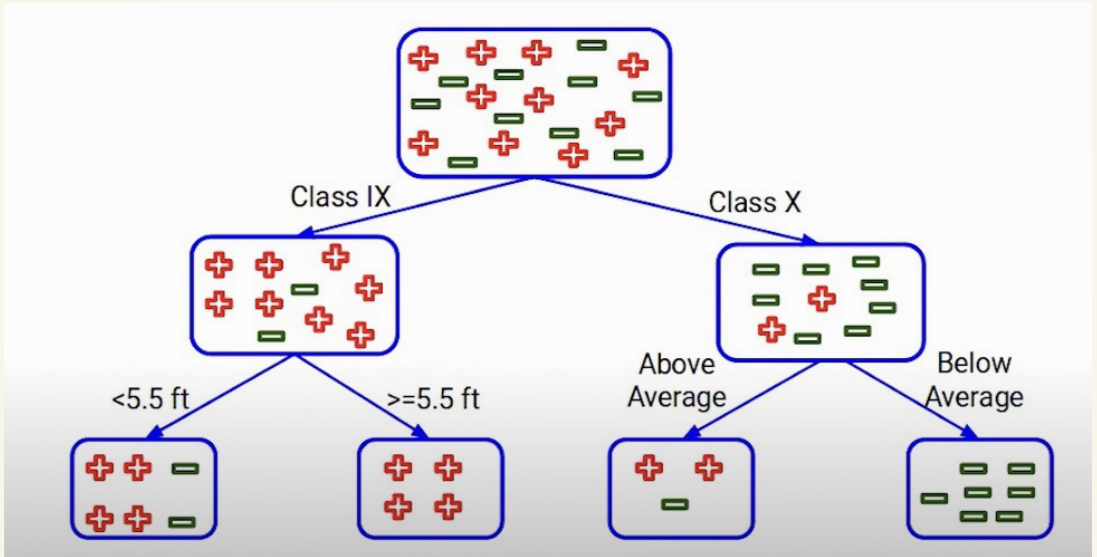
Ex 20 students in the class, 10 play cricket
fit a model to predict whether a student plays cricket
X Variables: Height, grades, class

Idea split on an X variable to classify the y variable



Which split is best? The rightmost split is best because the 80% in left bucket and 20% in right bucket separates the classes of cricket & not cricket as best as possible.

* decision trees allow variables to interact

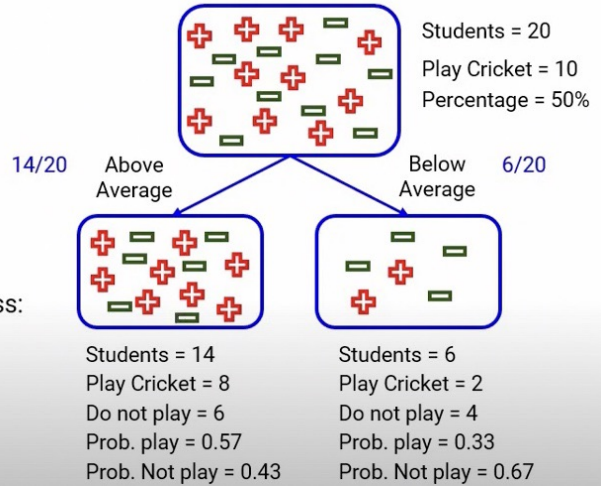


* How to actually select the best split point?

Gini 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:
 $(14/20)*0.49 + (6/20)*0.44 = 0.475$



Gini Impurity of the Above Avg Node = $1 - (P_+^2 + P_-^2) = 1 - (0.57)^2 - (0.43)^2 = .49$

Gini Impurity of the Below Avg Node = $1 - (P_+^2 + P_-^2) = 1 - (\frac{1}{3})^2 - (\frac{2}{3})^2 = .44$

fully informative: $P_+ = 1, P_- = 0$
 or
 $P_+ = 0, P_- = 1$

Gini Impurity = $1 - 1^2 - 0^2 = 0$

fully non-informative: $P_+ = \frac{1}{2}, P_- = \frac{1}{2}$

Gini Impurity = $1 - (\frac{1}{2})^2 - (\frac{1}{2})^2 = 0.5$

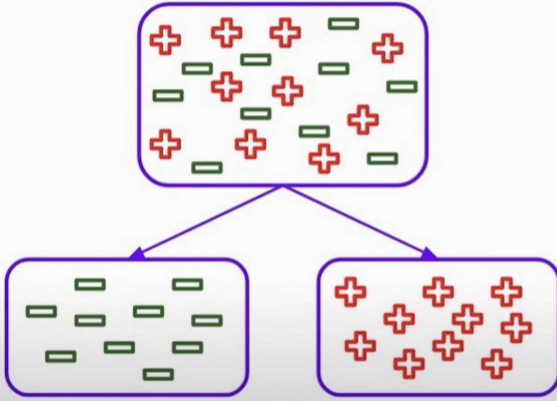
Weighted Gini Impurity = $\left(\frac{\# \text{ left bin}}{\text{total \#}}\right) (\text{Gini left}) + \left(\frac{\# \text{ right bin}}{\text{total \#}}\right) (\text{Gini right})$

$$= \left(\frac{14}{20}\right)(.49) + \left(\frac{6}{20}\right)(.44)$$
$$= 0.475$$

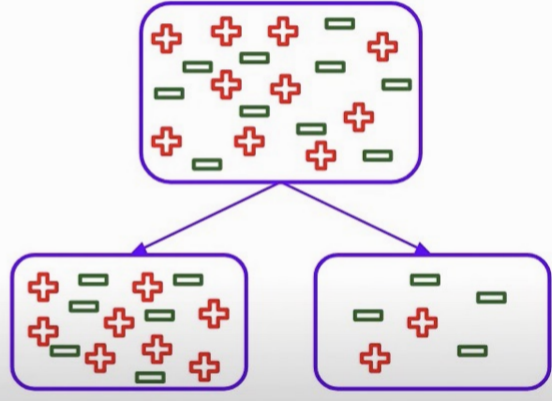
Split	Weighted Gini Impurity
Performance in Class	0.475
Class	0.32

→ split
on
Class

2. Information Gain / Entropy



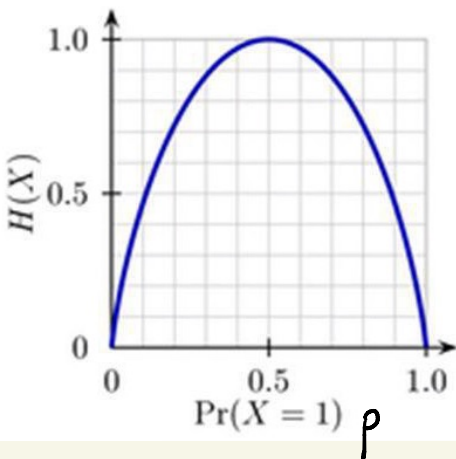
more info gained



less info gained

Information Gain = 1 - Entropy

$$\text{Entropy} = -p \log_2 p - (1-p) \log_2 (1-p)$$



$$p = P(+)=P(y=1)$$

$$p=0, \text{Entropy} = -0 \log_2 0 - 1 \cdot \log_2 1 = 0$$

$$p=1, \text{Entropy} = -1 \cdot \log_2 1 - 0 \cdot \log_2 0 = 0$$

$$p=\frac{1}{2}, \text{Entropy} = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) = 1$$



% Play = 0.50
% Not play = 0.50

$$\text{Entropy} = -(0.5) * \log_2(0.5) - (0.5) * \log_2(0.5)$$

$$= 1$$



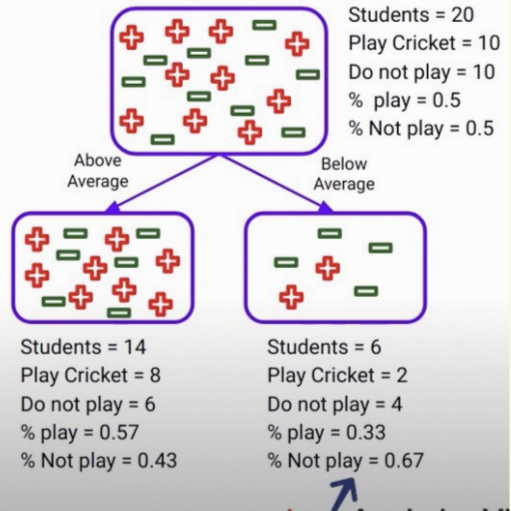
% Play = 0
% Not play = 1

$$\text{Entropy} = -(0) * \log_2(0) - (1) * \log_2(1)$$

$$= 0$$

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:
 $(14/20) * 0.98 + (6/20) * 0.91 = 0.959$



Left split bin :

$$-(0.57) \log_2(0.57) - (0.43) \log_2(0.43) = 0.98$$

Right split bin :

$$-(0.33) \log_2(0.33) - (0.67) \log_2(0.67) = 0.91$$

$$\text{weighted entropy} = \left(\frac{14}{20}\right)(0.98) + \left(\frac{6}{20}\right)(0.91) = 0.96$$

$$\text{info gain} = 1 - \text{weighted entropy} = 0.04$$

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

→ split on class!

* Our outcome variable $y = \{+, -\} = \left\{ \begin{array}{l} \text{Plays} \\ \text{doesn't play} \end{array} \right.$ is categorical, so $gini_{p^2}$ + entropy work well here.

But what if our outcome was continuous?

e.g. $y \in \mathbb{R}$ (height)

It breaks down.

Need another way to split,

3. Reduction in Variance

$$\text{Variance} = \frac{\sum [(X - \mu)^2]}{n}$$

2	6	7
4	7	9

Variance ~ 6

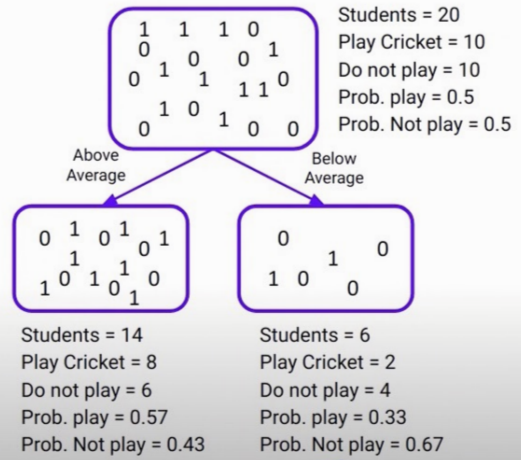
1	1	1
1	1	1

Variance = 0

- 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: $(14/20)*0.245 + (6/20)*0.222 = 0.238$

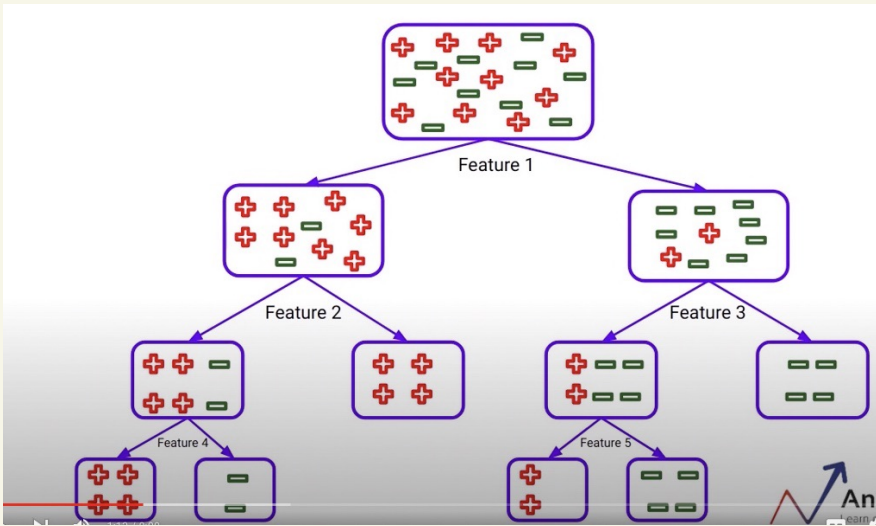


Split	Variance
Performance in Class	0.238
Class	0.16

↪ split on class!

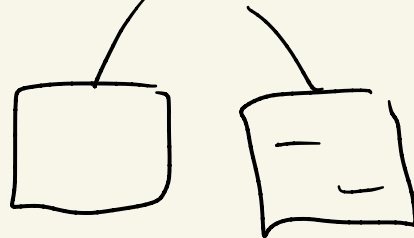
* 3+ ways to make splits,

* Could Iteratively make splits until all nodes are "pure" but this would overfit (memorize noise/randomness in the training dataset)



* Tuning Hyperparameters (to control overfitting)

- max depth of the tree
- min samples for a node split



HW: fit a decision tree WP model

* Random forests: How to further reduce overfitting in decision trees