

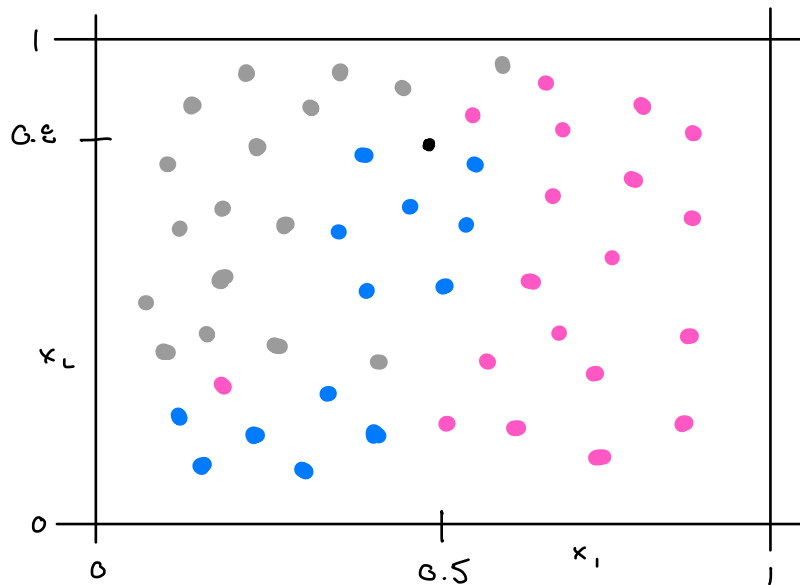
NON PARAMETRIC CLASSIFICATION

ESTIMATING $P[Y = k \mid X = x]$ WITH

- KNN
- TREES

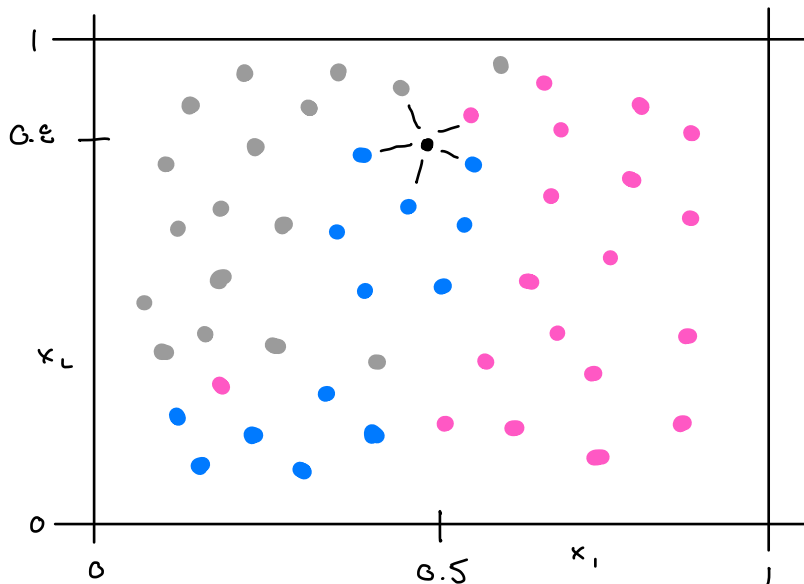
SETUP

y	x_1	x_2
A		
⋮		
A		
B		
⋮		
B		
C		
⋮		
C		
?	0.5	0.8



KNN

$$\hat{P}[Y=j | X=x] = \frac{1}{K} \sum_{\{i: x_i \in N_K(x, D)\}} I(y_i=j)$$



WITH $K=5$, AND $x=(0.5, 0.8)$

$$\hat{P}[Y=A | X=x] = 3/5$$

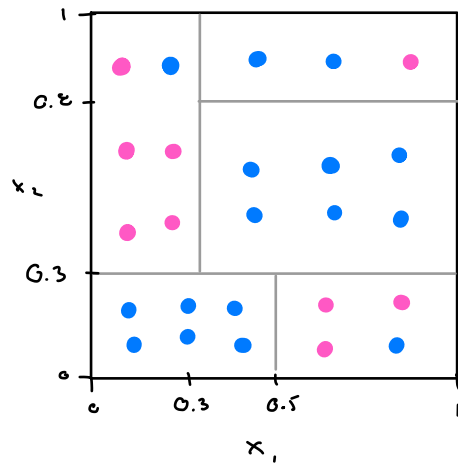
$$\hat{P}[Y=B | X=x] = 1/5$$

$$\hat{P}[Y=C | X=x] = 1/5$$

IF BINARY \rightarrow USE ODD K

\hookrightarrow AVOID TIES

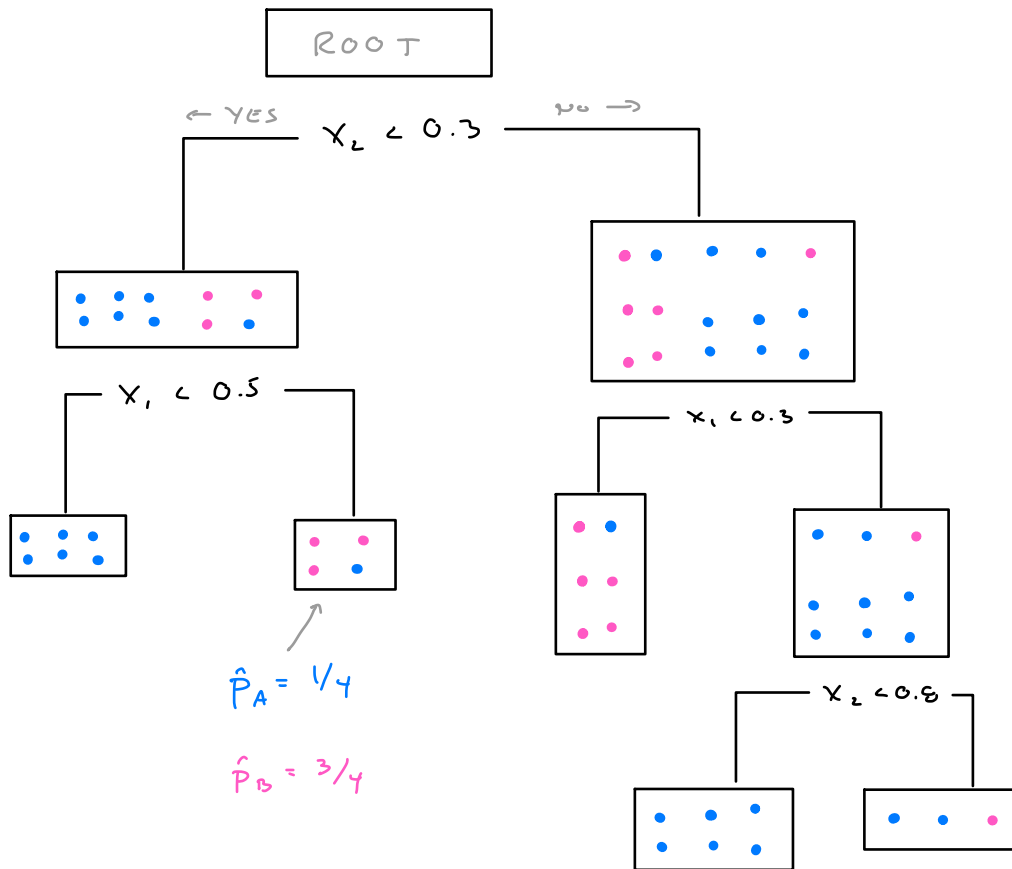
DECISION TREES



MIN SPLIT = 8

CP = 0

DIFFERENT INTERPRETATION
SAME USAGE



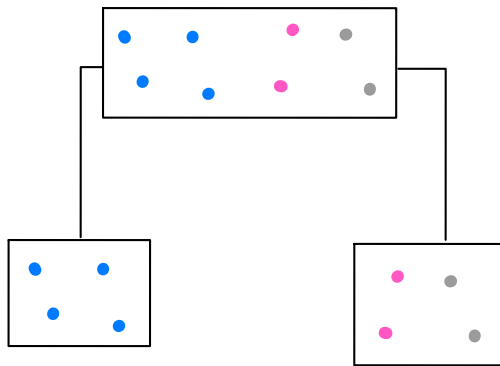
Node Probabilities

$$\hat{p}_k = \frac{\sum_i I(y_i = k) I(x_i \in A)}{\sum_i I(x_i \in A)}$$

$$\hat{p}_A = 4/8$$

$$\hat{p}_B = 2/8$$

$$\hat{p}_C = 2/8$$



$$\hat{p}_A = 4/4$$

$$\hat{p}_B = 0/4$$

$$\hat{p}_C = 0/4$$

$$\hat{p}_A = 0/4$$

$$\hat{p}_B = 2/4$$

$$\hat{p}_C = 2/4$$

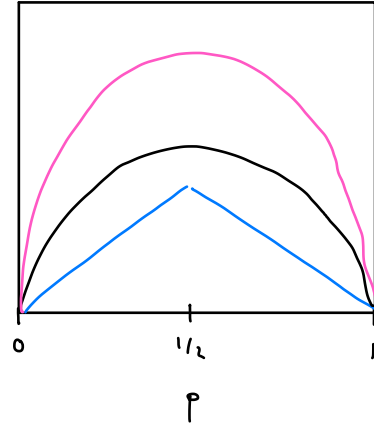
IMPURITY MEASURES FOR CATEGORICAL DATA

"VARIANCE"

$$Gini(A) = \sum_{k=1}^K \hat{p}_k (1 - \hat{p}_k) = 1 - \sum_{k=1}^K \hat{p}_k^2$$

$$\underline{Entropy(A)} = - \sum_{k=1}^K \hat{p}_k \log(\hat{p}_k)$$

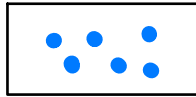
$$\underline{Error(A)} = 1 - \max_k (\hat{p}_k)$$



CALCULATING GINI

$$G_{INI}(A) = \sum_{k=1}^K \hat{p}_k (1 - \hat{p}_k) = 1 - \sum_{k=1}^K \hat{p}_k^2$$

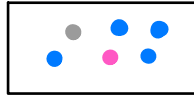
A:



$$\hat{p}_A = 6/6$$

$$\hat{p}_B = 0/6$$

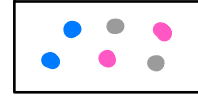
$$\hat{p}_C = 0/6$$



$$\hat{p}_A = 4/6$$

$$\hat{p}_B = 1/6$$

$$\hat{p}_C = 1/6$$



$$\hat{p}_A = 2/6$$

$$\hat{p}_B = 2/6$$

$$\hat{p}_C = 2/6$$

$$G_{INI}(A)$$

$$0$$

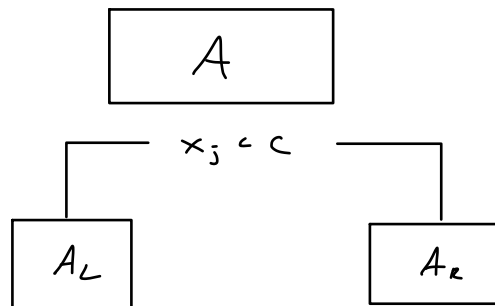
$$0.5$$

$$0.66\bar{6}$$

$$\hookrightarrow = 1 - \left[\left(\frac{4}{6} \right)^2 + \left(\frac{1}{6} \right)^2 + \left(\frac{1}{6} \right)^2 \right]$$

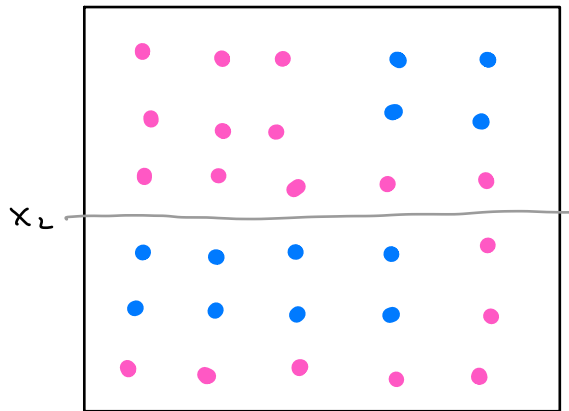
SPLITTING

FIND \rightarrow FEATURE x_j
 \rightarrow CUTOFF c

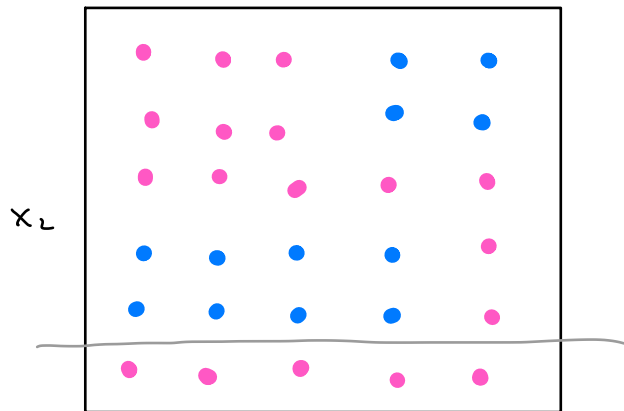


$$\min_{j,c} \left[\underbrace{\frac{|A_L|}{|A|}}_{\text{WEIGHTS}} \underbrace{G_{INI}(A_L)}_{\text{"VARIANCE"}} + \underbrace{\frac{|A_R|}{|A|}}_{\text{WEIGHTS}} \underbrace{G_{INI}(A_R)}_{\text{"VARIANCE"}} \right]$$

WHICH SPLIT?



0.44



SMALLER GINI

0.416

ln R

KNN

caret :: knn3()

TREES

rpart :: rpart()