

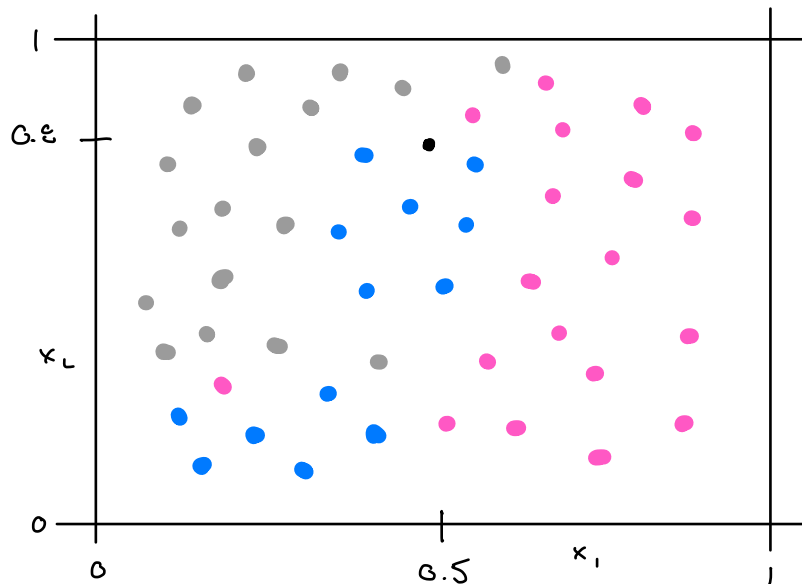
# NON PARAMETRIC CLASSIFICATION

ESTIMATING  $P[Y = k | X = x]$  WITH

- KNN
- TREES

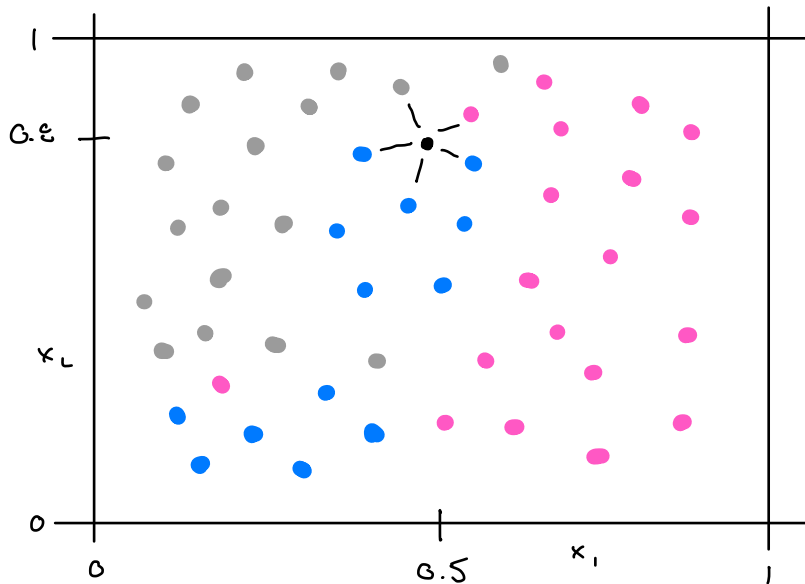
# SETUP

$y$	$x_1$	$x_2$
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  $x=(0.5, 0.8)$

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

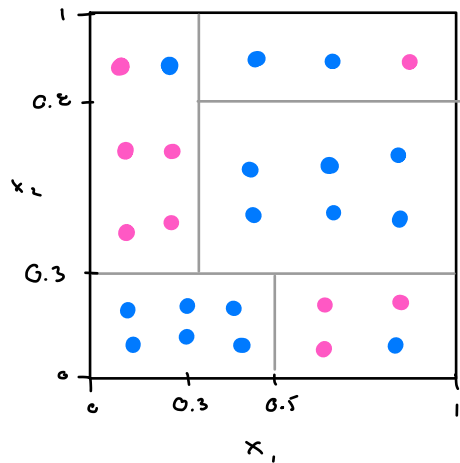
$$\hat{P}[Y=B | X=x] = 1/5$$

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

(F BINARY  $\rightarrow$  USE ODD  $K$ )

$\hookrightarrow$  AVOID TIES

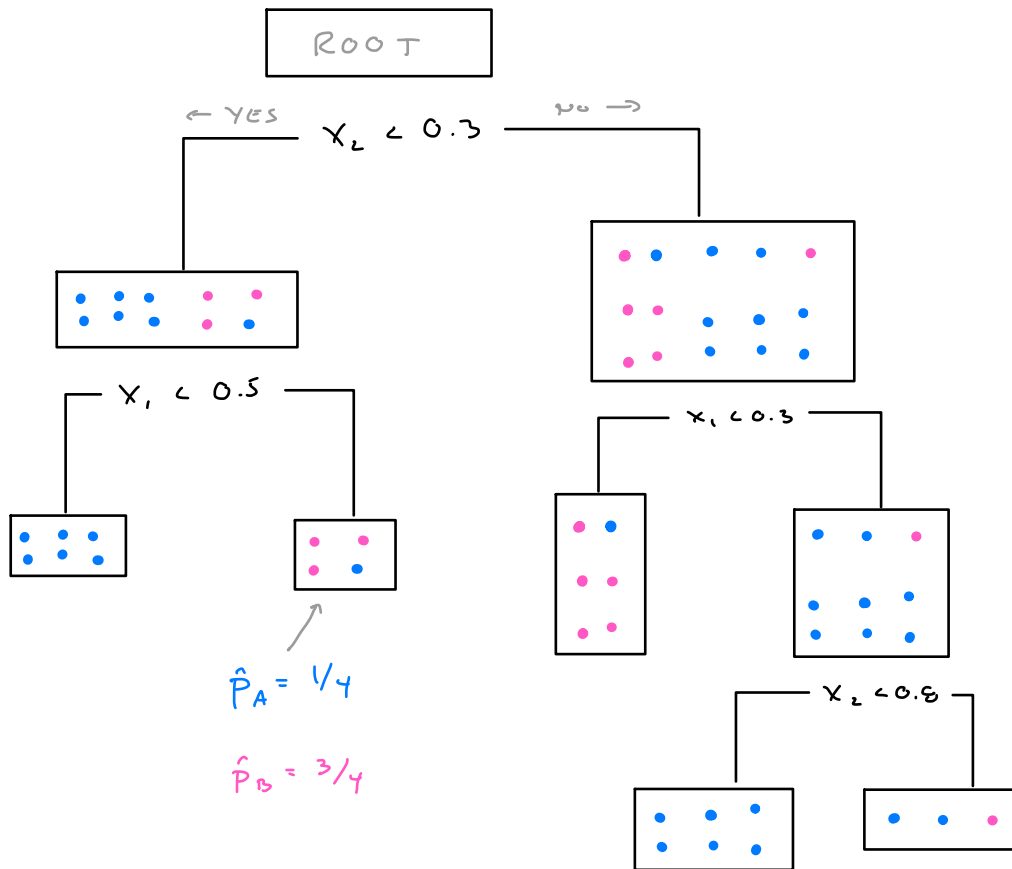
# DECISION TREES



MIN SPLIT = 8

CP = 0

DIFFERENT INTERPRETATION  
SAME USAGE



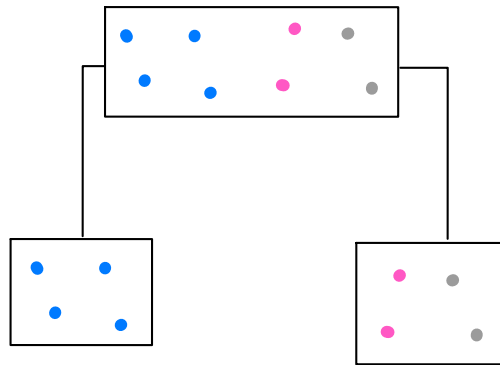
# NODE PROBABILITIES

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

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

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

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



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

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

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

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

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

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

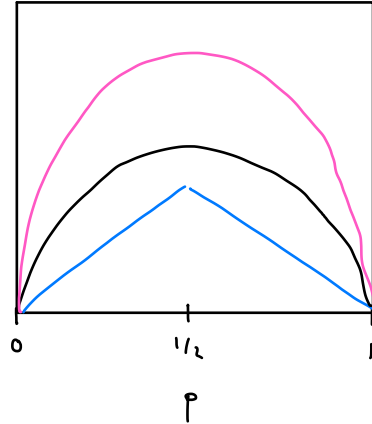
# IMPURITY MEASURES FOR CATEGORICAL DATA

"VARIANCE"

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

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

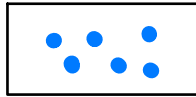
$$\underline{Error(A)} = 1 - \max_k (\hat{p}_k)$$



# CALCULATING GINI

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

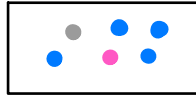
A:



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

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

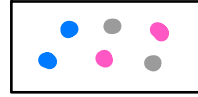
$$\hat{p}_C = 0/6$$



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

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

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



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

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

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

$G_{INI}(A)$

0

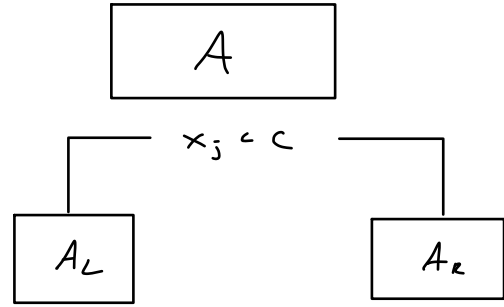
0.5

0.666

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

# SPLITTING

FIND  $\begin{cases} \rightarrow \text{FEATURE } x_j \\ \rightarrow \text{CUTOFF } c \end{cases}$



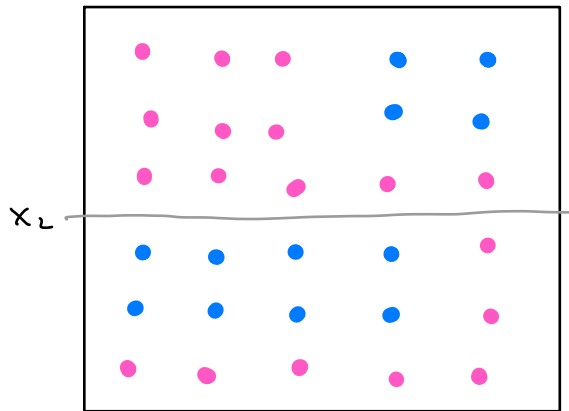
$$\min_{j,c} \left[ \frac{|A_L|}{|A|} \text{GINI}(A_L) + \frac{|A_R|}{|A|} \text{GINI}(A_R) \right]$$

WEIGHTS (blue arrows pointing to  $\frac{|A_L|}{|A|}$  and  $\frac{|A_R|}{|A|}$ )

"VARIANCE" (pink arrows pointing to  $\text{GINI}(A_L)$  and  $\text{GINI}(A_R)$ )

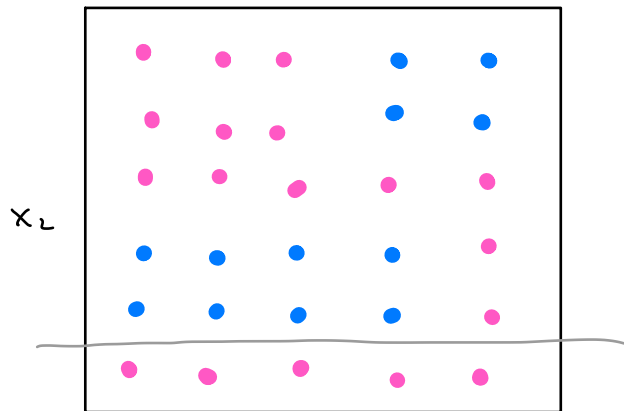


WHICH SPLIT?



$x_1$

$0.4\bar{4}$



$x_2$

$x_1$

SMALLER GINI

$0.416$

ln R

KNN

caret :: knn3()

TREES

caret :: rpart()