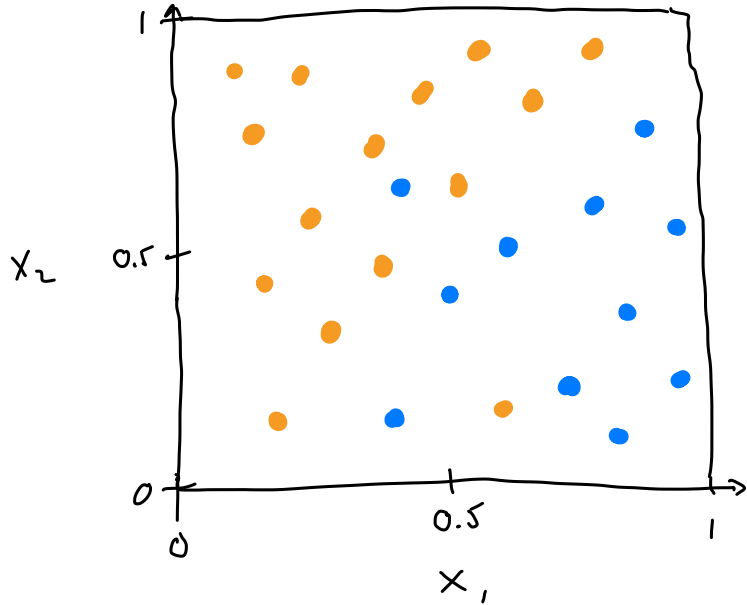


NONPARAMETRIC CLASSIFICATION

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

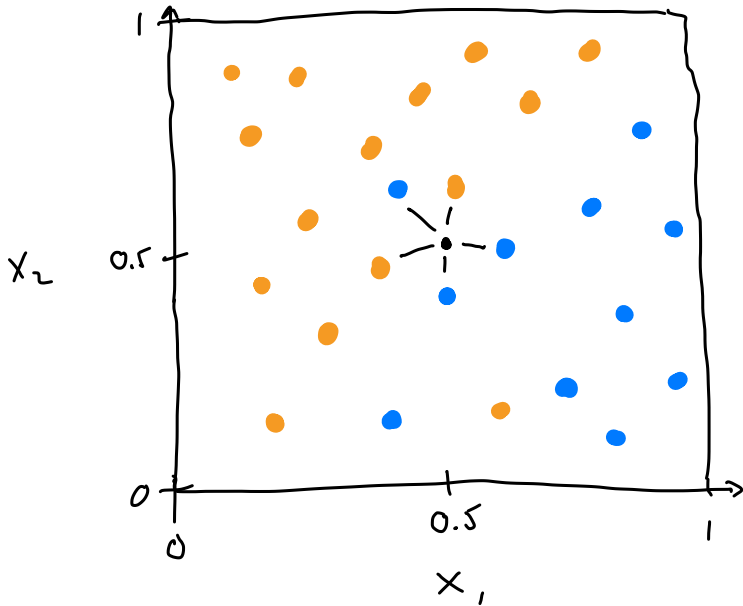
- KNN
- DECISION TREES

SETUP



| y | x_1 | x_2 |
|-----|-------|-------|
| A | ⋮ | ⋮ |
| A | ⋮ | ⋮ |
| A | ⋮ | ⋮ |
| A | ⋮ | ⋮ |
| B | ⋮ | ⋮ |
| B | ⋮ | ⋮ |
| B | ⋮ | ⋮ |
| B | ⋮ | ⋮ |
| B | ⋮ | ⋮ |
| ? | ⋮ | 0.5 |
| ? | ⋮ | 0.5 |

KNN



DISTANCE ?

- EUCLIDEAN !
- WHATEVER !

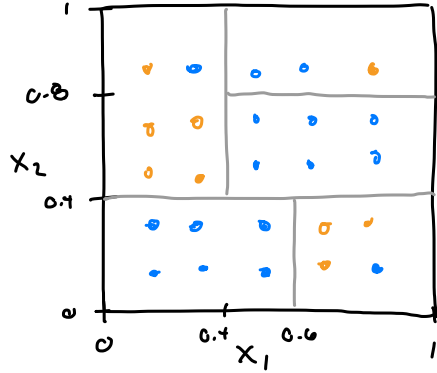
$$\hat{P}[x=j | X=x] = \frac{1}{K} \sum_{i \in N_k(x, D)} I(y_i=j)$$

$$KNN=5 \begin{cases} \hat{P}[y=\bullet | x_1=0.5, x_2=0.5] = 2/5 \\ \hat{P}[y=\circ | x_1=0.5, x_2=0.5] = 3/5 \end{cases}$$

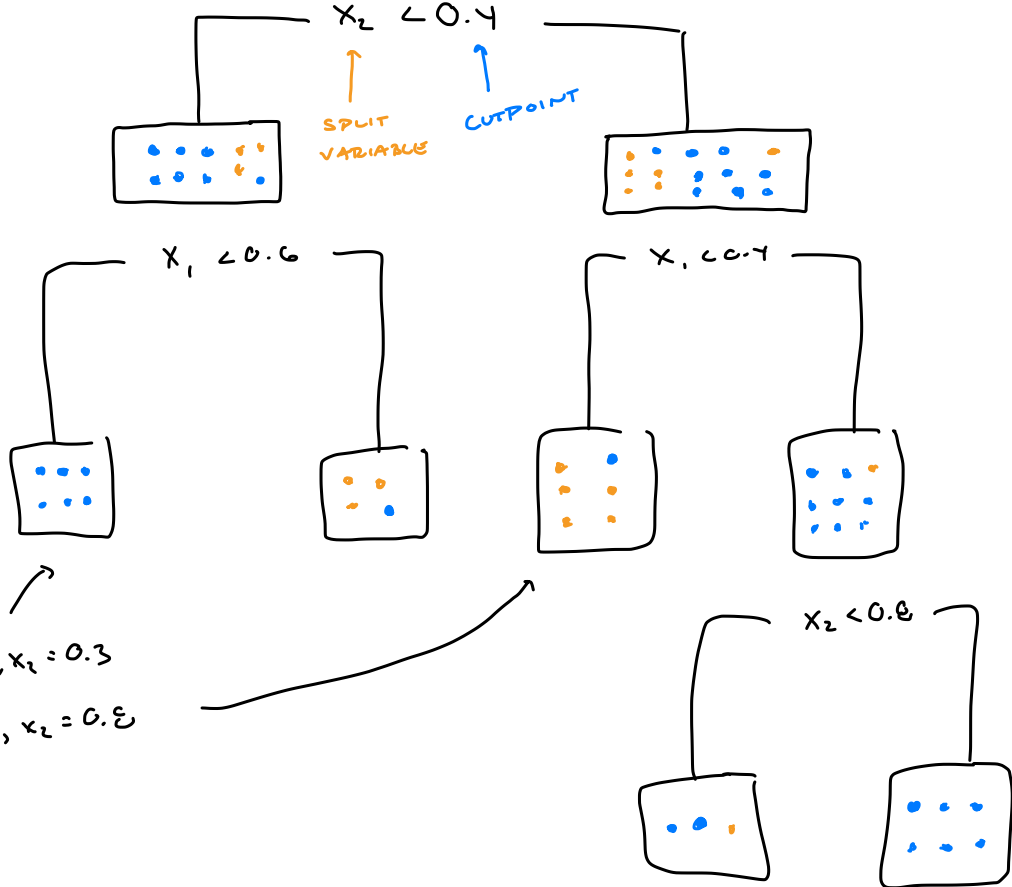
BINARY ? → USE ODD K

↳ AVOID TIES

DECISION TREES



ROOT



- MIN SPLIT = ϵ
- $C_p = 0$

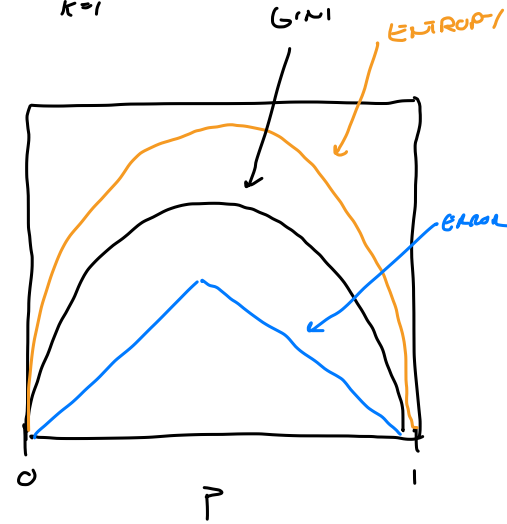
- $X_1 = 0.3, X_2 = 0.3$
- $X_1 = 0.2, X_2 = 0.8$

VARIANCE (IMPURITY) MEASURES IN CATEGORICAL DATA

$$\rightarrow \bullet \text{GINI}(A) = \sum_{k=1}^K \hat{p}_k (1 - \hat{p}_k) = 1 - \sum_{k=1}^K p_k^2$$

$$\bullet \text{ENTROPY}(A) = - \sum_{k=1}^K \hat{p}_k \log(\hat{p}_k)$$

$$\bullet \text{ERROR}(A) = 1 - \max_k (\hat{p}_k)$$

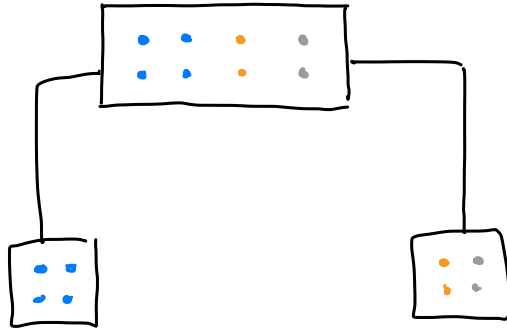


$$\hat{p}_k = \frac{\sum_i I(y_i = k) I(x_i \in A)}{\sum_i I(x_i \in A)}$$

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

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

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



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

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

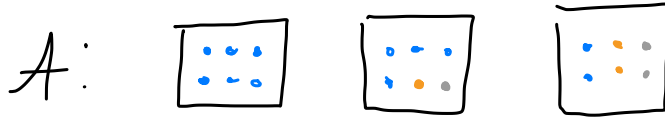
$$\hat{p}_C = 0/4$$

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

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

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

$$G_{IN_1}(A) = \sum_{k=1}^K \hat{p}_k (1 - \hat{p}_k) = 1 - \sum_{k=1}^K p_k^2$$



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

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

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

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

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

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

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

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

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

$$G_{IN}(A)$$

$$0$$

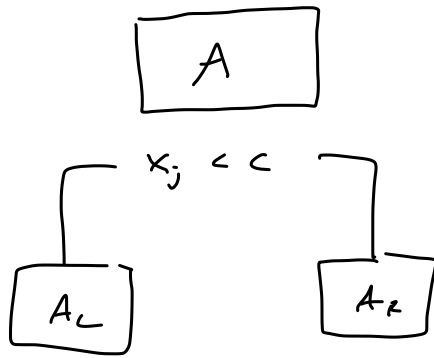
$$0.5$$

$$0.66\bar{6}$$

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

SPLITTING

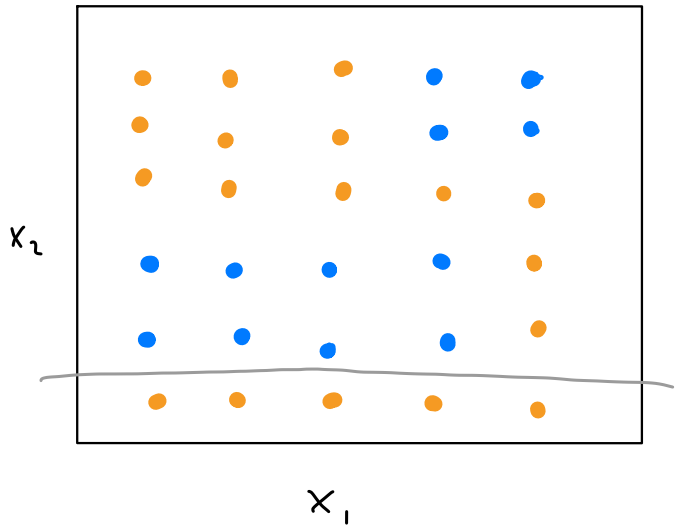
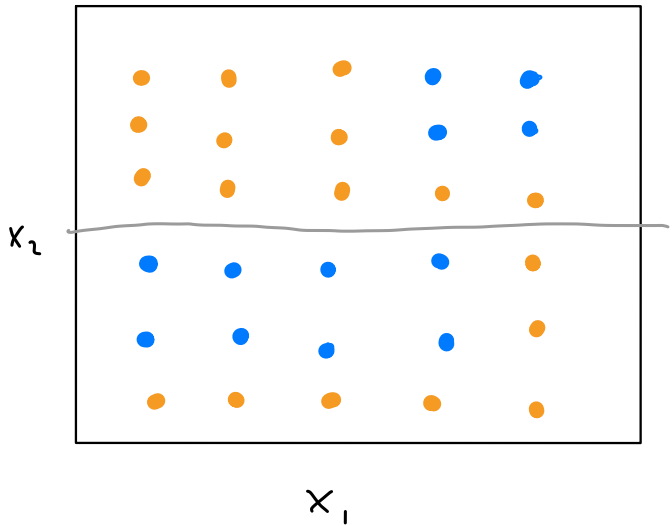
FIND \rightarrow VARIABLE x_j
 \rightarrow CUTOFF c



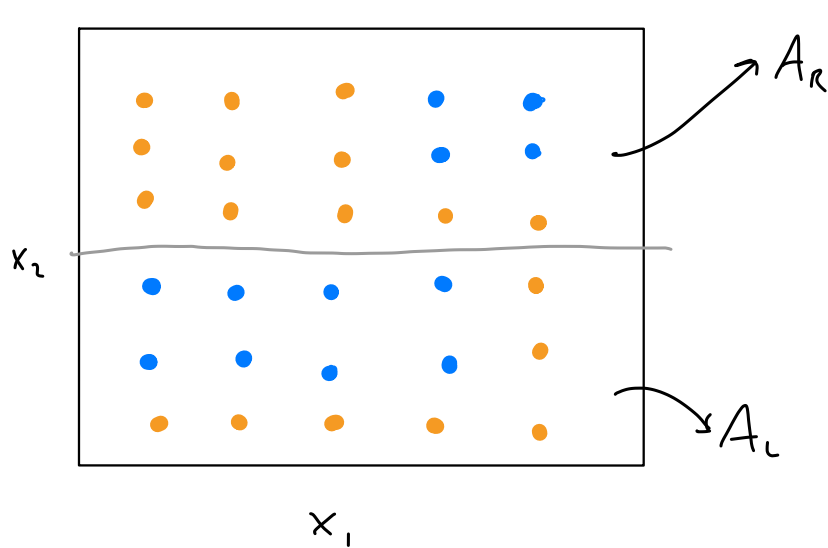
$$\min_{j, c} \left[\frac{|A_L|}{|A|} \text{GINI}(A_L) + \frac{|A_R|}{|A|} \text{GINI}(A_R) \right]$$

Annotations:
- Blue arrows labeled "WEIGHTS" point to the fractions $\frac{|A_L|}{|A|}$ and $\frac{|A_R|}{|A|}$.
- Orange arrows labeled "VARIANCE" point to the $\text{GINI}(A_L)$ and $\text{GINI}(A_R)$ terms.

WHICH SPLIT?



↑
SMALLER GINI



$$\hat{P}_A = 11/15$$

$$\hat{P}_B = 4/15$$

$$|A_R| = 15$$

$$\hat{P}_A = 7/15$$

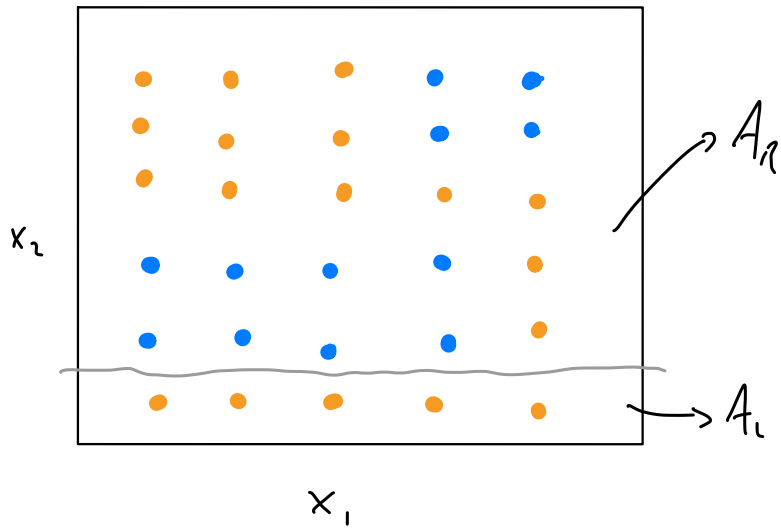
$$\hat{P}_B = 8/15$$

$$|A_L| = 15$$

$$G_{INI}(A_R) = 1 - \left[\left(\frac{11}{15} \right)^2 + \left(\frac{4}{15} \right)^2 \right] = \frac{28}{225}$$

$$G_{INI}(A_L) = 1 - \left[\left(\frac{7}{15} \right)^2 + \left(\frac{8}{15} \right)^2 \right] = \frac{112}{225}$$

$$\frac{|A_L|}{|A|} G_{INI}(A_L) + \frac{|A_R|}{|A|} G_{INI}(A_R) = \frac{15}{30} \left(\frac{112}{225} \right) + \frac{15}{30} \left(\frac{28}{225} \right) = 0.44\bar{4}$$



$$\hat{p}_A = 13/25$$

$$\hat{p}_B = 12/25$$

$$|A_U| = 25$$

$$\hat{p}_A = 5/5$$

$$\hat{p}_B = 0/5$$

$$|A_L| = 5$$

$$G_{INI}(A_U) = 1 - \left[\left(\frac{13}{25} \right)^2 + \left(\frac{12}{25} \right)^2 \right] = 312/625$$

$$G_{INI}(A_L) = 1 - \left[\left(\frac{5}{5} \right)^2 + \left(\frac{0}{5} \right)^2 \right] = 0$$

$$\frac{|A_L|}{|A|} G_{INI}(A_L) + \frac{|A_U|}{|A|} G_{INI}(A_U) = \frac{5}{30} (0) + \frac{25}{30} \left(\frac{312}{625} \right) = 0.416$$

TREE QUESTIONS

CATEGORICAL VARIABLES ?

MISSING DATA ?

STOPPING RULES ?

PRUNING ? COMPLEXITY ?