CS 307

Spring 2024

Dalpiaz

Decision Trees

Classification
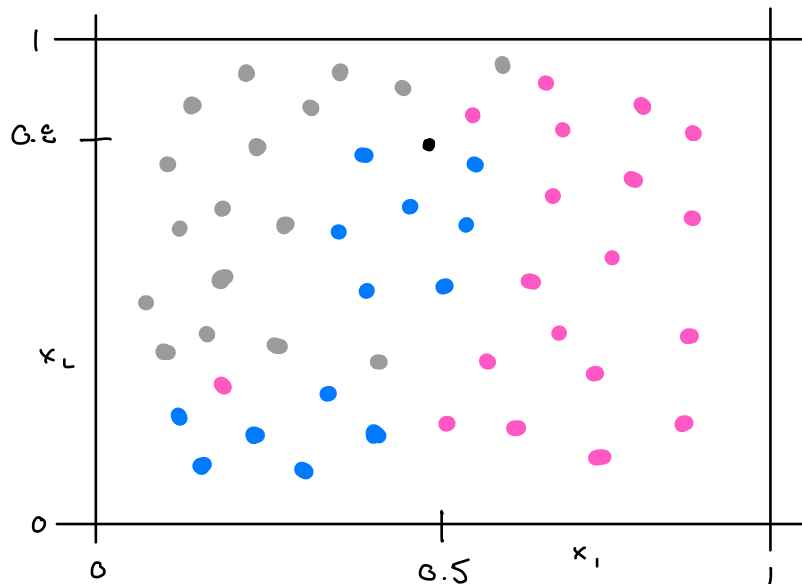
# Non Parametric Classification

Estimating $P[Y = k \mid X = x]$ with

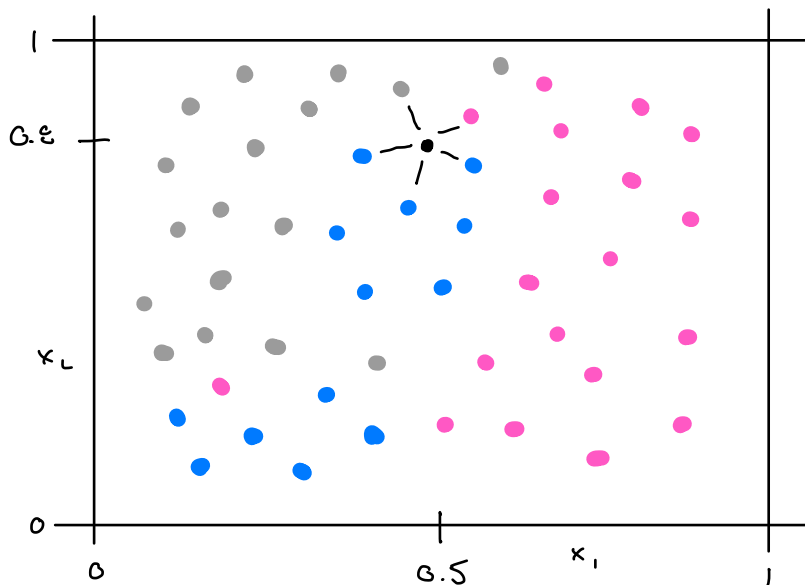- KNN
- TREES

# $\underline{S_{ETUP}}$

| y | $x_1$ | $x_2$ |
|---|---|---|
| A | | |
| ⋮ | | |
| A | | |
| B | | |
| ⋮ | | |
| B | | |
| C | | |
| ⋮ | | |
| C | | |
| ? | 0.5 | 0.8 |

# KNN

$$\hat{P}\left[Y = j \mid X = x\right] = \frac{1}{K} \sum_{\{i \,:\, x_i \in N_k(x, D)\}} \mathbb{I}(y_i = j)$$

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



$$\hat{P}\left[Y = A \mid X = x\right] = 3/5 \leftarrow$$

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

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

IF BINARY → USE ODD $k$
   ↳ AVOID TIES

# Decision Trees



ROOT

← YES          $X_2 < 0.3$          NO →

$X_1 < 0.5$

$X_1 < 0.3$

$X_2 < 0.8$

$$\hat{P}\left[Y = \bullet \mid x_1 = 0.9, \; x_2 = 0.1\right] = 1/4$$

# $N_{ODE}$ $P_{ROBABILITIES}$

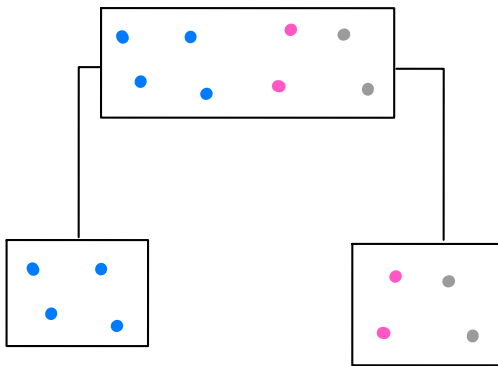$$\hat{P}_k = \frac{\sum_i I(y_i = k) I(x_i \in A)}{\sum_i I(x_i \in A)}$$

$$\downarrow$$

$$\hat{P}\left[Y = k \mid x \in N\right]$$

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

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

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



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

$$\hat{P}_B = 0 \qquad\qquad\quad \hat{P}_B = 2/4$$

$$\hat{P}_C = 0 \qquad\qquad\quad \hat{P}_C = 2/4$$
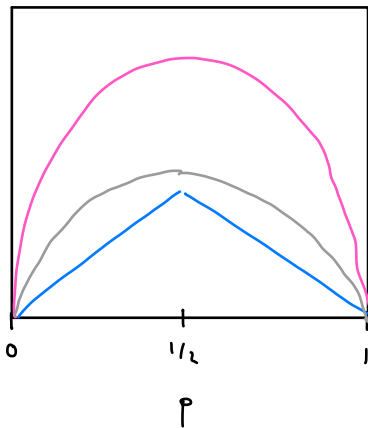
# Impurity Measures For Categorical Data

"VARIANCE"

# CATEGORIES

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

$$\text{Entropy}\,(A) = -\sum_{k=1}^{K} \hat{p}_k\,\log\!\left(\hat{p}_k\right)$$
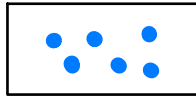
$$\text{Error}\,(A) = 1 - \max_k\left(\hat{p}_k\right)$$

# CALCULATING GINI

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

$A$ :



| | | |
|---|---|---|
| $\hat{p}_A = 6/6$ | $\hat{p}_A = 4/6$ | $\hat{p}_A = 2/6$ |
| $\hat{p}_B = 0/6$ | $\hat{p}_B = 1/6$ | $\hat{p}_B = 2/6$ |
| $\hat{p}_C = 0/6$ | $\hat{p}_C = 1/6$ | $\hat{p}_C = 2/6$ |

$\text{GINI}(A)$      $0$      $0.5$      $0.6\bar{6}$

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

# Splitting

Find
- Feature $x_j$
- Cutoff $c$

$N$

$x_j < c$

$N_L$   $N_R$

$$\min_{j,c} \left[ \frac{|N_L|}{|N|} \, \text{Gini}(N_L) + \frac{|N_R|}{|N|} \, \text{Gini}(N_R) \right]$$

(weights)

Variances

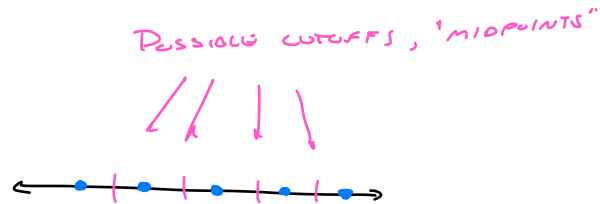# WHICH SPLIT?



$x_2$

$x_1$

0.44

$x_2$

$x_1$

0.416

- FOR EACH CURRENT TERMINAL NODE :
  - IF STOP CONDITIONS MET:
    - SKIP NODE
  - IF NO NODES CAN BE SPLIT:
    - STOP ALGORITHM
  - INIT STORAGE FOR IMPURITIES
  - FOR EACH FEATURE VARIABLE :
    - FOR EACH MIDPOINT, C :
      - CALCULATE AND STORE IMPURITY
  - FIND j AND C THAT GIVE LOWEST IMPURITY
  - SPLIT CURRENT NODE GIVEN J AND C

POSSIBLE CUTOFFS, 'MIDPOINTS'

# TREES ARE INDIFFERENT TO SCALING



The figure shows two horizontal number lines. The top line is labeled $X_1$ with tick marks at values 1000, 2000, and 3000. The bottom line is labeled $X_1$ -SCALED with tick marks at values 1, 2, and 3.