

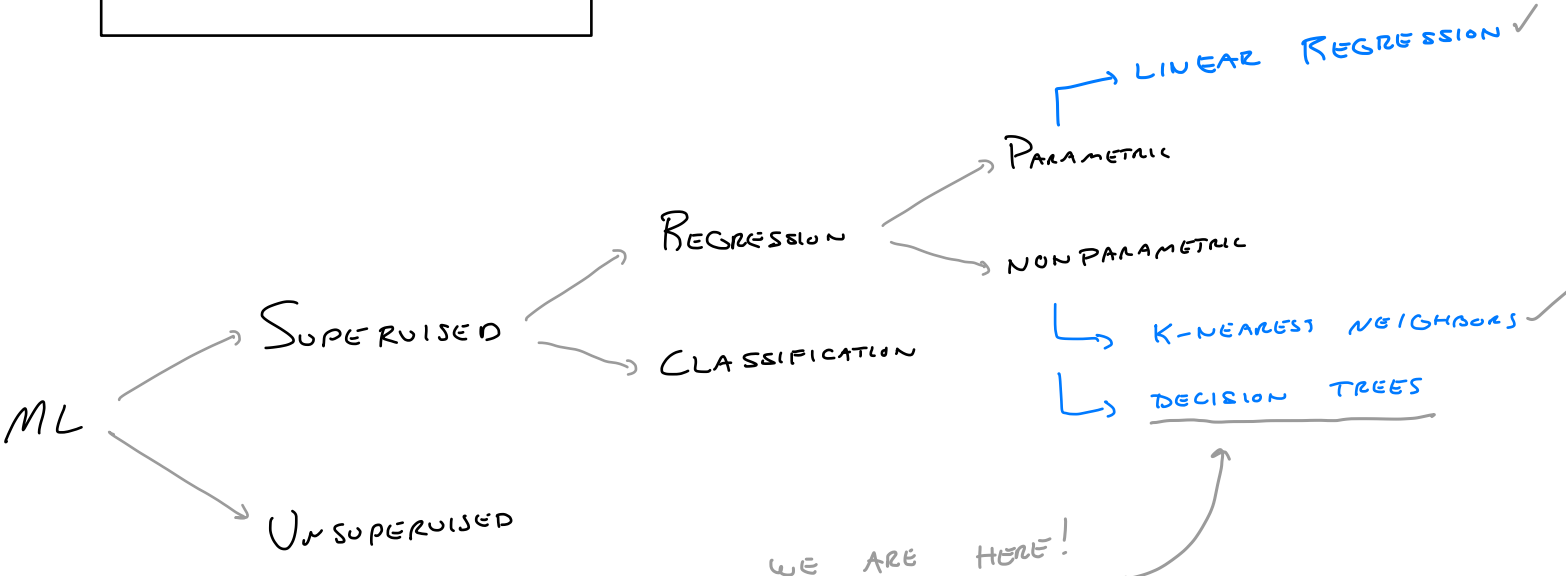
CS 307

FALL 2023

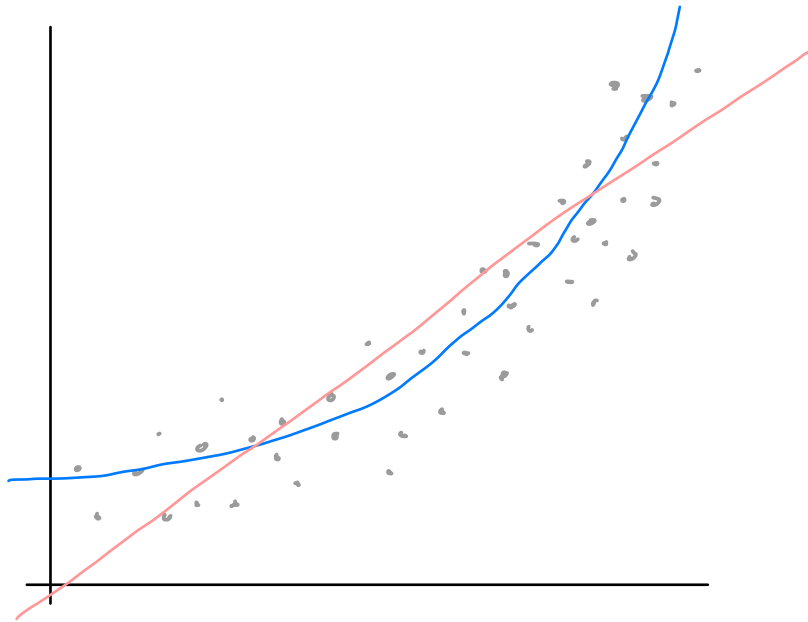
DALPIAZ

WEEK 03

NONPARAMETRIC  
REGRESSION



x	y
·	·
·	·
·	·
·	·
·	·
·	·
·	·



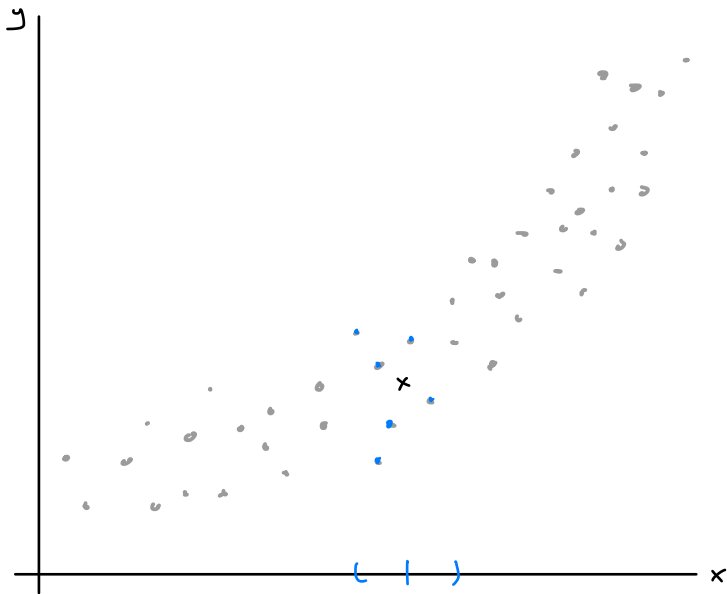
WANT  
 $E[Y | X=x]$

ASSUME

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

ASSUME

$$Y = \beta_0 + \beta_1 x + \epsilon$$



WANT

---


$$\mathbb{E}[Y | X=x]$$

•  $\hat{\mathbb{E}}[Y | X=x] = \text{AVE} \left( \left\{ y_i \text{ WHERE } x_i = x \right\} \right)$  ← WON'T WORK

•  $\hat{\mathbb{E}}[Y | X=x] = \text{AVE} \left( \left\{ y_i \text{ WHERE } x_i \text{ "CLOSE" TO } x \right\} \right)$



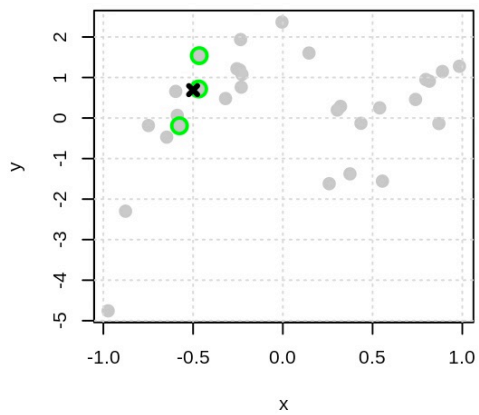
# k-NEAREST NEIGHBORS

To estimate  $\mu(x) = \mathbb{E}[Y | X=x]$

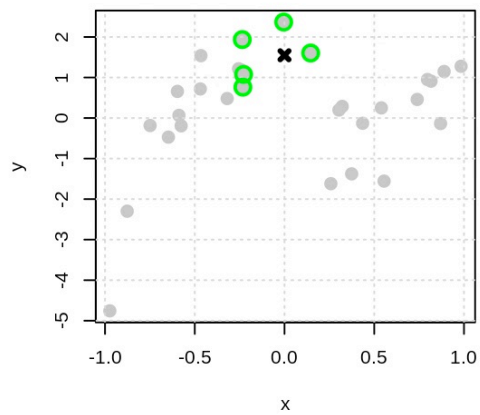
USE  $\hat{\mu}_k(x) = \frac{1}{k} \sum_{\{i: x_i \in N_h(x, D)\}} y_i$

k OBSERVATIONS WITH  $x_i$  NEAREST TO  $x$

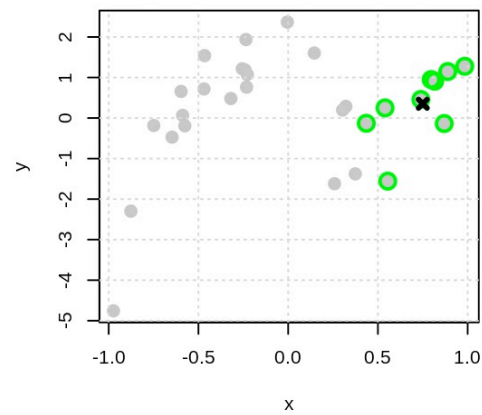
**k = 3, x = -0.5**



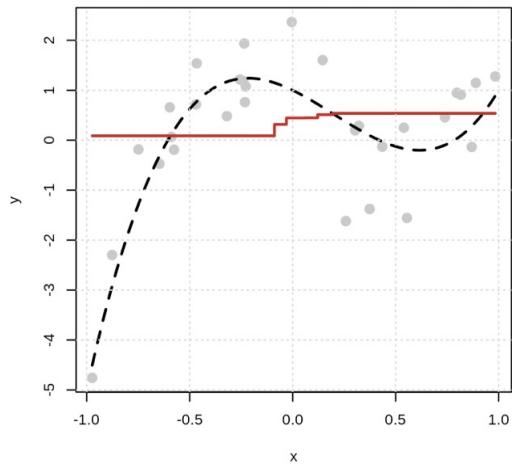
**k = 5, x = 0**



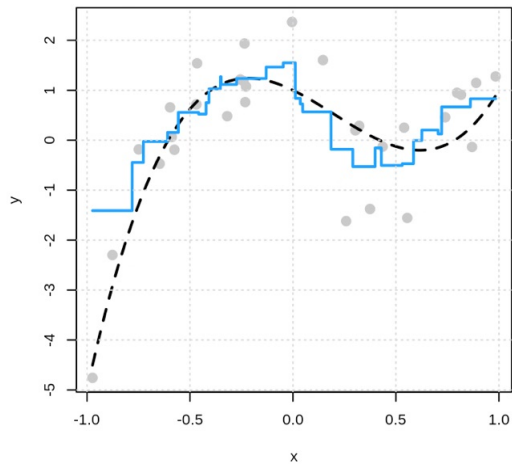
**k = 9, x = 0.75**



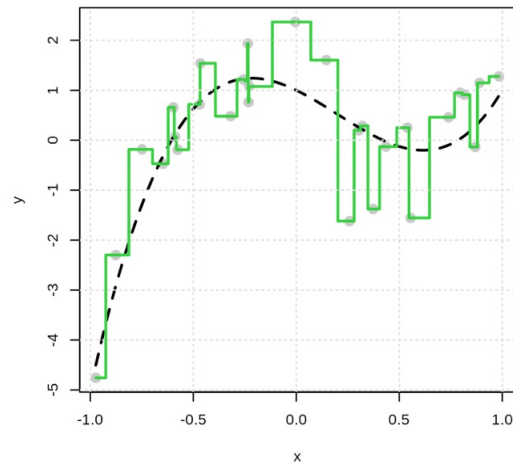
**k = 25**



**k = 5**



**k = 1**



# TUNING PARAMETERS

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

↑   ↑   ↑   ↑  
MODEL PARAMETERS

↪ LEARNED FROM DATA

K IN KNN

↑  
TUNING PARAMETER

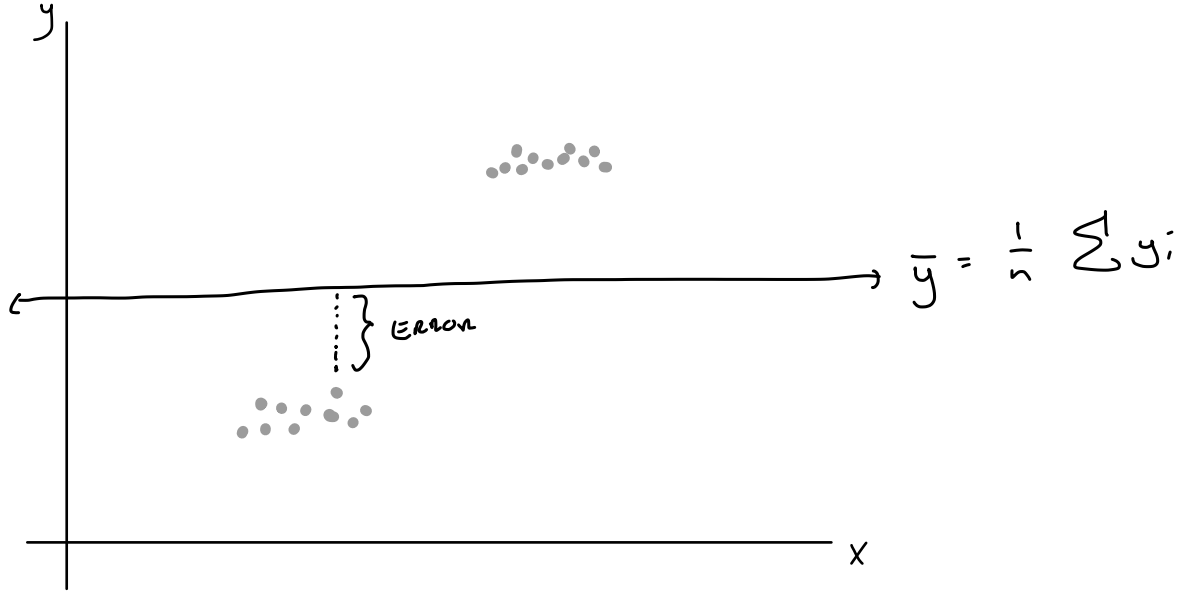
↪ DEFINES HOW TO LEARN FROM DATA

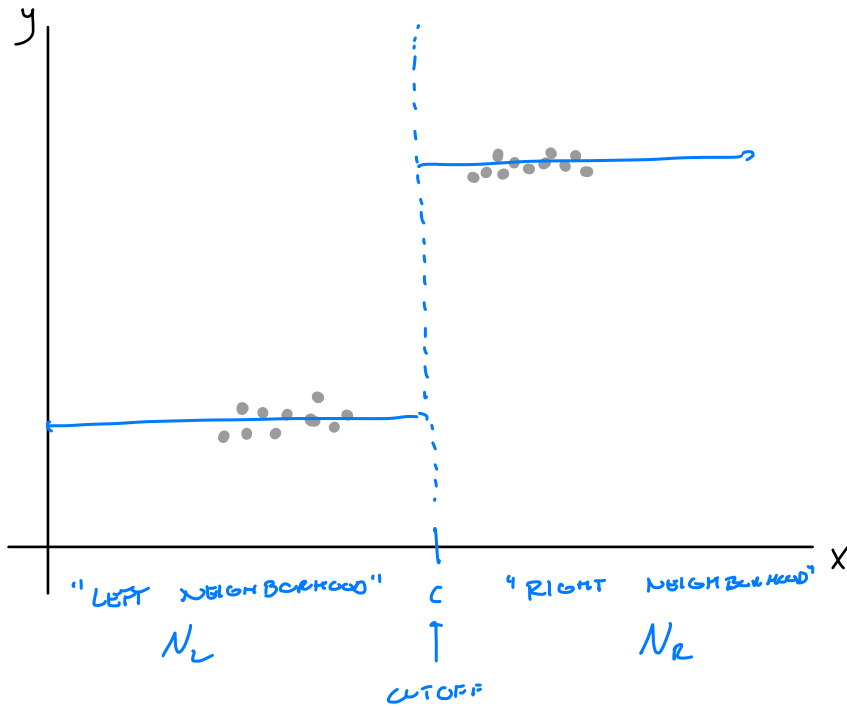


# OTHER KNN NOTES

- 'FAST' TO TRAIN, 'SLOW' TO PREDICT LAZY!
- WHICH FEATURES SHOULD BE USED? ???
- CATEGORICAL FEATURES? DUMMY ENCODING
- HOW TO CALCULATE DISTANCE? YOU PICK!
- FEATURE SCALING? !!!

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$





IDEA: FIND NEIGHBORHOODS, PREDICT AVERAGE OF  $y$ ; IN NEIGHBORHOODS

# DECISION TREES

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

→ FEATURE + CUTOFF  
FIND "SPLIT" THAT  
MINIMIZES

$c$  = "CUTOFF"

$$\sum_{i \in N_L} (y_i - \hat{\mu}_{N_L})^2 + \sum_{i \in N_R} (y_i - \hat{\mu}_{N_R})^2$$

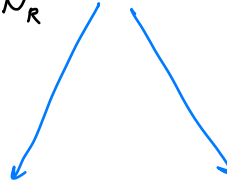
→ AVG  $y_i$  IN  $N_L$       → AVG  $y_i$  IN  $N_R$

↓  $x < c$       ↓  $x > c$

# RECURSIVE BINARY PARTITIONING

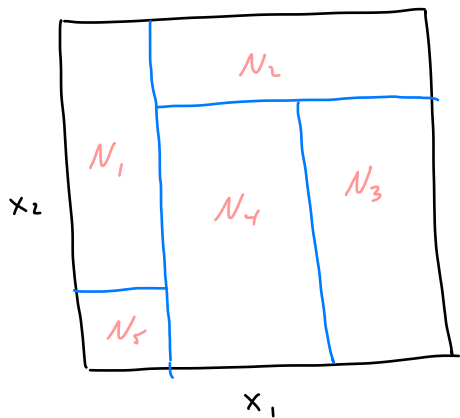
DO IT AGAIN!

$$\sum_{i \in N_L} (y_i - \hat{\mu}_{N_L})^2 + \sum_{i \in N_R} (y_i - \hat{\mu}_{N_R})^2$$



$$\sum_{i \in N_{R1}} (y_i - \hat{\mu}_{N_{R1}})^2 + \sum_{i \in N_{R2}} (y_i - \hat{\mu}_{N_{R2}})^2$$

# RECURSIVE PARTITIONING



$\#$  NEIGHBORHOODS

$$SSE = \sum_{j=1}^J \sum_{i \in N_j} (y_i - \hat{\mu}_j)^2$$

↑  
AVE  $y_i$  IN  $N_j$

$$R^2 = 1 - \frac{SSE}{SST}$$

How TO STOP ?

sklearn.tree.DecisionTreeRegressor

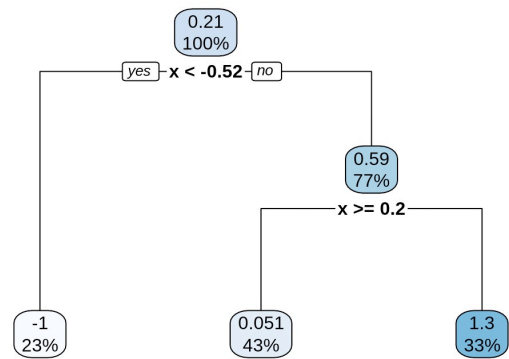
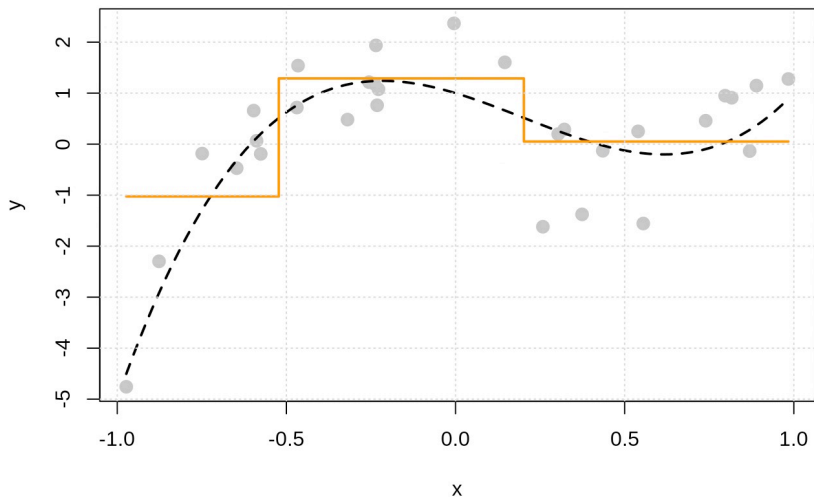
max\_depth

min\_samples\_split

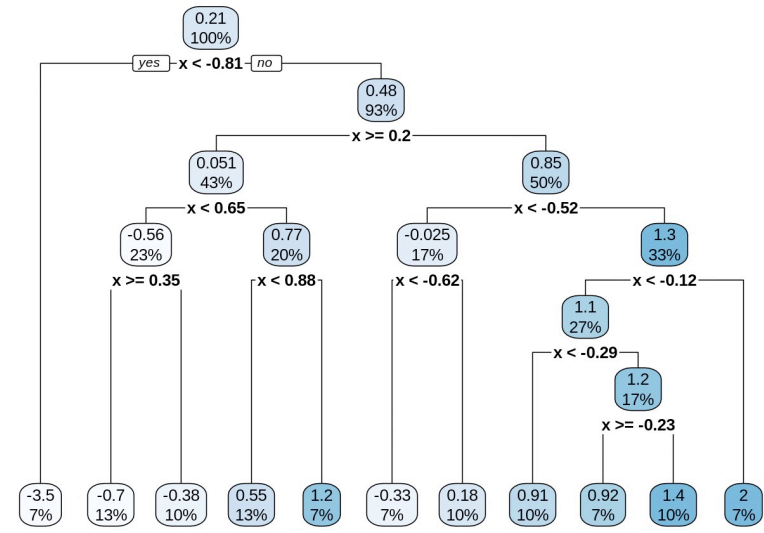
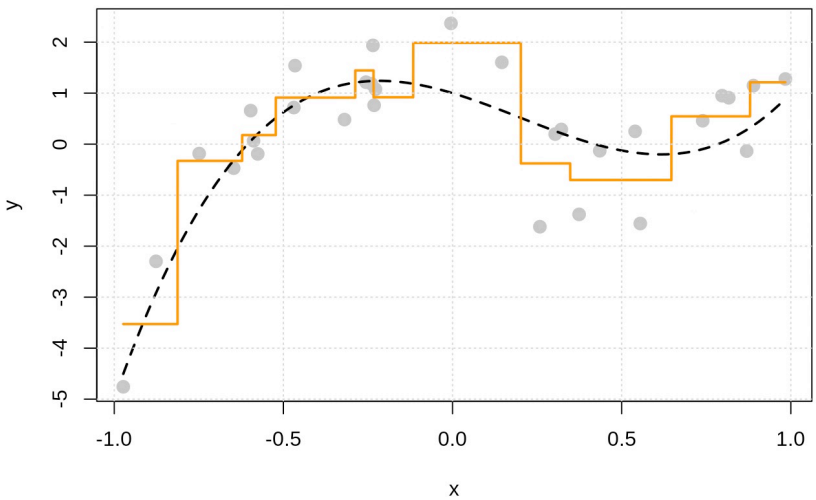
⋮

How "deep" (till) CAN THE TREE GROW?

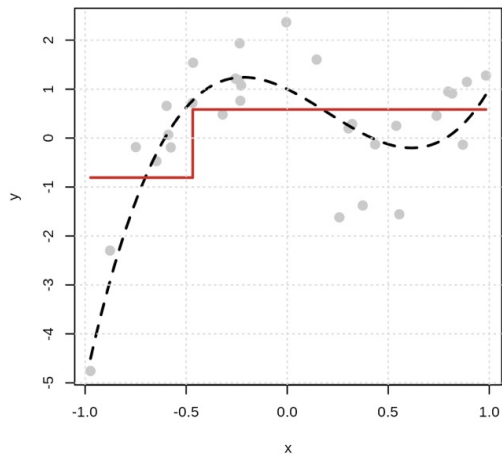
How BIG MUST A NODE BE TO CONSIDER A SPLIT?



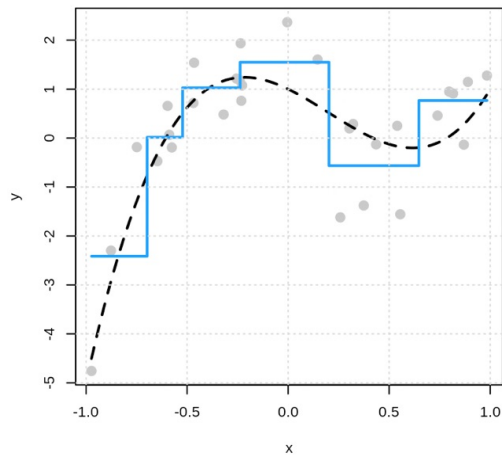




cp = 0.01, minsplit = 25



cp = 0.01, minsplit = 10



cp = 0.01, minsplit = 2

