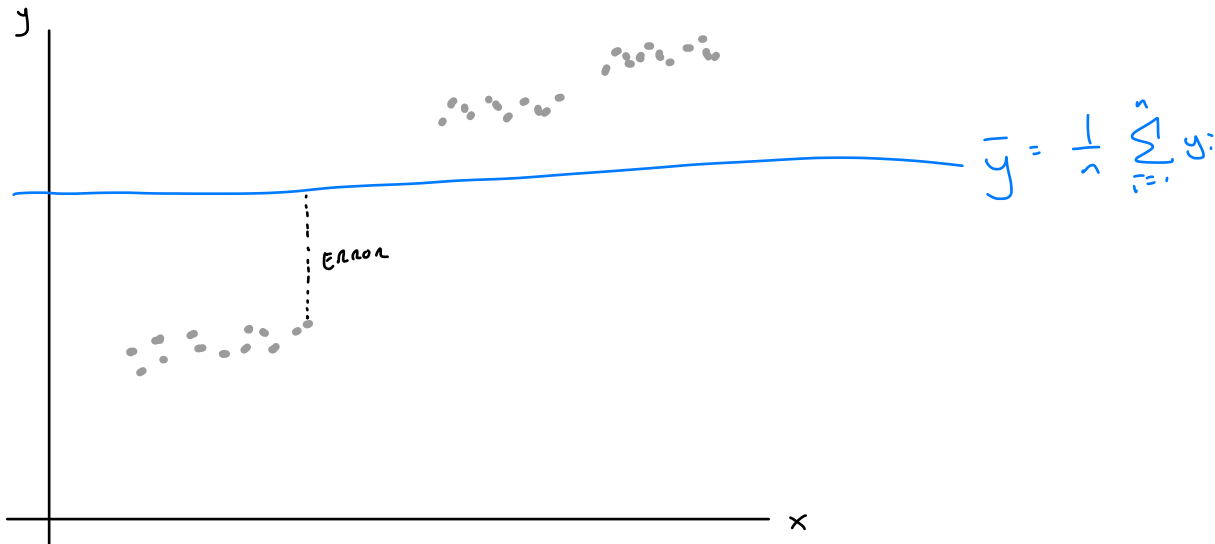


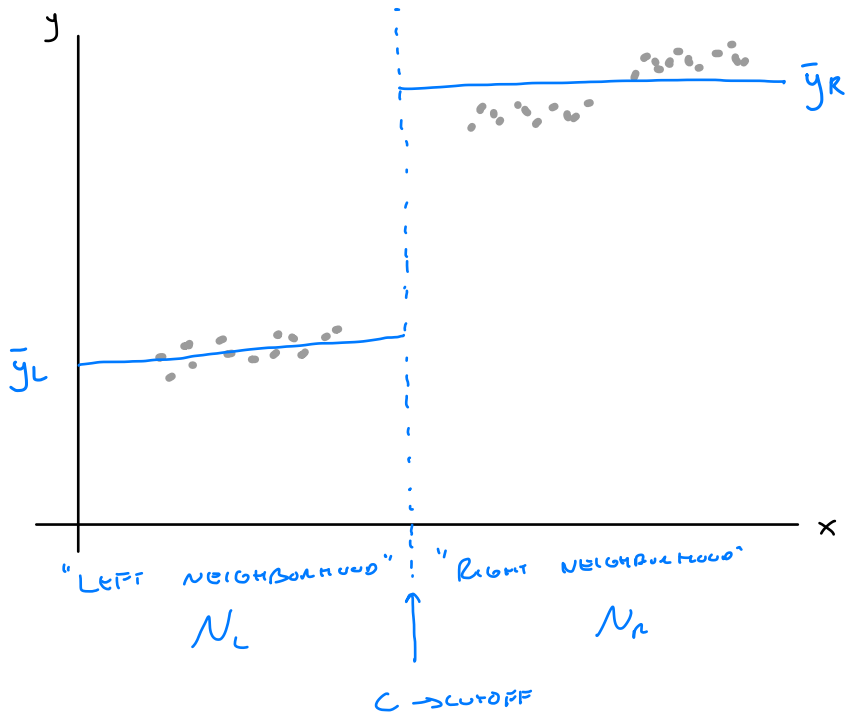
Decision Trees For Regression

CS 307



"OVERALL ERROR" → $SST = \sum_{i=1}^n (y_i - \bar{y})^2$
 ↳ "VARIANCE"

SUM SQUARES TOTAL



IDEA: CREATE NEIGHBORHOODS, PREDICT MEAN IN NEIGHBORHOOD

Decision Tree

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

$$SSE = \sum_{i \in N_L} (y_i - \bar{y}_L)^2 + \sum_{i \in N_R} (y_i - \bar{y}_R)^2$$

↑
error

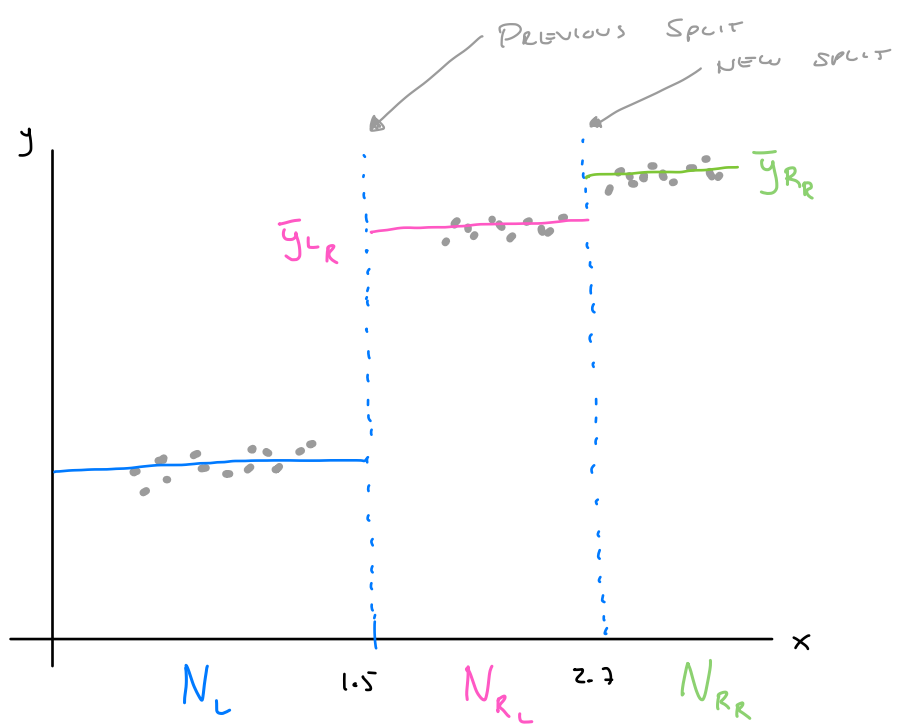
FIND c THAT MINIMIZES SSE

↳ CONSIDERING ALL FEATURES

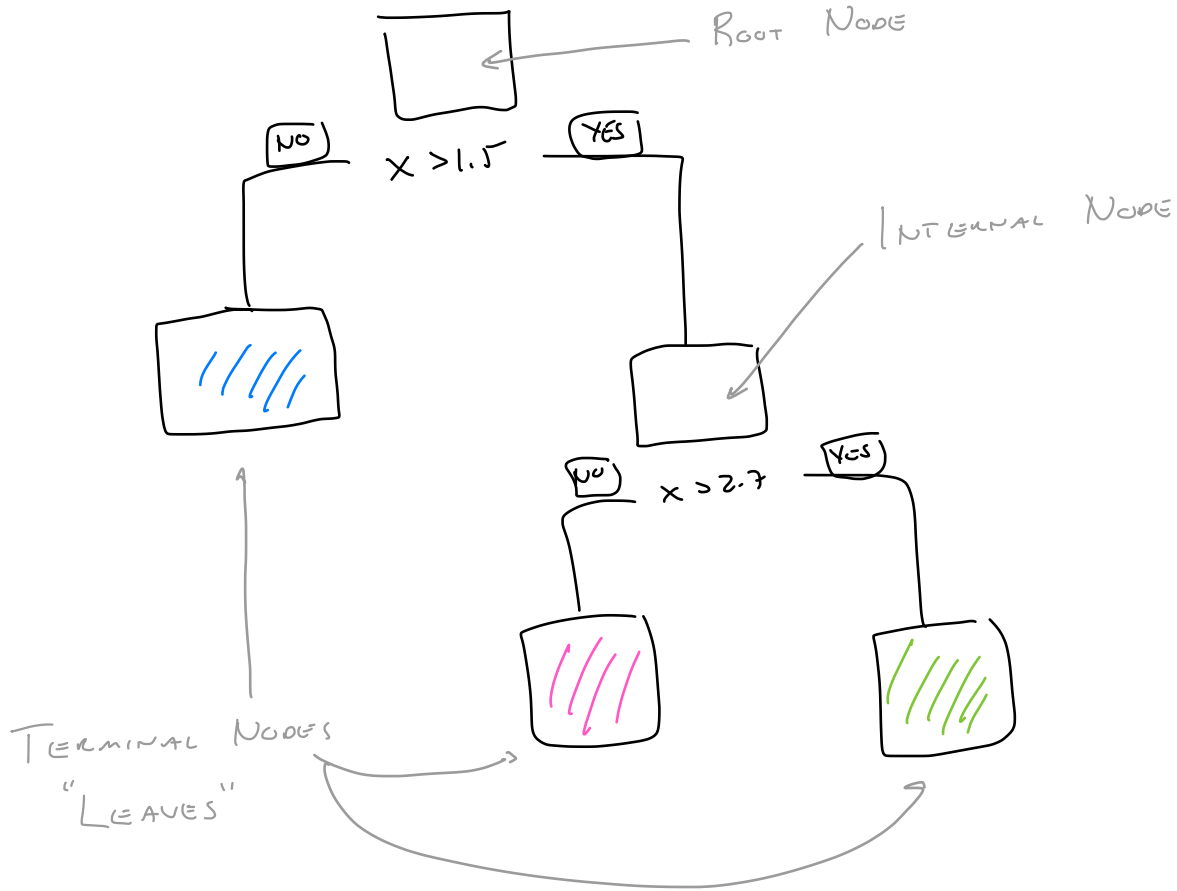
i	x	y
1	⋮	⋮
2	⋮	⋮
3	⋮	⋮
⋮	⋮	⋮
n	⋮	⋮

↑
IMPLIED

MEAN OF y IN N_L



$$SSE = \sum_{i \in N_L} (y_i - \bar{y}_L)^2 + \sum_{i \in N_{R_L}} (y_i - \bar{y}_{R_L})^2 + \sum_{i \in N_{R_R}} (y_i - \bar{y}_{R_R})^2$$



NOTE: THIS TREE IS "FLIPPED" COMPARED TO SKLEARN, WHICH GOES LEFT FOR "YES."

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

ONE SPLIT

$$SSE = \sum_{i \in N_L} (y_i - \bar{y}_L)^2 + \sum_{i \in N_R} (y_i - \bar{y}_R)^2$$

TWO SPLITS

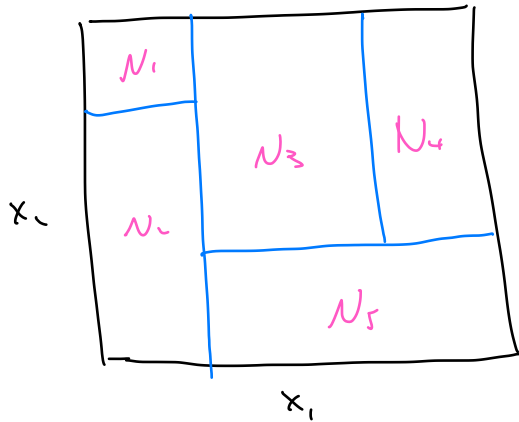
$$SSE = \sum_{i \in N_L} (y_i - \bar{y}_L)^2 + \sum_{i \in N_{R_1}} (y_i - \bar{y}_{R_1})^2 + \sum_{i \in N_{R_2}} (y_i - \bar{y}_{R_2})^2$$

NOTE :

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

RECURSIVE PARTITIONING

KEEP FINDING NEW SPLITS IN
CURRENT TERMINAL NODES



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

↑
MEAN OF POINTS
IN NODE j , N_j

WHEN TO STOP?

- DON'T ← BAD IDEA FOR A SINGLE TREE
 - MIN SAMPLES SPLIT
 - MAX DEPTH
- } TUNING PARAMETERS

SKETCH OF ALGORITHM

- INIT ROOT NODE, SET STOP CONDITIONS
- FOR EACH CURRENT TERMINAL NODE:
 - IF STOP CONDITIONS MET:
 - SKIP NODE
 - IF NO NODES CAN BE SPLIT:
 - STOP ALGORITHM, RETURN TREE
- FOR EACH FEATURE:
 - FOR EACH MIDPOINT c :
 - CALCULATE AND STORE ERROR
 - FIND j AND c COMBINATION WITH LOWEST ERROR
 - SPLIT CURRENT NODE AT USING $X_j < c$

