Nonparametric
Regression

ML

Supervised

Unsupervised

Regression

Classification

Parametric

Nonparametric

↳ K-Nearest Neighbors

↳ Decision Trees
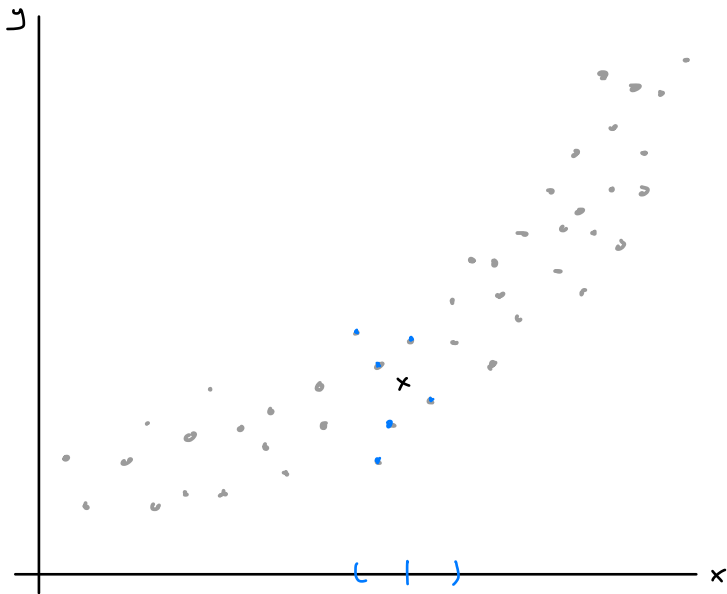
| x | y |
|---|---|
|   |   |

WANT

$$E\left[Y \mid X=x\right]$$

ASSUME

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2 + \varepsilon$$

ASSUME

$$Y = \beta_0 + \beta_1 x_1 + \varepsilon$$

$y$

WANT

$$\hat{\mathbb{E}}\left[Y \mid X = x\right]$$

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

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

# k — NEAREST NEIGHBORS

TO ESTIMATE $\mu(x) = \mathbb{E}\left[Y \mid X = x\right]$

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

$k$ OBSERVATIONS WITH $x_i$ NEAREST TO $x$
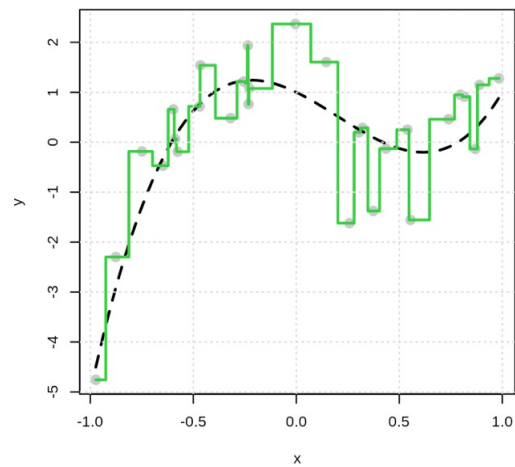
k = 3, x = -0.5    k = 5, x = 0    k = 9, x = 0.75

# Tuning Parameters

$$Y = B_0 + B_1 x_1 + B_2 x_2 + B_3 x_3 + \epsilon$$

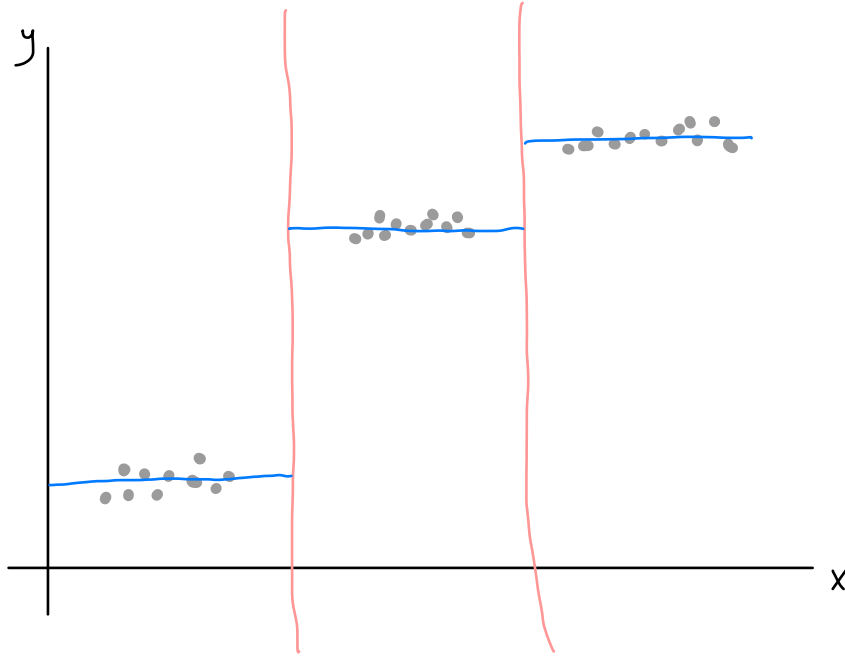MODEL PARAMETERS

↳ 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?    !!!

IDEA: FIND NEIGHBORHOODS, PREDICT AVERAGE OF $y_i$ IN NEIGHBORHOODS

# DECISION TREES

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

FEATURE + CUTOFF

FIND "SPLIT" THAT
MINIMIZES

AVG $y_i$ IN $N_L$

AVG $y_i$ IN $N_R$

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

$i \in N_L$

$x < c$

$i \in N_R$

$x > c$

# Recursive Partitioning

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

$$\sum_{i \in N_{R_1}} \left( y_i - \hat{\mu}_{N_{R_1}} \right)^2 \; + \; \sum_{i \in N_{R_2}} \left( y_i - \hat{\mu}_{N_{R_2}} \right)^2$$

# Recursive Partitioning



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

\# NEIGHBORHOODS

AVE $y_i$ IN $N_j$

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

# How to stop?

rpart :: rpart in R

**minsplit**    ONLY CONSIDER SPLIT IN NEIGHBORHOOD
IF IT HAS AT LEAST THIS MANY OBSERVATIONS
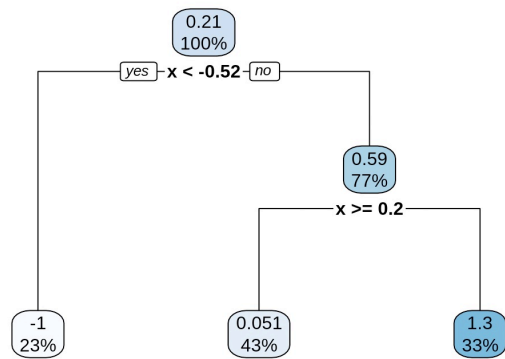
MINSPLIT = 2 ⟹ CAN ALWAYS SPLIT

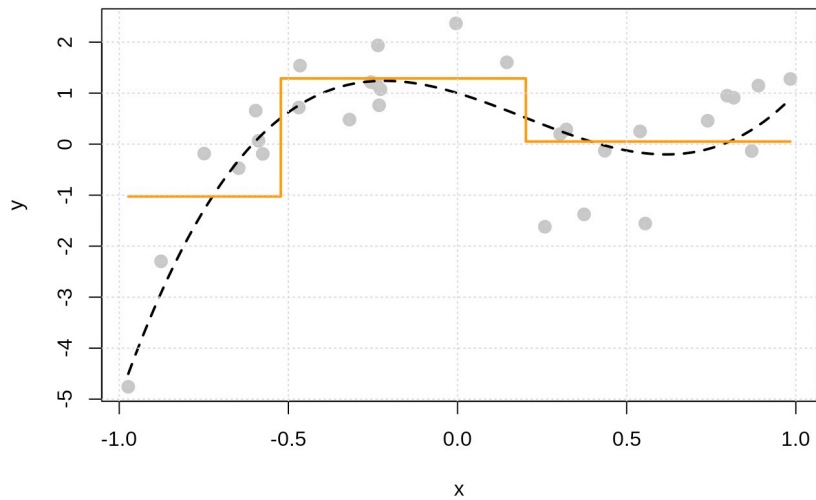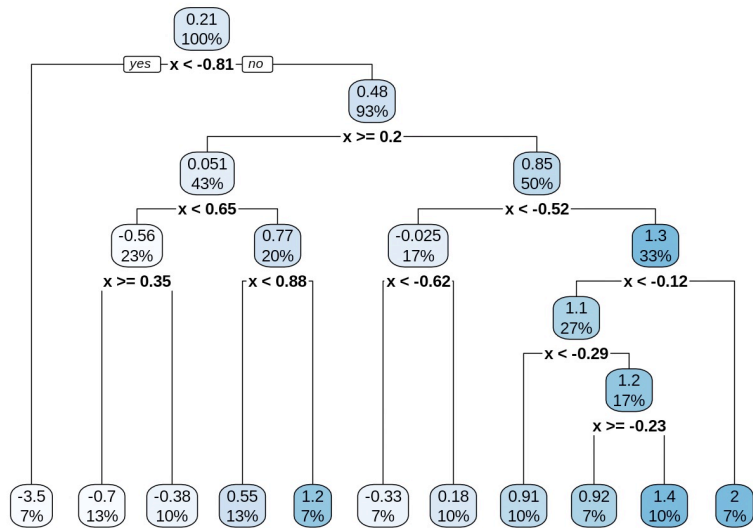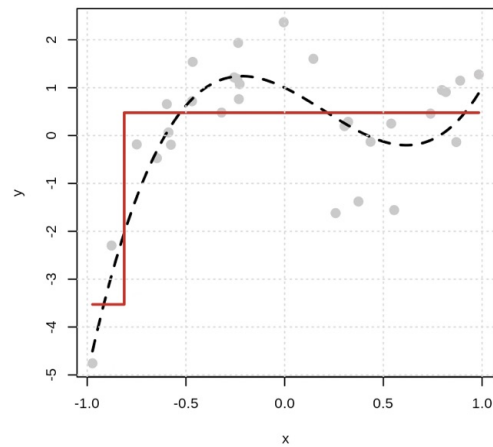**CP**    ONLY ACCEPT A SPLIT IF IT INCREASES
$R^2$ BY THIS AMOUNT OR MORE

CP = 0 ⟹ ANY SPLIT WILL BE
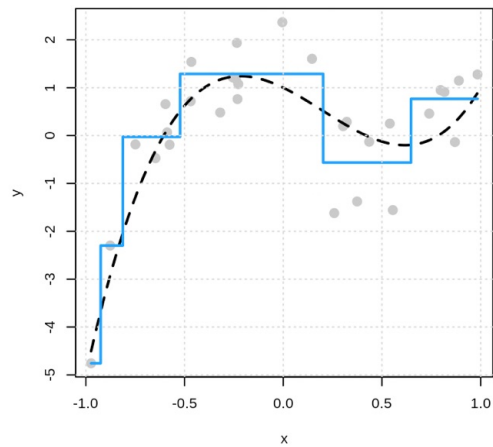ACCEPTED

↑

"COMPLEXITY PARAMETER"

**cp = 0.10, minsplit = 2**  **cp = 0.05, minsplit = 2**  **cp = 0.00, minsplit = 2**