# Regression Trees
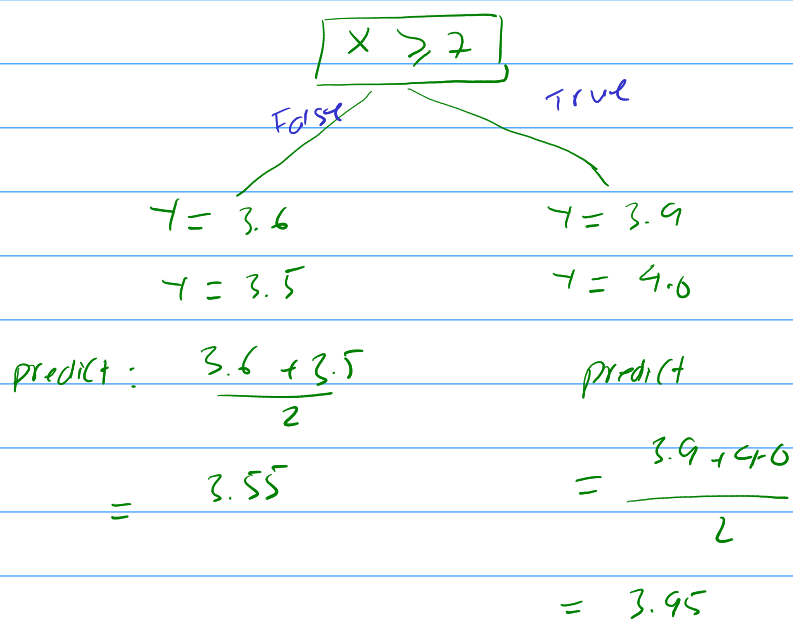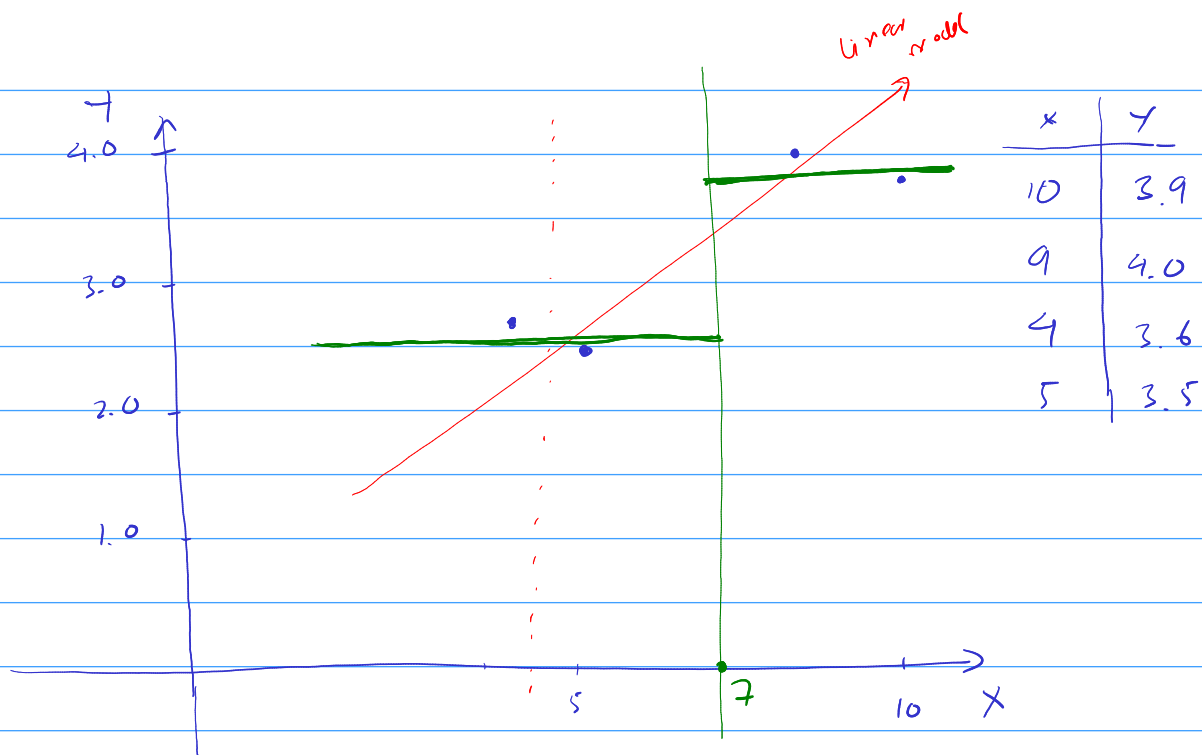
# Regression Trees

# Regression Trees

▶ The tree will search for all combination of predictors and cutoff value to decide the best split

▶ In Regression tree, the best split is the split that minimizes

$$\underbrace{\sum_{i:\mathbf{x}_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2}_{\text{RSS of obs. in left branch}} \quad + \quad \underbrace{\sum_{i:\mathbf{x}_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2}_{\text{RSS of obs. in right branch}}$$

▶ $\hat{y}_{R_1}$ and $\hat{y}_{R_2}$ are the means of the responses falling in to the left branch and right branch, respectively.
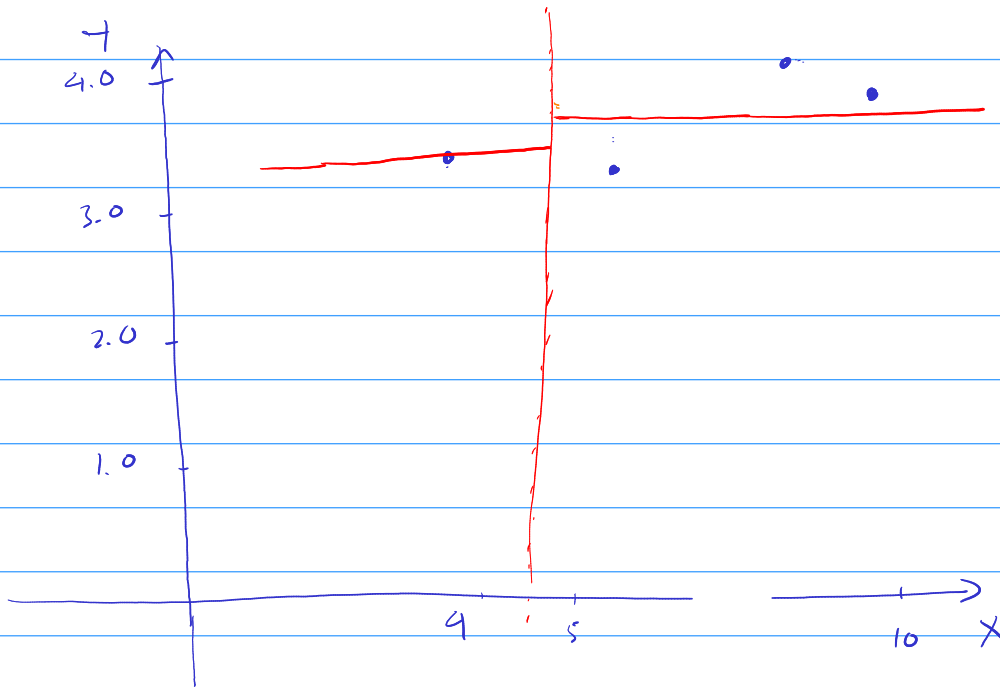
| x | Y |
|---|---|
| 10 | 3.9 |
| 9 | 4.0 |
| 4 | 3.6 |
| 5 | 3.5 |

linear model

$$\boxed{x \geq 7}$$

False                                    True

$Y = 3.6$                                $Y = 3.9$

$Y = 3.5$                                $Y = 4.0$

predict: $\dfrac{3.6 + 3.5}{2}$          predict

$= 3.55$                                 $= \dfrac{3.9 + 4.0}{2}$

$= 3.95$

Sum Square Errors of this Regression tree.

| x | Y | $\hat{Y}$ (tree predicts) | Squared errors |
|---|---|---|---|
| 10 | 3.9 | 3.95 | $(3.95 - 3.9)^2$ |
| 9 | 4.0 | 3.95 | $(3.95 - 4)^2$ |
| 4 | 3.6 | 3.55 | $(3.55 - 3.6)^2$ + |
| 5 | 3.5 | 3.55 | $(3.55 - 3.5)^2$ |

$$SSE = .01$$

$\boxed{x \geq 4.5}$

True

$Y = 3.6$

$-1 = 3.9, 4.0, 3.5 \Rightarrow predict = \dfrac{3.9 + 4.0 + 3.5}{3} = 3.8$



| x | Y | $\hat{Y}$ |
|---|---|---|
| 10 | 3.9 | 3.8 |
| 9 | 4.0 | 3.8 |
| 4 | 3.6 | 3.6 |
| 5 | 3.5 | 3.8 |

$SSE = (3.9 - 3.8)^2 + (4.0 - 3.8)^2 + (3.6 - 3.6)^2 + (3.5 - 3.8)^2$

$= .14$



$x > 7$

True

$x > 4.5$      $x > 9.5$

$Y = 3.6$    $Y = 3.5$    $Y = 4.0$    $Y = 3.5$

$SSE = 0$

# Example

| $X_1$ | $X_2$ | $Y$ |
|:-----:|:-----:|:---:|
| 1 | 0 | 1.2 |
| 2 | 1 | 2.1 |
| 3 | 2 | 1.5 |
| 4 | 1 | 3.0 |
| 2 | 2 | 2.0 |
| 1 | 1 | 1.6 |

Using the RSS to decide the best split among

▶ Split 1: Region 1 $X_1 < 4$, Region 2 $X_1 \geq 4$
▶ Split 2: Region 1 $X_2 < 2$, Region 2 $X_2 \geq 2$

# Example

$RSS = SSE$

| | $X_1$ | $X_2$ | $Y$ |
|---|---|---|---|
| A | 1 | 0 | 1.2 |
| B | 2 | 1 | 2.1 |
| C | 3 | 2 | 1.5 |
| D | 4 | 1 | 3.0 |
| E | 2 | 2 | 2.0 |
| F | 1 | 1 | 1.6 |

Using the RSS to decide the best split among

▶ Split 1: Region 1 $X_1 < 4$, Region 2 $X_1 \geq 4$
▶ Split 2: Region 1 $X_2 < 2$, Region 2 $X_2 \geq 2$

**Split 1 :**

$$X_1 < 4$$

| | $X_1$ | $X_2$ | $Y$ |
|---|---|---|---|
| A | 1 | 0 | 1.2 |
| B | 2 | 1 | 2.1 |
| C | 3 | 2 | 1.5 |
| D | 4 | 1 | 3.0 |
| E | 2 | 2 | 2.0 |
| F | 1 | 1 | 1.6 |

True → A, B, C, E F

D

$Y = 3.0$

$\hat{Y} = 3.0$

$Y = \{1.2, 2.1, 1.5, 2.0, 1.6\}$

$$\hat{Y} = \frac{1.2 + 2.1 + 1.5 + 2.0 + 1.6}{5} = 1.68$$

$$RSS = SSE = \underbrace{\sum (Y - \hat{Y})^2}_{branch\,1} + \underbrace{\sum (Y - \hat{Y})^2}_{branch\,2}$$

$$= (3 - 3.0)^2 + (1.2 - 1.68)^2 + (2.1 - 1.68)^2 + (1.5 - 1.68)^2 + (2.0 - 1.68)^2 + (1.6 - 1.68)^2$$

$$= .548$$

**Split 2 :**

$$X_2 < 2$$

| | $X_1$ | $X_2$ | $Y$ |
|---|---|---|---|
| A | 1 | 0 | 1.2 |
| B | 2 | 1 | 2.1 |
| C | 3 | 2 | 1.5 |
| D | 4 | 1 | 3.0 |
| E | 2 | 2 | 2.0 |
| F | 1 | 1 | 1.6 |

true

C, E                     A, B, D, F

$Y = 1.5, 2.0$           $Y = 1.2, 2.1, 3.0, 1.6$

$$\hat{Y} = \frac{1.5 + 2.0}{2}$$          $$\hat{Y} = \frac{1.2 + 2.1 + 3.0 + 1.6}{4}$$

$$= 1.75$$                              $$= 1.975$$

$$SSE = \underbrace{(1.5 - 1.75)^2 + (2.0 - 1.75)^2}_{left\ branch} + \underbrace{(1.2 - 1.975)^2 + (2.1 - 1.975)^2 + (3.0 - 1.975)^2 + (1.6 - 1.975)^2}_{right\ branch}$$
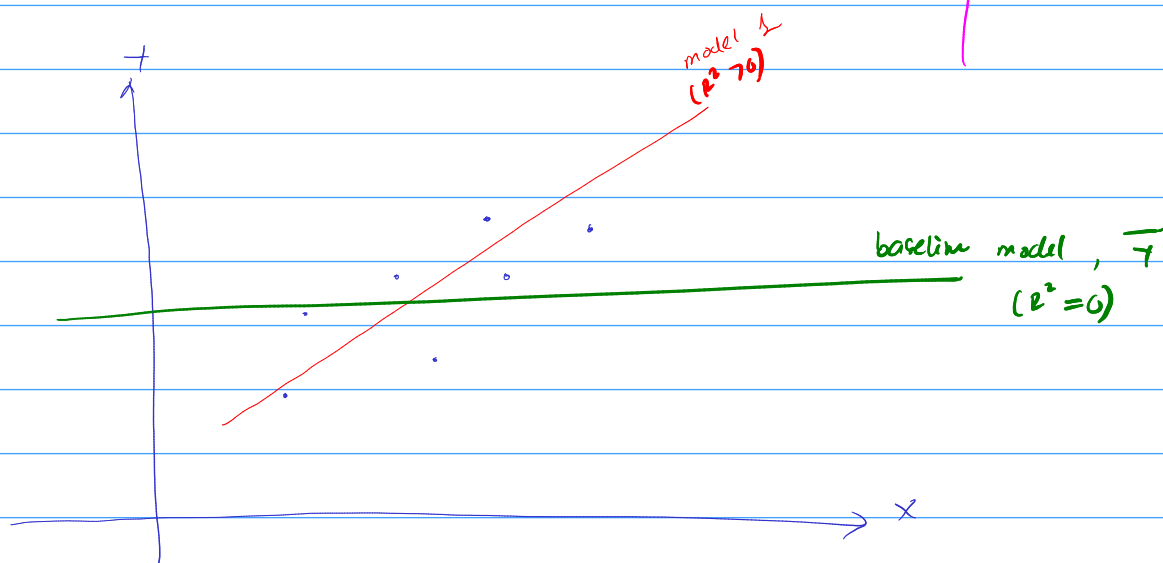
$$\boxed{SSE = 1.9325}$$

# Example

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 1 | 0 | 1.2 |
| 2 | 1 | 2.1 |
| 3 | 2 | 1.5 |
| 4 | 1 | 3.0 |
| 2 | 2 | 2.0 |
| 1 | 1 | 1.6 |

Using the RSS to decide the best split among

▶ Split 1: Region 1 $X_1 < 4$, Region 2 $X_1 \geq 4$
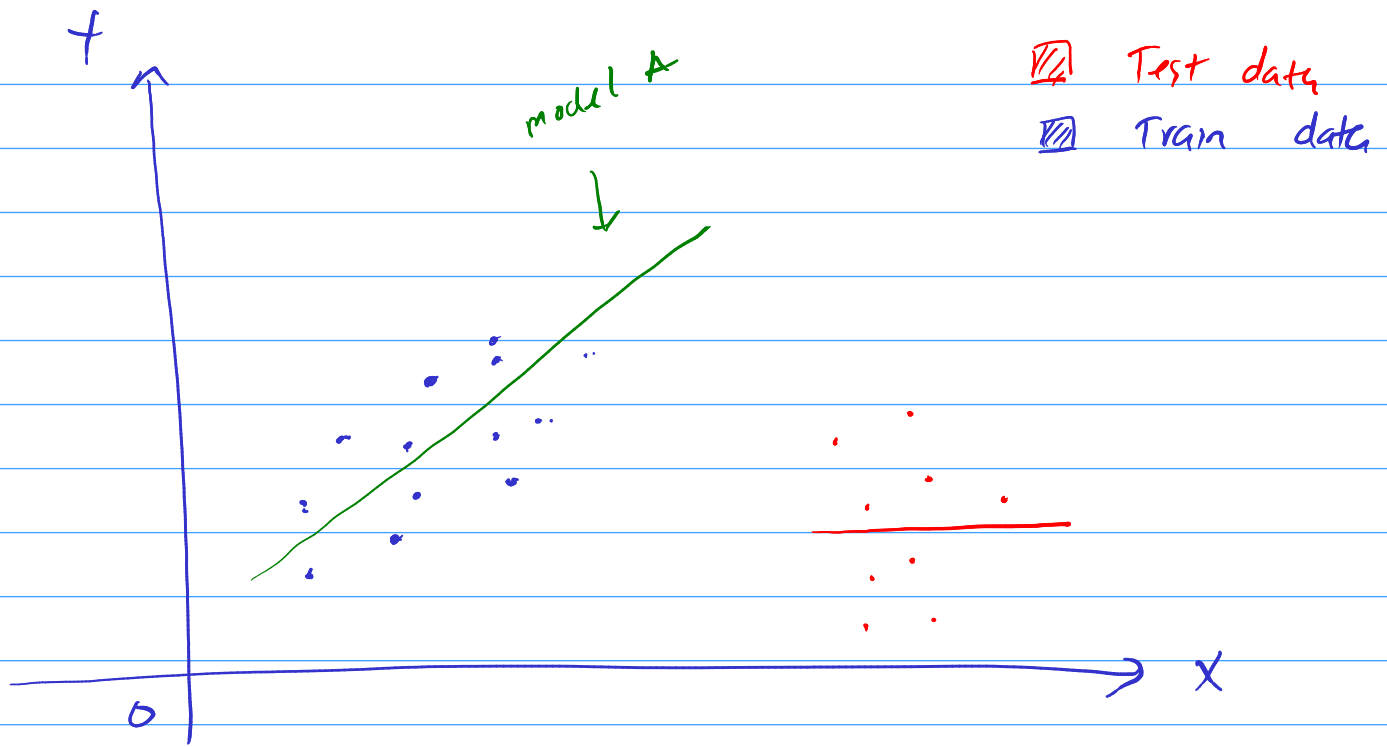▶ Split 2: Region 1 $X_2 < 2$, Region 2 $X_2 \geq 2$

Ⓐ $R^2$

model 2 $(R^2 < 0)$
↓

model 1
$(R^2 > 0)$

baseline model, $\overline{7}$
$(R^2 = 0)$

$t$

$x$

$$R^2 = 1 - \frac{\text{SSE of the model}}{\text{SSE of the baseline model, } \overline{7}}$$

① when the model A is the baseline model, $R^2$ of A $= 0$

② when the model A is "perfect", $R^2 = 1 - \dfrac{0}{\text{SSE of baseline}}$

$(\Rightarrow)$ $\boxed{R^2 = 1}$

③ when the model is "worse" than the baseline model, $R^2 < 0$

model A

Test data
Train data

$R^2$ of A on training is positive but on testing is negative

(Back to the example)

Split 1:

$X_1 < 4$

True

$A, B, C, E F$

$Y = \{ 1.2, 2.1, 1.5, 2.0, 1.6 \}$

$\hat{Y} = \dfrac{1.2 + 2.1 + 1.5 + 2.0 + 1.6}{5} = 1.68$

$Y = 3.0$

$\hat{Y} = 3.0$

| | $X_1$ | $X_2$ | $Y$ |
|---|---|---|---|
| A | 1 | 0 | 1.2 |
| B | 2 | 1 | 2.1 |
| C | 3 | 2 | 1.5 |
| D | 4 | 1 | 3.0 |
| E | 2 | 2 | 2.0 |
| F | 1 | 1 | 1.6 |

Let calculate the $R^2$ of this model / split

SSE of the baseline model $= \sum (Y - \bar{Y})^2$

$$\bar{Y} = \frac{1.2 + 2.1 + 1.5 + 3.0 + 2.0 + 1.6}{6} = 1.9$$

$$SSE \text{ of } \bar{Y} = (1.2 - 1.9)^2 + (2.1 - 1.9)^2$$
$$+ (1.5 - 1.9)^2 + (3.0 - 1.9)^2$$
$$+ (2 - 1.9)^2 + (1.6 - 1.9)^2$$
$$= 2$$

$$R^2 = 1 - \frac{.548}{2} = .726.$$