# 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

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

▶ $\hat{y}_{R_1}$ and $\hat{y}_{R_2}$ are the means of the responses falling in to the left branch and right branch, respectively.

# Example

*continuous*

| | $X_1$ | $X_2$ | $Y$ |
|---|---|---|---|
| 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 $X_2 \geq 2$

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

Split 1

$X_1 < 4$

no → $Y = 3.0$ → $\bar{Y}_2 = 3.0$

→ $RSS_2 = 0$

yes ↓

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

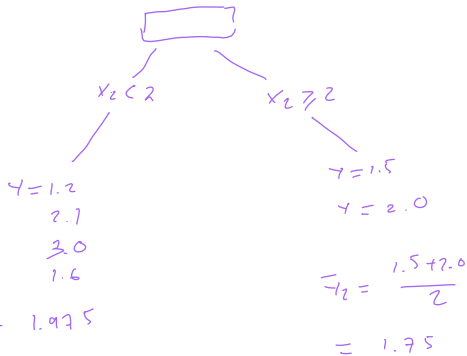$$\bar{Y}_1 = \frac{1.2 + 2.1 + 1.5 + 2.0 + 1.6}{5} = 1.68$$

$$RSS_1 = (1.2 - \bar{Y}_1)^2 + (2.1 - \bar{Y}_1)^2 + (1.5 - \bar{Y}_1)^2$$
$$+ (2.0 - \bar{Y}_1)^2 + (1.6 - \bar{Y}_1)^2$$

$$= 1.548$$

$$RSS = RSS_1 + RSS_2$$
$$= .548$$

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

split 2



$X_2 < 2$                     $X_2 \geq 2$

$Y = 1.2$                      $Y = 1.5$
$\quad\; 2.1$                   $Y = 2.0$
$\quad\; 3.0$
$\quad\; 1.6$                   $\bar{Y_2} = \dfrac{1.5 + 2.0}{2}$

$\bar{Y_1} = 1.975$                        $= 1.75$

$RSS_1 = \sum (Y_j - \bar{Y_1})^2 = 1.8075$

$RSS_2 =$

$(1.5 - 1.75)^2 + (2 - 1.75)^2$

$= 0.125$

Total $RSS = RSS_1 + RSS_2$

$= 1.9325 >$ RSS of split 1 $\Rightarrow$ Split 1 is better.

Split 1

$Y_1 < 4$

yes ← → no

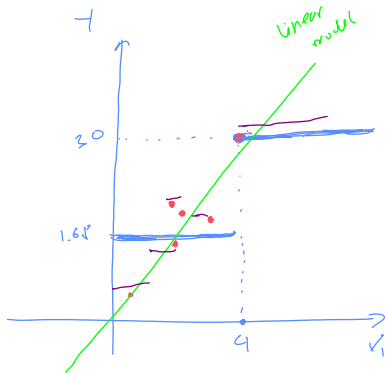$Y = 3.0 \rightarrow \bar{Y}_2 = 3.0$

predict by the mean

$\bar{Y}_2 = 3.0$

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

predict by the mean $\bar{Y}_1 = 1.68$



For example, if $x_1 = 0$, $x_2 = 100$, then this tree will predict $y = 1.68$

Clacsification: Mis. classification, ROC, Sensitivity....

Model

Evaluation

Regression:

$R^2 =$

$$RSS = \sum \left( t_i - \hat{t}_i \right)^2$$

$$\sum \left| t_i - \hat{t}_i \right|$$

$$MAE = \sum \frac{\left| t_i - \hat{t}_i \right|}{n}$$

( mean absolute

error )

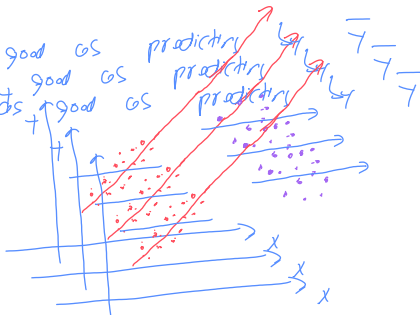| True value | predicted value |
|---|---|
| $t_1$ | $\hat{t}_1$ |
| $t_2$ | $\hat{t}_2$ |
| $\vdots$ | $\vdots$ |

# Linear model and regression model

$$R^2 = 1 - \frac{RSS}{Total\ SS} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

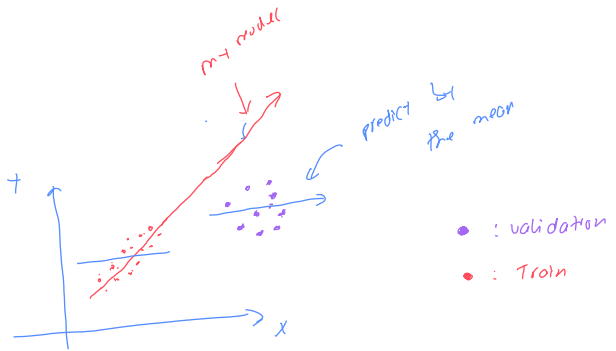(*)    $R^2 = 1 \iff RSS = 0$

(*)    $R^2 = 0 \iff RSS = TSS$

The model is as good as predicting $y_i$ by $\bar{y}$

The model is as good as predicting $y_i$ by $\bar{y}$

The model is as good as predicting $y_i$ by $\bar{y}$

(*)    $R^2 < 0$    $(?)$

my model

predict by the mean

$t$

$x$

• : validation

• : Train

(a) In training

$R^2 > 0$

(b) In validation :

$R^2 < 0$

classify        0        1        :        2%

predict 1    →    0

regression
              2.5

$R^2 = -100$