

# Multiple Linear Regression

# Linear Regression is an Example of a Linear Model

Input instance – feature vector:  $\mathbf{x} = (x_0, x_1, \dots, x_n)$

Predicted  
output:

$$\hat{y} = \hat{w}_0 x_0 + \hat{w}_1 x_1 + \dots + \hat{w}_n x_n + \hat{b}$$

Parameters  
to estimate:

$\hat{\mathbf{w}} = (\hat{w}_0, \dots, \hat{w}_n)$ : feature weights/  
model coefficients  
 $\hat{\mathbf{b}}$ : constant bias term / intercept

# Linear Models

- A linear model is a sum of weighted variables that predicts a target output value given an input data instance. Example: *predicting housing prices*

- House features: taxes per year ( $X_{TAX}$ ), age in years ( $X_{AGE}$ )

$$\widehat{Y_{PRICE}} = 212000 + 109 X_{TAX} - 2000 X_{AGE}$$

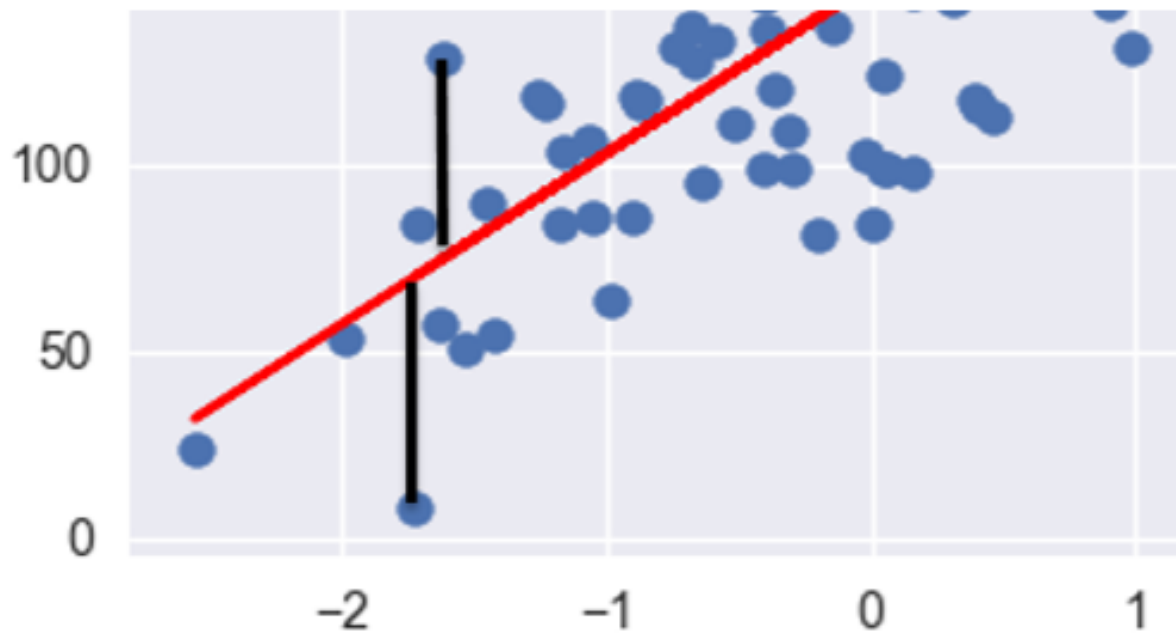
- A house with feature values ( $X_{TAX}, X_{AGE}$ ) of (10000, 75) would have a predicted selling price of:

$$\widehat{Y_{PRICE}} = 212000 + 109 \cdot \mathbf{10000} - 2000 \cdot \mathbf{75} = \mathbf{1,152,000}$$

# Least-Squares Linear Regression

## ("Ordinary least-squares")

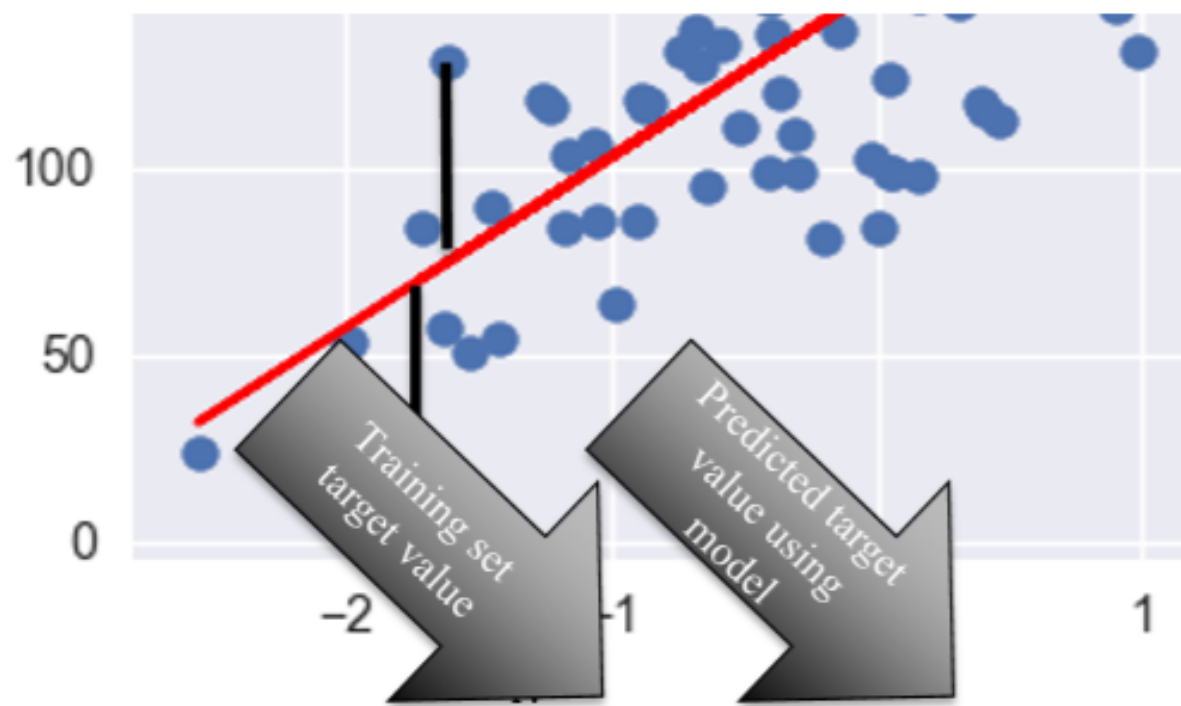
- Finds  $w$  and  $b$  that minimizes the sum of squared differences (RSS) over the training data between predicted target and actual target values.
- a.k.a. mean squared error of the linear model
- No parameters to control model complexity.



$$RSS(\mathbf{w}, b) = \sum_{i=1}^N (y_i - (\mathbf{w} \cdot \mathbf{x}_i + b))^2$$

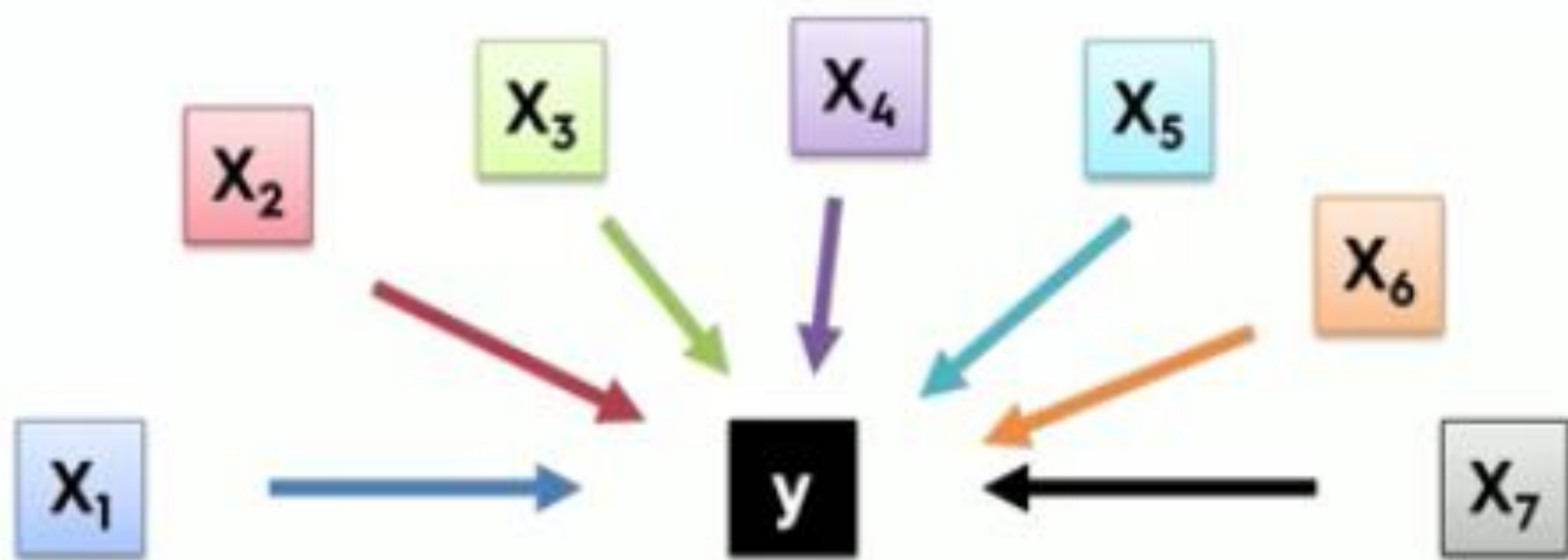
# Least-Squares Linear Regression ("Ordinary Least-Squares")

- Finds  $w$  and  $b$  that minimizes the sum of squared differences (RSS) over the training data between predicted target and actual target values.
- a.k.a. mean squared error of the linear model
- No parameters to control model complexity.

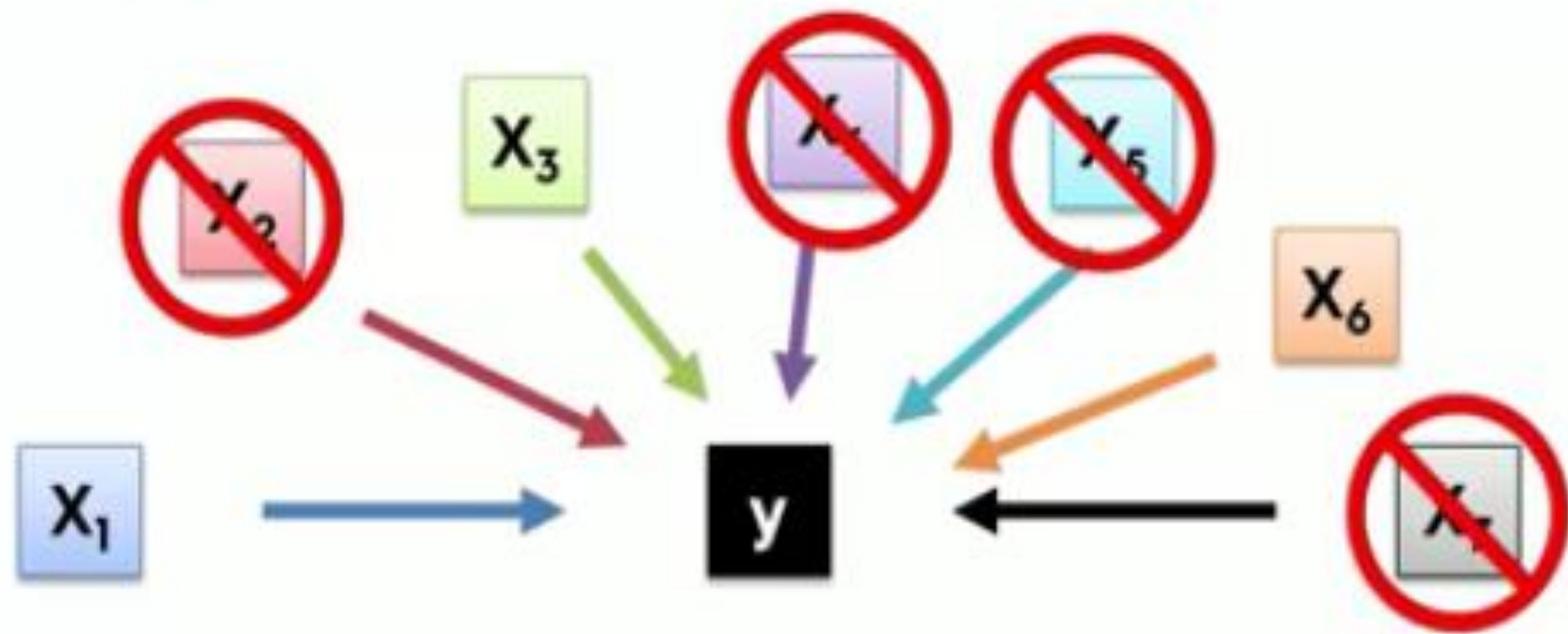


$$RSS(\mathbf{w}, b) = \sum_{\{i=1\}} (y_i - (\mathbf{w} \cdot \mathbf{x}_i + b))^2$$

# Building A Model



# Building A Model



# P-Value



$H_0$ : This is a fair coin

$H_1$ : This is not a fair coin



0.5



0.25



0.12



0.06



0.03



0.01





$H_0$ : This is a fair coin

$H_1$ : This is not a fair coin



0.5



0.25



0.12



0.06



0.03



0.01

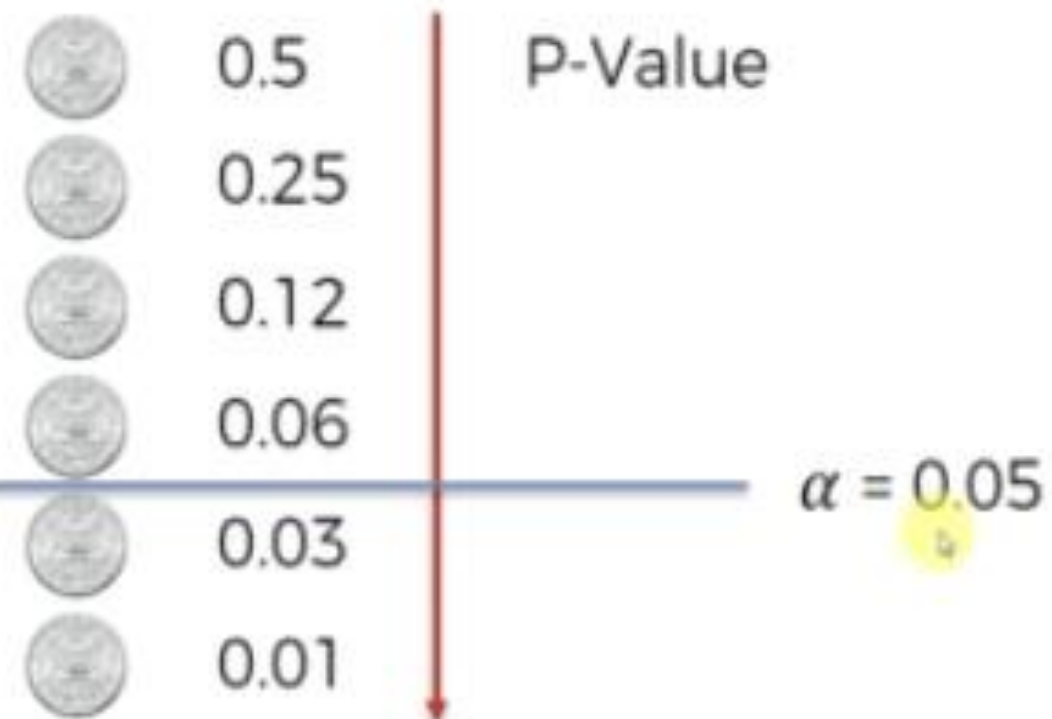
P-Value





$H_0$ : This is a fair coin

$H_1$ : This is not a fair coin



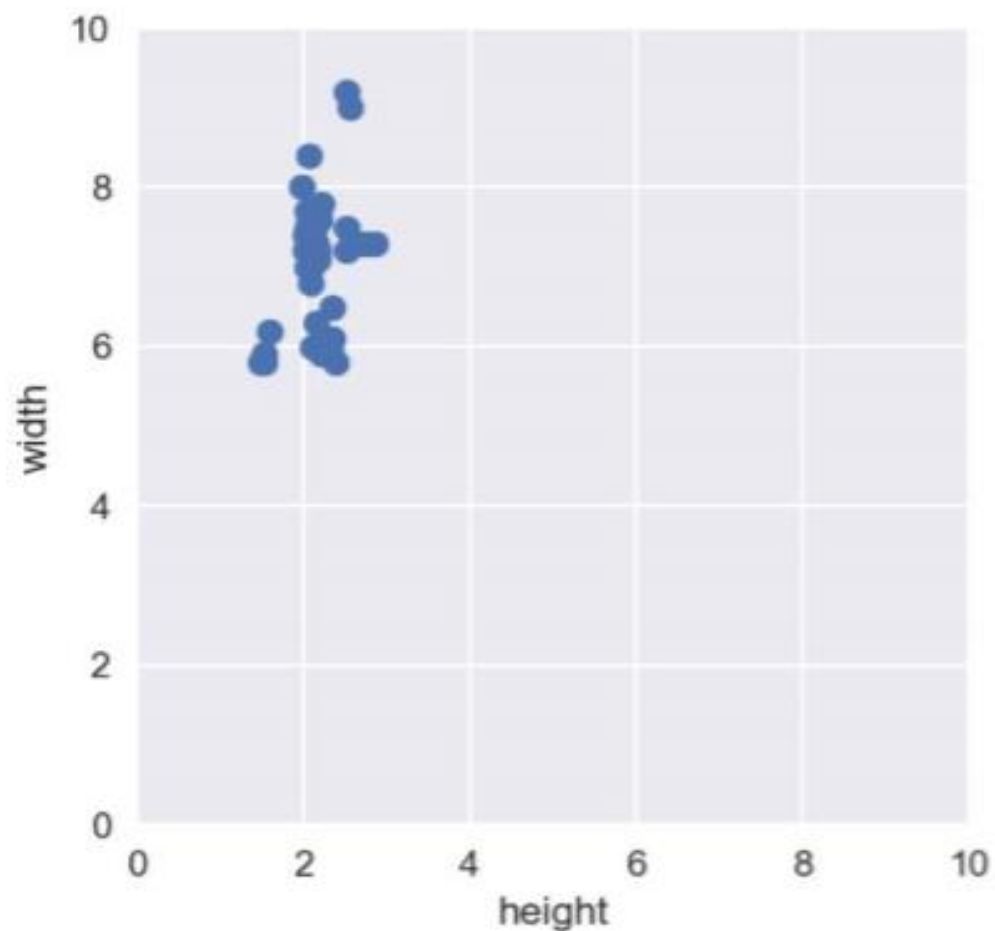
# Ridge Regression

- Ridge regression learns  $w$ ,  $b$  using the same least-squares criterion but adds a penalty for large variations in  $w$  parameters

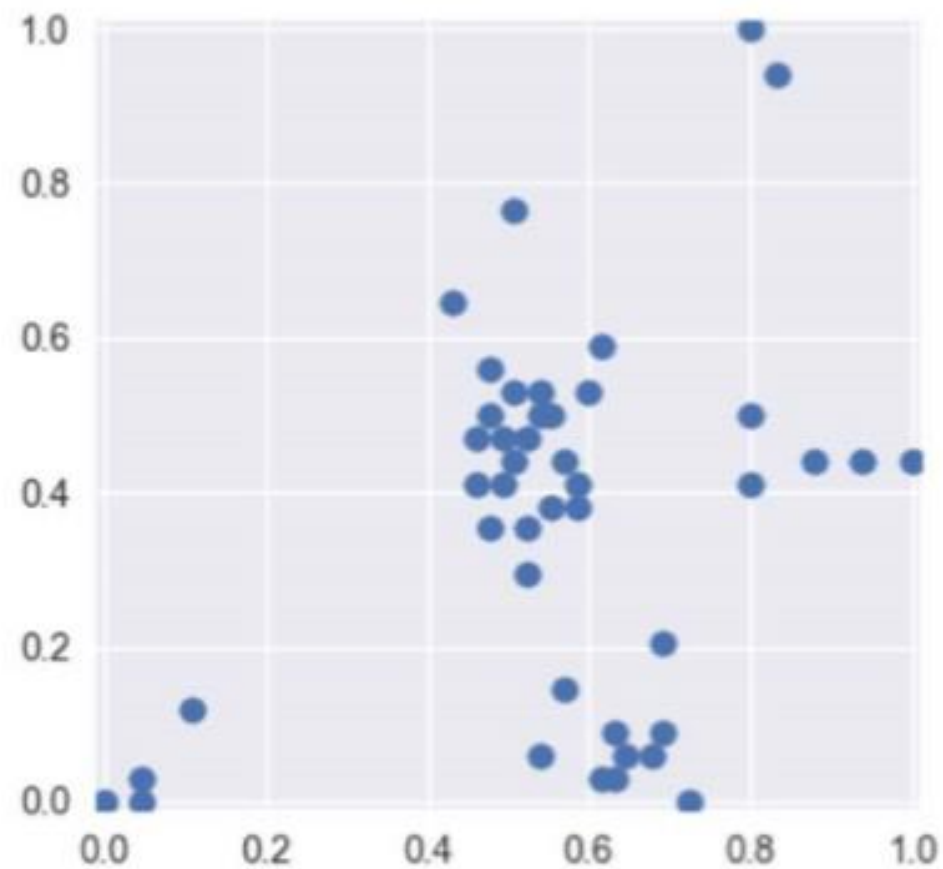
$$RSS_{RIDGE}(w, b) = \sum_{\{i=1\}}^N (y_i - (w \cdot x_i + b))^2 + \alpha \sum_{\{j=1\}}^p w_j^2$$

- Once the parameters are learned, the ridge regression prediction formula is the same as ordinary least-squares.
- The addition of a parameter penalty is called regularization. Regularization prevents overfitting by restricting the model, typically to reduce its complexity.
- Ridge regression uses L2 regularization: minimize sum of squares of  $w$  entries
- The influence of the regularization term is controlled by the  $\alpha$  parameter.
- Higher alpha means more regularization and simpler models.

# Feature Normalization with MinMaxScaler



Unnormalized data points



Normalized with MinMaxScaler

## Lasso regression is another form of regularized linear regression that uses an L1 regularization penalty for training (instead of ridge's L2 penalty)

- **L1 penalty:** Minimize the sum of the absolute values of the coefficients

$$RSS_{LASSO}(w, b) = \sum_{\{i=1\}}^N (y_i - (w \cdot x_i + b))^2 + \alpha \sum_{\{j=1\}}^p |w_j|$$

- This has the effect of setting parameter weights in  $w$  to zero for the least influential variables. This is called a sparse solution: a kind of feature selection
- The parameter  $\alpha$  controls amount of L1 regularization (default = 1.0).
- The prediction formula is the same as ordinary least-squares.
- **When to use ridge vs lasso regression:**
  - *Many small/medium sized effects: use ridge.*
  - *Only a few variables with medium/large effect: use lasso.*