

# Linear Regression

Supervised Learning

# Simple Linear Regression



~



# Simple Linear Regression

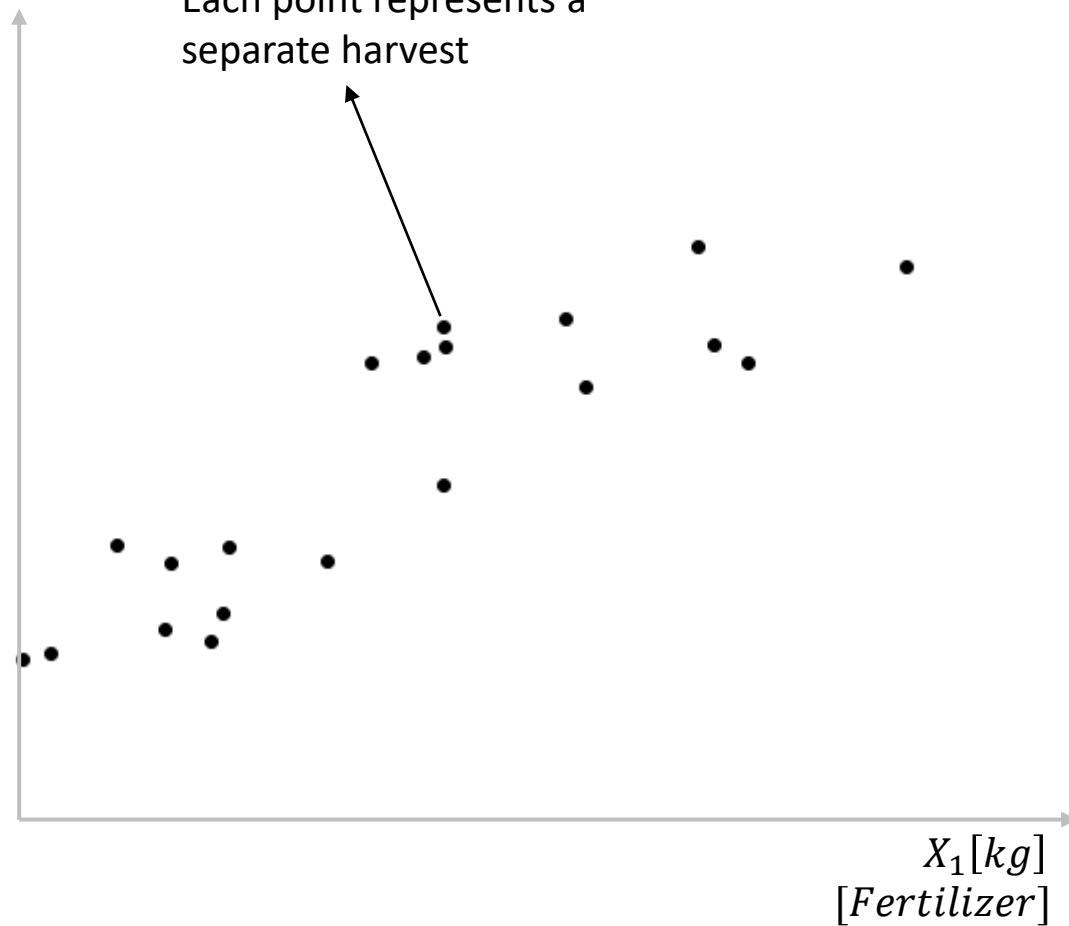


~



$y[lb]$   
[Wheat yield]

Each point represents a  
separate harvest



# Simple Linear Regression

$$\hat{y} = b_0 + b_1 X_1$$



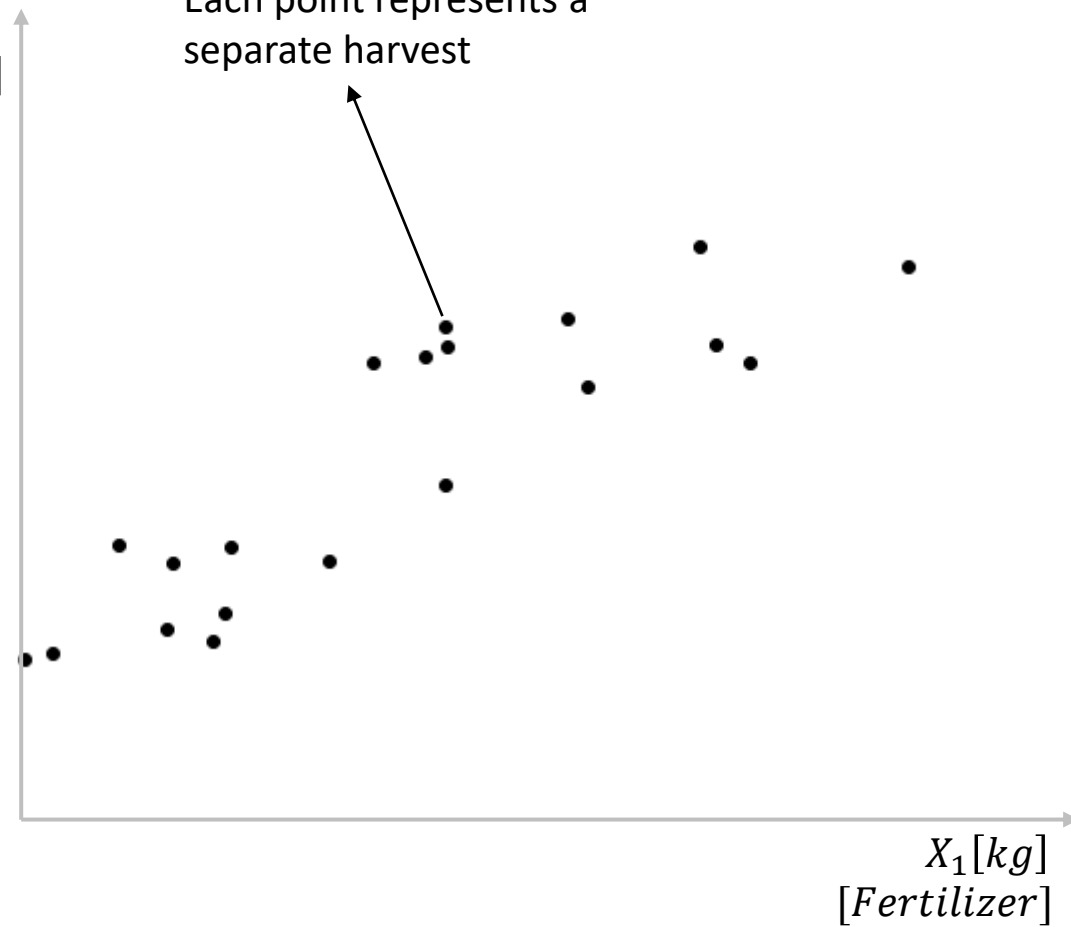
~



$$wheat[lb] = b_0 + b_1 * Fertilizer[kg]$$

$y[lb]$   
[Wheat yield]

Each point represents a  
separate harvest



# Simple Linear Regression

$$\hat{y} = b_0 + b_1 X_1$$



~



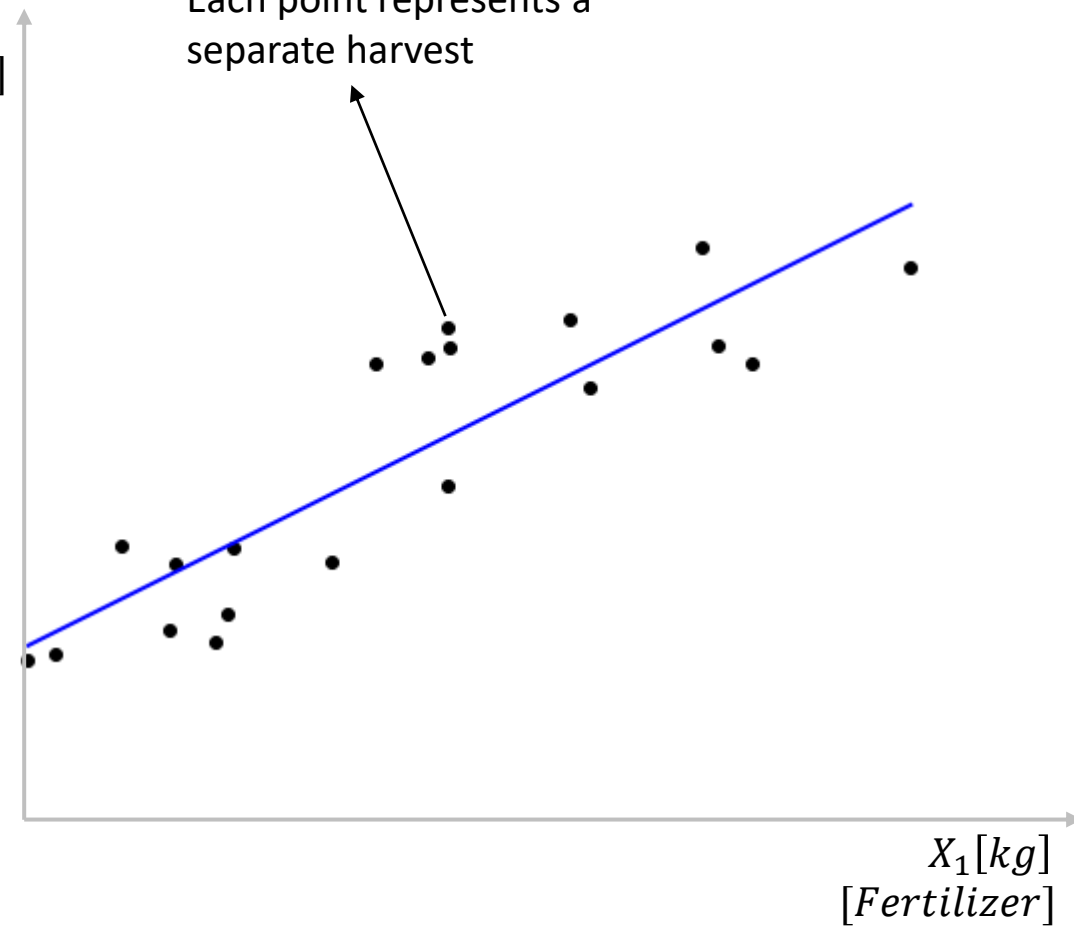
$$wheat[lb] = b_0 + b_1 * Fertilizer[kg]$$

$$b_0 = 10[lb]$$

$$b_1 = 4\left[\frac{lb}{kg}\right]$$

$y[lb]$   
[Wheat yield]

Each point represents a  
separate harvest



# Simple Linear Regression

$$\hat{y} = b_0 + b_1 X_1$$



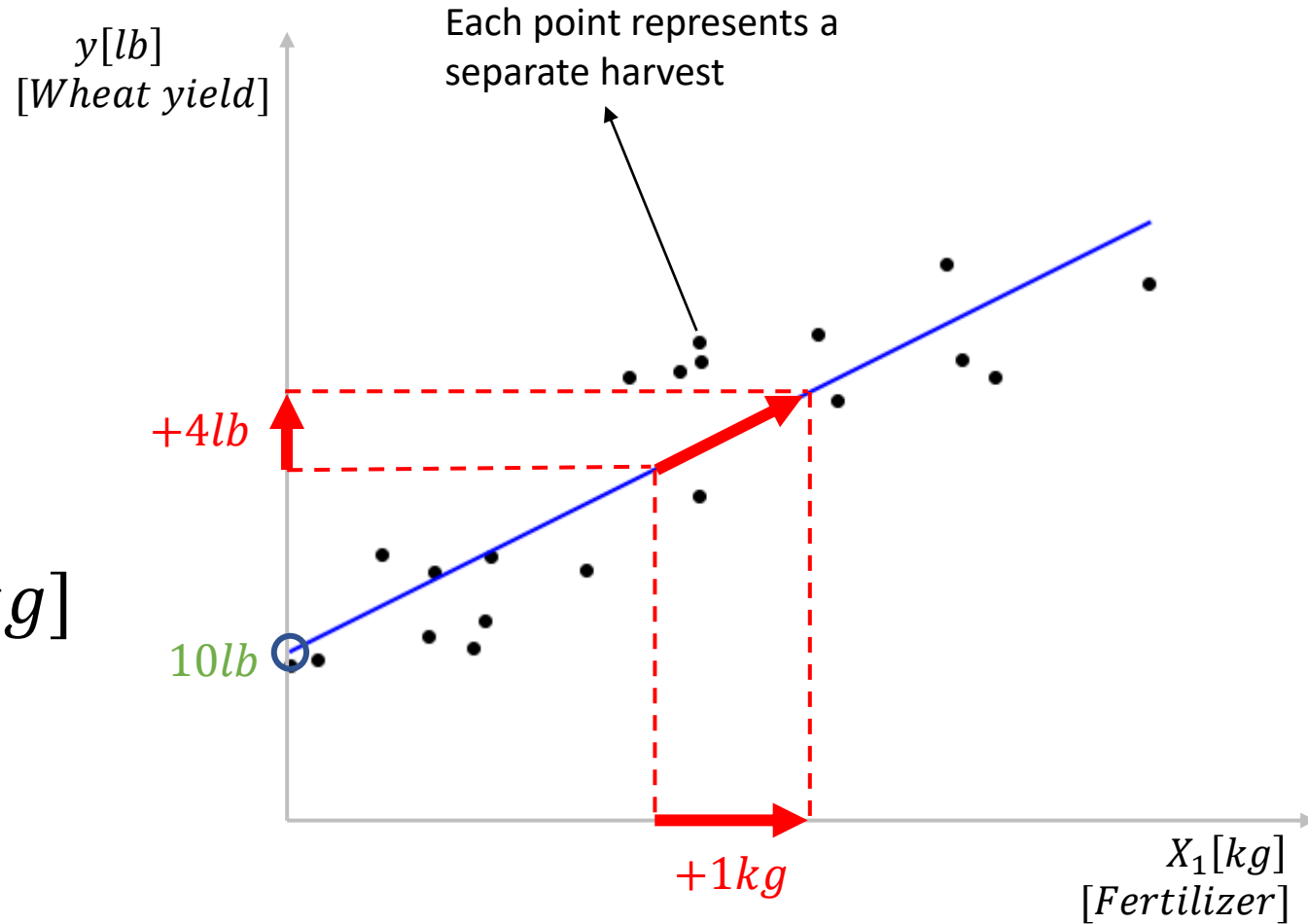
~



$$wheat[lb] = b_0 + b_1 * Fertilizer[kg]$$

$$b_0 = 10[lb]$$

$$b_1 = 4\left[\frac{lb}{kg}\right]$$



# Simple Linear Regression

$$\hat{y} = b_0 + b_1 X_1$$

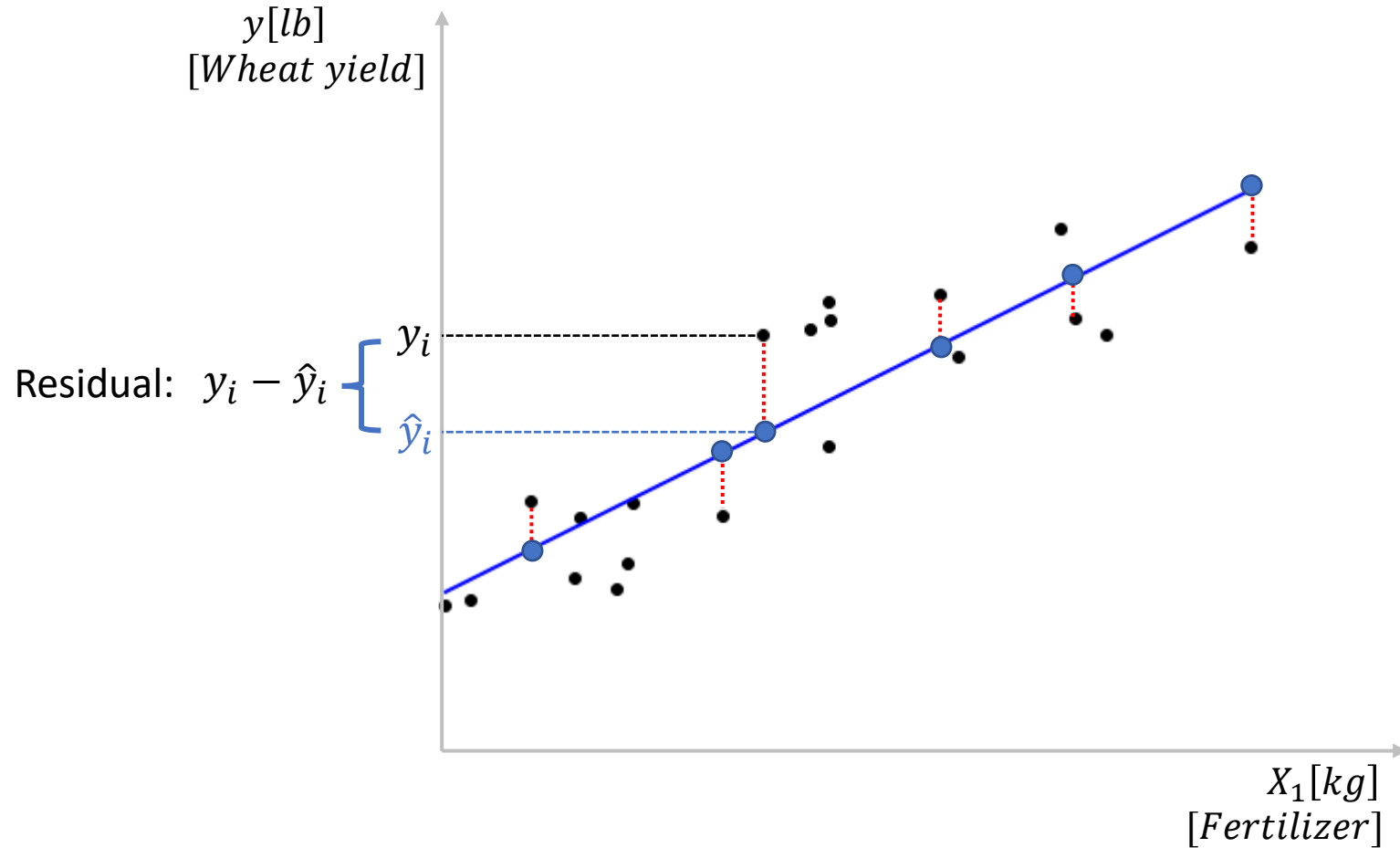
Ordinary Least Squares:

Estimates the values of

$b_0, b_1$  such that:

$$\sum_i (y_i - \hat{y}_i)^2$$

is minimized.



# Simple Linear Regression

$$\hat{y} = b_0 + b_1 X_1$$

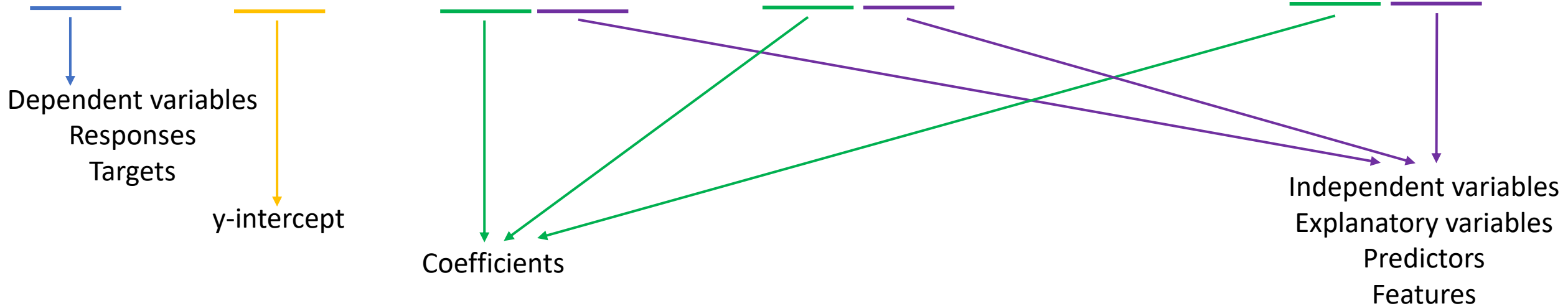
The diagram illustrates the components of the Simple Linear Regression equation  $\hat{y} = b_0 + b_1 X_1$ . Each term is underlined with a colored line, and a corresponding colored arrow points down to its label:

- $\hat{y}$  (blue underline) points to: Dependent variable, Response, Target
- $b_0$  (orange underline) points to: y-intercept
- $b_1$  (green underline) points to: Slope coefficient
- $X_1$  (purple underline) points to: Independent variable, Explanatory variable, Predictor, Feature



# Multiple Linear Regression

$$\hat{y} = b_0 + b_1X_1 + b_2X_2 + \cdots + b_nX_n$$



# Multiple Linear Regression

$$\hat{y} = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3$$

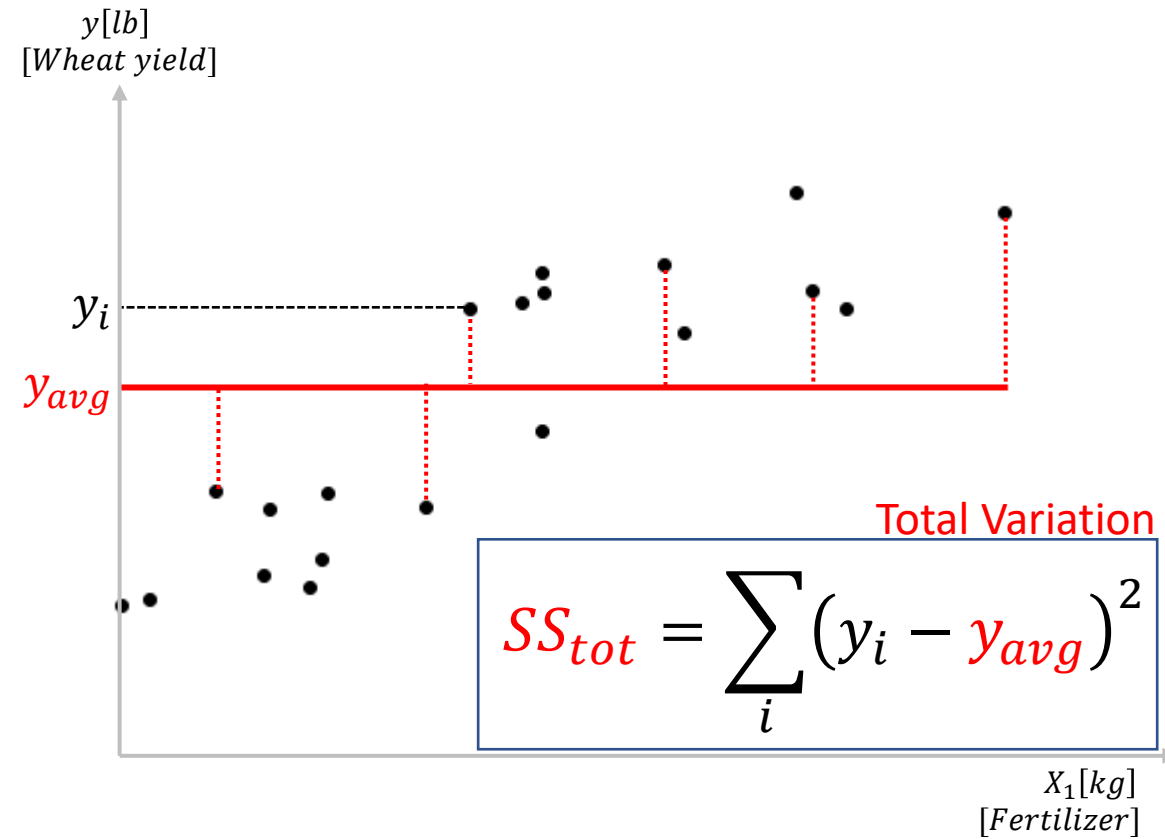
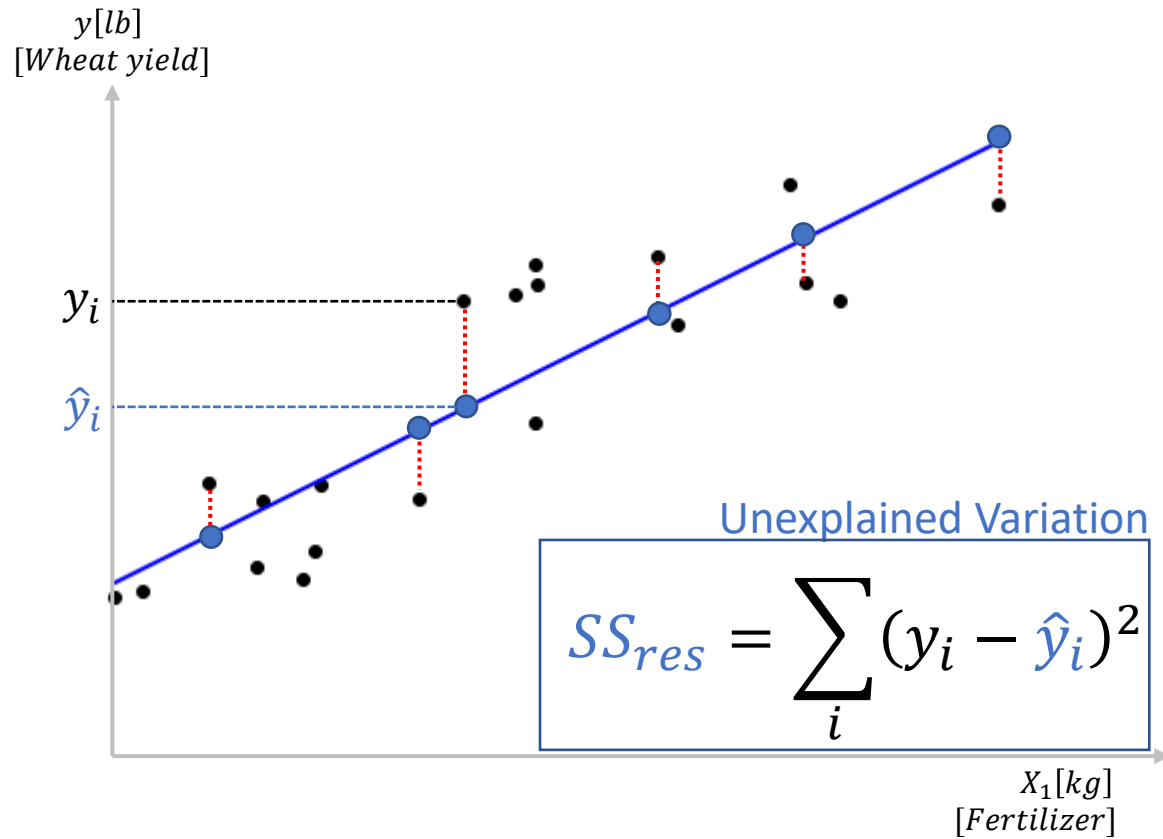


~



$$wheat[lb] = 7lb + 3 \frac{lb}{kg} * Fertilizer[kg] - 0.5 \frac{lb}{F} * Temp[F] + 1.2 \frac{lb}{mm} * Rain[mm]$$

# Coefficient of Determination ( $R^2$ )



$$R^2 = 1 - \frac{\text{Unexplained Variation}}{\text{Total Variation}} = 1 - \frac{SS_{res}}{SS_{tot}}$$

- Measures the % of variance in the target variable explained by the model
- Values between [0, 1]. Higher the better in general.

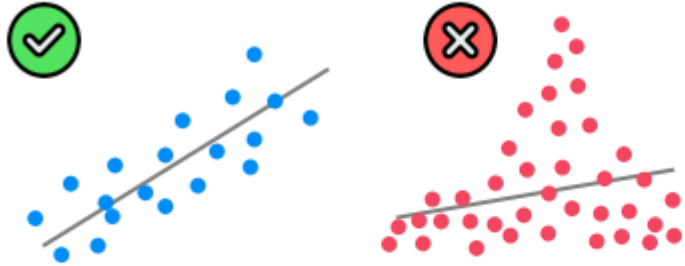
# Evaluation metrics

| $R^2$   | Adjusted $R^2$   | Mean Absolute Error (MAE)  | Root Mean Squared Error (RMSE)  |
|---|--|--|---|
| $R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$ <ul style="list-style-type: none"> <li>Measures the % of variance in the target variable explained by the model</li> <li>Values between [0, 1]</li> <li>Higher the better in general</li> </ul> | $Adj. R^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - k - 1}$ <ul style="list-style-type: none"> <li>Similar to <math>R^2</math>, but penalizes for the addition of too many variables</li> <li>Adding more variables always increase <math>R^2</math>, but not <math>Adj. R^2</math></li> <li>Higher the better in general</li> </ul> | $MAE = \frac{1}{n} \sum_{i=1}^n  y_i - \hat{y}_i $ <ul style="list-style-type: none"> <li>Easy to understand and interpret</li> <li>Same unit as target variable</li> <li>Not sensitive to outliers</li> <li>Lower the better</li> </ul> | $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$ <ul style="list-style-type: none"> <li>Same unit as target variable</li> <li>Sensitive to outliers – errors will be magnified due to the squared term</li> <li>Lower the better</li> </ul> |

# Assumptions of Linear Regression

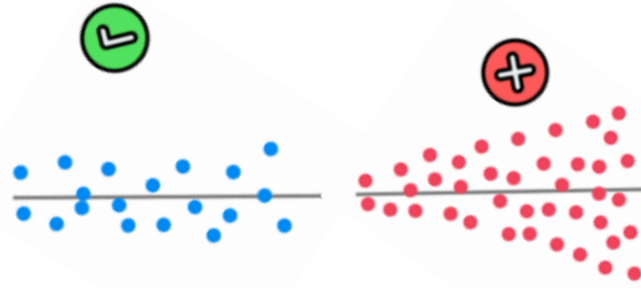
## 1. Linearity

(Linear relationship between Y and each X)



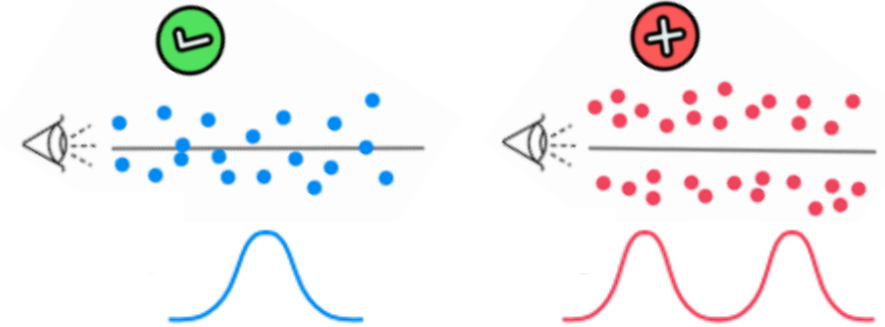
## 2. Homoscedasticity

(Equal variance)



## 3. Multivariate Normality

(Normality of error distribution)



## 4. Independence

(of observations. Includes "no autocorrelation")



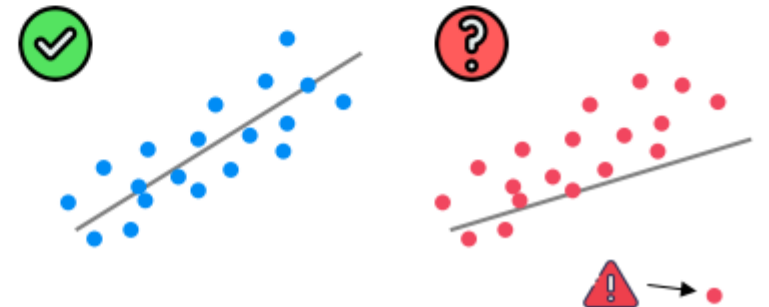
## 5. Lack of Multicollinearity

(Predictors are not correlated with each other)

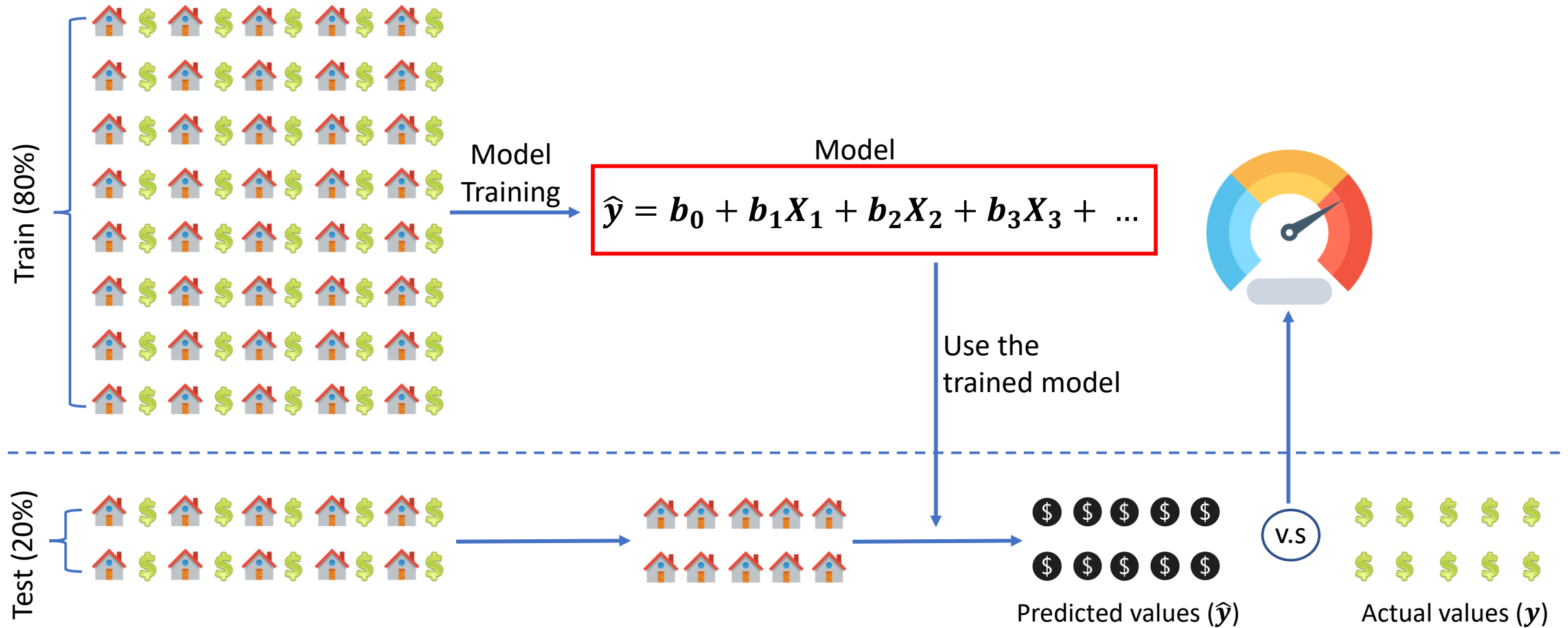
Green checkmark:  $X_1 \not\sim X_2$       Red 'X':  $X_1 \sim X_2$

## 6. The Outlier Check

(This is not an assumption, but an "extra")



# Linear Regression in action



# Pros and cons of Linear Regression

## **Pros**

- Easy to implement and understand.
- Computationally efficient.

## **Cons**

- Too simple to capture many of real-world complexities.
- Sensitive to outliers.
- Too many statistical assumptions which are usually hard to satisfy in real world data.

# Categorical Encodings



# Categorical Encodings

- Real world data is usually a combination of numerical columns (features) and categorical (strings) columns.
- Most implementations of machine learning algorithms require the input data to be numeric.
- Need to convert categorical columns to some numeric counterpart so that the machine learning algorithms can process them.
- Categorical Encoding - Process of transforming a categorical column into one or more numeric column(s).
- E.g. OneHotEncoding, OrdinalEncoding, CountEncoding, HashingEncoding, TargetEncoding, etc.

# OneHotEncoding (Dummy variables)

| $y$             | $X_1$           | $X_2$             | $X_3$     |         |
|-----------------|-----------------|-------------------|-----------|---------|
| Life expectancy | Adult Mortality | Total expenditure | Schooling | Country |
| 76.8            | 114             | 13.7              | 12.3      | USA     |
| 76.9            | 115             | 13.73             | 12.3      | USA     |
| 77              | 115             | 14.55             | 12.3      | USA     |
| 71              | 179             | 7.19              | 14.3      | Brazil  |
| 71.4            | 176             | 7.13              | 14.6      | Brazil  |
| 71.8            | 172             | 6.94              | 14.8      | Brazil  |
| 68.6            | 16              | 4.18              | 8         | Morocco |
| 69              | 155             | 4.44              | 8.5       | Morocco |
| 69.5            | 15              | 5.31              | 8.8       | Morocco |

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + ?$$

# OneHotEncoding (Dummy variables)

| $y$             | $X_1$           | $X_2$             | $X_3$     | Categorical variable |
|-----------------|-----------------|-------------------|-----------|----------------------|
| Life expectancy | Adult Mortality | Total expenditure | Schooling | Country              |
| 76.8            | 114             | 13.7              | 12.3      | USA                  |
| 76.9            | 115             | 13.73             | 12.3      | USA                  |
| 77              | 115             | 14.55             | 12.3      | USA                  |
| 71              | 179             | 7.19              | 14.3      | Brazil               |
| 71.4            | 176             | 7.13              | 14.6      | Brazil               |
| 71.8            | 172             | 6.94              | 14.8      | Brazil               |
| 68.6            | 16              | 4.18              | 8         | Morocco              |
| 69              | 155             | 4.44              | 8.5       | Morocco              |
| 69.5            | 15              | 5.31              | 8.8       | Morocco              |



How should we handle this?

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + ?$$

# OneHotEncoding (Dummy variables)

| $y$             | $X_1$           | $X_2$             | $X_3$     | Categorical variable | Dummy variables |       |       |
|-----------------|-----------------|-------------------|-----------|----------------------|-----------------|-------|-------|
| Life expectancy | Adult Mortality | Total expenditure | Schooling | Country              | $D_1$           | $D_2$ | $D_3$ |
| 76.8            | 114             | 13.7              | 12.3      | USA                  | 1               | 0     | 0     |
| 76.9            | 115             | 13.73             | 12.3      | USA                  | 1               | 0     | 0     |
| 77              | 115             | 14.55             | 12.3      | USA                  | 1               | 0     | 0     |
| 71              | 179             | 7.19              | 14.3      | Brazil               | 0               | 1     | 0     |
| 71.4            | 176             | 7.13              | 14.6      | Brazil               | 0               | 1     | 0     |
| 71.8            | 172             | 6.94              | 14.8      | Brazil               | 0               | 1     | 0     |
| 68.6            | 16              | 4.18              | 8         | Morocco              | 0               | 0     | 1     |
| 69              | 155             | 4.44              | 8.5       | Morocco              | 0               | 0     | 1     |
| 69.5            | 15              | 5.31              | 8.8       | Morocco              | 0               | 0     | 1     |

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 D_1 + b_5 D_2 + b_6 D_3$$

# OneHotEncoding (Dummy variables)

| $y$             | $X_1$           | $X_2$             | $X_3$     | Categorical variable | Dummy variables |       |       |
|-----------------|-----------------|-------------------|-----------|----------------------|-----------------|-------|-------|
| Life expectancy | Adult Mortality | Total expenditure | Schooling | Country              | $D_1$           | $D_2$ | $D_3$ |
| 76.8            | 114             | 13.7              | 12.3      | USA                  | 1               | 0     | 0     |
| 76.9            | 115             | 13.73             | 12.3      | USA                  | 1               | 0     | 0     |
| 77              | 115             | 14.55             | 12.3      | USA                  | 1               | 0     | 0     |
| 71              | 179             | 7.19              | 14.3      | Brazil               | 0               | 1     | 0     |
| 71.4            | 176             | 7.13              | 14.6      | Brazil               | 0               | 1     | 0     |
| 71.8            | 172             | 6.94              | 14.8      | Brazil               | 0               | 1     | 0     |
| 68.6            | 16              | 4.18              | 8         | Morocco              | 0               | 0     | 1     |
| 69              | 155             | 4.44              | 8.5       | Morocco              | 0               | 0     | 1     |
| 69.5            | 15              | 5.31              | 8.8       | Morocco              | 0               | 0     | 1     |

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 D_1 + b_5 D_2 + b_6 D_3$$

but,  $D_3 = 1 - D_1 - D_2$

# OneHotEncoding (Dummy variables)

| $y$             | $X_1$           | $X_2$             | $X_3$     | Categorical variable | Dummy variables |       |       |
|-----------------|-----------------|-------------------|-----------|----------------------|-----------------|-------|-------|
| Life expectancy | Adult Mortality | Total expenditure | Schooling | Country              | $D_1$           | $D_2$ | $D_3$ |
| 76.8            | 114             | 13.7              | 12.3      | USA                  | 1               | 0     | 0     |
| 76.9            | 115             | 13.73             | 12.3      | USA                  | 1               | 0     | 0     |
| 77              | 115             | 14.55             | 12.3      | USA                  | 1               | 0     | 0     |
| 71              | 179             | 7.19              | 14.3      | Brazil               | 0               | 1     | 0     |
| 71.4            | 176             | 7.13              | 14.6      | Brazil               | 0               | 1     | 0     |
| 71.8            | 172             | 6.94              | 14.8      | Brazil               | 0               | 1     | 0     |
| 68.6            | 16              | 4.18              | 8         | Morocco              | 0               | 0     | 1     |
| 69              | 155             | 4.44              | 8.5       | Morocco              | 0               | 0     | 1     |
| 69.5            | 15              | 5.31              | 8.8       | Morocco              | 0               | 0     | 1     |

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 D_1 + b_5 D_2 + b_6 D_3$$

- Always need to omit one dummy variable (doesn't matter which one is omitted)
- If there are  $n$  levels in a categorical variable, we only need  $n - 1$  dummy variables

# OneHotEncoding (Dummy variables)

Original data

| Life expectancy | Adult Mortality | Total expenditure | Schooling | Country |
|-----------------|-----------------|-------------------|-----------|---------|
| 76.8            | 114             | 13.7              | 12.3      | USA     |
| 76.9            | 115             | 13.73             | 12.3      | USA     |
| 77              | 115             | 14.55             | 12.3      | USA     |
| 71              | 179             | 7.19              | 14.3      | Brazil  |
| 71.4            | 176             | 7.13              | 14.6      | Brazil  |
| 71.8            | 172             | 6.94              | 14.8      | Brazil  |
| 68.6            | 16              | 4.18              | 8         | Morocco |
| 69              | 155             | 4.44              | 8.5       | Morocco |
| 69.5            | 15              | 5.31              | 8.8       | Morocco |

Modified data  
(Dummified)

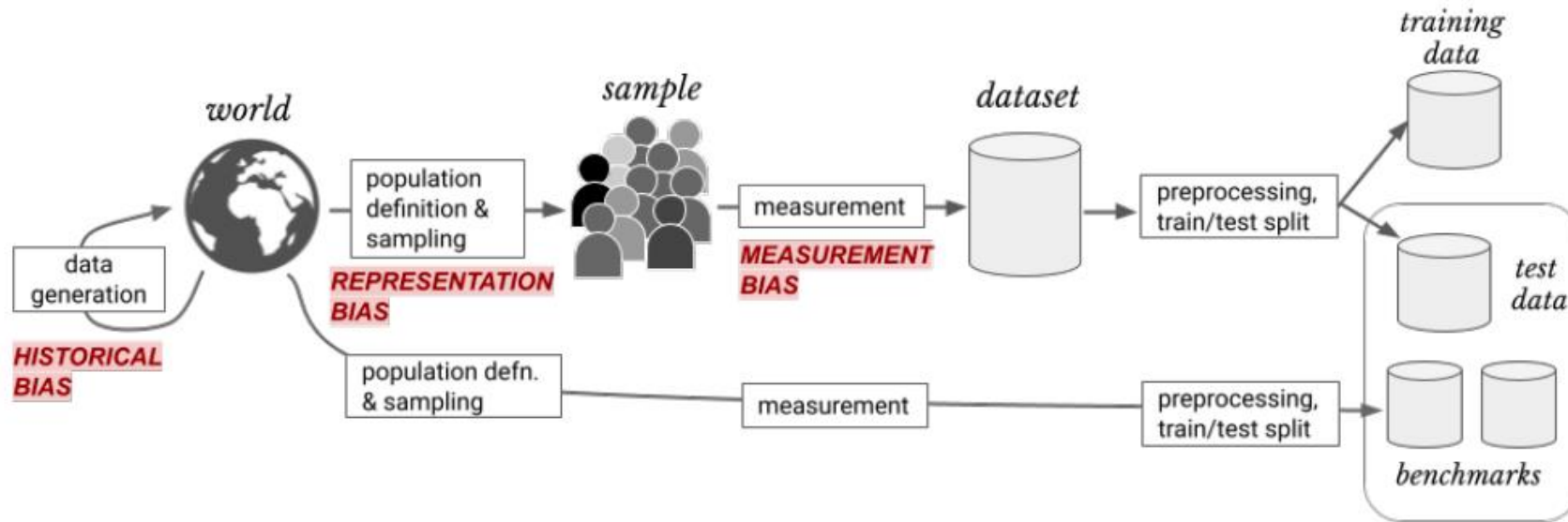
| Life expectancy | Adult Mortality | Total expenditure | Schooling | USA | Brazil |
|-----------------|-----------------|-------------------|-----------|-----|--------|
| 76.8            | 114             | 13.7              | 12.3      | 1   | 0      |
| 76.9            | 115             | 13.73             | 12.3      | 1   | 0      |
| 77              | 115             | 14.55             | 12.3      | 1   | 0      |
| 71              | 179             | 7.19              | 14.3      | 0   | 1      |
| 71.4            | 176             | 7.13              | 14.6      | 0   | 1      |
| 71.8            | 172             | 6.94              | 14.8      | 0   | 1      |
| 68.6            | 16              | 4.18              | 8         | 0   | 0      |
| 69              | 155             | 4.44              | 8.5       | 0   | 0      |
| 69.5            | 15              | 5.31              | 8.8       | 0   | 0      |

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + b_4 D_1 + b_5 D_2$$

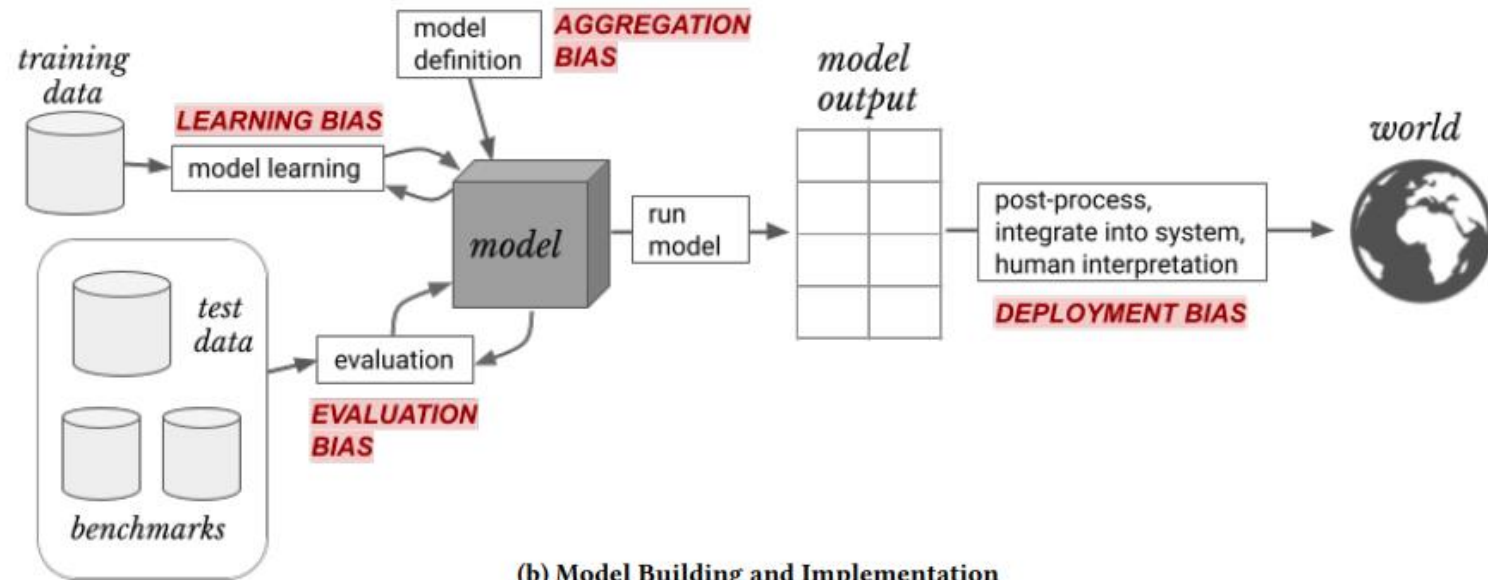
# Bias and Variance



# Bias-Variance

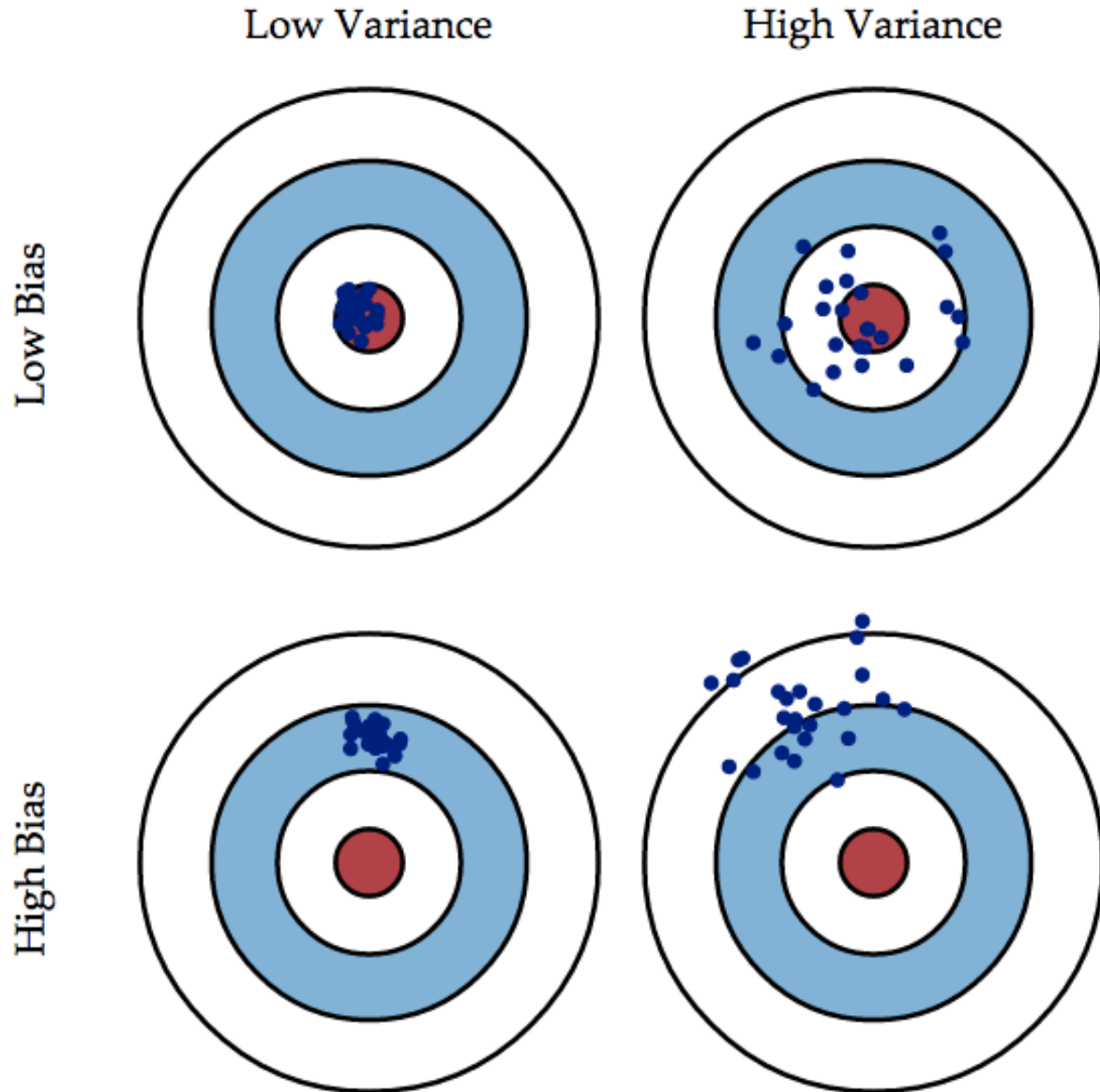


(a) Data Generation



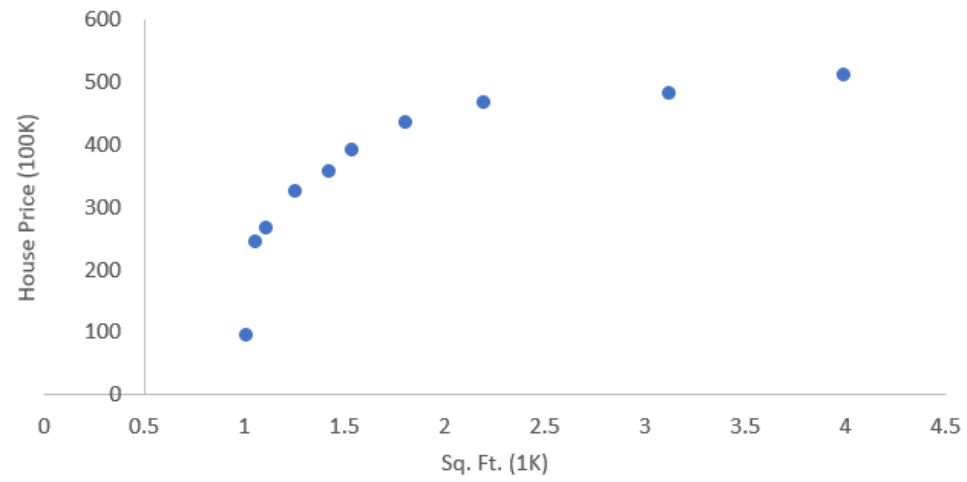
(b) Model Building and Implementation

# Bias-Variance

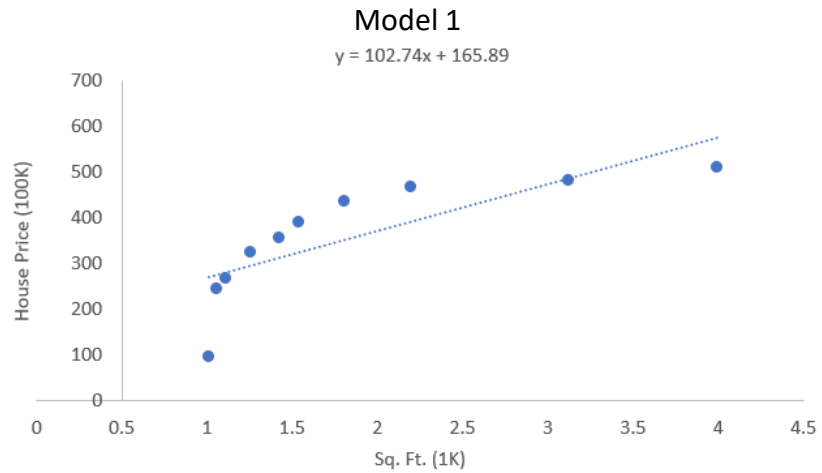
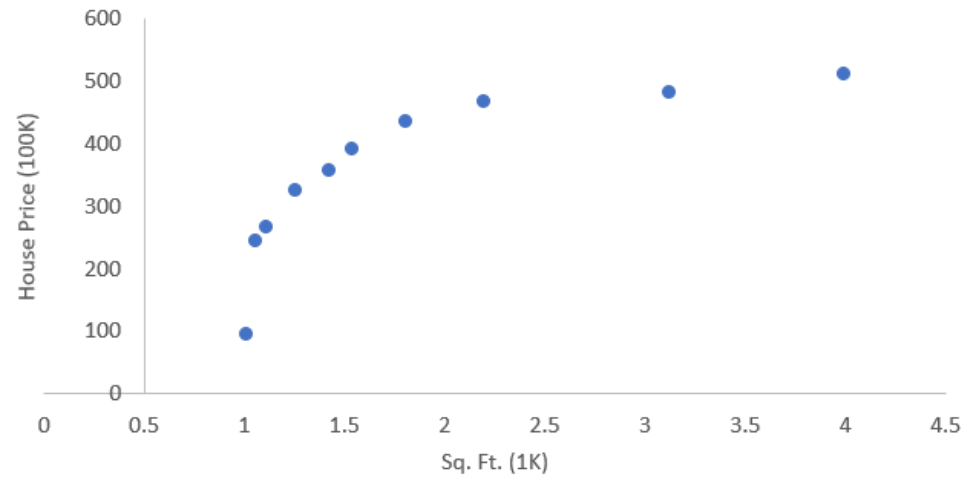


- The ***error due to bias*** is the amount by which the expected model prediction differs from the true value or target, over the training data.
- If these average prediction values are substantially different than the true value, bias will be high.
- The ***error due to variance*** is the amount by which the prediction, over one training set, differs from the expected predicted value, over all the training sets.
- Variance measures how inconsistent are the predictions from one another, over different training sets, not whether they are accurate or not.

# Bias-Variance Tradeoff

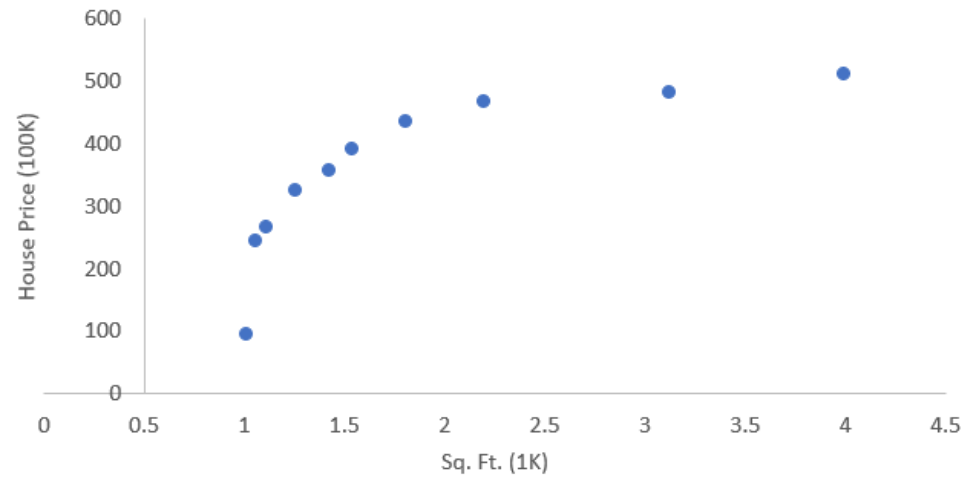


# Bias-Variance Tradeoff



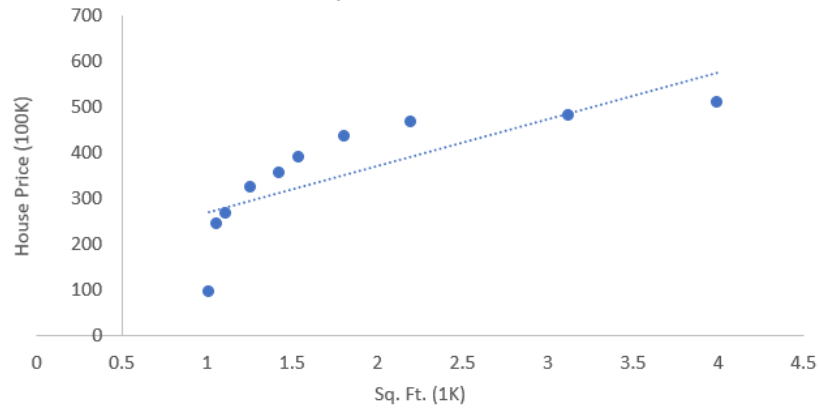
$$y = b_0 + b_1X_1$$

# Bias-Variance Tradeoff



Model 1

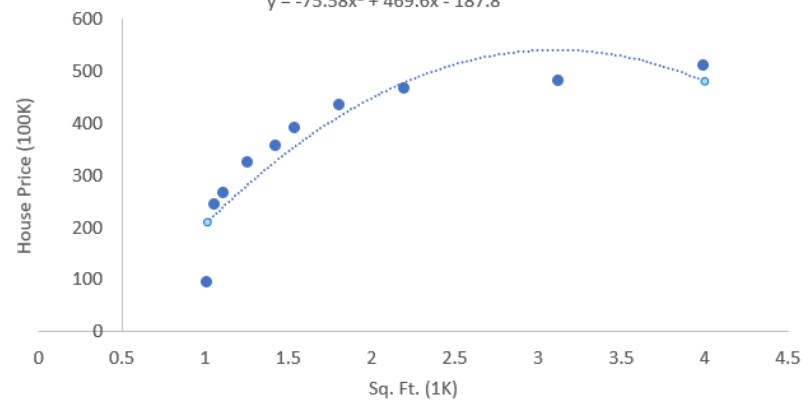
$$y = 102.74x + 165.89$$



$$y = b_0 + b_1X_1$$

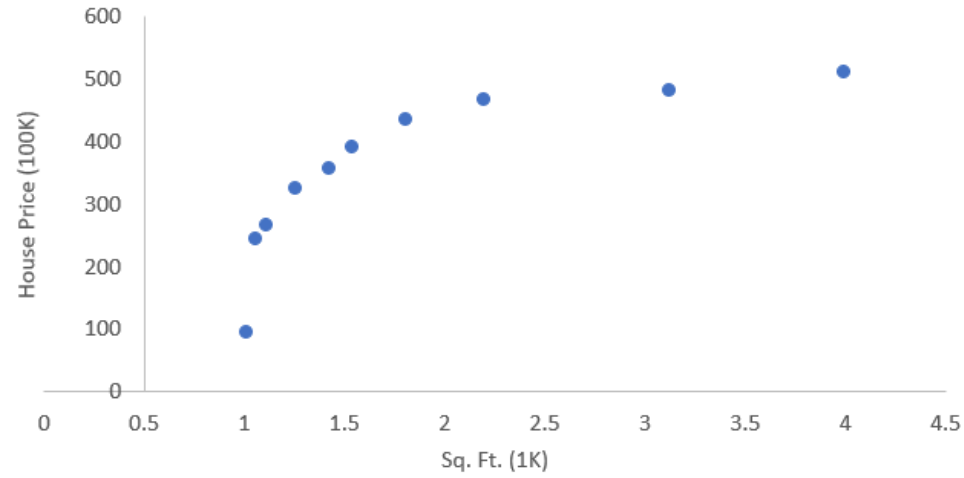
Model 2

$$y = -75.58x^2 + 469.6x - 187.8$$



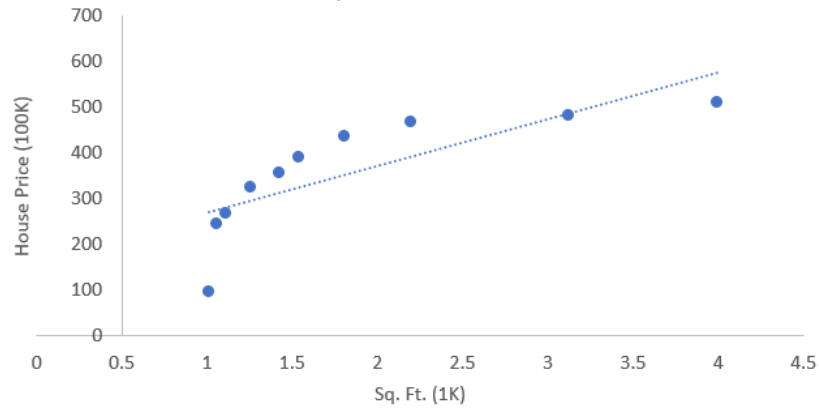
$$y = b_0 + b_1X_1 + b_2X_1^2$$

# Bias-Variance Tradeoff



Model 1

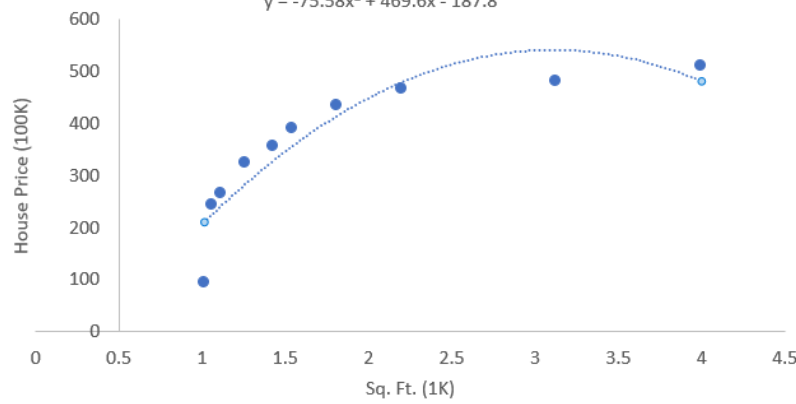
$$y = 102.74x + 165.89$$



$$y = b_0 + b_1X_1$$

Model 2

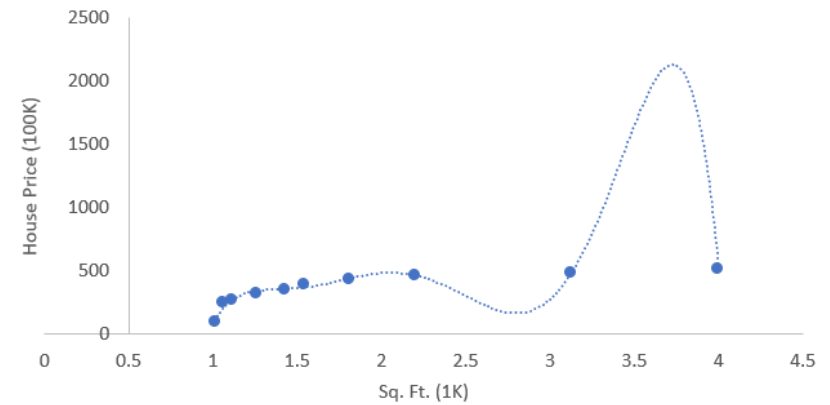
$$y = -75.58x^2 + 469.6x - 187.8$$



$$y = b_0 + b_1X_1 + b_2X_1^2$$

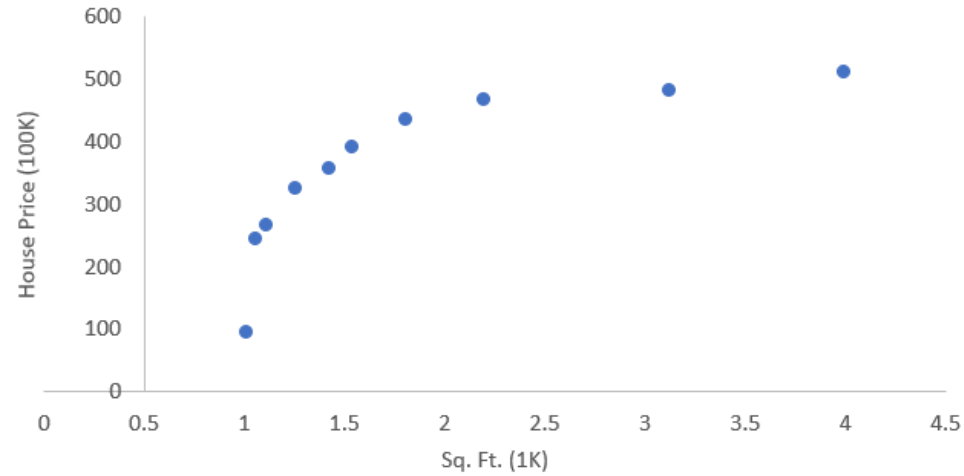
Model 3

$$y = -541.34x^6 + 7355.4x^5 - 39997x^4 + 111513x^3 - 168501x^2 + 131353x - 41098$$



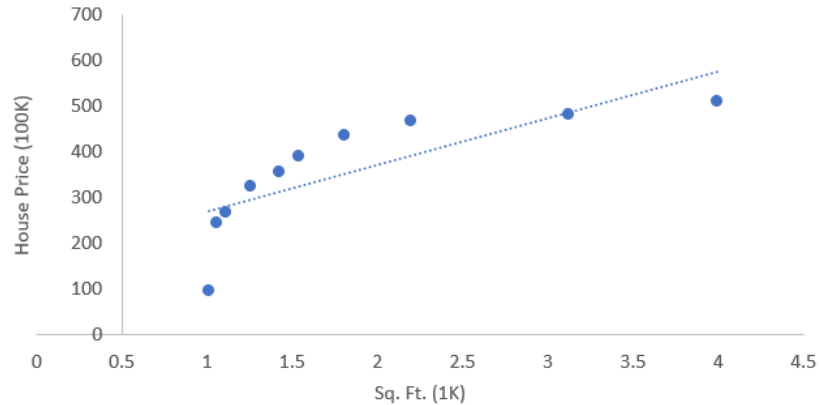
$$y = b_0 + b_1X_1 + b_2X_1^2 + b_3X_1^3 + b_4X_1^4 + b_5X_1^5 + b_6X_1^6$$

# Bias-Variance Tradeoff



Model 1

$$y = 102.74x + 165.89$$

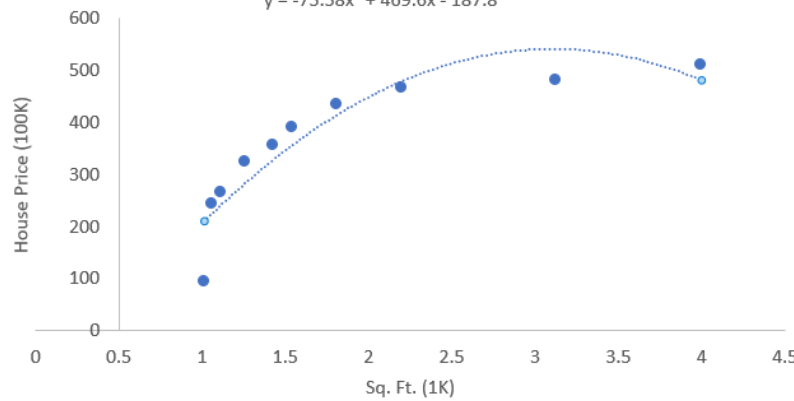


$$y = b_0 + b_1X_1$$

- Underfitting
- Biased

Model 2

$$y = -75.58x^2 + 469.6x - 187.8$$

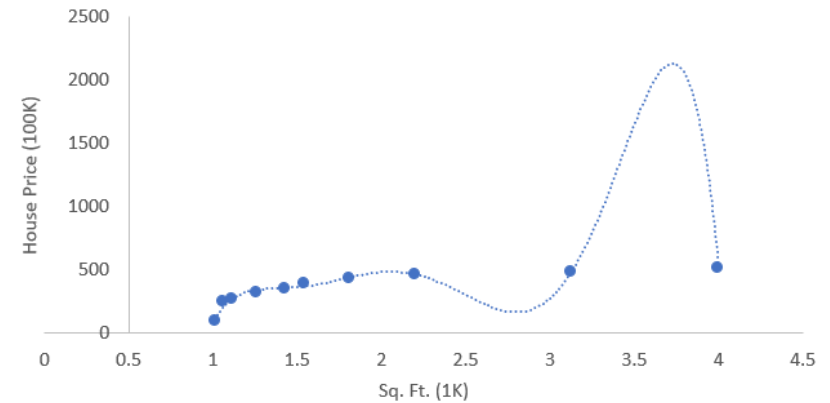


$$y = b_0 + b_1X_1 + b_2X_1^2$$

- Just right

Model 3

$$y = -541.34x^6 + 7355.4x^5 - 39997x^4 + 111513x^3 - 168501x^2 + 131353x - 41098$$



$$y = b_0 + b_1X_1 + b_2X_1^2 + b_3X_1^3 + b_4X_1^4 + b_5X_1^5 + b_6X_1^6$$

- Overfitting
- High variance

# Bias-Variance Tradeoff

