# Identify Data-level Issues in Machine Learning Models

**Ravikiran Srinivasulu**
SOFTWARE CONSULTANT

ravikirans.com | ravikirans.com/YouTube

# Imbalanced Dataset for Classification Problems

# Imbalanced Dataset for Classification Problems

75%

25%

< $50K    > $50K

Model

Predicted    Actual

75% accuracy

# Imbalanced Dataset for Classification Problems

# Undersampling



Protein 2 (vertical axis), Protein 1 (horizontal axis)

● Normal condition    ● Patient with a rare disease

# Random Oversampling

# Synthetic Minority Oversampling Technique



Protein 2

Protein 1

- ● Observation of interest
- ● Nearest neighbor
- ● Synthetic data points

Synthetic Minority Oversampling Technique

# Demo

Use SMOTE to increase minority samples in Census dataset

# Data Scale Issues in Distance-based Models
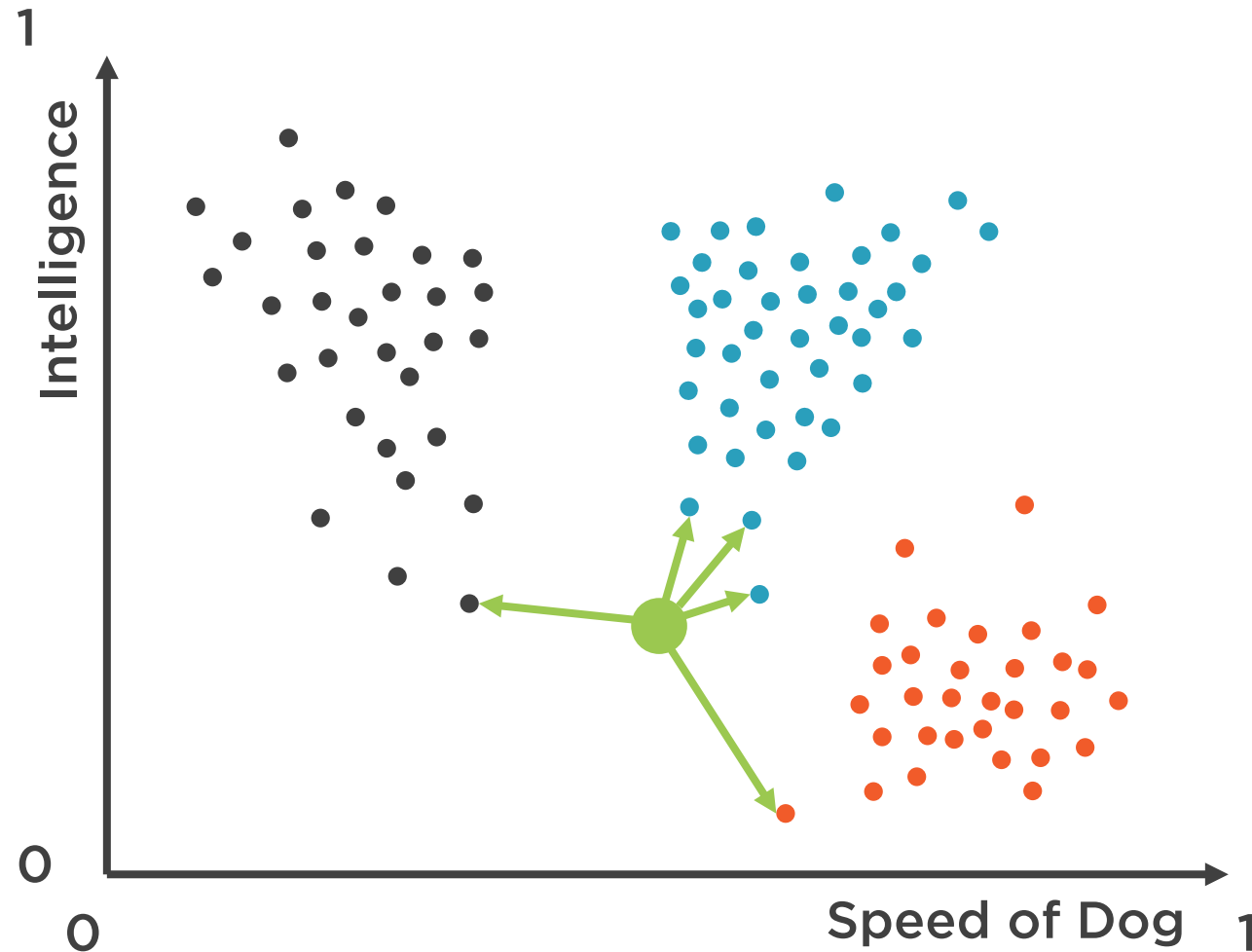
# Data Scale Issues in Distance-based Models
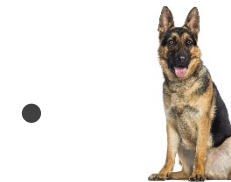
**Greyhound**

**German Shepherd**
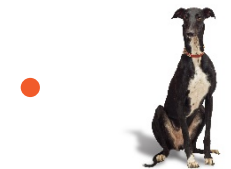
**Doberman**

# Data Scale Issues in Distance-based Models

# Data Scale Issues in Distance-based Models

# Data Scale Issues in Distance-based Models



$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

$$\sqrt{(30 - 34)^2 + (y_1 - y_2)^2}$$

$$\sqrt{(30 - 34)^2 + (7000 - 5400)^2}$$

**Intelligence dominates distance calculation**

(30, 7000)

(34, 5400)

# Multicollinearity in Multiple Regression

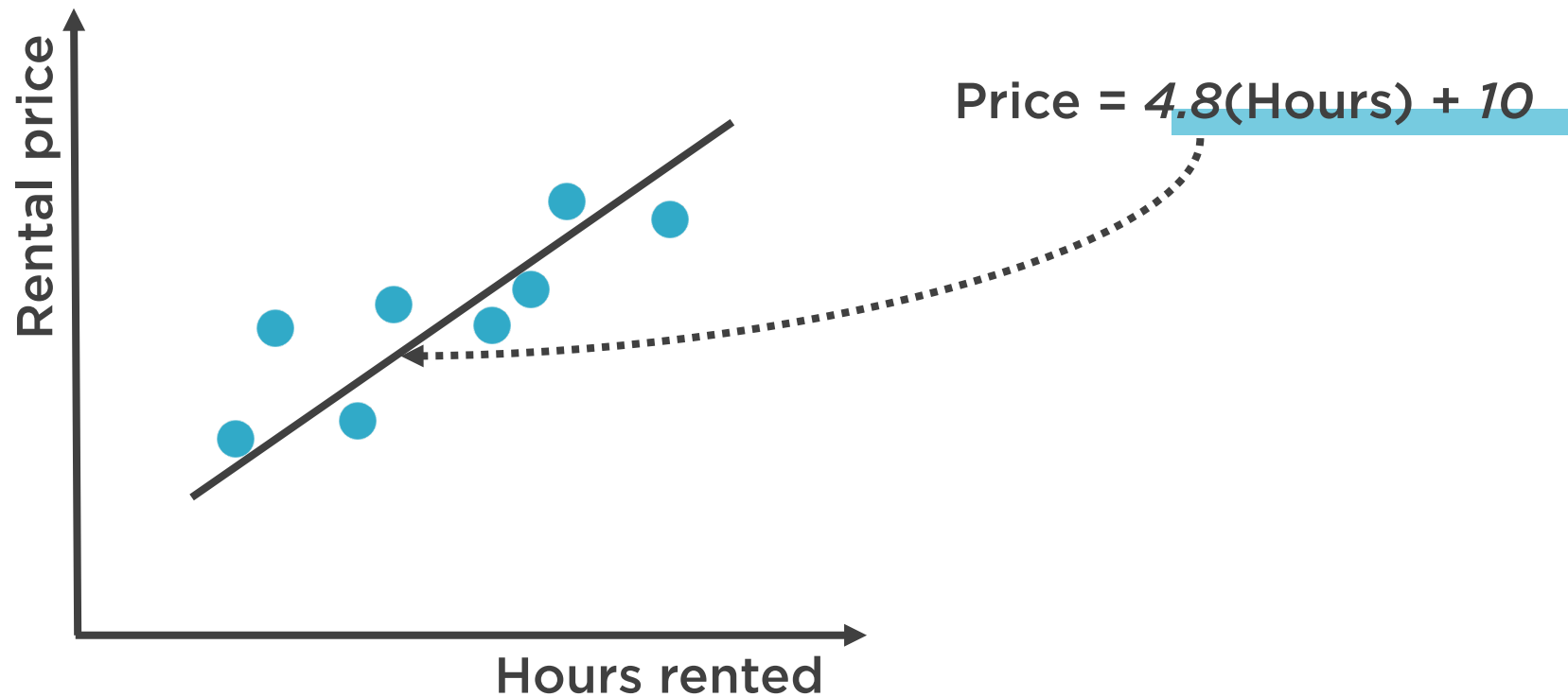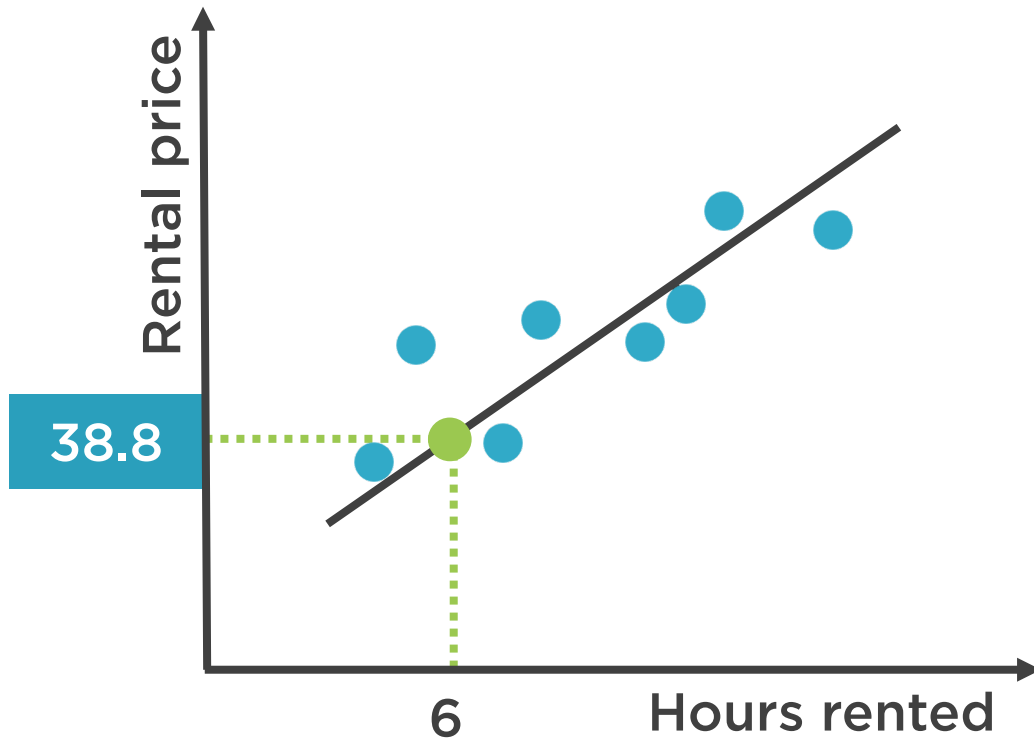# Multicollinearity

When one predictor variable in multiple regression can be linearly predicted from the others with a substantial degree of accuracy
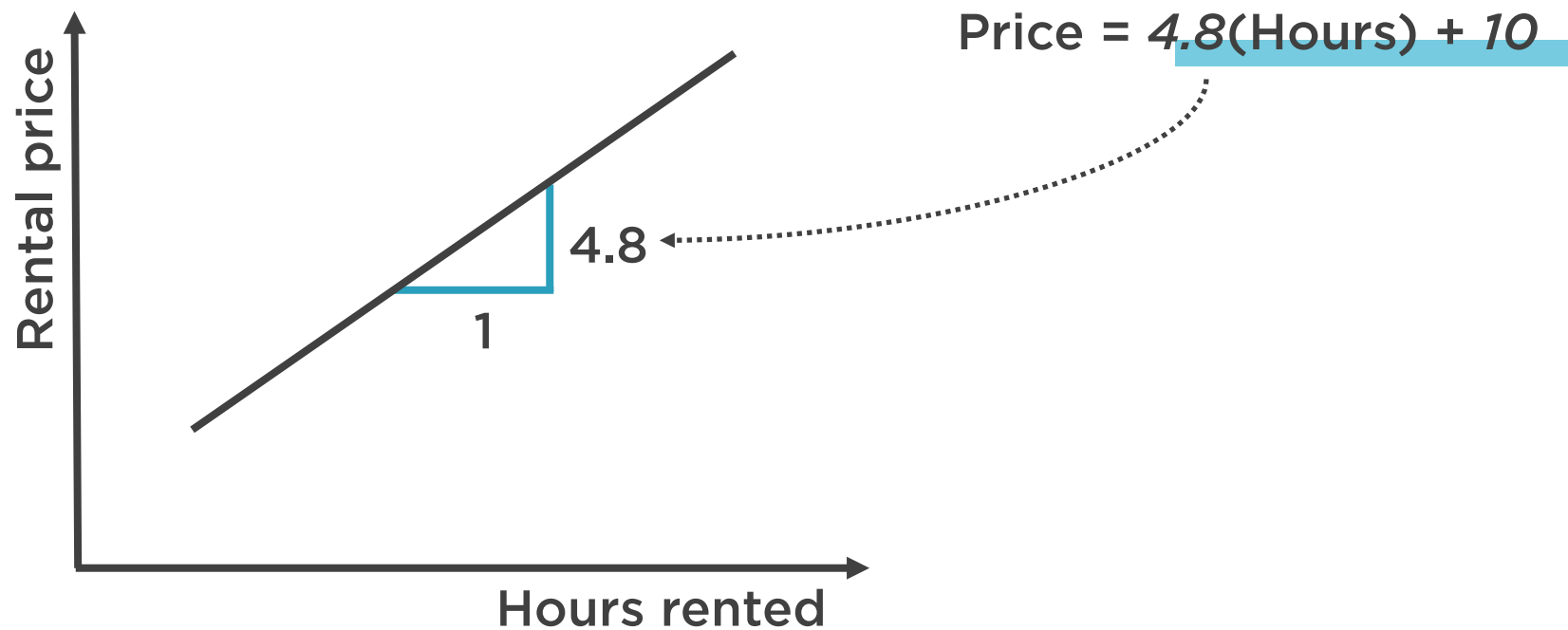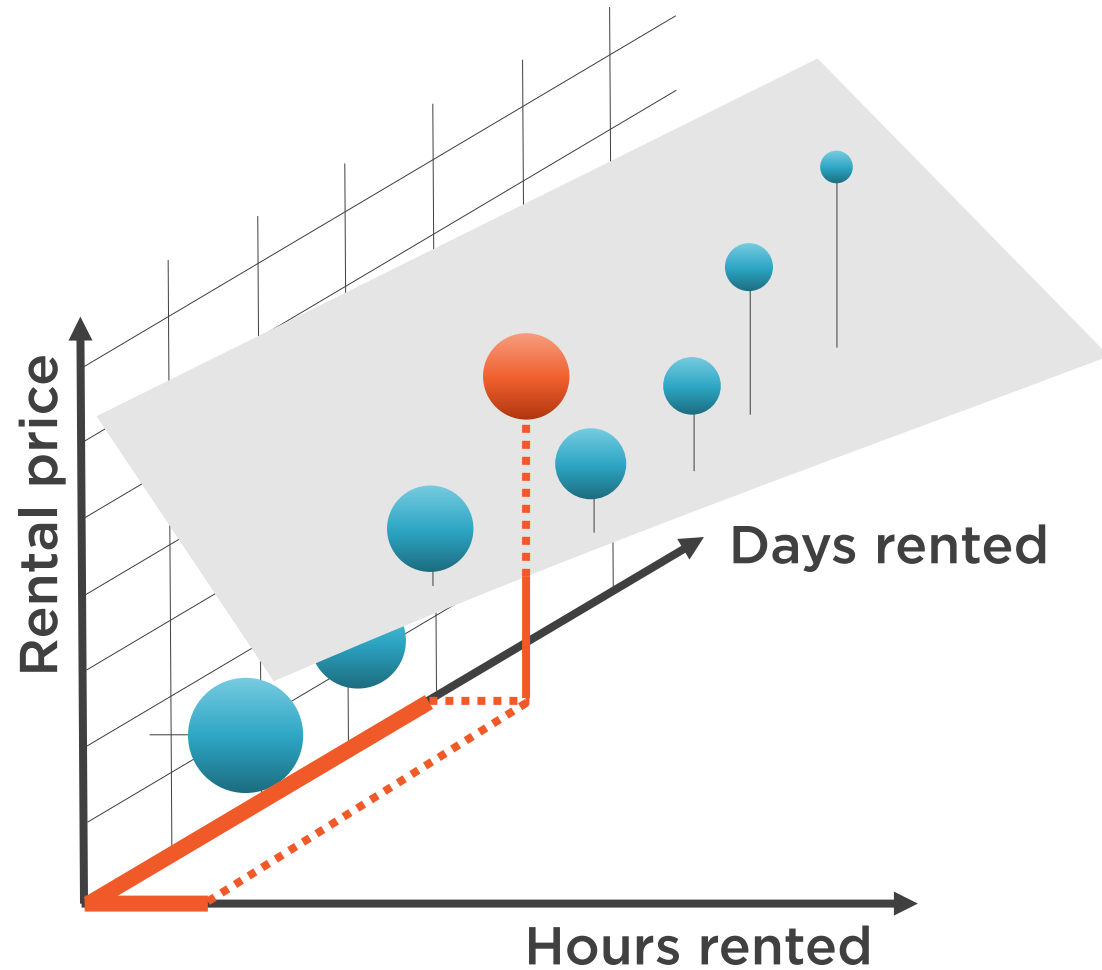
# Multicollinearity in Multiple Regression

Rental price

Hours rented

Price = *4.8*(Hours) + *10*

# Multicollinearity in Multiple Regression



$38.8 = 4.8(6) + 10$

# Multicollinearity in Multiple Regression

Rental price

Hours rented

4.8

1

Price = 4.8(Hours) + 10

# Multicollinearity in Multiple Regression

$$Price = .12(Hours) + 100(days) + C$$

Rental price

Days rented

Hours rented

# Multicollinearity in Multiple Regression

Price = .12(Hours) + 100(days) + C

Rental price
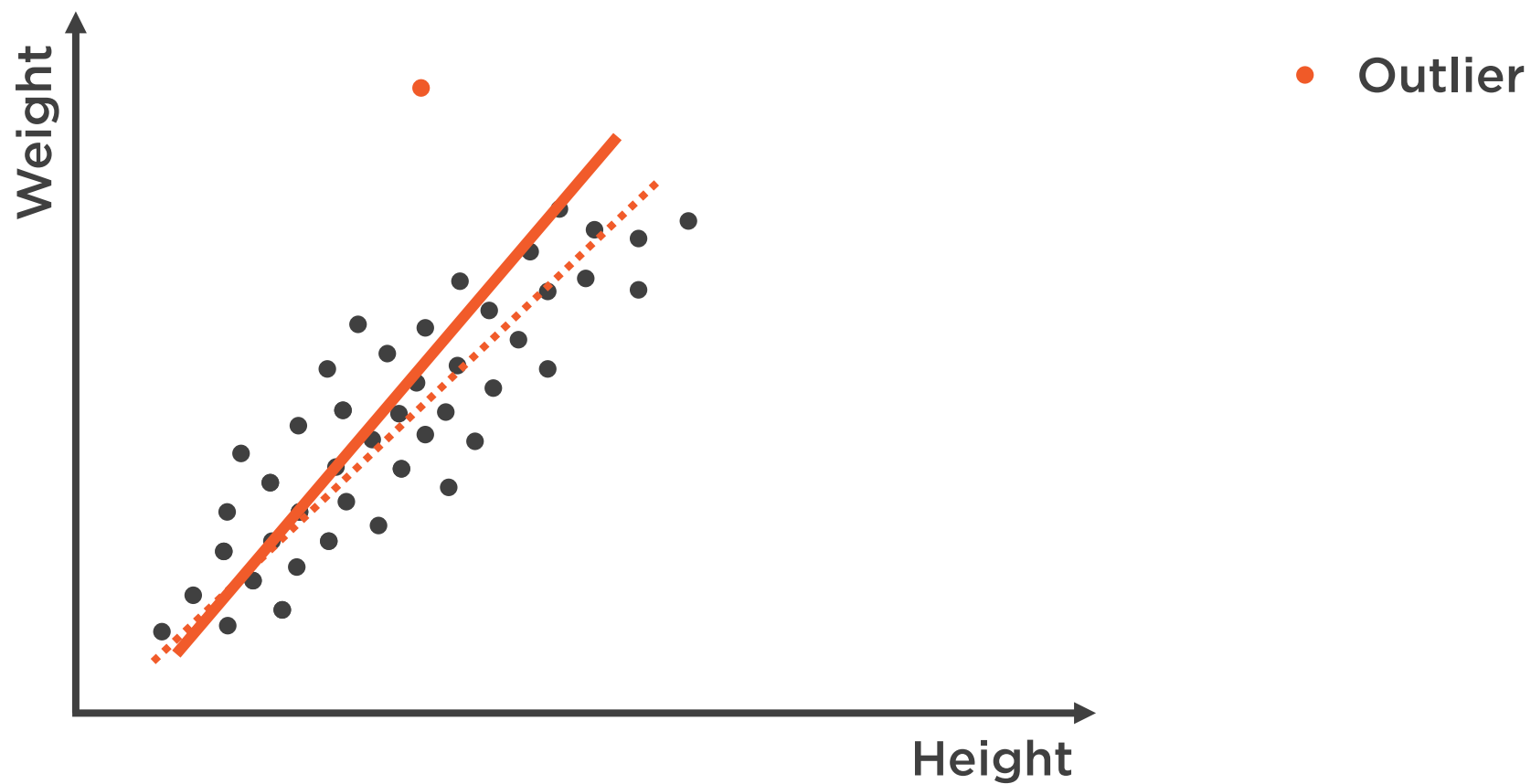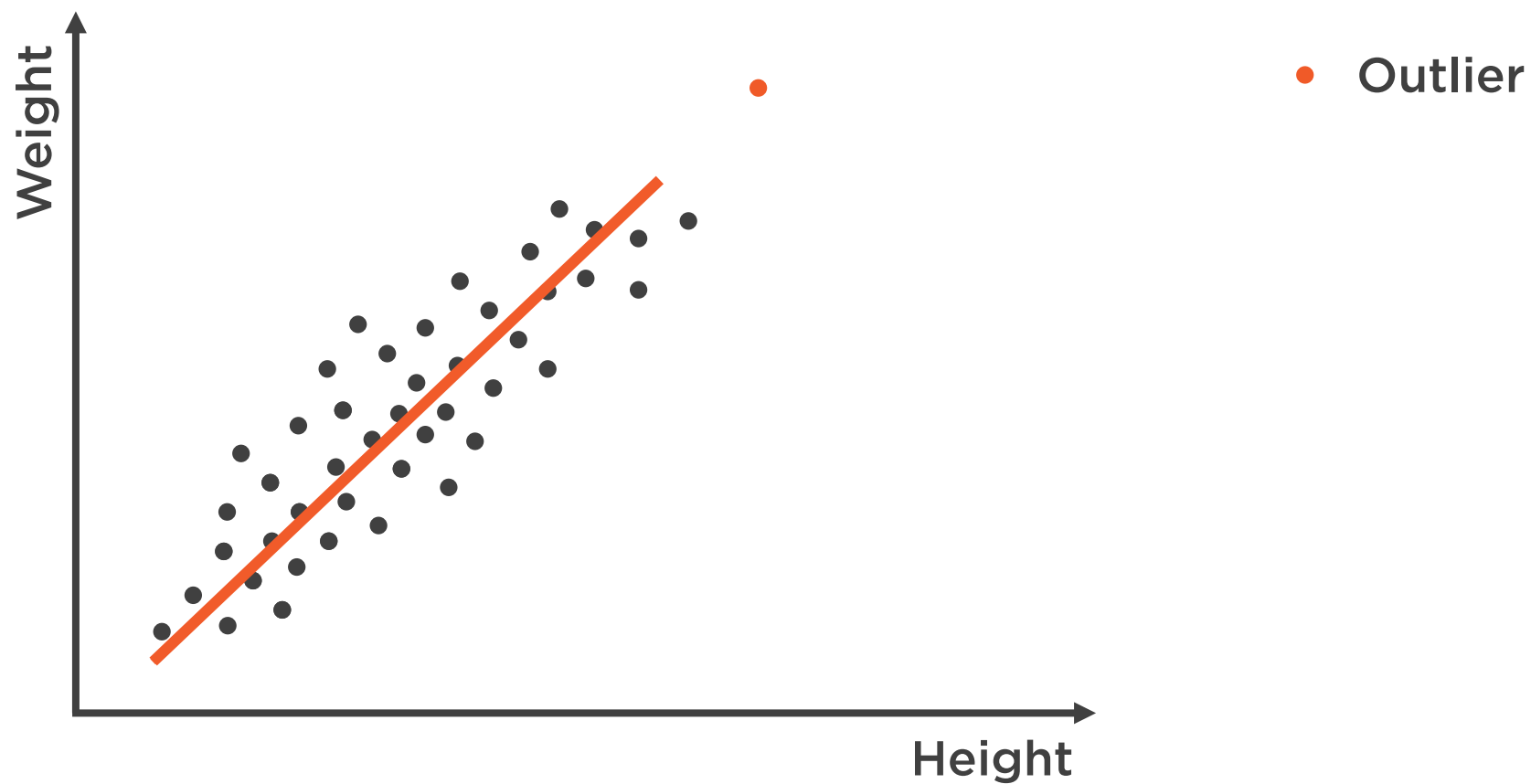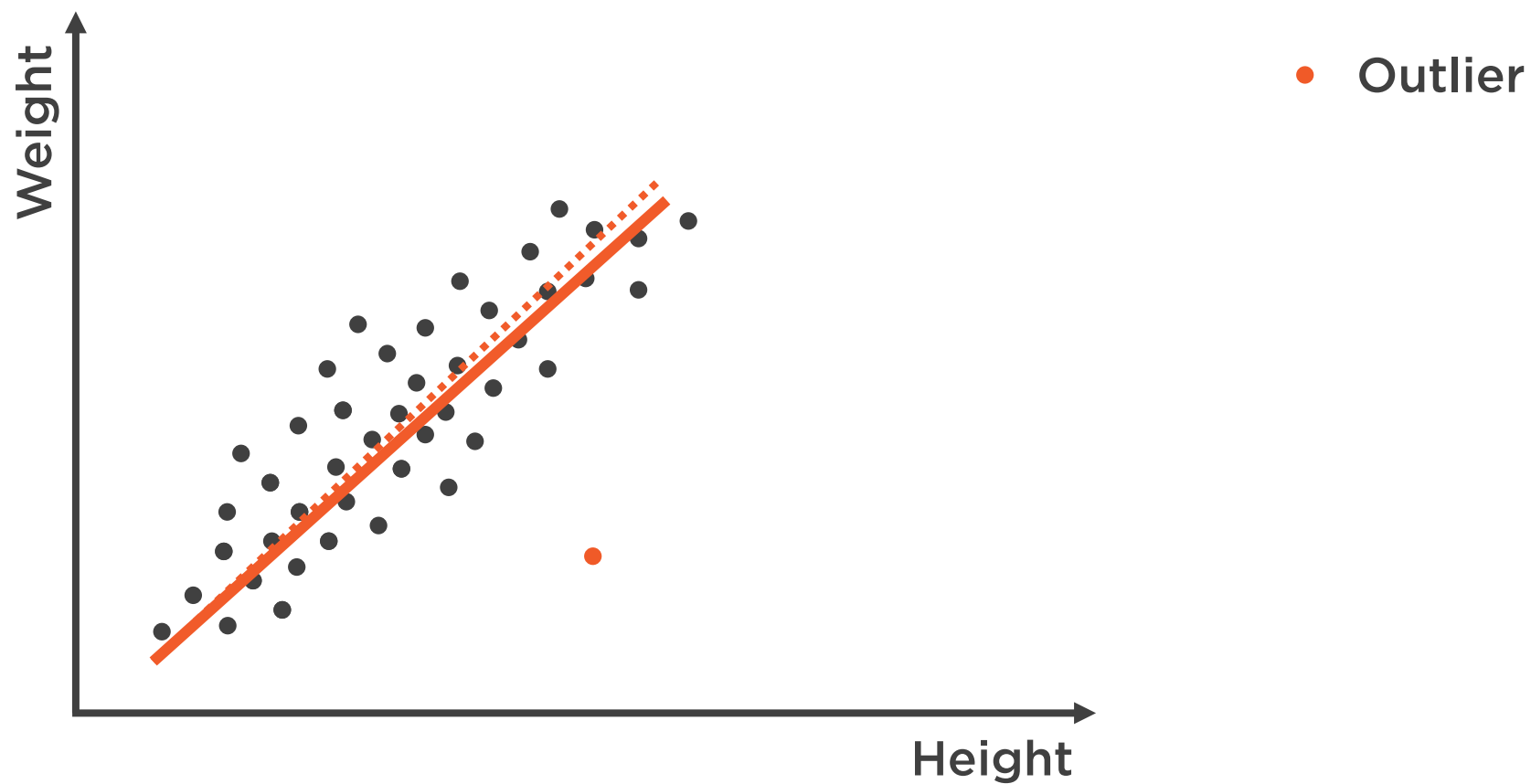
Days rented

Hours rented

# Outliers in Regression Models

# Outliers in Regression Models
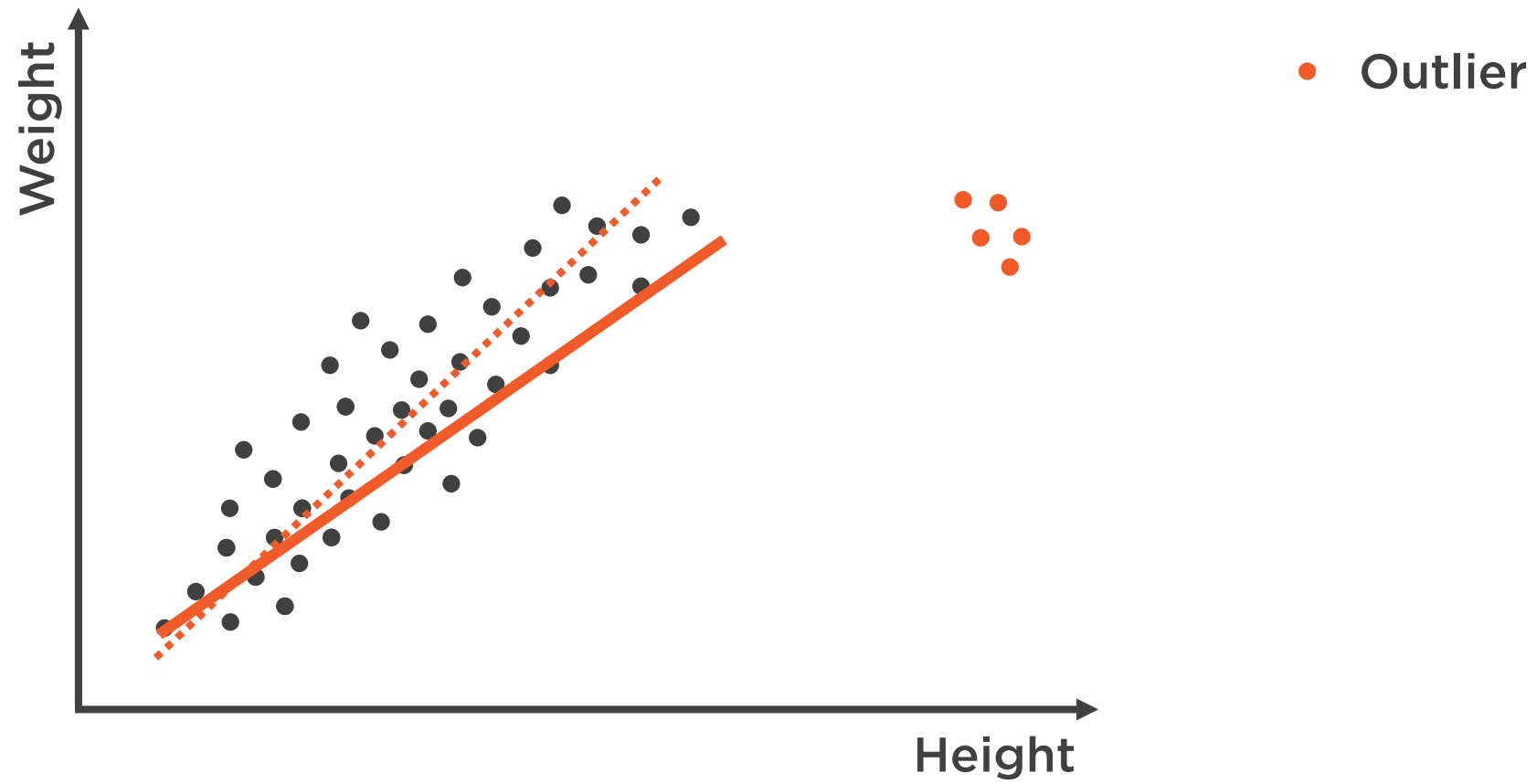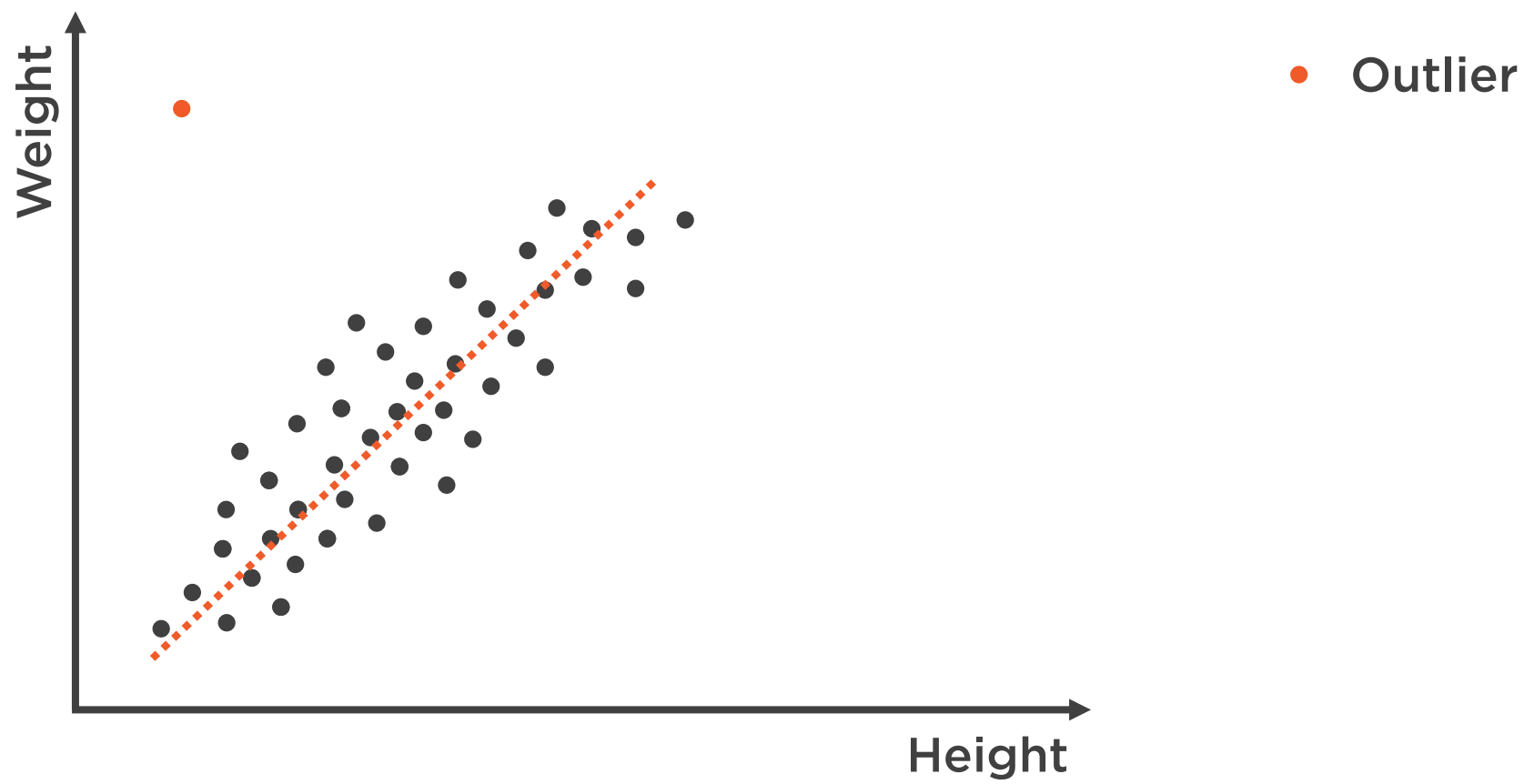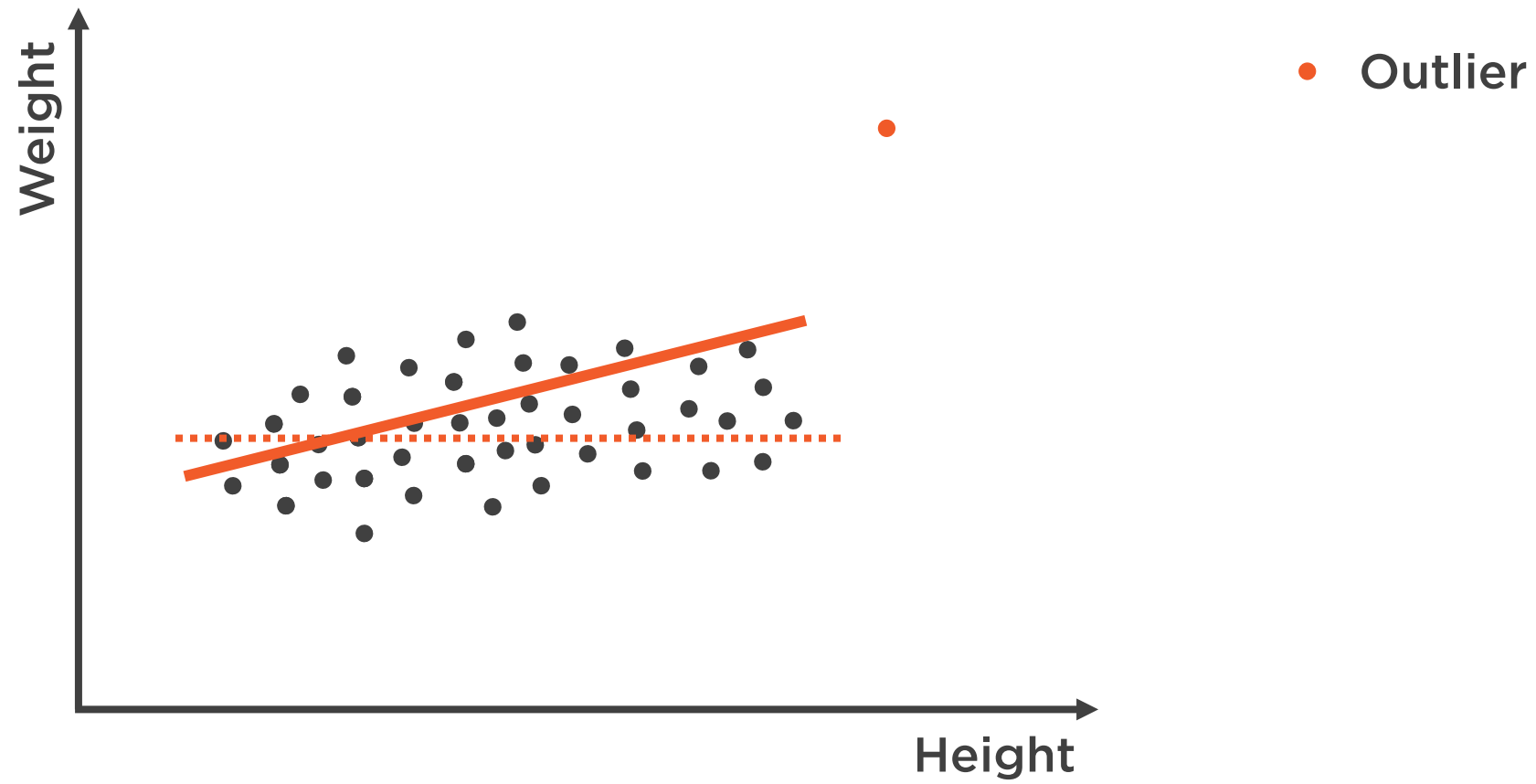
# Outliers in Regression Models

# Outliers in Regression Models

# Outliers in Regression Models

# Outliers in Regression Models
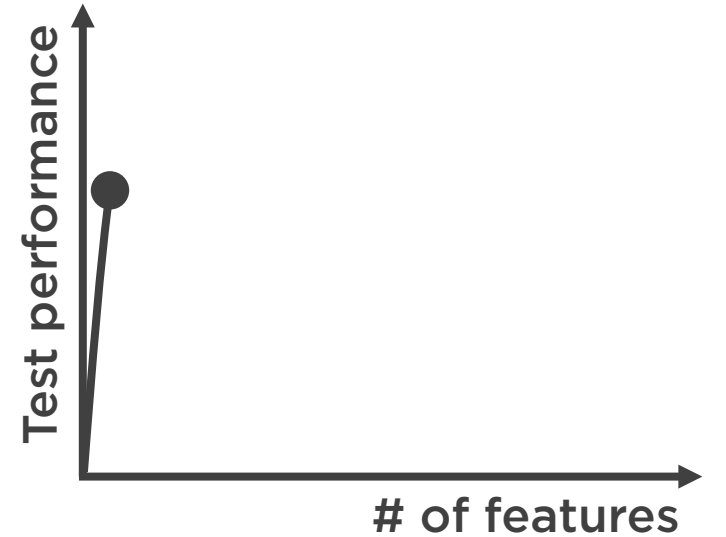
# Outliers in Regression Models

# Outliers in Regression Models
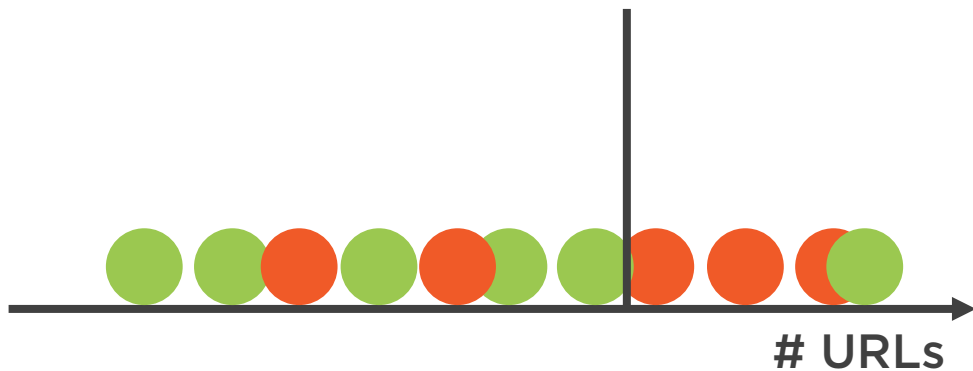
# Problem with High-dimensional Datasets
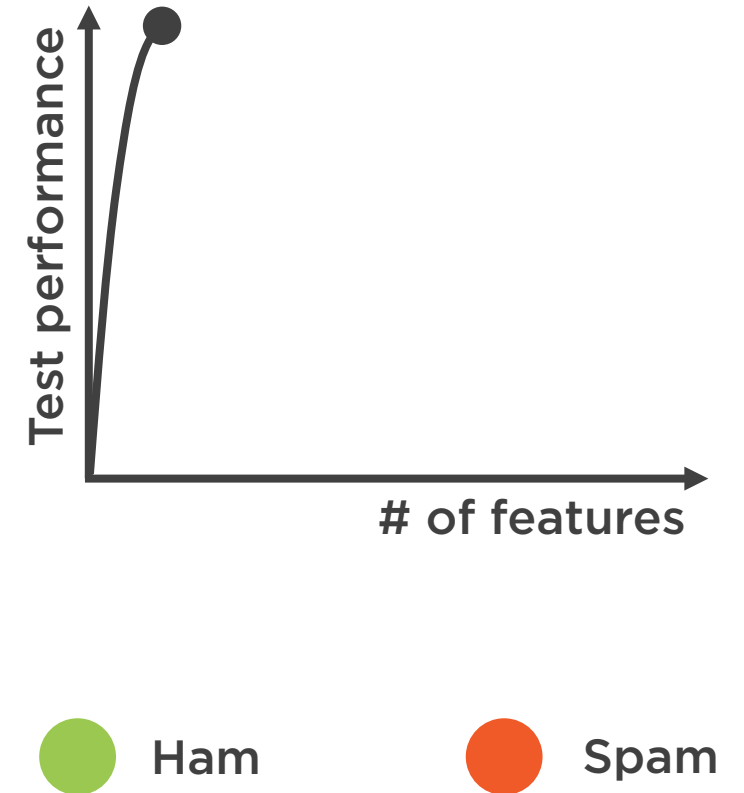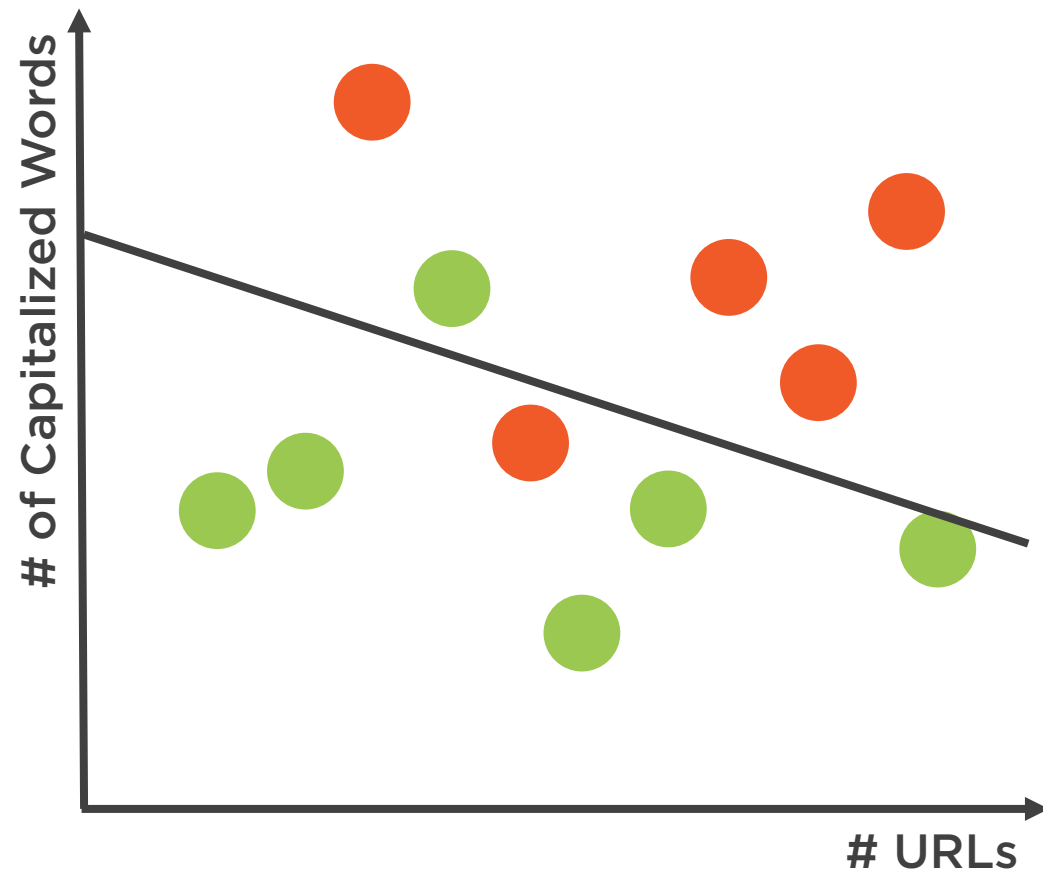
# Problem with High-dimensional Datasets
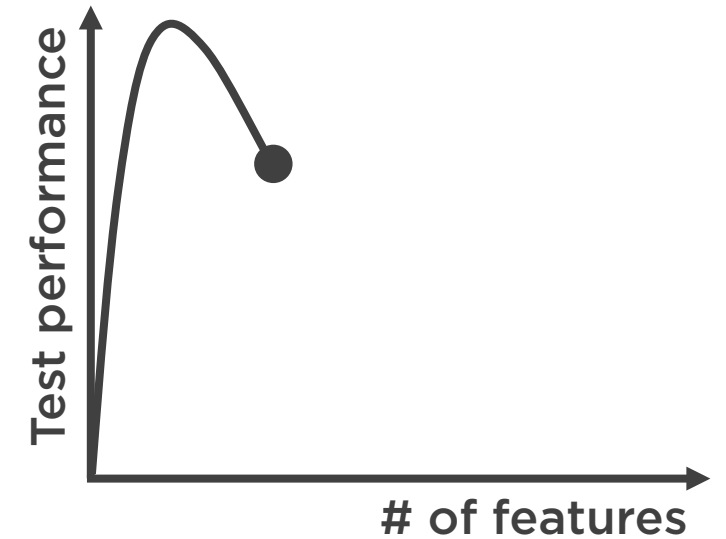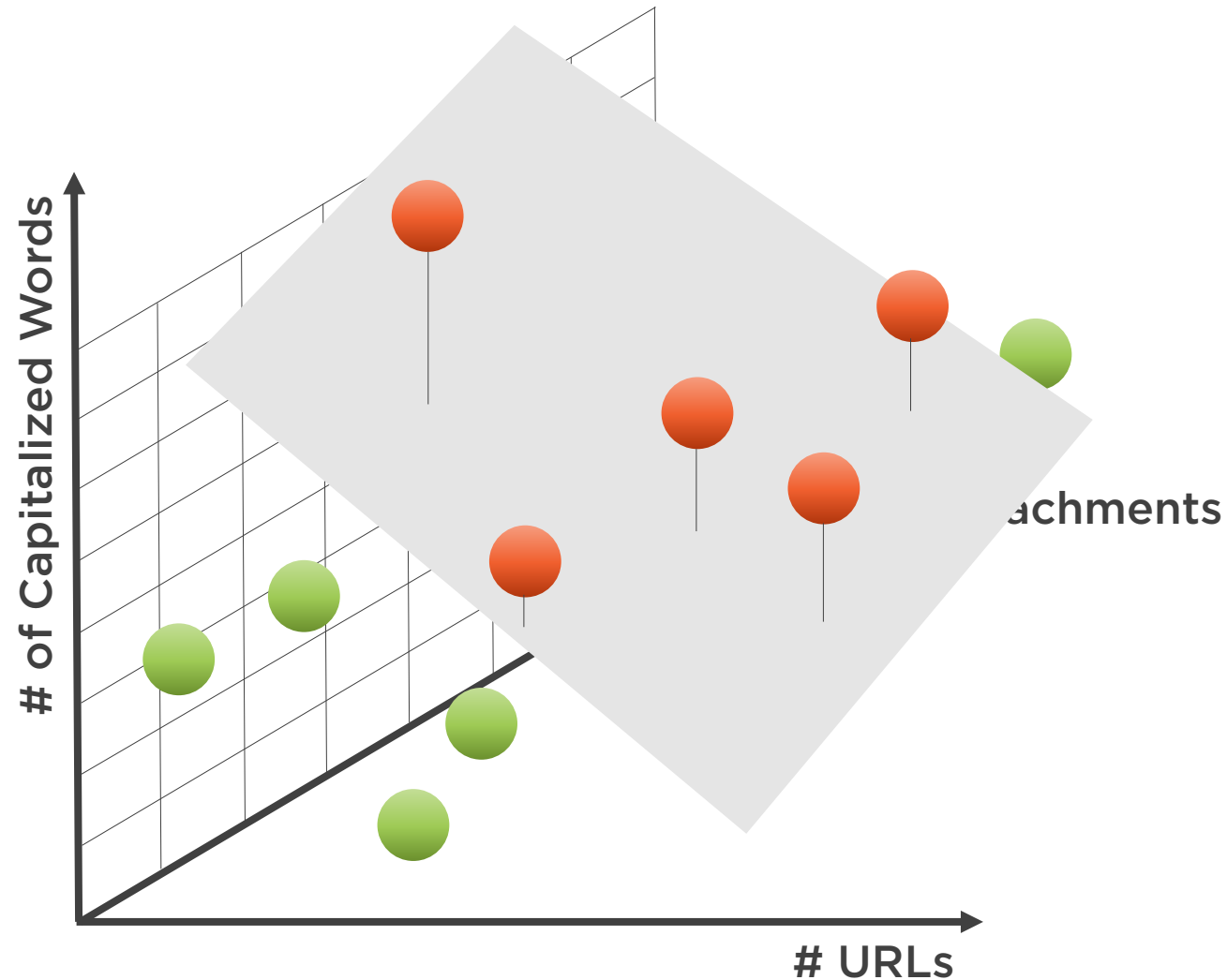
# Problem with High-dimensional Datasets

# Problem with High-dimensional Datasets
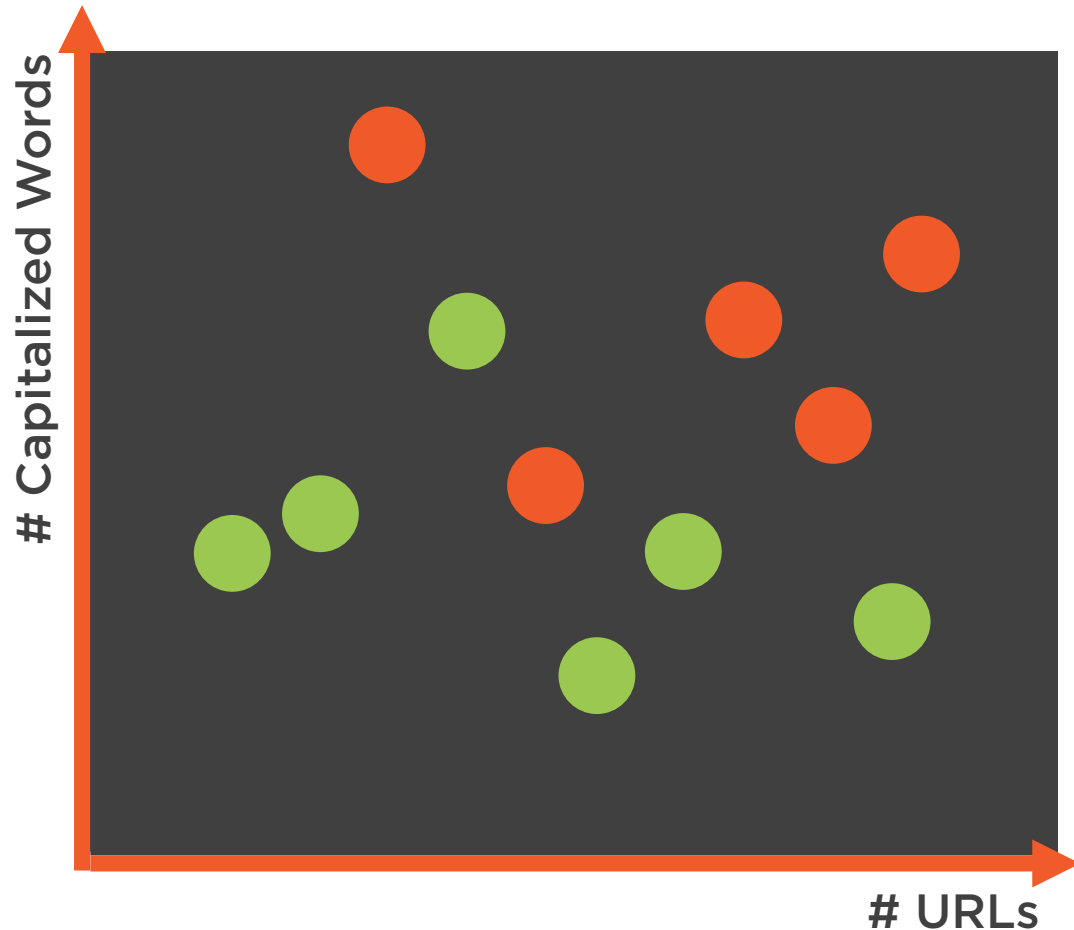
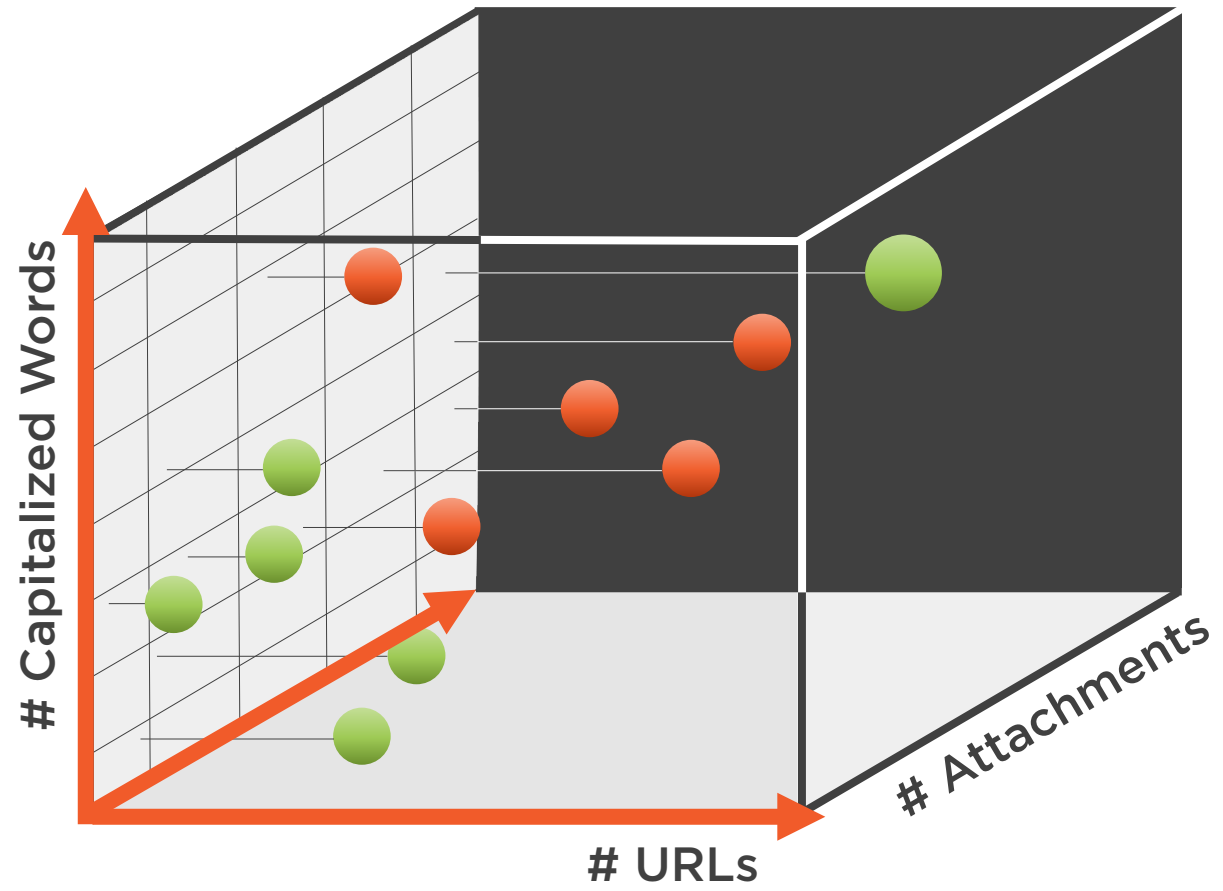# Problem with High-dimensional Datasets



# URLs

Data density:

**10/10 = 1**

# Problem with High-dimensional Datasets



**# Capitalized Words** (y-axis)

**# URLs** (x-axis)

**Data density:**

**10/100 = 0.1**

# Problem with High-dimensional Datasets



- Increase the number of observations
- Remove unnecessary features
- Use PCA

# Summary

**Data-level issues indicate the importance of data transformation**

**Play with SMOTE to improve model performance**

**PCA can help solve a variety of issues**