Lesson 4.3: Machine Learning Fundamentals

Slide 2: Welcome Back!

# Welcome Back!

Machine learning in a business context

Fundamentals of machine learning

Slide 3: Learning Objectives
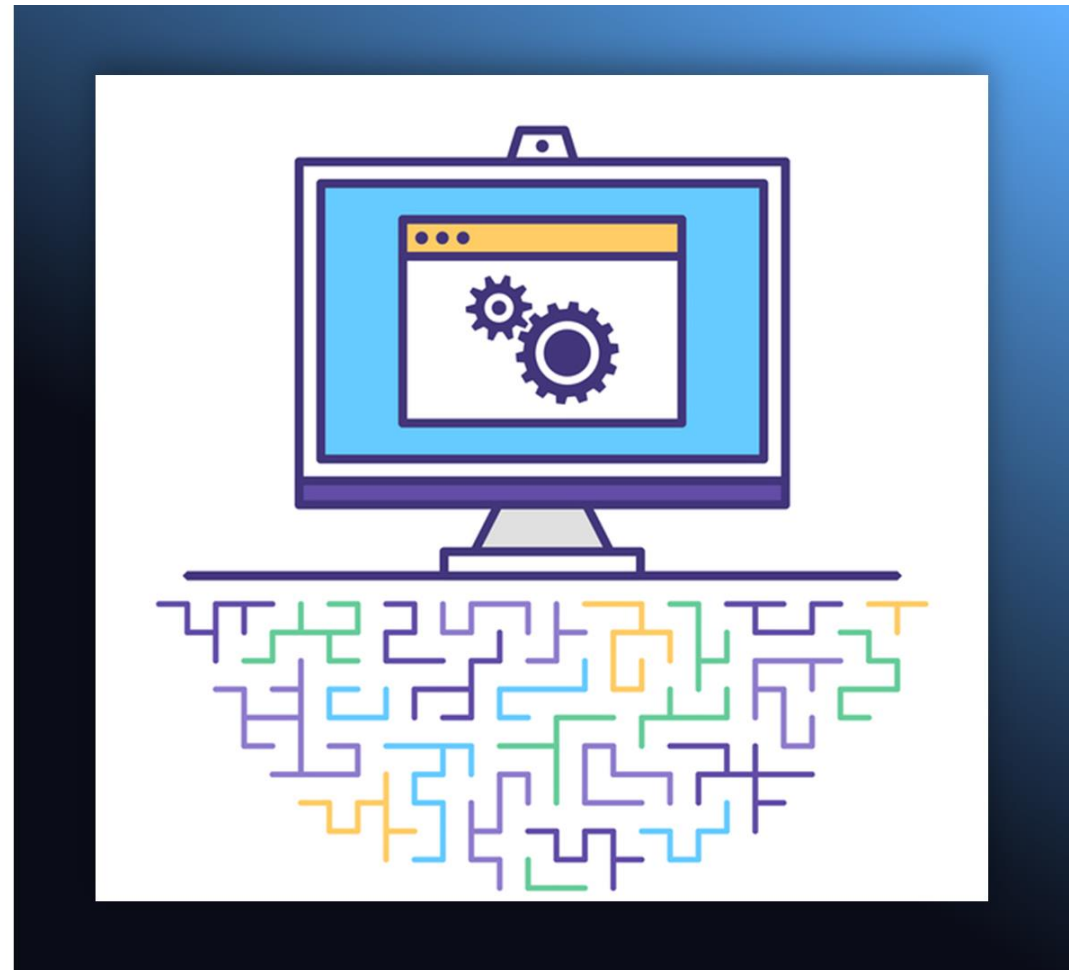
## Learning Objectives

Differentiate between regression and classification use cases

Quantify what is model success

Slide 4: What is Machine Learning?

# What Is Machine Learning?

A broad array of techniques that learns patterns and data, without being explicitly programmed
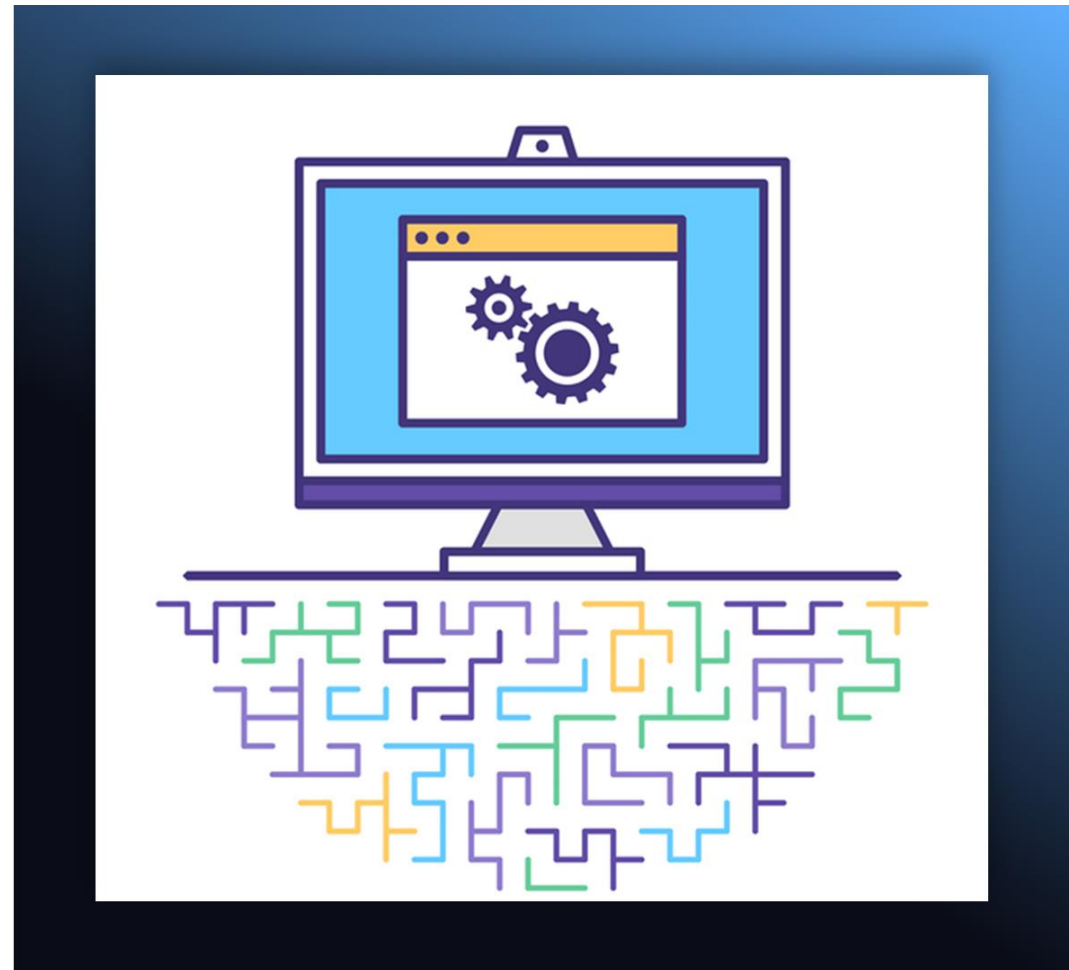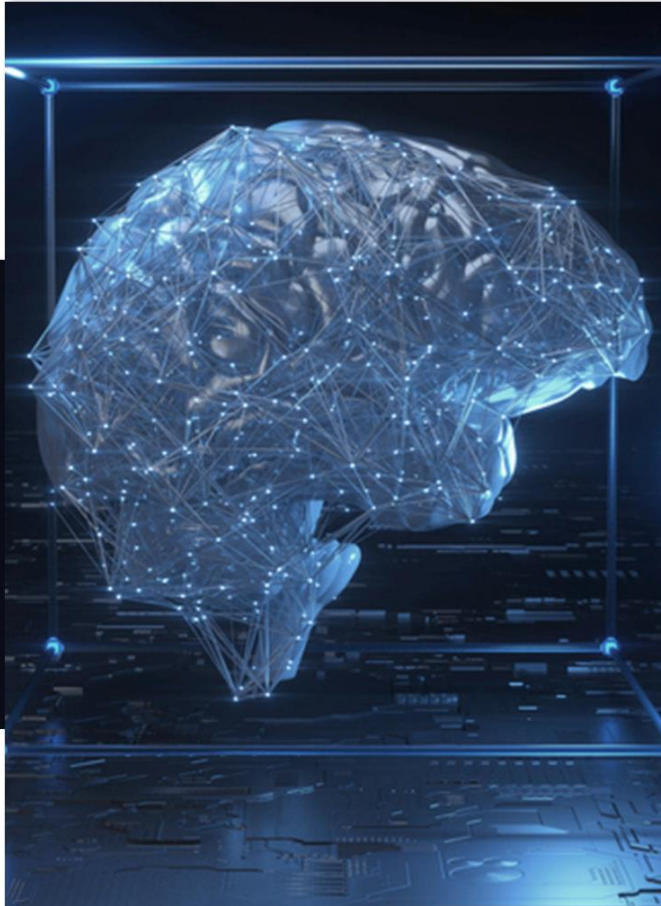
Slide 5: Churn Analysis

Slide 6: What is Machine Learning?

# What is Machine Learning?

A function that maps features to an output

Slide 7: Types of Machine Learning



Types of Machine Learning
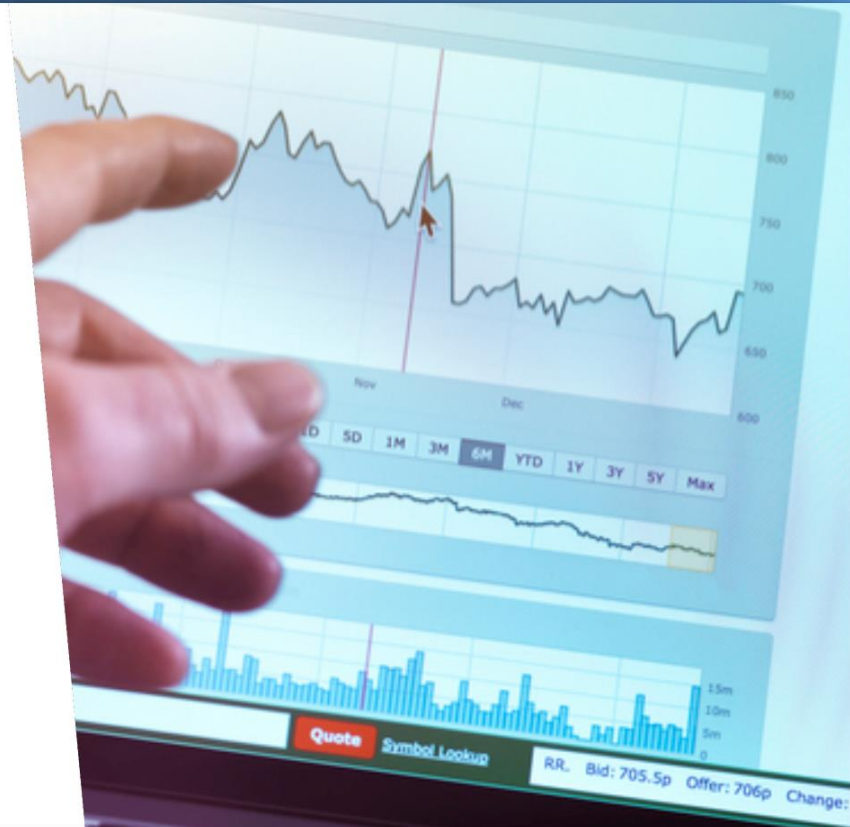
**Supervised**

**Unsupervised**

Reinforcement

Semi-supervised

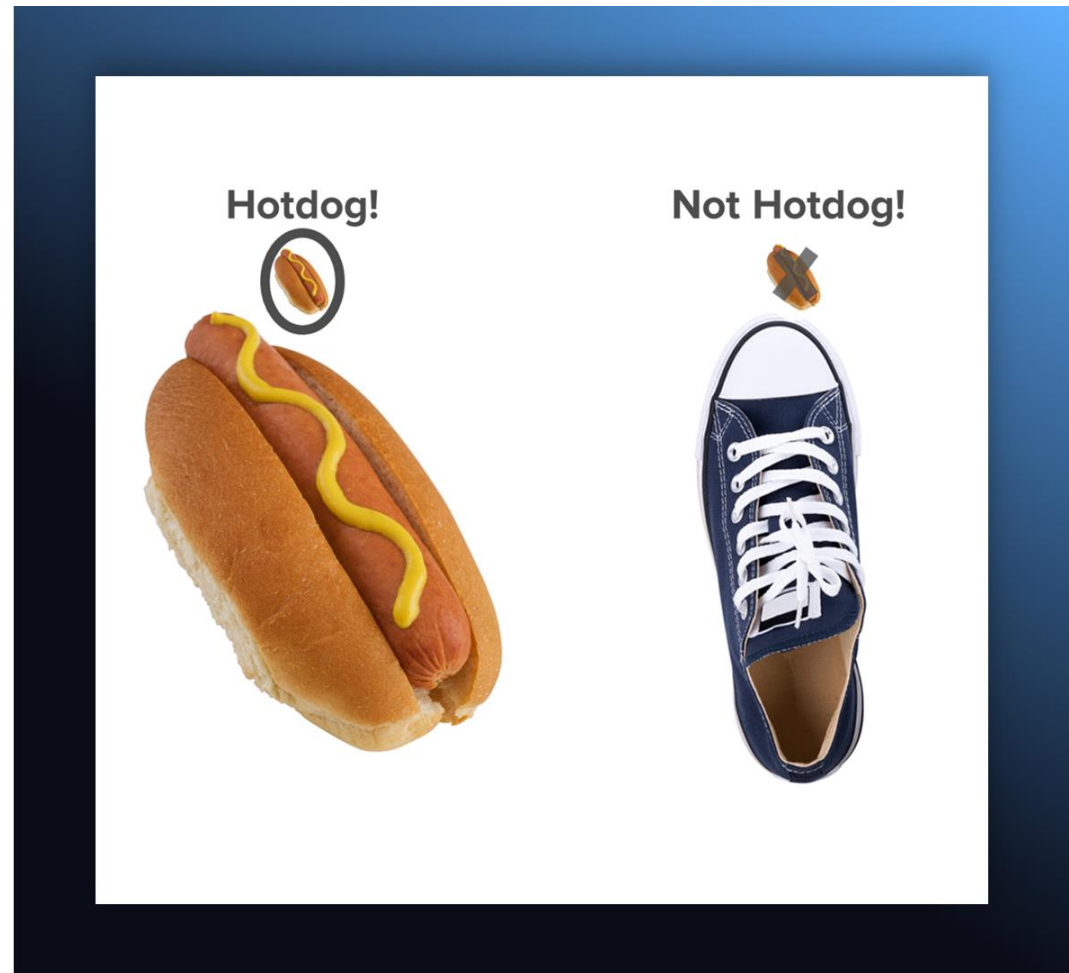Slide 8: Supervised Machine Learning

Slide 9: Classification Tasks

# Classification Tasks

Predicts a discrete set of categories

Binary classification
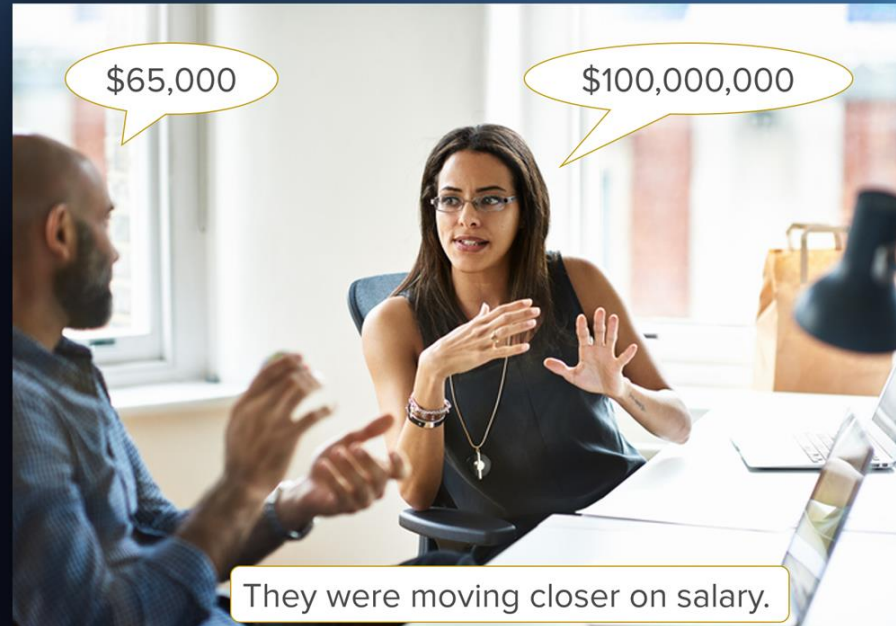
Multiclass classification
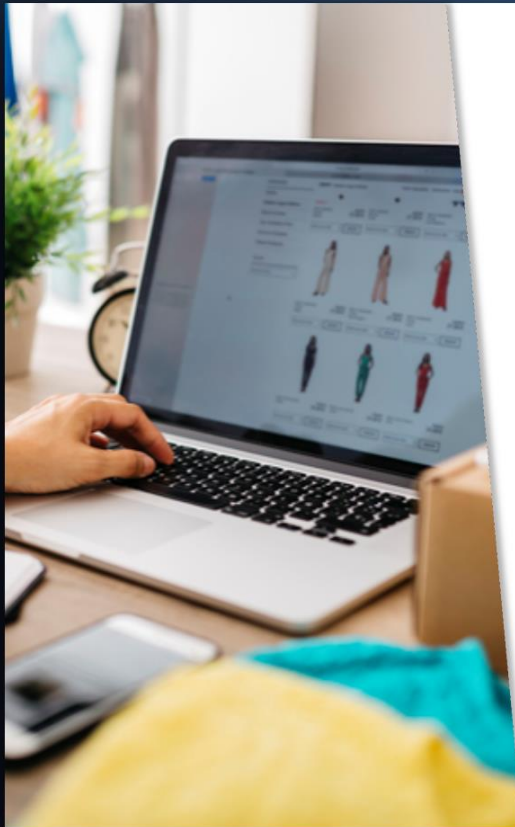
Slide 10: Regression Tasks

# Regression Tasks

Predict a continuous value

Financial forecasting

Unbounded number rather than a category

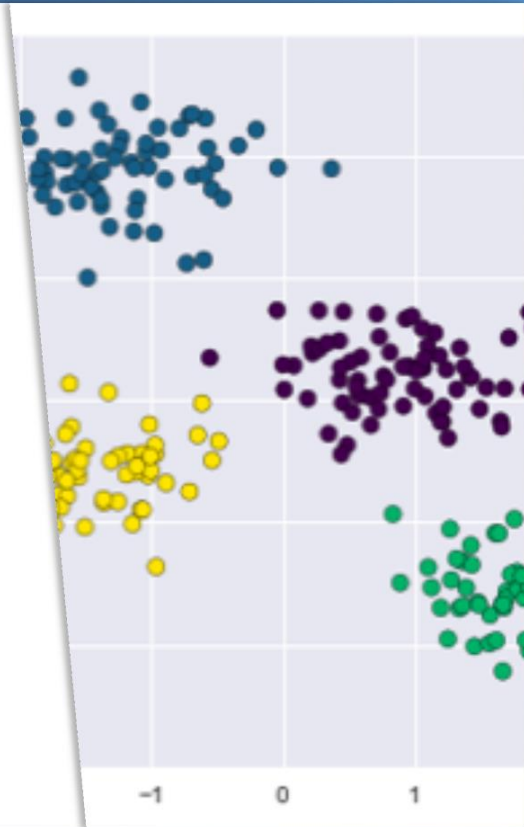Slide 11: Unsupervised Machine Learning

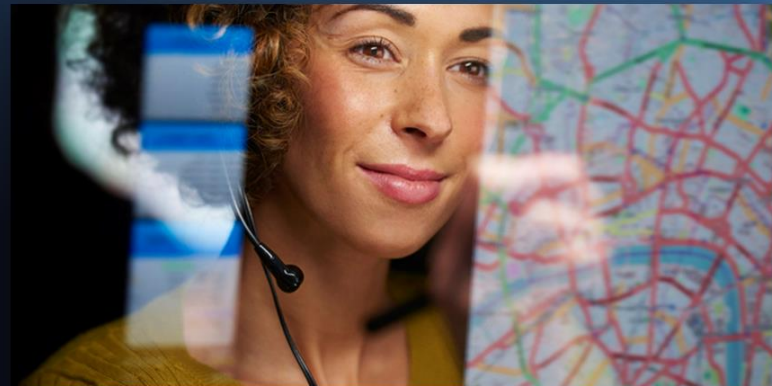Slide 12: Applying Machine Learning – Fire Call Dataset

Slide 13: Calculating Error

## Calculating Error

Predict response times

Look at the difference between predicted and true values

$$Error = (y_i - \hat{y}_i)$$
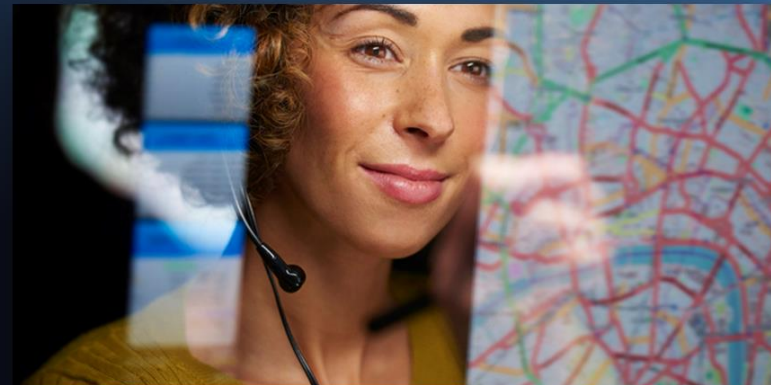
Slide 14: Root Means Squared Error

Slide 15: Compute the Sum of the Squared Error



## Compute the Sum of the Squared Error

The lower the RMSE, the better

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

Slide 16: Compute the Sum of the Squared Error

$$SSE = \sum_{i=1}^{n} (y_i - \widehat{y_i})^2$$

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \widehat{y_i})^2$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \widehat{y_i})^2}$$

Slide 17: Applying Machine Learning – Fire Call Dataset

Slide 18: Why Have a Baseline Model for Comparison?

Slide 19: Coming Up

## Coming Up

How to build a regression model