



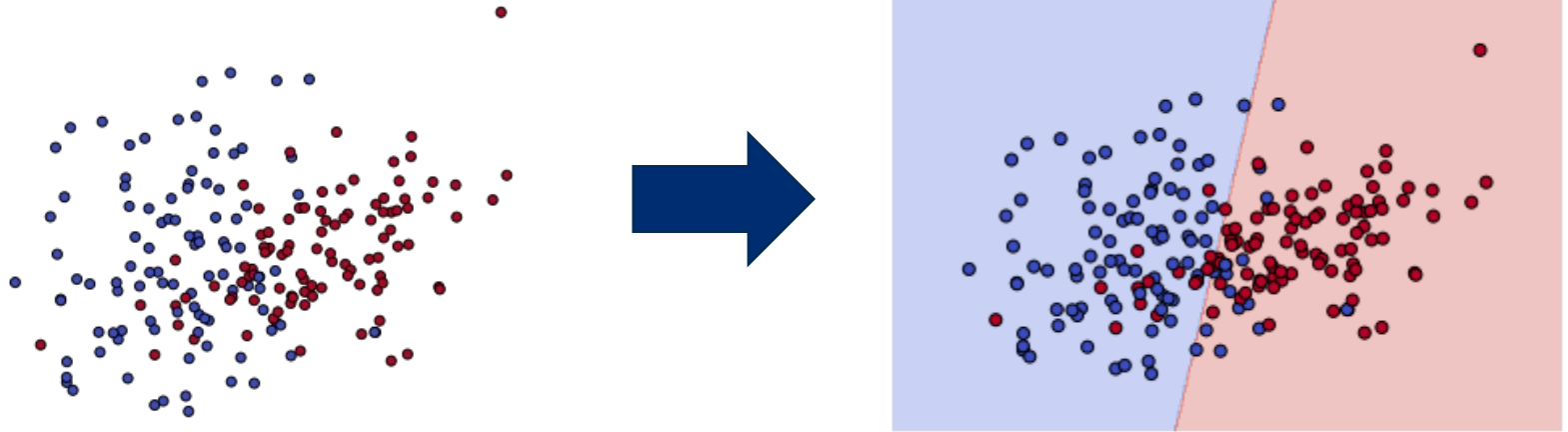
JOHNS HOPKINS

WHITING SCHOOL  
of ENGINEERING

# 685.621 Algorithms for Data Science

Supervised Learning: Classification Pipeline

# The Classification Pipeline



# Step 1: Preparing the Data for Classification

## Structured

- CSV, TSV, Databases

## Unstructured

- Text, Images, Video

## Preprocessing Tasks

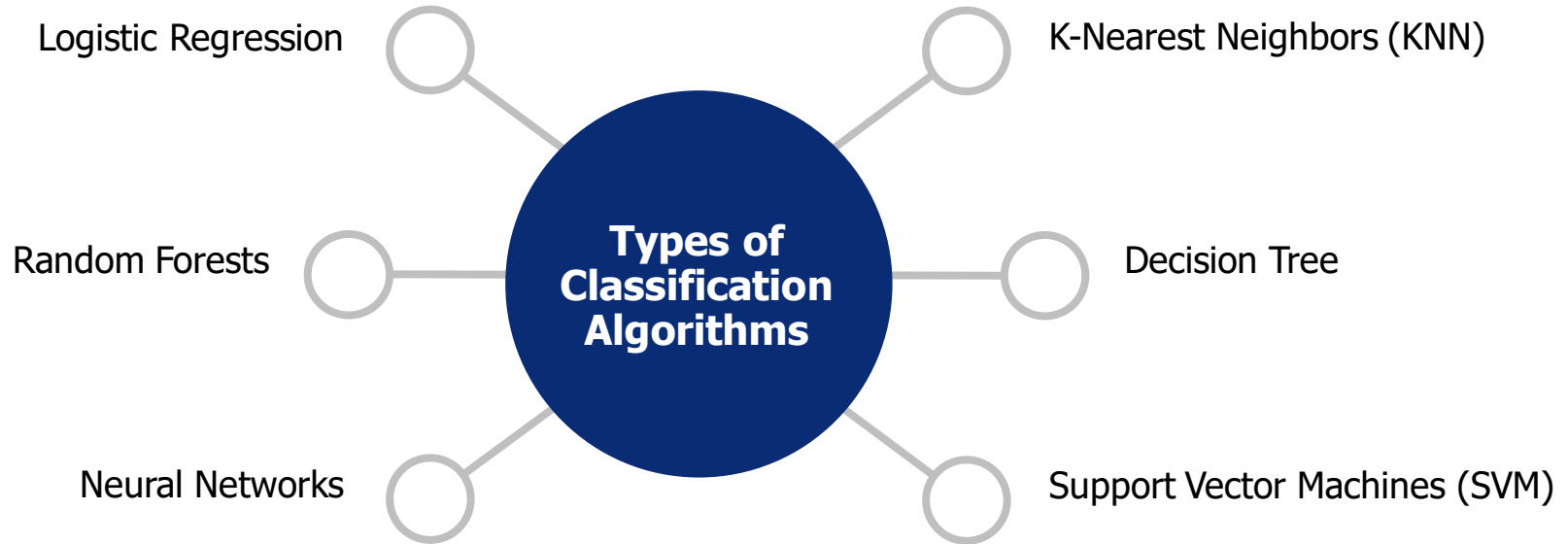
- Handling **missing values** (imputation, removal)
- Encoding **categorical values** (Label encoding, One-Hot Encoding)
- Scaling **numerical features** (Standardization, Normalization)
- Balancing **imbalanced datasets** (SMOTE)

# Step 2: Extract Meaningful Features

**Feature Engineering** – Transforming raw data into useful input features.

- **Types of Features:**
  - **Numerical Features** (e.g., Age, Salary)
  - **Categorical Features** (e.g., Gender, Product Category)
  - **Derived Features** (e.g., Total Purchase / Year, BMI)
- **Dimensionality Reduction:**
  - **Principal Component Analysis** (PCA)
  - **Feature Selection** (Eigenvalue Decomposition, Feature Importance, Fisher's Linear Discriminant Ratio)

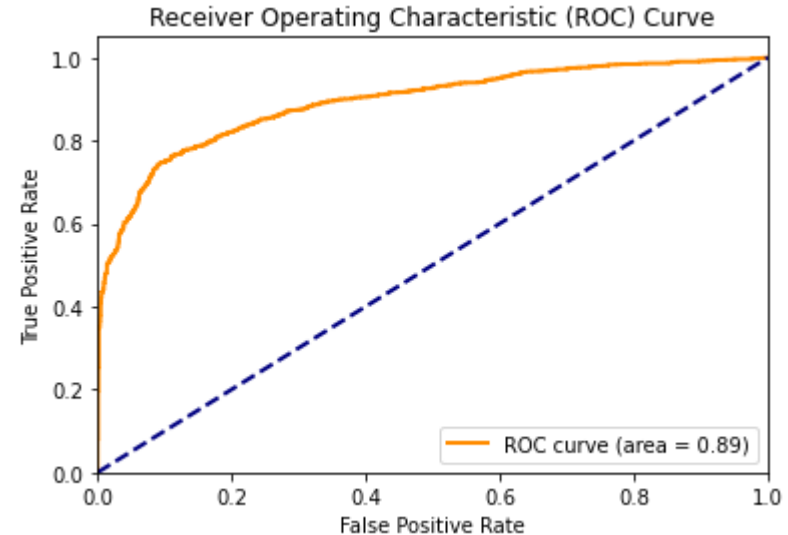
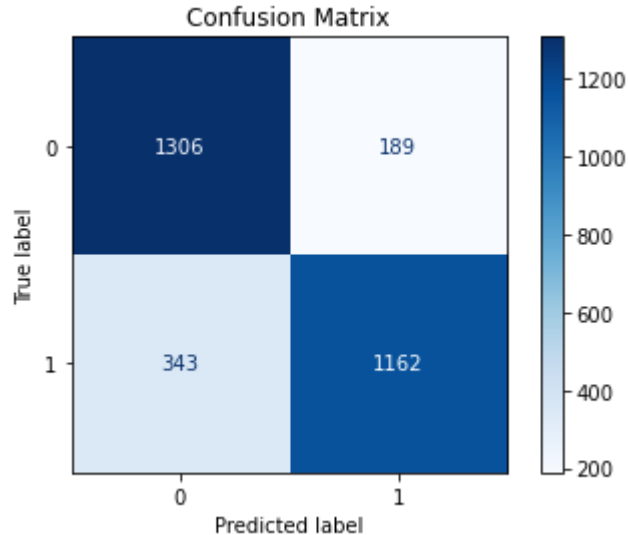
# Step 3: Choosing the Right Algorithm



# Step 4: Training & Validating the Model

<b>K-FOLD CROSS VALIDATION</b>	<b>Fold 1</b>	<b>Fold 2</b>	<b>Fold 3</b>	<b>Fold 4</b>	<b>Fold 5</b>
<b>Experiment 1</b>	TRAIN	TRAIN	TRAIN	TRAIN	<b>TEST</b>
<b>Experiment 2</b>	TRAIN	TRAIN	TRAIN	<b>TEST</b>	TRAIN
<b>Experiment 3</b>	TRAIN	TRAIN	<b>TEST</b>	TRAIN	TRAIN
<b>Experiment 4</b>	TRAIN	<b>TEST</b>	TRAIN	TRAIN	TRAIN
<b>Experiment 5</b>	<b>TEST</b>	TRAIN	TRAIN	TRAIN	TRAIN

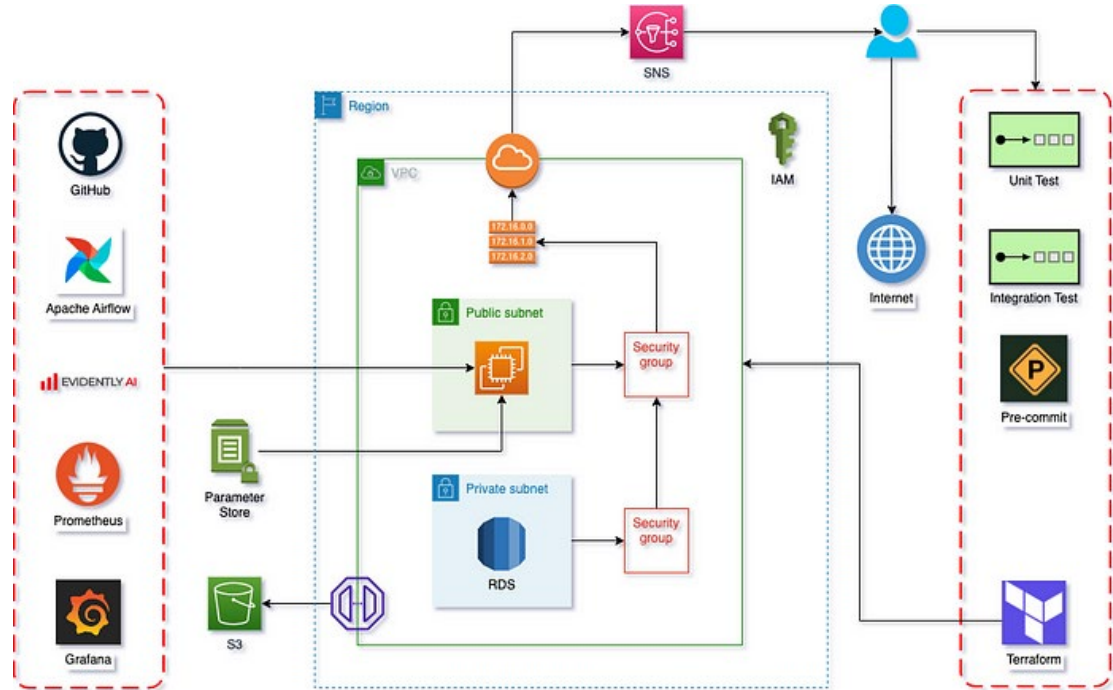
# Step 5: Evaluation Model Performance



# Step 6: Making Predictions & Deployment

## Deployment

- Production-ready system
- Deploy Models
- Monitor Model Performance



How to Put an ML Model into Production





# JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

© The Johns Hopkins University 2024, All Rights Reserved.