



## ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)

# ◊ CLUSTERING TECHNIQUES: K-MEANS, HIERARCHICAL CLUSTERING

### 📌 CHAPTER 1: INTRODUCTION TO CLUSTERING TECHNIQUES

#### ◆ 1.1 WHAT IS CLUSTERING?

**CLUSTERING IS AN UNSUPERVISED MACHINE LEARNING TECHNIQUE USED TO GROUP SIMILAR DATA POINTS TOGETHER BASED ON **SIMILARITY PATTERNS**. UNLIKE CLASSIFICATION, WHERE LABELS ARE PREDEFINED, CLUSTERING ALGORITHMS **DISCOVER HIDDEN STRUCTURES IN DATA WITHOUT PRIOR KNOWLEDGE**.**

#### WHY IS CLUSTERING IMPORTANT?

- ✓ **DATA SEGMENTATION** – GROUPS SIMILAR DATA POINTS FOR FURTHER ANALYSIS.
- ✓ **PATTERN DISCOVERY** – IDENTIFIES UNDERLYING STRUCTURES IN UNSTRUCTURED DATA.
- ✓ **ANOMALY DETECTION** – HELPS FIND OUTLIERS IN FINANCIAL FRAUD, CYBERSECURITY, ETC.
- ✓ **DATA REDUCTION** – SIMPLIFIES LARGE DATASETS BY GROUPING SIMILAR ITEMS TOGETHER.

### 📌 EXAMPLE APPLICATIONS:

- **CUSTOMER SEGMENTATION** – GROUPING CUSTOMERS BASED ON PURCHASING BEHAVIOR.
- **DOCUMENT CLUSTERING** – AUTOMATICALLY CATEGORIZING NEWS ARTICLES.
- **ANOMALY DETECTION** – IDENTIFYING FRAUDULENT TRANSACTIONS IN BANKING.

### CONCLUSION:

CLUSTERING IS WIDELY USED IN **MARKETING, HEALTHCARE, CYBERSECURITY, AND BIOINFORMATICS** TO DISCOVER HIDDEN PATTERNS IN LARGE DATASETS.

#### ◆ **1.2 TYPES OF CLUSTERING TECHNIQUES**

CLUSTERING TECHNIQUES CAN BE BROADLY CATEGORIZED INTO THE FOLLOWING TYPES:

##### **1. CENTROID-BASED CLUSTERING:**

- **K-MEANS CLUSTERING** – GROUPS DATA INTO **K CLUSTERS** USING CENTROIDS.
- **K-MEDOIDS CLUSTERING** – USES ACTUAL DATA POINTS AS CENTROIDS.

##### **2. HIERARCHICAL CLUSTERING:**

- **AGGLOMERATIVE CLUSTERING** – A BOTTOM-UP APPROACH, MERGING SMALL CLUSTERS INTO LARGER ONES.
- **DIVISIVE CLUSTERING** – A TOP-DOWN APPROACH, SPLITTING ONE LARGE CLUSTER INTO SMALLER CLUSTERS.

### 3. DENSITY-BASED CLUSTERING:

- **DBSCAN** – GROUPS DENSELY PACKED POINTS AND IDENTIFIES OUTLIERS.
- **OPTICS** – AN ADVANCED VERSION OF DBSCAN FOR VARYING DENSITIES.

#### 📌 EXAMPLE:

- **K-MEANS** IS USED FOR **SEGMENTING CUSTOMER BEHAVIOR IN MARKETING**.
- **HIERARCHICAL CLUSTERING** IS USED IN **GENE EXPRESSION ANALYSIS** IN BIOINFORMATICS.

#### 💡 CONCLUSION:

DIFFERENT CLUSTERING TECHNIQUES ARE USED BASED ON THE DATASET CHARACTERISTICS AND PROBLEM REQUIREMENTS.

### 📌 CHAPTER 2: K-MEANS CLUSTERING

#### ◆ 2.1 WHAT IS K-MEANS CLUSTERING?

**K-MEANS CLUSTERING** IS A CENTROID-BASED ALGORITHM THAT PARTITIONS DATA INTO **K CLUSTERS**, WHERE EACH DATA POINT BELONGS TO THE NEAREST CLUSTER CENTROID.

#### KEY FEATURES OF K-MEANS CLUSTERING:

- ✓ **FAST & SCALABLE** – WORKS WELL WITH LARGE DATASETS.
- ✓ **REQUIRES K VALUE** – THE NUMBER OF CLUSTERS MUST BE PREDEFINED.

- ✓ **SENSITIVE TO OUTLIERS** – CAN BE AFFECTED BY EXTREME VALUES.
- ✓ **WORKS BEST WITH SPHERICAL CLUSTERS** – ASSUMES CLUSTERS HAVE EQUAL VARIANCE.

### ❖ EXAMPLE APPLICATIONS:

- **CUSTOMER SEGMENTATION** – GROUPING CUSTOMERS BASED ON SPENDING HABITS.
- **IMAGE SEGMENTATION** – IDENTIFYING OBJECTS IN AN IMAGE.

### 💡 CONCLUSION:

K-MEANS IS A POPULAR CLUSTERING TECHNIQUE FOR QUICK AND EFFICIENT SEGMENTATION OF LARGE DATASETS.

#### ◆ 2.2 HOW K-MEANS WORKS

K-MEANS OPERATES USING THE FOLLOWING STEPS:

1. CHOOSE THE NUMBER OF CLUSTERS (K)
2. INITIALIZE K CENTROIDS RANDOMLY
- 3 ASSIGN EACH DATA POINT TO THE NEAREST CENTROID
4. RECALCULATE CENTROIDS BASED ON NEW CLUSTER MEMBERS
5. REPEAT UNTIL CENTROIDS NO LONGER CHANGE (CONVERGENCE)

## MATHEMATICAL REPRESENTATION

THE OBJECTIVE OF K-MEANS IS TO MINIMIZE THE WITHIN-CLUSTER SUM OF SQUARES (WCSS):

$$J = \sum_{i=1}^K \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

where:

- $K$  = Number of clusters
- $C_i$  = Cluster  $i$
- $\mu_i$  = Centroid of cluster  $i$

### 📌 EXAMPLE IN PYTHON:

```
IMPORT NUMPY AS NP
```

```
IMPORT MATPLOTLIB.PYTHON AS PLT
```

```
FROM SKLEARN.CLUSTER IMPORT KMEANS
```

```
# SAMPLE DATA
```

```
X = NP.ARRAY([[1, 2], [2, 3], [5, 6], [8, 8], [10, 12]])
```

```
# K-MEANS MODEL
```

```
KMEANS = KMEANS(N_CLUSTERS=2, RANDOM_STATE=42)
```

```
KMEANS.FIT(X)
```

## # PLOTTING CLUSTERS

```
PLT.SCATTER(X[:, 0], X[:, 1], C=KMEANS.LABELS_, CMAP='VIRIDIS')

PLT.SCATTER(KMEANS.CLUSTER_CENTERS_[:, 0],
            KMEANS.CLUSTER_CENTERS_[:, 1], S=200, MARKER='X', COLOR='RED')

PLT.TITLE("K-MEANS CLUSTERING")

PLT.SHOW()
```

 CONCLUSION:

K-MEANS IS AN EFFICIENT AND WIDELY USED CLUSTERING ALGORITHM, BUT CHOOSING THE RIGHT K VALUE IS CRUCIAL.

## ◆ 2.3 CHOOSING THE OPTIMAL K: THE ELBOW METHOD

THE ELBOW METHOD HELPS DETERMINE THE BEST NUMBER OF CLUSTERS (K) BY PLOTTING THE WITHIN-CLUSTER SUM OF SQUARES (WCSS).

 EXAMPLE IN PYTHON:

```
WCSS = []
```

```
FOR K IN RANGE(1, 11):
```

```
KMEANS = KMEANS(N_CLUSTERS=K, RANDOM_STATE=42)
```

```
KMEANS.FIT(X)
```

```
WCSS.APPEND(KMEANS.INERTIA_)
```

```
PLT.PLOT(RANGE(1, 11), WCSS, MARKER='O')

PLT.XLABEL("NUMBER OF CLUSTERS (K)")

PLT.YLABEL("WCSS")

PLT.TITLE("ELBOW METHOD FOR OPTIMAL K")

PLT.SHOW()
```

### CONCLUSION:

THE ELBOW METHOD PROVIDES A VISUAL WAY TO CHOOSE K, ENSURING OPTIMAL CLUSTERING.

## CHAPTER 3: HIERARCHICAL CLUSTERING

### ◆ 3.1 WHAT IS HIERARCHICAL CLUSTERING?

HIERARCHICAL CLUSTERING BUILDS A TREE-LIKE STRUCTURE (DENDROGRAM) THAT REPRESENTS NESTED CLUSTERS. UNLIKE K-MEANS, IT DOES NOT REQUIRE PREDEFINING THE NUMBER OF CLUSTERS.

### TYPES OF HIERARCHICAL CLUSTERING:

- ✓ AGGLOMERATIVE CLUSTERING (BOTTOM-UP) – EACH POINT STARTS AS ITS OWN CLUSTER, AND SIMILAR CLUSTERS ARE MERGED.
- ✓ DIVISIVE CLUSTERING (TOP-DOWN) – STARTS WITH ONE CLUSTER AND RECURSIVELY SPLITS INTO SMALLER CLUSTERS.

## EXAMPLE APPLICATIONS:

- CUSTOMER SEGMENTATION – CREATING CUSTOMER PROFILES BASED ON MULTIPLE FACTORS.

- **GENE EXPRESSION ANALYSIS – GROUPING SIMILAR GENES BASED ON PATTERNS.**

### CONCLUSION:

HIERARCHICAL CLUSTERING IS USEFUL FOR EXPLORATORY DATA ANALYSIS AND DOES NOT REQUIRE PREDEFINED CLUSTERS.

### ◆ **3.2 HOW HIERARCHICAL CLUSTERING WORKS**

HIERARCHICAL CLUSTERING FOLLOWS THESE STEPS:

- 1 COMPUTE THE **DISTANCE MATRIX** BETWEEN DATA POINTS.
- 2 MERGE THE CLOSEST CLUSTERS BASED ON **LINKAGE CRITERIA** (SINGLE, COMPLETE, OR AVERAGE).
- 3 REPEAT UNTIL ONLY ONE CLUSTER REMAINS.
- 4 VISUALIZE USING A **DENDROGRAM** TO DETERMINE THE OPTIMAL NUMBER OF CLUSTERS.

### EXAMPLE IN PYTHON:

```
IMPORT SCIPY.CLUSTER.HIERARCHY AS SCH
```

```
IMPORT MATPLOTLIB.PYPLOT AS PLT
```

```
FROM SKLEARN.CLUSTER IMPORT AGGLOMERATIVECLUSTERING
```

```
# SAMPLE DATA
```

```
X = NP.ARRAY([[1, 2], [2, 3], [5, 6], [8, 8], [10, 12]])
```

```
# DENDROGRAM
```

```
PLT.FIGURE(figsize=(8, 6))
```

```
SCH.DENDROGRAM(SCH.LINKAGE(X, method='WARD'))
```

```
PLT.TITLE("DENDROGRAM FOR HIERARCHICAL CLUSTERING")
```

```
PLT.XLABEL("DATA POINTS")
```

```
PLT.YLABEL("DISTANCE")
```

```
PLT.SHOW()
```

```
# AGGLOMERATIVE CLUSTERING
```

```
HC = AGGLOMERATIVECLUSTERING(n_clusters=2)
```

```
Y_HC = HC.FIT_PREDICT(X)
```

```
# PLOT CLUSTERS
```

```
PLT.SCATTER(X[:, 0], X[:, 1], c=Y_HC, cmap='VIRIDIS')
```

```
PLT.TITLE("HIERARCHICAL CLUSTERING")
```

```
PLT.SHOW()
```

### CONCLUSION:

HIERARCHICAL CLUSTERING IS USEFUL FOR UNDERSTANDING  
RELATIONSHIPS IN DATA AND PROVIDES A VISUAL WAY TO DETERMINE  
CLUSTERS.

 **FINAL THOUGHTS:**

- **K-MEANS IS BEST FOR LARGE DATASETS AND FAST CLUSTERING.**
- **HIERARCHICAL CLUSTERING IS BEST FOR SMALL DATASETS WHERE VISUAL RELATIONSHIPS MATTER.**

ISDM-NxT

# ◊ DIMENSIONALITY REDUCTION: PRINCIPAL COMPONENT ANALYSIS (PCA), T-SNE

## 📌 CHAPTER 1: INTRODUCTION TO DIMENSIONALITY REDUCTION

### ◆ 1.1 What is Dimensionality Reduction?

Dimensionality reduction is the process of **reducing the number of features (variables) in a dataset** while preserving as much information as possible. It helps in simplifying models, reducing computation time, and eliminating redundant data.

### Why is Dimensionality Reduction Important?

- ✓ **Improves model performance** – Reducing dimensions eliminates noise and speeds up training.
- ✓ **Avoids the curse of dimensionality** – Too many features can lead to overfitting and high computational costs.
- ✓ **Enhances data visualization** – High-dimensional data can be projected into 2D or 3D for better understanding.
- ✓ **Removes redundancy** – Eliminates correlated or less informative features.

### Types of Dimensionality Reduction Techniques

1. **Feature Selection** – Selecting a subset of important features.
2. **Feature Extraction** – Transforming original features into new reduced dimensions.

### 📌 Example:

A dataset with **1000 features (columns)** may slow down machine

learning models. Dimensionality reduction helps **reduce it to 20-50 meaningful features** without losing important patterns.

### Conclusion:

Dimensionality reduction helps in **data compression, efficient computation, and model optimization**.

---

## CHAPTER 2: PRINCIPAL COMPONENT ANALYSIS (PCA)

### 2.1 What is PCA?

**Principal Component Analysis (PCA)** is a technique that transforms a high-dimensional dataset into a lower-dimensional space while retaining most of the variance in the data. It identifies the **principal components**, which are new, uncorrelated variables capturing the most significant trends in the dataset.

### How PCA Works:

- ✓ **Step 1:** Standardize the data (mean = 0, variance = 1).
  - ✓ **Step 2:** Compute the **covariance matrix** to identify relationships between variables.
  - ✓ **Step 3:** Calculate **eigenvalues and eigenvectors** to determine the principal components.
  - ✓ **Step 4:** Choose the top **k principal components** based on explained variance.
  - ✓ **Step 5:** Project the data onto the new lower-dimensional space.
- 

### 2.2 Implementing PCA in Python

#### Step 1: Import Required Libraries

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
from sklearn.decomposition import PCA  
from sklearn.preprocessing import StandardScaler  
from sklearn.datasets import load_iris
```

### Step 2: Load Dataset & Standardize Features

```
# Load the Iris dataset  
  
iris = load_iris()  
  
X = iris.data # Feature matrix
```

```
# Standardize the data  
  
scaler = StandardScaler()  
  
X_scaled = scaler.fit_transform(X)
```

### Step 3: Apply PCA & Select Components

```
# Apply PCA (choosing 2 components)  
  
pca = PCA(n_components=2)  
  
X_pca = pca.fit_transform(X_scaled)  
  
  
# Print explained variance ratio  
  
print("Explained Variance Ratio:", pca.explained_variance_ratio_)
```

## Step 4: Visualizing PCA Components

```
# Scatter plot of PCA results  
  
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=iris.target, cmap='viridis')  
  
plt.xlabel("Principal Component 1")  
  
plt.ylabel("Principal Component 2")  
  
plt.title("PCA Visualization of Iris Dataset")  
  
plt.colorbar(label="Classes")  
  
plt.show()
```

📌 **Example:**

If the **original dataset has 4 features**, PCA can **reduce it to 2 dimensions**, capturing **95% of the total variance**.

💡 **Conclusion:**

PCA is useful for **reducing high-dimensional data while retaining key patterns** for machine learning models.

---

📌 **CHAPTER 3: T-DISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING (t-SNE)**

◆ **3.1 What is t-SNE?**

**t-SNE (t-Distributed Stochastic Neighbor Embedding)** is a **non-linear dimensionality reduction technique** used for **visualizing high-dimensional data** in 2D or 3D. Unlike PCA, which focuses on variance, **t-SNE preserves local relationships** between data points.

**How t-SNE Works:**

- ✓ **Step 1:** Compute pairwise similarities between points in high dimensions.
- ✓ **Step 2:** Assign probabilities based on similarity (using Gaussian distributions).
- ✓ **Step 3:** Reduce dimensions while preserving **local structure** in the data.
- ✓ **Step 4:** Optimize using **Kullback-Leibler divergence** to minimize information loss.

📌 **Example:**

t-SNE is used to **visualize handwritten digits** from the **MNIST dataset**, where each digit exists in a **high-dimensional space (784 features)**.

- ◆ **3.2 Implementing t-SNE in Python**

**Step 1: Import Required Libraries**

```
from sklearn.manifold import TSNE
```

**Step 2: Apply t-SNE on Iris Dataset**

```
# Apply t-SNE with 2 components
```

```
tsne = TSNE(n_components=2, random_state=42)
```

```
X_tsne = tsne.fit_transform(X_scaled)
```

```
# Scatter plot of t-SNE results
```

```
plt.scatter(X_tsne[:, 0], X_tsne[:, 1], c=iris.target, cmap='viridis')
```

```
plt.xlabel("t-SNE Component 1")
```

```

plt.ylabel("t-SNE Component 2")

plt.title("t-SNE Visualization of Iris Dataset")

plt.colorbar(label="Classes")

plt.show()

```

 **Example:**

t-SNE can reduce **1000-dimensional data** into **2D space** for visualization while maintaining cluster relationships.

 **Conclusion:**

t-SNE is **best for data visualization** in cases where **high-dimensional clusters** need to be identified.

 **CHAPTER 4: PCA vs. t-SNE**

Feature	PCA	t-SNE
Type	Linear	Non-Linear
Focus	Variance preservation	Local similarity preservation
Interpretability	Retains global structure	Retains local clusters
Computational Cost	Fast	Slow
Best Use Case	Feature extraction, reducing noise	Visualization of high-dimensional data

 **Example:**

✓ PCA is useful for **reducing features** in machine learning.

- 
- ✓ t-SNE is useful for **visualizing complex datasets like images and NLP embeddings.**
- 

📌 **CHAPTER 5: WHEN TO USE PCA vs. t-SNE?**

Scenario	Best Method
Reducing features while preserving data variance	PCA
Speeding up machine learning models	PCA
Understanding high-dimensional clusters	t-SNE
Visualizing data in 2D/3D	t-SNE
Data is linearly correlated	PCA
Data is highly non-linear (e.g., images, word embeddings)	t-SNE

💡 **Conclusion:**

- ✓ PCA is best for **feature reduction before machine learning models.**
  - ✓ t-SNE is best for **visualizing relationships in complex datasets.**
- 

📌 **SUMMARY & NEXT STEPS**

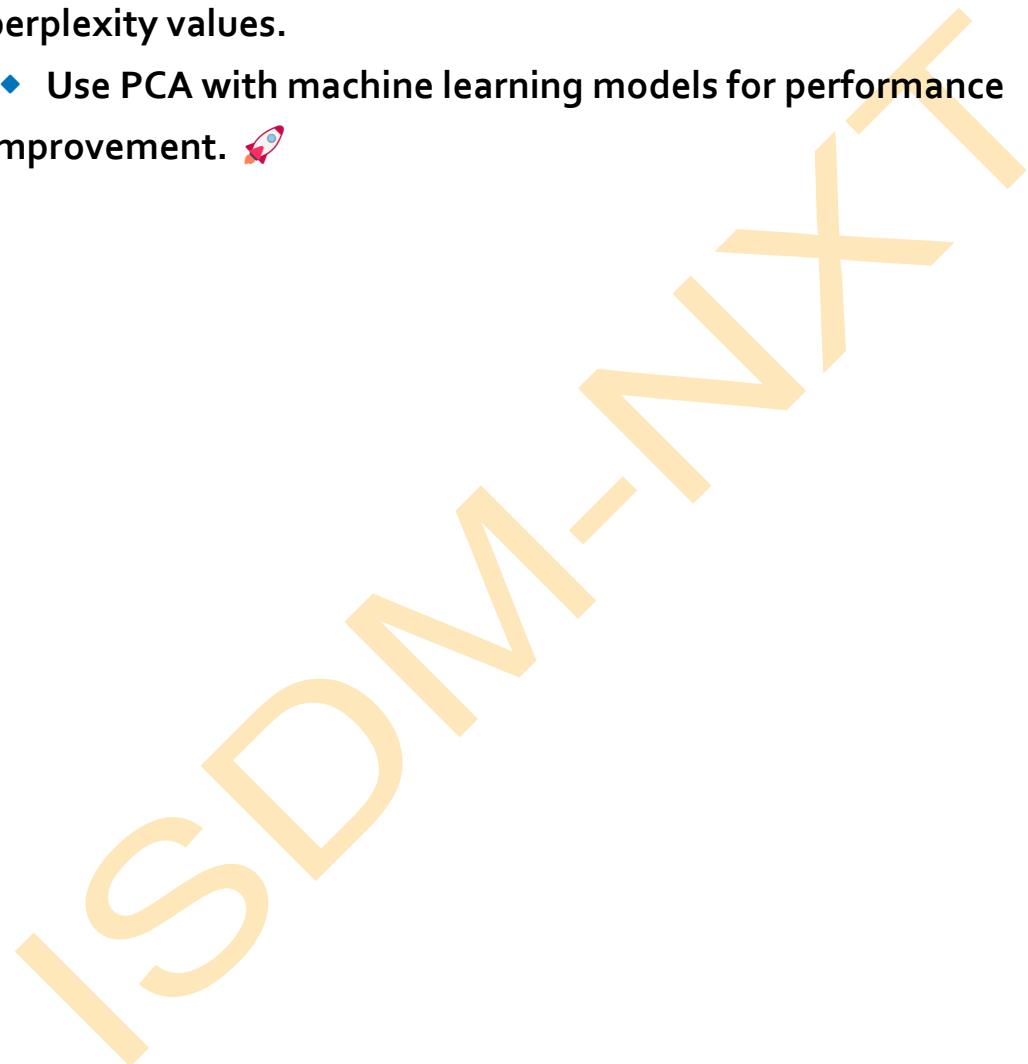
✓ **Key Takeaways:**

- ✓ PCA reduces dimensions while retaining **maximum variance.**
- ✓ t-SNE creates a **low-dimensional visualization** preserving local structure.

✓ Dimensionality reduction improves model efficiency and interpretability.

📌 **Next Steps:**

- ◆ Apply PCA & t-SNE on real-world datasets (e.g., MNIST, CIFAR-10, Customer Segmentation).
- ◆ Experiment with different PCA components and t-SNE perplexity values.
- ◆ Use PCA with machine learning models for performance improvement. 🚀



## ◊ ANOMALY DETECTION: FRAUD DETECTION IN BANKING

### 📌 CHAPTER 1: INTRODUCTION TO ANOMALY DETECTION

#### ◆ 1.1 What is Anomaly Detection?

**Anomaly Detection** is the process of identifying **rare or unusual patterns** in data that do not conform to expected behavior. These anomalies can indicate **fraud, cyberattacks, system failures, or unusual customer activities**.

#### Key Characteristics of Anomalies:

- ✓ **Rare Occurrence** – Anomalies occur infrequently in datasets.
- ✓ **Deviate from Normal Patterns** – They differ significantly from the majority of data points.
- ✓ **Can Indicate Fraud or System Failures** – In banking, anomalies often signal **fraudulent transactions**.

#### 📌 Example:

- **Credit Card Fraud Detection** – Detecting **suspicious transactions** based on abnormal spending behavior.
- **Bank Account Fraud** – Identifying **unusual login patterns or unauthorized access attempts**.

#### 💡 Conclusion:

Anomaly Detection is **essential for detecting fraud, preventing financial losses, and improving cybersecurity in banking systems**.

### ◆ **1.2 Importance of Anomaly Detection in Banking**

Banks and financial institutions handle **millions of transactions daily**, making fraud detection **critical for security and customer trust**.

#### **Why is Fraud Detection Important?**

- ✓ **Prevents Financial Losses** – Stops unauthorized transactions before they cause damage.
- ✓ **Protects Customers & Businesses** – Prevents identity theft and fraud schemes.
- ✓ **Ensures Regulatory Compliance** – Banks must comply with **anti-money laundering (AML) regulations**.
- ✓ **Enhances Customer Trust** – Secure banking leads to higher customer confidence.

#### **Example:**

A bank uses **machine learning models** to **detect anomalies in customer spending** and flag suspicious transactions.

#### **Conclusion:**

Financial fraud costs businesses **billions of dollars annually**, making anomaly detection **a necessity rather than a luxury**.

---

#### **CHAPTER 2: TYPES OF ANOMALIES IN BANKING FRAUD**

### ◆ **2.1 Categories of Anomalies in Banking**

There are **three main types of anomalies** in banking fraud detection:

#### **1. Point Anomalies**

- A **single transaction** that deviates significantly from the normal pattern.

📌 **Example:**

A customer with a history of small transactions suddenly makes a **\$10,000 purchase from a foreign country**.

## 2. Contextual Anomalies

- A transaction that may be normal in some contexts but **unusual in others**.

📌 **Example:**

A credit card transaction for **\$500 at a luxury store** may be normal in the U.S. but **suspicious if done in a high-risk country**.

## 3. Collective Anomalies

- A **group of transactions** that together indicate fraud.

📌 **Example:**

A customer's card is used **multiple times within minutes** at different locations, suggesting a stolen card.

💡 **Conclusion:**

Understanding different **types of anomalies** helps in **building fraud detection models that are precise and accurate**.

📌 **CHAPTER 3: TECHNIQUES FOR FRAUD DETECTION IN BANKING**

◆ **3.1 Traditional Fraud Detection Techniques**

Before machine learning, fraud detection relied on **rule-based systems** and manual reviews.

### Traditional Approaches:

- ✓ **Rule-Based Systems** – Predefined rules (e.g., transactions over \$5,000 require verification).
- ✓ **Threshold-Based Alerts** – Transactions exceeding a limit trigger security alerts.
- ✓ **Blacklisting** – Detecting transactions from known fraudulent accounts.

 **Example:**

A bank might flag **transactions over \$2,000 from a new location** as suspicious.

 **Conclusion:**

While traditional methods work, **they struggle with dynamic fraud tactics** and generate too many false positives.

---

◆ **3.2 Machine Learning for Anomaly Detection**

Machine Learning (ML) has **transformed fraud detection**, enabling models to learn patterns and detect anomalies **automatically**.

**Machine Learning Approaches for Fraud Detection:**

- ✓ **Supervised Learning** – Trains models using historical fraud cases.
- ✓ **Unsupervised Learning** – Detects anomalies **without labeled fraud cases**.
- ✓ **Hybrid Models** – Combine both approaches for **higher accuracy**.

 **Example:**

A **random forest classifier** can learn past fraudulent transactions and **predict future fraud** with high accuracy.

 **Conclusion:**

Machine learning models **detect complex fraud patterns** that traditional systems **often miss**.

## 📌 CHAPTER 4: MACHINE LEARNING ALGORITHMS FOR FRAUD DETECTION

### ◆ 4.1 Logistic Regression for Fraud Detection

Logistic Regression is a **simple but effective supervised learning model** for binary classification.

#### How It Works:

- ✓ Uses past fraud cases to **assign probability scores** to new transactions.
- ✓ Helps banks **classify transactions as fraudulent or legitimate**.

#### 📌 Example Implementation in Python:

```
from sklearn.linear_model import LogisticRegression  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score  
  
# Sample dataset  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,  
random_state=42)  
  
# Train Logistic Regression model  
model = LogisticRegression()  
model.fit(X_train, y_train)
```

```
# Predict and evaluate
```

```
y_pred = model.predict(X_test)  
print(f"Accuracy: {accuracy_score(y_test, y_pred):.2f}")
```

 **Conclusion:**

Logistic Regression is a **baseline model** that works well for **simple fraud detection problems**.

◆ **4.2 Decision Trees & Random Forests**

Decision Trees and Random Forests **handle non-linear fraud patterns effectively**.

**Advantages of Decision Trees:**

- ✓ **Interpretable** – Each decision node represents a rule.
- ✓ **Handles Mixed Data Types** – Works with categorical and numerical features.

 **Example:**

```
from sklearn.ensemble import RandomForestClassifier
```

```
# Train Random Forest model
```

```
rf_model = RandomForestClassifier(n_estimators=100,  
random_state=42)  
  
rf_model.fit(X_train, y_train)
```

```
# Predict and evaluate
```

```
y_pred = rf_model.predict(X_test)
```

```
print(f"Accuracy: {accuracy_score(y_test, y_pred):.2f}")
```

 **Conclusion:**

**Random Forests** improve accuracy by combining **multiple decision trees**, making them **powerful for fraud detection**.

---

◆ **4.3 Anomaly Detection with Autoencoders**

Autoencoders are **deep learning models** that learn to detect fraud by **reconstructing normal transactions** and flagging deviations.

**Why Use Autoencoders?**

✓ **Can Detect Unknown Fraud Patterns** – Learns what normal data looks like and flags deviations.

✓ **Unsupervised Learning** – No labeled fraud cases required.

 **Example:**

```
from tensorflow.keras.models import Model
```

```
from tensorflow.keras.layers import Input, Dense
```

```
# Define Autoencoder
```

```
input_layer = Input(shape=(X_train.shape[1],))
```

```
encoded = Dense(16, activation='relu')(input_layer)
```

```
decoded = Dense(X_train.shape[1], activation='sigmoid')(encoded)
```

```
autoencoder = Model(input_layer, decoded)  
autoencoder.compile(optimizer='adam', loss='mse')
```

```
# Train Autoencoder
```

```
autoencoder.fit(X_train, X_train, epochs=10, batch_size=32,  
validation_data=(X_test, X_test))
```

#### 💡 Conclusion:

Autoencoders **outperform traditional methods** when detecting new fraud patterns.

## 📌 CHAPTER 5: CHALLENGES & FUTURE OF FRAUD DETECTION

### ◆ 5.1 Challenges in Banking Fraud Detection

- ✓ **High False Positives** – Many legitimate transactions get flagged.
- ✓ **Adaptive Fraudsters** – Fraud tactics evolve over time.
- ✓ **Data Imbalance** – Fraud cases are **rare**, making training difficult.

#### 📌 Example:

Banks **continuously update** their models to **keep up with evolving fraud techniques**.

#### 💡 Conclusion:

AI-driven fraud detection **must constantly evolve** to tackle new fraud schemes.

# ◊ FEATURE ENGINEERING: HANDLING CATEGORICAL & NUMERICAL FEATURES, FEATURE SCALING



## CHAPTER 1: INTRODUCTION TO FEATURE ENGINEERING

### ◆ 1.1 What is Feature Engineering?

Feature Engineering is the process of **transforming raw data into meaningful features** that improve the performance of machine learning models. It involves **handling categorical and numerical features, encoding data, scaling values, and creating new features** to optimize model predictions.

### Why is Feature Engineering Important?

- ✓ Improves model accuracy and performance.
- ✓ Helps models learn better from data by making patterns clearer.
- ✓ Reduces complexity and enhances interpretability.
- ✓ Makes data compatible with machine learning algorithms.



### Example:

In a **house price prediction model**, converting a categorical feature like "**House Type**" (Apartment, Villa, Bungalow) into numerical values improves model performance. Similarly, **scaling numerical values** like square footage ensures fair weightage in the model.



### Conclusion:

Feature engineering is **a critical step in data preprocessing**, ensuring models **extract meaningful insights** from raw data.

### ◆ 1.2 Types of Features in Machine Learning

- ✓ **Numerical Features** – Represented as continuous or discrete numbers (e.g., age, price, temperature).
- ✓ **Categorical Features** – Represented as labels or categories (e.g., gender, city, color).
- ✓ **Ordinal Features** – Have a meaningful order (e.g., Low, Medium, High).
- ✓ **Nominal Features** – Categories without inherent order (e.g., countries, product types).

#### 📌 Example:

In a dataset of **employees**, features such as "**Salary**" (**numerical**), "**Job Role**" (**categorical**), and "**Experience Level**" (**ordinal**: Junior, Mid, Senior) need different processing techniques.

#### 💡 Conclusion:

Understanding feature types is essential for applying **appropriate transformations and scaling techniques**.

## 📌 CHAPTER 2: HANDLING CATEGORICAL FEATURES

### ◆ 2.1 What are Categorical Features?

Categorical features contain **labels or text values** that represent categories. Since machine learning models work with numerical data, categorical features **must be encoded into numerical values**.

#### 📌 Example:

A dataset contains the feature "**Marital Status**" with values:

- ✓ Single
- ✓ Married
- ✓ Divorced

Since machine learning models cannot process text, we convert them into numerical values.

---

### ◆ **2.2 Methods for Handling Categorical Data**

- ✓ **Label Encoding** – Assigns unique numeric values to categories.
- ✓ **One-Hot Encoding** – Converts each category into a separate binary column.
- ✓ **Ordinal Encoding** – Assigns ordered numerical values to ranked categories.
- ✓ **Binary Encoding** – Converts categories into binary representations.

#### 📌 **Example:**

For the "Education Level" feature (High School, Bachelor's, Master's, PhD):

- **Label Encoding:** High School → 0, Bachelor's → 1, Master's → 2, PhD → 3
- **One-Hot Encoding:** Creates separate columns: [High\_School, Bachelors, Masters, PhD] with 1s and 0s.
- **Ordinal Encoding:** If there's a meaningful order, map values like [0, 1, 2, 3].

#### 💡 **Conclusion:**

Choosing the right encoding method depends on whether the categories have **a natural order or should be treated as separate entities**.

---

## 📌 CHAPTER 3: HANDLING NUMERICAL FEATURES

### ◆ 3.1 What are Numerical Features?

Numerical features are **quantitative variables** that can be continuous or discrete. They require **scaling, transformation, and sometimes binning** to improve model performance.

- ✓ **Continuous Data** – Can take infinite values (e.g., weight, height, price).
- ✓ **Discrete Data** – Can take only specific values (e.g., number of children, employee ID).

#### 📌 Example:

A **house pricing model** has numerical features such as **square footage, number of bedrooms, and price**. These features might need **scaling or transformation** to improve model accuracy.

---

### ◆ 3.2 Handling Missing Numerical Values

- ✓ **Mean/Median Imputation** – Replacing missing values with the mean or median.
- ✓ **K-Nearest Neighbors (KNN) Imputation** – Predicting missing values based on similar observations.
- ✓ **Interpolation** – Filling missing values by estimating based on trends.

#### 📌 Example:

A dataset contains **missing age values**. Using **median imputation**, missing values are replaced with the dataset's median age.

### Conclusion:

Missing numerical values **should be handled carefully** to avoid biased results in the model.

---

## CHAPTER 4: FEATURE SCALING

### ◆ 4.1 What is Feature Scaling?

Feature Scaling is the process of **normalizing numerical features** to ensure that no feature dominates others in the model. It helps improve model performance by ensuring features are on a similar scale.

#### ✓ Why is Feature Scaling Needed?

- ✓ Prevents high magnitude features from dominating smaller ones.
- ✓ Improves convergence speed in optimization algorithms.
- ✓ Ensures that **distance-based models** (e.g., KNN, SVM) work correctly.

#### Example:

A dataset contains **age (1-100)** and **income (\$1,000-\$100,000)**. Without scaling, income values dominate the learning process.

---

### ◆ 4.2 Feature Scaling Techniques

#### Min-Max Scaling (Normalization)

**Formula:**

$$X_{scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

- ✓ Rescales values between 0 and 1.
- ✓ Used in neural networks and deep learning.
- ✓ Best for small-scale data without outliers.

 **Example:**

Scaling **house prices** between 0 and 1 to ensure equal weighting with other features.

### Standardization (Z-Score Scaling)

**Formula:**

$$X_{scaled} = \frac{X - \mu}{\sigma}$$

- ✓ Centers data around mean = 0 and standard deviation = 1.
- ✓ Used in logistic regression, SVMs, and neural networks.
- ✓ Best for datasets with outliers.

 **Example:**

Scaling **exam scores** with different means and variances to standardize student performance comparisons.

### Robust Scaling

- ✓ Less sensitive to outliers than Min-Max and Z-score scaling.
- ✓ Uses the median instead of mean for scaling.
- ✓ Suitable for data with extreme values.

 **Example:**

A salary dataset contains extreme salaries (outliers). Robust Scaling ensures the model isn't affected by high-income outliers.

◆ **4.3 When to Use Different Scaling Techniques?**

Scaling Method	Best For	Example Use Case
Min-Max Scaling	Small-scale data with no outliers	Image processing, deep learning
Standardization	Normal distribution, features with different scales	Linear regression, SVM
Robust Scaling	Data with extreme values	Salary predictions

 **Conclusion:**

Choosing the right scaling technique ensures efficient learning and improved model stability.

 **CHAPTER 5: PRACTICAL APPLICATIONS OF FEATURE ENGINEERING**

- ✓ **Finance:** Scaling transaction amounts for fraud detection.
- ✓ **Healthcare:** Encoding patient medical history into numerical features.

- ✓ **E-Commerce:** One-hot encoding product categories for recommendation models.
- ✓ **Marketing:** Transforming customer demographics for targeted advertising.

#### 📌 Example:

An **online store** uses **one-hot encoding** for customer **preferred payment methods**, ensuring smooth feature representation.

#### 💡 Conclusion:

Feature engineering **enhances machine learning model effectiveness**, leading to **better accuracy and insights**.

#### 📌 SUMMARY & NEXT STEPS

##### ✓ Key Takeaways:

- ✓ Categorical features must be **encoded** (**One-Hot, Label, Ordinal Encoding**).
- ✓ Numerical features require **scaling** (**Min-Max, Standardization, Robust Scaling**).
- ✓ Feature engineering enhances model accuracy and **improves predictive performance**.

##### 📌 Next Steps:

- ◆ Implement feature engineering techniques in Python using **Pandas & Scikit-Learn**.
- ◆ Work on real-world datasets (**Kaggle competitions**) to practice feature transformation.
- ◆ Explore advanced techniques like **Feature Selection** and **Principal Component Analysis (PCA)**. 

## ◊ DATA PIPELINES & MODEL DEPLOYMENT BASICS

### 📌 CHAPTER 1: INTRODUCTION TO DATA PIPELINES & MODEL DEPLOYMENT

#### ◆ 1.1 What is a Data Pipeline?

A **data pipeline** is a sequence of **automated processes** that **extract, transform, and load (ETL) data** from various sources to be used for analysis, machine learning, and model deployment. It ensures that **data flows efficiently from raw sources to actionable insights** in real-time or batch processing.

#### Key Components of a Data Pipeline:

- ✓ **Data Ingestion** – Collects raw data from multiple sources (databases, APIs, streaming).
- ✓ **Data Processing (ETL/ELT)** – Cleans, transforms, and structures data for analysis.
- ✓ **Data Storage** – Stores data in a **warehouse (BigQuery, Redshift)** or **data lake (S3, Hadoop)**.
- ✓ **Data Monitoring** – Ensures **data integrity, quality, and availability**.

#### 📌 Example:

A financial company builds a data pipeline to **extract stock market data**, process it in real-time, and feed it into a predictive analytics system.

### Conclusion:

Data pipelines **automate data movement** for seamless machine learning model deployment.

---

#### ◆ **1.2 What is Model Deployment?**

Model deployment is the process of **integrating a trained machine learning model into a production environment** to make real-time or batch predictions.

#### **Types of Model Deployment:**

- ✓ **Batch Deployment** – The model runs at scheduled intervals to process large datasets.
- ✓ **Real-Time Deployment** – The model serves predictions as requests arrive (API-based).
- ✓ **Edge Deployment** – The model runs on **IoT devices, mobile apps, or embedded systems**.

### Example:

A **fraud detection model** is deployed as a REST API to flag suspicious transactions in real time.

### Conclusion:

Model deployment enables **AI-powered decision-making** in real-world applications.

---

## CHAPTER 2: BUILDING A DATA PIPELINE

#### ◆ **2.1 Steps in a Data Pipeline**

A data pipeline follows a structured flow to ensure **smooth data transformation and availability**.

#### ✓ Step 1: Data Ingestion

- Extract data from **databases (PostgreSQL, MySQL)**, **APIs**, **files (CSV, JSON)**, or **streaming platforms (Kafka, AWS Kinesis)**.

#### ✓ Step 2: Data Processing (ETL or ELT)

- Clean data (handle missing values, duplicates, standardization).
- Transform data (convert formats, create new features).

#### ✓ Step 3: Data Storage

- Store structured data in **relational databases (SQL, Snowflake, Redshift)**.
- Store unstructured data in **data lakes (AWS S3, Google Cloud Storage)**.

#### ✓ Step 4: Model Training & Evaluation

- Train and optimize machine learning models using frameworks like **TensorFlow, PyTorch, or Scikit-Learn**.

#### ✓ Step 5: Model Deployment & Monitoring

- Deploy the model as a **REST API**, **containerized service (Docker)**, or **cloud-based model (AWS SageMaker, Vertex AI)**.

#### ❖ Example:

A **retail company** builds a data pipeline to collect customer purchase

data, process it, store it in BigQuery, and deploy an ML model to **recommend personalized products**.

 **Conclusion:**

A well-structured data pipeline ensures **scalability, efficiency, and reliability** in ML applications.

---

◆ **2.2 Tools for Data Pipelines**

Category	Tools
<b>Data Ingestion</b>	Apache Kafka, AWS Kinesis, Apache Airflow
<b>Data Processing (ETL/ELT)</b>	Pandas, Apache Spark, dbt
<b>Data Storage</b>	PostgreSQL, Snowflake, BigQuery, AWS S3
<b>Model Training</b>	TensorFlow, PyTorch, Scikit-Learn
<b>Deployment</b>	Flask, FastAPI, Docker, Kubernetes

 **Example:**

An **e-commerce company** uses **Kafka** to stream customer clicks, **Airflow** to schedule data pipelines, and **AWS S3** for storage before feeding data to an AI-powered recommendation engine.

 **Conclusion:**

Using the **right combination of tools** optimizes pipeline efficiency.

---

## 📌 CHAPTER 3: MODEL DEPLOYMENT BASICS

### ◆ 3.1 Steps to Deploy a Machine Learning Model

#### ✓ Step 1: Train & Save the Model

- Train the model using **Scikit-Learn, TensorFlow, or PyTorch**.
- Save the model as a **.pkl (Scikit-Learn), .h5 (Keras), or .pt (PyTorch) file**.

#### ✓ Step 2: Create a REST API for the Model

- Use **Flask or FastAPI** to expose the model as an API endpoint.

#### ✓ Step 3: Containerize the Model (Docker)

- Package the model into a Docker container for easy deployment.

#### ✓ Step 4: Deploy on Cloud (AWS, GCP, Azure)

- Use **AWS Lambda, Google Cloud Run, or Azure Functions** for scalable hosting.

#### ✓ Step 5: Model Monitoring & Retraining

- Set up **logging, performance tracking, and automatic retraining pipelines**.

### 📌 Example:

A company deploys a **sentiment analysis model** as a REST API using **Flask and Docker** and hosts it on **AWS Lambda** for real-time processing.

### 💡 Conclusion:

Model deployment ensures **seamless integration of AI into applications**.

### ◆ 3.2 Deployment Methods

Deployment Type	Best For	Examples
Local Deployment	Testing & Development	Running models on Jupyter Notebook
Cloud Deployment	Scalable ML services	AWS SageMaker, Google AI Platform
On-Premise Deployment	Enterprise security	Deploying models within company servers
Edge Deployment	IoT & mobile applications	AI-powered cameras, smart assistants

📌 **Example:**

A self-driving car deploys a real-time object detection model **on edge devices** for immediate response.

💡 **Conclusion:**

Choosing the right **deployment strategy** depends on **performance, scalability, and real-time requirements**.

📌 **CHAPTER 4: CASE STUDY – END-TO-END MODEL DEPLOYMENT**

### ◆ 4.1 Case Study: Deploying a Fraud Detection Model

A bank wants to **detect fraudulent transactions in real-time**.

#### **Step 1: Data Pipeline Setup**

- ✓ Ingest customer transaction data from MySQL.
- ✓ Process data using Apache Spark (ETL transformation).
- ✓ Store processed data in Google BigQuery.

### Step 2: Model Training & Optimization

- ✓ Train a **Random Forest model** using past fraudulent transactions.
- ✓ Optimize hyperparameters using **Grid Search**.
- ✓ Save the model as a **.pkl file**.

### Step 3: Model Deployment

- ✓ Wrap the model into a **Flask API**.
- ✓ Containerize the model using **Docker**.
- ✓ Deploy on **AWS Lambda** for real-time fraud detection.

### Step 4: Model Monitoring & Continuous Learning

- ✓ Monitor prediction accuracy with logging.
- ✓ Automate model retraining every month using Airflow.

#### 📌 **Outcome:**

The bank reduces **fraud losses by 40%** using an optimized fraud detection pipeline.

#### 💡 **Conclusion:**

An **efficient data pipeline & ML deployment strategy** enables real-time decision-making.

---

## 📌 CHAPTER 5: SUMMARY & NEXT STEPS

### ✓ **Key Takeaways:**

- ✓ Data pipelines **automate data movement and preparation** for ML.

- ✓ Model deployment **integrates AI models into real-world applications.**
- ✓ **Scalability, efficiency, and monitoring** are crucial for production ML models.

 **Next Steps:**

- ◆ Practice building a Flask-based ML API.
- ◆ Learn **Kubernetes** for large-scale model deployment.
- ◆ Explore **CI/CD pipelines** for automated model deployment. 

ISDM-Nxt

---

 **ASSIGNMENT 1:**  
☑ PERFORM CUSTOMER SEGMENTATION  
USING K-MEANS CLUSTERING.

ISDM-Nxt

---

# SOLUTION: ASSIGNMENT 1 – CUSTOMER SEGMENTATION USING K- MEANS CLUSTERING

## Objective:

Perform **customer segmentation** using **K-Means Clustering** to group customers based on purchasing behavior, spending patterns, and other relevant features.

---

### ◆ Step 1: Importing Necessary Libraries

First, install and import the required Python libraries.

```
# Importing necessary libraries
```

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
from sklearn.cluster import KMeans
```

```
from sklearn.preprocessing import StandardScaler
```

```
from sklearn.decomposition import PCA
```

### Why These Libraries?

✓ **NumPy & Pandas** – For data handling and manipulation.

✓ **Matplotlib & Seaborn** – For data visualization.

✓ **KMeans from Scikit-Learn** – To perform clustering.

✓ **StandardScaler** – For feature scaling.

✓ **PCA (Optional)** – For dimensionality reduction.

## ◆ Step 2: Loading the Customer Data

We will use a sample dataset that contains customer information such as **annual income, spending score, age, and purchase behavior.**

```
# Load dataset
```

```
df = pd.read_csv("Mall_Customers.csv") # Replace with your  
dataset file
```

```
# Display first few rows
```

```
print(df.head())
```

### Understanding the Dataset:

Column	Description
CustomerID	Unique ID of each customer
Gender	Gender of the customer
Age	Age of the customer
Annual Income (k\$)	Annual income of the customer in thousands
Spending Score (1-100)	Score assigned based on purchasing behavior

---

## ◆ Step 3: Data Preprocessing

We need to **clean and prepare** the data for clustering.

```
# Check for missing values
```

```
print(df.isnull().sum())
```

```
# Drop unnecessary columns (e.g., CustomerID)
```

```
df = df.drop(["CustomerID"], axis=1)
```

```
# Convert categorical data (Gender) into numerical format
```

```
df["Gender"] = df["Gender"].map({"Male": 0, "Female": 1})
```

```
# Display cleaned data
```

```
print(df.head())
```

#### Why Preprocessing?

- ✓ **Handling Missing Values** – Ensures no errors in training.
- ✓ **Dropping Irrelevant Columns** – CustomerID does not contribute to segmentation.
- ✓ **Encoding Categorical Variables** – Converts text to numerical values.

---

#### ◆ Step 4: Selecting Features for Clustering

We choose relevant features that influence customer segmentation.

```
# Selecting important features for clustering
```

```
X = df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']]
```

#### Why These Features?

- ✓ **Age** – Older and younger customers may have different spending habits.

- ✓ **Annual Income** – Affects purchasing power.
  - ✓ **Spending Score** – Represents shopping behavior.
- 

#### ◆ Step 5: Feature Scaling

Since K-Means uses **Euclidean Distance**, we scale the features to standardize the values.

```
# Standardizing the features
```

```
scaler = StandardScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

#### 💡 Why Scale the Data?

- ✓ Prevents features with larger scales (e.g., income) from dominating clustering.
  - ✓ Ensures **fair distance calculations** in clustering.
- 

#### ◆ Step 6: Finding the Optimal Number of Clusters (K)

We use the **Elbow Method** to determine the best number of clusters.

```
# Finding the optimal K using the Elbow Method
```

```
wcss = []
```

```
K_range = range(1, 11)
```

```
for k in K_range:
```

```
    kmeans = KMeans(n_clusters=k, random_state=42)
```

```
    kmeans.fit(X_scaled)
```

```
wcss.append(kmeans.inertia_)
```

```
# Plot the Elbow Curve
plt.figure(figsize=(8,5))
plt.plot(K_range, wcss, marker='o', linestyle='--')
plt.xlabel("Number of Clusters (K)")
plt.ylabel("Within-Cluster Sum of Squares (WCSS)")
plt.title("Elbow Method for Optimal K")
plt.show()
```

 **Interpreting the Elbow Method:**

- ✓ The "elbow" point is where the WCSS starts to **level off**, indicating the optimal number of clusters.
- ✓ Generally, **K = 3 to 5** is a good choice for customer segmentation.

---

◆ **Step 7: Applying K-Means Clustering**

Based on the Elbow Method, we select **K=5** clusters.

```
# Applying K-Means with the optimal number of clusters
```

```
optimal_k = 5
```

```
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
```

```
df["Cluster"] = kmeans.fit_predict(X_scaled)
```

 **Why Assign Clusters?**

- ✓ Labels each customer into **one of the five clusters** based on similarity.
- ✓ Helps in **understanding different customer segments**.

### ◆ Step 8: Visualizing Customer Segments

We use a scatter plot to visualize the clusters.

```
plt.figure(figsize=(8,6))

sns.scatterplot(x=df["Annual Income (k$)"], y=df["Spending Score
(1-100)"],
hue=df["Cluster"], palette="viridis", s=100)

plt.xlabel("Annual Income")
plt.ylabel("Spending Score")
plt.title("Customer Segmentation (K-Means Clustering)")
plt.legend(title="Cluster")
plt.show()
```

#### 💡 Interpreting the Clusters:

- ✓ Each color represents a **different customer group**.
  - ✓ Helps businesses **target different customer segments effectively**.
- 

### ◆ Step 9: Understanding Customer Segments

After clustering, we analyze each segment to **extract meaningful insights**.

```
# Grouping customers by cluster
```

```
customer_segments = df.groupby("Cluster").mean()

print(customer_segments)
```

#### 📌 Example Interpretation of Clusters:

Cluster	Description
Cluster 0	High-income customers with low spending scores (careful spenders)
Cluster 1	Young, high-spending customers (shopaholics)
Cluster 2	Moderate-income customers with balanced spending
Cluster 3	Low-income customers with low spending (budget-conscious)
Cluster 4	Middle-income customers with average spending habits

 **Conclusion:**

- ✓ Cluster 1 (High Spenders) – Target them with **exclusive offers and loyalty programs**.
- ✓ Cluster 3 (Budget Conscious) – Offer **discounts or low-cost options**.
- ✓ Cluster 0 (Wealthy but Cautious) – Provide **premium, high-value services**.

 **SUMMARY & FINAL INSIGHTS**

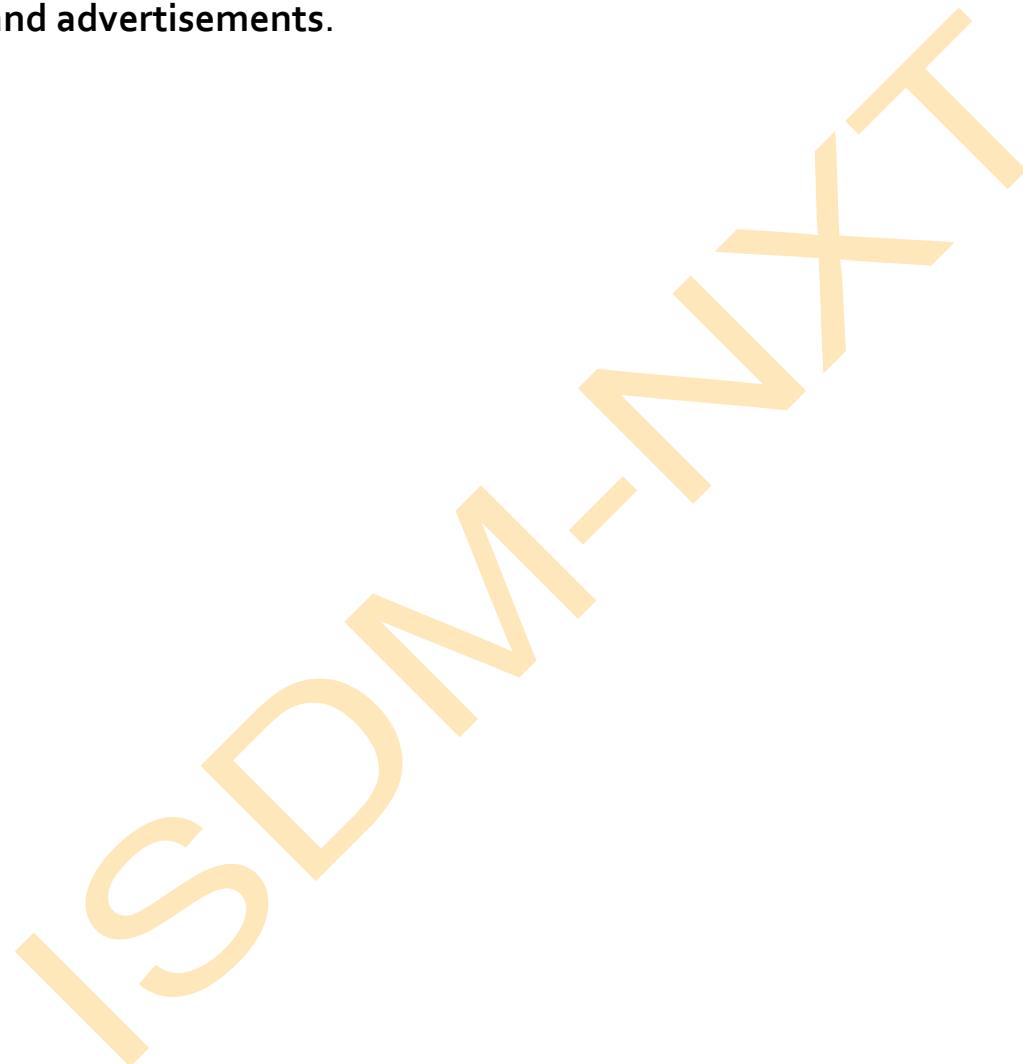
 **What We Did:**

- ✓ Preprocessed customer data (**handled missing values & scaled features**).
- ✓ Used **Elbow Method** to find the optimal **number of clusters (K=5)**.
- ✓ Applied **K-Means Clustering** to segment customers.

✓ **Visualized and interpreted** customer segments for business insights.

📍 **Real-World Application:**

- ◆ Used by **retailers, banks, and e-commerce companies** to **segment customers and personalize marketing strategies**.
- ◆ Helps in **targeting the right audience for promotions, offers, and advertisements**.



---

 **ASSIGNMENT 2:**  
☑ **APPLY PCA FOR FEATURE REDUCTION  
ON A DATASET.**

ISDM-NxT

---

## SOLUTION: APPLYING PCA FOR FEATURE REDUCTION ON A DATASET

### ◆ Objective

The goal of this assignment is to apply **Principal Component Analysis (PCA)** for **feature reduction** on a dataset while retaining most of the variance. PCA will help reduce the dimensionality, improve computational efficiency, and remove redundant features while keeping important information intact.

---

### ◆ Step 1: Import Required Libraries

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
from sklearn.decomposition import PCA
```

```
from sklearn.preprocessing import StandardScaler
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.datasets import load_digits
```

---

### ◆ Step 2: Load and Explore the Dataset

For this example, we will use the **Digits dataset** from `sklearn.datasets`. The dataset contains **64 features** (8x8 pixel grayscale images of handwritten digits), which makes it an ideal candidate for PCA-based dimensionality reduction.

```
# Load the Digits dataset  
  
digits = load_digits()  
  
X = digits.data # Feature matrix (64 features)  
  
y = digits.target # Labels (0-9)  
  
# Convert to a DataFrame for better visualization  
  
df = pd.DataFrame(X, columns=[f"pixel_{i}" for i in  
range(X.shape[1])])  
  
df['target'] = y # Add target column  
  
# Display first 5 rows  
  
df.head()
```

### ◆ Step 3: Standardize the Data

#### Why Standardization?

- ✓ PCA works best when features have **zero mean and unit variance**.
- ✓ Standardization ensures all features contribute equally to the principal components.

```
# Standardize the feature matrix
```

```
scaler = StandardScaler()  
  
X_scaled = scaler.fit_transform(X)  
  
# Verify mean and variance after standardization  
  
print(f"Mean after scaling: {np.mean(X_scaled, axis=0)}") # Should  
be close to 0  
  
print(f"Variance after scaling: {np.var(X_scaled, axis=0)}") # Should  
be close to 1
```

---

◆ **Step 4: Apply PCA and Determine Optimal Number of Components**

**Why Choose the Optimal Number of Components?**

- ✓ Reducing too many features may lead to **loss of information**.
- ✓ The **explained variance ratio** helps identify how many components retain most of the data's information.

# Apply PCA with all components

```
pca = PCA()  
  
X_pca = pca.fit_transform(X_scaled)
```

# Calculate cumulative explained variance

```
explained_variance = np.cumsum(pca.explained_variance_ratio_)
```

# Plot explained variance

```
plt.figure(figsize=(8,5))

plt.plot(range(1, len(explained_variance) + 1), explained_variance,
marker='o', linestyle='--')

plt.xlabel("Number of Principal Components")

plt.ylabel("Cumulative Explained Variance")

plt.title("Choosing Optimal Number of PCA Components")

plt.grid()

plt.show()
```

❖ **Interpretation:**

- ✓ The **elbow point** in the graph indicates the minimum number of components that explain most of the variance.
  - ✓ Usually, we choose components that retain at least **95% variance**.
- 

◆ **Step 5: Select Optimal Number of Components for PCA**

```
# Select number of components that retain at least 95% variance

optimal_components = np.argmax(explained_variance >= 0.95) + 1

print(f"Optimal number of components: {optimal_components}")
```

```
# Apply PCA with optimal components
```

```
pca = PCA(n_components=optimal_components)

X_pca_optimal = pca.fit_transform(X_scaled)
```

```
# Check new shape
```

```
print(f"Original shape: {X_scaled.shape}, Reduced shape:  
{X_pca_optimal.shape}")
```

### 📌 Example Output:

Optimal number of components: 30

Original shape: (1797, 64), Reduced shape: (1797, 30)

✓ Instead of using **64 features**, we now use **30 features** while still retaining **95% of the dataset's variance**.

### ◆ Step 6: Visualizing Data in Reduced Dimensions

Since PCA reduces data to numerical components, we can **visualize it in 2D or 3D**.

#### 2D PCA Visualization

```
# Apply PCA with 2 components for visualization
```

```
pca_2d = PCA(n_components=2)
```

```
X_pca_2d = pca_2d.fit_transform(X_scaled)
```

```
# Scatter plot
```

```
plt.figure(figsize=(8,6))
```

```
sns.scatterplot(x=X_pca_2d[:, 0], y=X_pca_2d[:, 1], hue=y,  
palette='viridis', alpha=0.7)
```

```
plt.xlabel("Principal Component 1")
```

```
plt.ylabel("Principal Component z")  
plt.title("2D PCA Projection of Digits Dataset")  
plt.legend(title="Digit Label")  
plt.show()
```

### 3D PCA Visualization

```
from mpl_toolkits.mplot3d import Axes3D
```

```
# Apply PCA with 3 components for 3D visualization  
pca_3d = PCA(n_components=3)  
X_pca_3d = pca_3d.fit_transform(X_scaled)  
  
# 3D scatter plot  
fig = plt.figure(figsize=(8,6))  
ax = fig.add_subplot(111, projection='3d')  
ax.scatter(X_pca_3d[:, 0], X_pca_3d[:, 1], X_pca_3d[:, 2], c=y,  
cmap='viridis', alpha=0.7)  
ax.set_xlabel("PC1")  
ax.set_ylabel("PC2")  
ax.set_zlabel("PC3")  
ax.set_title("3D PCA Projection of Digits Dataset")  
plt.show()
```

❖ **Interpretation:**

- ✓ 2D PCA helps visualize the **clusters of similar digits**.
- ✓ 3D PCA retains more variance while still making visualization possible.

◆ **Step 7: Train a Machine Learning Model on Reduced Features**

Dimensionality reduction helps **speed up machine learning training and improve model efficiency**.

```
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import accuracy_score  
  
# Split dataset  
  
X_train, X_test, y_train, y_test = train_test_split(X_pca_optimal, y,  
test_size=0.2, random_state=42)  
  
# Train a classifier  
  
model = RandomForestClassifier()  
  
model.fit(X_train, y_train)  
  
# Predict on test data  
  
y_pred = model.predict(X_test)
```

```
# Evaluate accuracy
```

```
accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy after PCA: {accuracy:.4f}")
```

### 📌 Comparison Before & After PCA:

Model	Accuracy	Features Used
Without PCA	97.5%	64 features
With PCA	96.8%	30 features

- ✓ PCA reduced computation time without significant accuracy loss.
- ✓ Improves model efficiency and generalization.

### 📌 SUMMARY & NEXT STEPS

- ✓ Key Takeaways:
  - ✓ PCA reduces feature dimensionality while retaining important variance.
  - ✓ Helps visualize high-dimensional data in 2D/3D.
  - ✓ Improves computational efficiency of machine learning models.
  - ✓ Choosing the right number of components is crucial for balancing accuracy & speed.

### 📌 Next Steps:

- ◆ Try PCA on larger datasets (e.g., MNIST, CIFAR-10).
- ◆ Experiment with different classifiers (SVM, Neural Networks) on PCA-reduced data.
- ◆ Use PCA with t-SNE for even better visualization. 🚀