**UNIT 4:**

**Distance-Based Models:**
1. **Neighbors and Examples**
2. **Nearest Neighbor Classification (KNN)**
   - **Introduction to KNN**:
     - KNN is used for both classification and regression.
     - For classification, it predicts the class based on the majority vote of the nearest neighbors.
     - For regression, it predicts the value based on the mean value of the nearest neighbors.
     - The principle is that "birds of a feather flock together."
   - **KNN Classification Process**:
     - Determine the value of $k$ (number of neighbors).
     - Calculate the distance to find the nearest neighbors.
     - Classify based on the majority class among the nearest neighbors.
   - **Finding the Optimal $k$ Value**:
     - Too small $k$ can lead to sensitivity to outliers.
     - Too large $k$ makes the model dominated by the majority class.
     - Use the square root of the number of data points or error plots to find an optimal $k$.
     - Always use an odd $k$ to avoid ties.
   - **Distance Measures**:
     - Euclidean Distance: Derived from Pythagoras theorem.
     - Manhattan Distance, Minkowski Distance.
   - **Example Dataset**:
     - Demonstrates KNN classification using height and weight to predict BMI categories.

**Kernel-Based Models:**
1. **Support Vector Machines (SVM)**
   - **Linear SVM**:
     - Classifies data by finding the best hyperplane that separates the classes.
   - **RBF SVM**:
     - Uses a radial basis function to handle non-linear classification.
   - **Sigmoid SVM**:
     - Uses a sigmoid kernel for the SVM algorithm.
   - **Polynomial SVM**:
     - Uses a polynomial kernel to handle complex relationships in the data.

**Probability-Based Models:**
1. **Conditional Probability**
2. **Bayes Theorem**
3. **Naive Bayes Classification**
   - Assumes independence between predictors and calculates the probability of each class.
4. **Bayesian Regression**
   - Combines Bayes' theorem with regression models to incorporate prior distributions.

**Case Studies:**
   - **Classification Algorithm for Student Learning Capacity**
     - Applying classification algorithms to predict student performance.

**Clustering Techniques (Related Topics):**
1. **K-Means Clustering**:

- **Elbow Method**: Determines the optimal number of clusters by plotting within-cluster sum of squares against the number of clusters.
- **Python Implementation**: Demonstrates how to preprocess data, apply K-Means, and analyze clusters.

2. **Hierarchical Risk Parity (HRP)**:
   - An innovative portfolio construction method that uses hierarchical clustering for asset allocation.
   - **Steps for Implementation**:
     - Calculate correlation and distances.
     - Perform hierarchical clustering.
     - Quasi-diagonalize the covariance matrix.
     - Recursive bisection for weight allocation.
   - **Python Implementation**: Uses numpy, pandas, scipy, and matplotlib.

**Pairs Trading Using K-Means Clustering:**

1. **Data Collection**:
   - Collect historical price data, volumes, and volatility measures.
2. **Feature Selection**:
   - Choose features like returns, volatility, trading volume, and other technical indicators.
3. **Preprocessing**:
   - Normalize data using z-score or min-max scaling.
   - Calculate distances for clustering.
4. **Apply K-Means Clustering**:
   - Choose the number of clusters (using elbow method).
   - Perform clustering and analyze clusters.
5. **Identify Potential Pairs**:
   - Select pairs with similar financial characteristics and test for correlation and cointegration.
6. **Trading Strategy Development**:
   - Define entry and exit points.
   - Backtest the strategy with historical data.
7. **Implementation and Monitoring**:
   - Deploy capital and monitor performance, adjusting as necessary.

**Tools and Technologies:**

- **Python Libraries**:
  - Scikit-learn for K-Means, Pandas for data manipulation, NumPy for numerical calculations, Matplotlib for plotting data.

This document provides a comprehensive overview of supervised learning techniques, focusing on KNN, SVM, and probability-based models, and extends into practical applications such as clustering for pairs trading and risk parity in portfolio management.

UNIT 5:

**Key Concepts:**
**1. Supervised vs. Unsupervised Learning**
- **Supervised Learning**:
  - The model learns using labeled data.
  - Example: Training a model to recognize cats and dogs using labeled images where each image is tagged as 'cat' or 'dog'.
- **Unsupervised Learning**:
  - The model learns using unlabeled data.
  - Example: Feeding a model with images of cats and dogs without labels. The model tries to find patterns and group similar images together.

**Unsupervised Learning Process:**
- The model analyzes the raw data to find hidden patterns.
- Suitable algorithms, such as K-Means clustering, are applied to group data into clusters based on similarities.

**K-Means Clustering:**
- **Objective**: To partition data points into $k$ clusters.
- **Example Data Points**:
  - A1(2,10), A2(2,5), A3(8,4)
  - B1(5,8), B2(7,5), B3(6,4)
  - C1(1,2), C2(4,9)
- **Distance Function**: Euclidean distance is used to measure the similarity between data points.

**Steps in K-Means Clustering:**
1. **Initial Cluster Centers**:
   - Start by assigning initial centers (centroids) for each cluster, e.g., A1, B1, and C1.
   - Calculate distances from each data point to the initial centroids.
   - Assign each data point to the nearest cluster based on the distance.
2. **Assigning Data Points to Clusters**:
   - For each data point, compute the distance to all centroids and assign it to the closest one.
   - Example:
     - Point A1 (2,10) is closest to centroid A1.
     - Point A2 (2,5) is closest to centroid C1.
     - And so on.
3. **Recomputing Centroids**:
   - After assigning all points to clusters, recalculate the centroids by averaging the points in each cluster.
   - Example:
     - If A1, B1, and C2 belong to Cluster 1, compute the new centroid by averaging their coordinates.
4. **Iterate**:
   - Repeat the assignment and centroid recalculation steps until the centroids no longer change significantly, indicating that the clusters are stable.

**Example Calculation:**
- **Initial Assignment**:
  - Calculate distances of each point from initial centroids A1, B1, and C1.
  - Assign points to the nearest centroids and form initial clusters.
- **Recompute Centroids**:
  - Example: For Cluster 1 with points (2,10), (5,8), (4,9):
    - New centroid = Average of coordinates = ((2+5+4)/3, (10+8+9)/3) = (3.67, 9)

- **Reassign Points**:
    - Calculate distances from each point to the new centroids.
    - Reassign points based on the new distances.
    - Continue until centroids stabilize.

**Final Clusters:**
- After several iterations, the data points stabilize into distinct clusters with minimal changes in centroids.

**Conclusion:**
- K-Means clustering is an iterative process that partitions data into clusters based on similarity.
- It involves selecting initial centroids, assigning points to clusters, recomputing centroids, and iterating until the clusters are stable.

This document provides a detailed explanation of unsupervised learning and a step-by-step guide to applying K-Means clustering to group data points into meaningful clusters.

UNIT 6:

## 1. Definition of NLP
- **Natural Language Processing (NLP)**: A field of computer science, artificial intelligence, and linguistics focused on the interaction between computers and human languages. It enables machines to understand, interpret, and respond to human language.

## 2. Applications of NLP
1. **Question Answering**: Systems that automatically answer questions posed by humans.
2. **Spam Detection**: Identifying and filtering unwanted emails.
3. **Sentiment Analysis**: Analyzing the attitude, emotions, and opinions expressed in text.
4. **Machine Translation**: Translating text or speech from one language to another (e.g., Google Translator).
5. **Spelling Correction**: Correcting spelling errors in text (e.g., Microsoft Word).
6. **Speech Recognition**: Converting spoken words into text, used in applications like virtual assistants and voice-activated systems.
7. **Chatbots**: Automated systems that engage in conversation with users, commonly used for customer service.
8. **Information Extraction**: Extracting structured information from unstructured or semi-structured text.
9. **Natural Language Understanding (NLU)**: Converting large text sets into structured, formal representations for easier manipulation by computers.

## 3. Challenges in NLP
1. **Contextual Understanding**: Words and phrases can have different meanings based on context.
2. **Synonyms**: Multiple words can express the same idea.
3. **Irony and Sarcasm**: Difficult for models to detect due to literal vs. intended meanings.
4. **Ambiguity**: Words or sentences can have multiple interpretations.
5. **Errors in Text and Speech**: Misspellings and misused words can cause issues.
6. **Colloquialisms and Slang**: Informal language varies by region and culture.
7. **Domain-Specific Language**: Different fields use specialized terminology.

## 4. NLP Pipeline
- **Steps to build an NLP Pipeline**:
    1. **Sentence Segmentation**: Breaking down text into sentences.
    2. **Word Tokenization**: Splitting sentences into words or tokens.
    3. **Stemming**: Reducing words to their root form.
    4. **Lemmatization**: Grouping different forms of a word to its base form.
    5. **Identifying Stop Words**: Removing common words like "is", "and", "the".
    6. **Dependency Parsing**: Analyzing the grammatical structure of a sentence.
    7. **POS Tagging**: Identifying parts of speech in text.
    8. **Named Entity Recognition (NER)**: Detecting and classifying named entities (e.g., people, organizations).
    9. **Chunking**: Grouping tokens into meaningful phrases.

## 5. Phases of NLP
1. **Lexical Analysis**: Breaking down text into words and sentences.
2. **Syntactic Analysis (Parsing)**: Checking grammar and structure.
3. **Semantic Analysis**: Understanding the meaning of words and sentences.
4. **Discourse Integration**: Understanding the context of sentences within a larger text.
5. **Pragmatic Analysis**: Interpreting the intended meaning based on context.

## 6. Text Summarization
- **Need for Text Summarization**: Managing the vast amount of data generated daily by creating concise summaries that convey the main points.

- **Approaches**:
  1. **Extraction-based Summarization**: Selecting key phrases from the text.
  2. **Abstraction-based Summarization**: Generating new sentences that capture the essence of the text.

**7. Advantages and Disadvantages of NLP**
- **Advantages**:
  - Efficient and direct responses to queries.
  - Improved communication between computers and humans.
  - Enhanced documentation processes and information retrieval.
- **Disadvantages**:
  - Contextual limitations.
  - Unpredictability and potential for requiring more input.
  - Inability to adapt to new domains without specific training.

**8. Components of NLP**
1. **Natural Language Understanding (NLU)**: Extracting meaningful information from text.
2. **Natural Language Generation (NLG)**: Generating natural language text from data.

The document provides a comprehensive overview of NLP, from its fundamental concepts and applications to the challenges and methodologies involved in processing natural language.