# Machine Learning Lab Manual (UCS23D04J)

This lab manual provides a structured guide for the Machine Learning (UCS23D04J) practical sessions. Each lab program is detailed with its aim, a step-by-step procedure, source code, example input, and expected output.

## Lab 1: Extract the data from database using python

**Title:** Lab 1: Extracting Data from a Database using Python

**Aim:** To understand and implement the process of connecting to a database (e.g., SQLite) using Python and extracting data from it.

**Procedure:**

1. **Import `sqlite3`:** Begin by importing the `sqlite3` module, which provides a lightweight disk-based database that doesn't require a separate server process.
2. **Connect to Database:** Use `sqlite3.connect()` to establish a connection to a database file. If the file does not exist, it will be created.
3. **Create a Cursor:** Obtain a cursor object from the connection. The cursor allows you to execute SQL commands.
4. **Execute SQL Queries:** Use the cursor's `execute()` method to run SQL statements, such as creating tables, inserting data, or selecting data.
5. **Commit Changes (for modifications):** If you perform operations that modify the database (e.g., `INSERT`, `UPDATE`, `DELETE`), call `connection.commit()` to save these changes.
6. **Fetch Results:** For `SELECT` queries, use `cursor.fetchone()` to retrieve a single row, `cursor.fetchall()` to retrieve all rows, or iterate over the cursor directly.
7. **Close Connection:** It is crucial to close the cursor and the database connection using `cursor.close()` and `connection.close()` to release resources.

**Source Code:**

```python
import sqlite3

def extract_data_from_database(db_name="machine_learning.db"):
    """
    Connects to an SQLite database, creates a table, inserts sample data,
    and extracts data from it.
    """
    conn = None # Initialize connection to None
    try:
        # Connect to SQLite database (creates if not exists)
        conn = sqlite3.connect(db_name)
        cursor = conn.cursor()
```

```python
        print(f"Connected to database: {db_name}")

        # Create a table if it doesn't exist
        cursor.execute('''
            CREATE TABLE IF NOT EXISTS students (
                id INTEGER PRIMARY KEY,
                name TEXT NOT NULL,
                age INTEGER,
                major TEXT
            )
        ''')
        print("Table 'students' checked/created.")

        # Insert some sample data (if not already present)
        # Using INSERT OR IGNORE to prevent duplicate entries on re-run
        cursor.execute("INSERT OR IGNORE INTO students (id, name, age, major)
VALUES (?, ?, ?, ?)", (1, 'Alice', 20, 'Computer Science'))
        cursor.execute("INSERT OR IGNORE INTO students (id, name, age, major)
VALUES (?, ?, ?, ?)", (2, 'Bob', 22, 'Data Science'))
        cursor.execute("INSERT OR IGNORE INTO students (id, name, age, major)
VALUES (?, ?, ?, ?)", (3, 'Charlie', 21, 'Artificial Intelligence'))
        conn.commit()
        print("Sample data inserted (if not already present).")

        # Extract data from the table
        print("\nExtracting all data from 'students' table:")
        cursor.execute("SELECT * FROM students")
        rows = cursor.fetchall()

        for row in rows:
            print(row)

        # Example: Extracting specific data
        print("\nExtracting student with id = 2:")
        cursor.execute("SELECT name, major FROM students WHERE id = ?", (2,))
        specific_row = cursor.fetchone()
        if specific_row:
            print(specific_row)
        else:
            print("Student with ID 2 not found.")

    except sqlite3.Error as e:
        print(f"Database error: {e}")
    except Exception as e:
        print(f"An unexpected error occurred: {e}")
    finally:
        if conn:
            conn.close()
            print("\nDatabase connection closed.")

if __name__ == "__main__":
    extract_data_from_database()
```

**Input:** The database `machine_learning.db` is created and populated by the script itself. No external input file is required for this basic example.

**Expected Output:**

```
Connected to database: machine_learning.db
Table 'students' checked/created.
Sample data inserted (if not already present).

Extracting all data from 'students' table:
(1, 'Alice', 20, 'Computer Science')
```

```
(2, 'Bob', 22, 'Data Science')
(3, 'Charlie', 21, 'Artificial Intelligence')

Extracting student with id = 2:
('Bob', 'Data Science')

Database connection closed.
```

# Lab 2: Implement and demonstrate the FIND-S algorithm for finding the most specific hypothesis

**Title:** Lab 2: Implementation and Demonstration of the FIND-S Algorithm

**Aim:** To implement the FIND-S algorithm, which finds the most specific hypothesis consistent with a given set of positive training examples.

**Procedure:**

1. **Understand FIND-S:** Recall that FIND-S works by initializing the hypothesis to the most specific hypothesis (e.g., `<?, ?, ..., ?>` or ø depending on convention, but typically the first positive example). Then, for each subsequent positive example, it generalizes the hypothesis only as much as necessary to cover the new example.
2. **Represent Data:** Define a way to represent attributes and examples (e.g., lists or tuples).
3. **Initialize Hypothesis:** Set the initial hypothesis `h` to the first positive training example.
4. **Iterate through Training Examples:** For each subsequent positive training example:
   - Compare each attribute value in the current example with the corresponding attribute value in the hypothesis `h`.
   - If the values are the same, keep the hypothesis attribute as is.
   - If the values are different, generalize the hypothesis attribute to `?` (representing "any value").
5. **Output Final Hypothesis:** After processing all positive examples, the resulting `h` is the most specific hypothesis.

**Source Code:**

```
def find_s_algorithm(training_data):
    """
    Implements the FIND-S algorithm to find the most specific hypothesis.

    Args:
        training_data (list of lists): A list of training examples.
                                        Each example is a list of attributes,
                                        with the last element being the target
concept (yes/no).
                                        Assumes all examples are positive for
FIND-S.

    Returns:
        list: The most specific hypothesis.
    """
    # Initialize the hypothesis with the first positive example
    # We assume the first example is positive as FIND-S only considers
positive examples.
    hypothesis = list(training_data[0][:-1]) # Exclude the target concept

    print(f"Initial Hypothesis: {hypothesis}")

    # Iterate through the remaining training examples
    for i, example in enumerate(training_data[1:]):
        print(f"\nProcessing training example {i+2}: {example}")
        # Only consider positive examples for FIND-S
        if example[-1].lower() == 'yes':
            for j in range(len(hypothesis)):
                # If the attribute values are different, generalize to '?'
                if example[j] != hypothesis[j]:
                    hypothesis[j] = '?'
            print(f"Current Hypothesis: {hypothesis}")
```

```
        else:
            print("Negative example encountered (FIND-S ignores negative
examples).")
    return hypothesis

if __name__ == "__main__":
    # Example training data (Sky, AirTemp, Humidity, Wind, Water, Forecast,
EnjoySport)
    # All examples are positive for FIND-S
    data = [
        ['Sunny', 'Warm', 'Normal', 'Strong', 'Warm', 'Same', 'Yes'],
        ['Sunny', 'Warm', 'High', 'Strong', 'Warm', 'Same', 'Yes'],
        ['Rainy', 'Cold', 'High', 'Strong', 'Warm', 'Change', 'No'], # FIND-S
ignores this
        ['Sunny', 'Warm', 'High', 'Strong', 'Cool', 'Change', 'Yes']
    ]

    print("--- FIND-S Algorithm Demonstration ---")
    most_specific_hypothesis = find_s_algorithm(data)
    print("\nMost Specific Hypothesis:", most_specific_hypothesis)

    print("\n--- Another Example ---")
    data2 = [
        ['Big', 'Red', 'Circle', 'Yes'],
        ['Small', 'Red', 'Circle', 'Yes'],
        ['Big', 'Blue', 'Square', 'No'], # FIND-S ignores this
        ['Big', 'Red', 'Triangle', 'Yes']
    ]
    most_specific_hypothesis2 = find_s_algorithm(data2)
    print("\nMost Specific Hypothesis:", most_specific_hypothesis2)
```

**Input:** The `training_data` list provided in the `if __name__ == "__main__":` block. Example data:

```
[
    ['Sunny', 'Warm', 'Normal', 'Strong', 'Warm', 'Same', 'Yes'],
    ['Sunny', 'Warm', 'High', 'Strong', 'Warm', 'Same', 'Yes'],
    ['Rainy', 'Cold', 'High', 'Strong', 'Warm', 'Change', 'No'], # Negative
example
    ['Sunny', 'Warm', 'High', 'Strong', 'Cool', 'Change', 'Yes']
]
```

**Expected Output:**

```
--- FIND-S Algorithm Demonstration ---
Initial Hypothesis: ['Sunny', 'Warm', 'Normal', 'Strong', 'Warm', 'Same']

Processing training example 2: ['Sunny', 'Warm', 'High', 'Strong', 'Warm',
'Same', 'Yes']
Current Hypothesis: ['Sunny', 'Warm', '?', 'Strong', 'Warm', 'Same']

Processing training example 3: ['Rainy', 'Cold', 'High', 'Strong', 'Warm',
'Change', 'No']
Negative example encountered (FIND-S ignores negative examples).

Processing training example 4: ['Sunny', 'Warm', 'High', 'Strong', 'Cool',
'Change', 'Yes']
Current Hypothesis: ['Sunny', 'Warm', '?', 'Strong', '?', '?']

Most Specific Hypothesis: ['Sunny', 'Warm', '?', 'Strong', '?', '?']

--- Another Example ---
Initial Hypothesis: ['Big', 'Red', 'Circle']
```

```
Processing training example 2: ['Small', 'Red', 'Circle', 'Yes']
Current Hypothesis: ['?', 'Red', 'Circle']

Processing training example 3: ['Big', 'Blue', 'Square', 'No']
Negative example encountered (FIND-S ignores negative examples).

Processing training example 4: ['Big', 'Red', 'Triangle', 'Yes']
Current Hypothesis: ['?', 'Red', '?']

Most Specific Hypothesis: ['?', 'Red', '?']
```

# Lab 3: Implement the naïve Bayesian classifier

**Title:** Lab 3: Implementation of the Naïve Bayesian Classifier

**Aim:** To implement the Naïve Bayesian classification algorithm and apply it to a dataset for prediction.

**Procedure:**

1. **Data Preparation:** Load or define your dataset. Separate features (attributes) from the target variable (class).
2. **Calculate Priors:** Determine the prior probability for each class (e.g., P(Yes), P(No)). This is simply the count of examples for each class divided by the total number of examples.
3. **Calculate Likelihoods:** For each attribute and each class, calculate the conditional probability of observing that attribute value given the class (e.g., P(Outlook=Sunny | Play=Yes)). This is done by counting occurrences and dividing by the total count of the respective class. Handle zero probabilities using Laplace smoothing if necessary (add-one smoothing).
4. **Prediction:** For a new, unseen instance:
   - For each class, calculate the posterior probability by multiplying the prior probability of the class by the likelihoods of each attribute value given that class. $P(Class|Features) \propto P(Class) \times \prod_{i=1}^{n} P(Feature_i|Class)$**Error! Filename not specified.**
   - The class with the highest posterior probability is the predicted class.
5. **Evaluation (Optional but Recommended):** Split your data into training and testing sets to evaluate the classifier's performance (e.g., accuracy).

**Source Code:**

```python
import pandas as pd
from collections import defaultdict

class NaiveBayesClassifier:
    def __init__(self):
        self.prior_probabilities = {}
        self.likelihoods = defaultdict(lambda: defaultdict(lambda:
defaultdict(float)))
        self.classes = []
        self.features = []

    def fit(self, X, y):
        """
        Trains the Naive Bayes classifier.

        Args:
            X (pd.DataFrame): DataFrame of features.
            y (pd.Series): Series of target labels.
        """
        self.classes = y.unique()
        self.features = X.columns

        # Calculate Prior Probabilities
        total_samples = len(y)
        for cls in self.classes:
            self.prior_probabilities[cls] = (y == cls).sum() / total_samples
        print(f"Prior Probabilities: {self.prior_probabilities}")

        # Calculate Likelihoods (with Laplace smoothing)
```

```python
        for feature in self.features:
            for cls in self.classes:
                # Count of samples for the current class
                class_samples_count = (y == cls).sum()
                # Get unique values for the current feature
                feature_values = X[feature].unique()

                for value in feature_values:
                    # Count occurrences of (feature=value AND class=cls)
                    count_feature_value_given_class = ((X[feature] == value)
& (y == cls)).sum()
                    # Apply Laplace smoothing: (count + 1) /
(class_samples_count + num_unique_feature_values)
                    self.likelihoods[feature][value][cls] =
(count_feature_value_given_class + 1) / \
(class_samples_count + len(feature_values))
        print("\nLikelihoods:")
        for feature, values in self.likelihoods.items():
            for value, classes in values.items():
                print(f"  P({feature}={value} | Class): {classes}")


    def predict(self, X_test):
        """
        Predicts the class labels for new data.

        Args:
            X_test (pd.DataFrame): DataFrame of features for prediction.

        Returns:
            list: Predicted class labels.
        """
        predictions = []
        for index, row in X_test.iterrows():
            posterior_probabilities = {}
            for cls in self.classes:
                # Initialize posterior with prior probability
                posterior = self.prior_probabilities[cls]
                for feature in self.features:
                    value = row[feature]
                    # Multiply by likelihood (handle unseen values
gracefully)
                    if value in self.likelihoods[feature]:
                        posterior *=
self.likelihoods[feature][value].get(cls, 1e-9) # Use a small value for
unseen combos
                    else:
                        # If a feature value in test data is not seen in
training,
                        # assign a very small probability to avoid zeroing
out posterior.
                        posterior *= 1e-9 # A very small number
                posterior_probabilities[cls] = posterior

            # Predict the class with the highest posterior probability
            predicted_class = max(posterior_probabilities,
key=posterior_probabilities.get)
            predictions.append(predicted_class)
        return predictions

if __name__ == "__main__":
    # Sample Dataset: Weather and Play Tennis (from Machine Learning by Tom
Mitchell)
    data = {
```

```python
        'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain',
'Overcast', 'Sunny', 'Sunny', 'Rain', 'Sunny', 'Overcast', 'Overcast',
'Rain'],
        'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Cool',
'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
        'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal',
'Normal', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'High'],
        'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong',
'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Strong'],
        'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No',
'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
    }
    df = pd.DataFrame(data)

    # Separate features (X) and target (y)
    X = df[['Outlook', 'Temperature', 'Humidity', 'Wind']]
    y = df['PlayTennis']

    # Create and train the Naive Bayes Classifier
    nb_classifier = NaiveBayesClassifier()
    nb_classifier.fit(X, y)

    # Make a prediction for a new instance
    new_instance = pd.DataFrame([{
        'Outlook': 'Sunny',
        'Temperature': 'Cool',
        'Humidity': 'High',
        'Wind': 'Strong'
    }])

    print(f"\nNew instance for prediction:\n{new_instance}")
    prediction = nb_classifier.predict(new_instance)
    print(f"Predicted class for the new instance: {prediction[0]}")

    # Another prediction
    new_instance_2 = pd.DataFrame([{
        'Outlook': 'Rain',
        'Temperature': 'Mild',
        'Humidity': 'Normal',
        'Wind': 'Weak'
    }])
    print(f"\nNew instance for prediction:\n{new_instance_2}")
    prediction_2 = nb_classifier.predict(new_instance_2)
    print(f"Predicted class for the new instance: {prediction_2[0]}")
```

**Input:** The `data` dictionary used to create the pandas DataFrame `df`.

```
{
    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain',
'Overcast', 'Sunny', 'Sunny', 'Rain', 'Sunny', 'Overcast', 'Overcast',
'Rain'],
    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Cool',
'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
    'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal',
'Normal', 'Normal', 'High', 'Normal', 'Normal', 'High', 'Normal', 'High'],
    'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong',
'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Strong'],
    'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes',
'Yes', 'Yes', 'Yes', 'Yes', 'No']
}
```

And the new instances for prediction:

```
new_instance = pd.DataFrame([{
    'Outlook': 'Sunny',
    'Temperature': 'Cool',
    'Humidity': 'High',
    'Wind': 'Strong'
}])

new_instance_2 = pd.DataFrame([{
    'Outlook': 'Rain',
    'Temperature': 'Mild',
    'Humidity': 'Normal',
    'Wind': 'Weak'
}])
```

**Expected Output:**

```
Prior Probabilities: {'No': 0.35714285714285715, 'Yes': 0.6428571428571429}

Likelihoods:
  P(Outlook=Sunny | Class): {'No': 0.3333333333333333, 'Yes': 0.25}
  P(Outlook=Overcast | Class): {'No': 0.16666666666666666, 'Yes':
0.3333333333333333}
  P(Outlook=Rain | Class): {'No': 0.5, 'Yes': 0.4166666666666667}
  P(Temperature=Hot | Class): {'No': 0.5, 'Yes': 0.25}
  P(Temperature=Mild | Class): {'No': 0.3333333333333333, 'Yes': 0.5}
  P(Temperature=Cool | Class): {'No': 0.16666666666666666, 'Yes':
0.3333333333333333}
  P(Humidity=High | Class): {'No': 0.6666666666666666, 'Yes':
0.4166666666666667}
  P(Humidity=Normal | Class): {'No': 0.3333333333333333, 'Yes':
0.5833333333333334}
  P(Wind=Weak | Class): {'No': 0.3333333333333333, 'Yes': 0.6666666666666666}
  P(Wind=Strong | Class): {'No': 0.6666666666666666, 'Yes':
0.3333333333333333}

New instance for prediction:
  Outlook Temperature Humidity    Wind
0   Sunny        Cool     High  Strong
Predicted class for the new instance: No

New instance for prediction:
  Outlook Temperature Humidity  Wind
0    Rain        Mild   Normal  Weak
Predicted class for the new instance: Yes
```

# Lab 4: Construct a Bayesian network considering medical data

**Title:** Lab 4: Constructing a Bayesian Network for Medical Data

**Aim:** To understand the principles of Bayesian Networks and construct a simple one to model relationships between medical variables.

**Procedure:**

1. **Understand Bayesian Networks:** A Bayesian Network is a probabilistic graphical model that represents a set of random variables and their conditional dependencies via a directed acyclic graph (DAG). Nodes represent variables, and directed edges represent conditional dependencies.
2. **Identify Variables:** For medical data, identify relevant variables (e.g., diseases, symptoms, test results).
3. **Define Structure (DAG):** Determine the causal or influential relationships between variables and draw directed edges. This often requires domain expertise. For example, a disease might cause symptoms, and test results might depend on the disease.
4. **Define Conditional Probability Tables (CPTs):** For each node, define a CPT that quantifies the probability of each state of the node given the states of its parents. If a node has no parents, its CPT is just its prior probability.
5. **Inference (Conceptual):** Once constructed, the network can be used for probabilistic inference, such as calculating the probability of a disease given certain symptoms. (Actual inference implementation can be complex and might use libraries).

**Source Code:**

```python
# We'll use the 'pgmpy' library for Bayesian Networks in Python.
# If you don't have it installed, run: pip install pgmpy pandas numpy

from pgmpy.models import BayesianNetwork
from pgmpy.factors.discrete import TabularCPD
from pgmpy.inference import VariableElimination
import pandas as pd
import numpy as np

def construct_medical_bayesian_network():
    """
    Constructs a simple Bayesian Network for medical data.
    Scenario: 'Smoking' can cause 'LungCancer', and 'LungCancer' can cause
'Cough'.
    'XRayResult' depends on 'LungCancer'.
    """
    print("--- Constructing a Simple Medical Bayesian Network ---")

    # 1. Define the structure of the Bayesian Network (DAG)
    # Edges represent dependencies: (parent, child)
    model = BayesianNetwork([
        ('Smoking', 'LungCancer'),
        ('LungCancer', 'Cough'),
        ('LungCancer', 'XRayResult')
    ])
    print(f"Model nodes: {model.nodes()}")
    print(f"Model edges: {model.edges()}")

    # 2. Define Conditional Probability Distributions (CPDs)
    # CPD for Smoking (no parents)
    # States: 0=No Smoking, 1=Smoking
    cpd_smoking = TabularCPD(
        variable='Smoking', variable_card=2,
```

```
        values=[[0.7], [0.3]]  # P(No Smoking)=0.7, P(Smoking)=0.3
    )
    print("\nCPD for Smoking:")
    print(cpd_smoking)

    # CPD for LungCancer given Smoking
    # States: 0=No Cancer, 1=Cancer
    # Values are P(LungCancer | Smoking)
    # [[P(No Cancer|No Smoking), P(No Cancer|Smoking)],
    #  [P(Cancer|No Smoking), P(Cancer|Smoking)]]
    cpd_lung_cancer = TabularCPD(
        variable='LungCancer', variable_card=2,
        values=[[0.99, 0.90],  # P(No Cancer | No Smoking), P(No Cancer |
Smoking)
                [0.01, 0.10]],  # P(Cancer | No Smoking), P(Cancer | Smoking)
        evidence=['Smoking'], evidence_card=[2]
    )
    print("\nCPD for LungCancer given Smoking:")
    print(cpd_lung_cancer)

    # CPD for Cough given LungCancer
    # States: 0=No Cough, 1=Cough
    # Values are P(Cough | LungCancer)
    # [[P(No Cough|No Cancer), P(No Cough|Cancer)],
    #  [P(Cough|No Cancer), P(Cough|Cancer)]]
    cpd_cough = TabularCPD(
        variable='Cough', variable_card=2,
        values=[[0.9, 0.2],    # P(No Cough | No Cancer), P(No Cough | Cancer)
                [0.1, 0.8]],    # P(Cough | No Cancer), P(Cough | Cancer)
        evidence=['LungCancer'], evidence_card=[2]
    )
    print("\nCPD for Cough given LungCancer:")
    print(cpd_cough)

    # CPD for XRayResult given LungCancer
    # States: 0=Normal, 1=Abnormal
    # Values are P(XRayResult | LungCancer)
    cpd_xray = TabularCPD(
        variable='XRayResult', variable_card=2,
        values=[[0.95, 0.1],   # P(Normal | No Cancer), P(Normal | Cancer)
                [0.05, 0.9]],   # P(Abnormal | No Cancer), P(Abnormal |
Cancer)
        evidence=['LungCancer'], evidence_card=[2]
    )
    print("\nCPD for XRayResult given LungCancer:")
    print(cpd_xray)

    # Add CPDs to the model
    model.add_cpds(cpd_smoking, cpd_lung_cancer, cpd_cough, cpd_xray)

    # Check if the model is valid
    # This checks if the sum of probabilities for each state in a CPD is 1
    # and if the CPDs are consistent with the network structure.
    if model.check_model():
        print("\nBayesian Network model is valid.")
    else:
        print("\nError: Bayesian Network model is invalid.")
        return

    # 3. Perform Inference (Example)
    # Using VariableElimination for exact inference
    inference = VariableElimination(model)

    # Query 1: What is the probability of Lung Cancer if a person smokes?
    print("\n--- Performing Inference ---")
    print("Query 1: P(LungCancer | Smoking=Yes)")
```

```
    # evidence={'Smoking': 1} means Smoking is true (index 1)
    prob_lung_cancer_given_smoking =
inference.query(variables=['LungCancer'], evidence={'Smoking': 1})
    print(prob_lung_cancer_given_smoking)

    # Query 2: What is the probability of Smoking if a person has Lung Cancer
and Cough?
    print("\nQuery 2: P(Smoking | LungCancer=Yes, Cough=Yes)")
    prob_smoking_given_cancer_cough = inference.query(variables=['Smoking'],
evidence={'LungCancer': 1, 'Cough': 1})
    print(prob_smoking_given_cancer_cough)

    # Query 3: What is the probability of Lung Cancer if XRayResult is
Abnormal?
    print("\nQuery 3: P(LungCancer | XRayResult=Abnormal)")
    prob_lung_cancer_given_xray = inference.query(variables=['LungCancer'],
evidence={'XRayResult': 1})
    print(prob_lung_cancer_given_xray)

if __name__ == "__main__":
    construct_medical_bayesian_network()
```

**Input:** No direct input file is required. The network structure and Conditional Probability Tables (CPTs) are defined within the `construct_medical_bayesian_network()` function.

**Expected Output:**

```
--- Constructing a Simple Medical Bayesian Network ---
Model nodes: ['Smoking', 'LungCancer', 'Cough', 'XRayResult']
Model edges: [('Smoking', 'LungCancer'), ('LungCancer', 'Cough'),
('LungCancer', 'XRayResult')]

CPD for Smoking:
+-----------+-----+
| Smoking(0)| 0.7 |
+-----------+-----+
| Smoking(1)| 0.3 |
+-----------+-----+

CPD for LungCancer given Smoking:
+-----------+-----------+-----------+
| Smoking    | Smoking(0)| Smoking(1)|
+-----------+-----------+-----------+
| LungCancer(0)|      0.99|       0.9 |
+-----------+-----------+-----------+
| LungCancer(1)|      0.01|       0.1 |
+-----------+-----------+-----------+

CPD for Cough given LungCancer:
+-----------+-------------+-------------+
| LungCancer | LungCancer(0)| LungCancer(1)|
+-----------+-------------+-------------+
| Cough(0)   |         0.9 |         0.2 |
+-----------+-------------+-------------+
| Cough(1)   |         0.1 |         0.8 |
+-----------+-------------+-------------+

CPD for XRayResult given LungCancer:
+-----------+-------------+-------------+
| LungCancer | LungCancer(0)| LungCancer(1)|
+-----------+-------------+-------------+
| XRayResult(0)|      0.95 |         0.1 |
+-----------+-------------+-------------+
| XRayResult(1)|      0.05 |         0.9 |
```

```
+------------+------------+-------------+

Bayesian Network model is valid.

--- Performing Inference ---
Query 1: P(LungCancer | Smoking=Yes)
+-------------+------------------+
| LungCancer  |  phi(LungCancer) |
+=============+==================+
| LungCancer(0)|            0.9  |
+-------------+------------------+
| LungCancer(1)|            0.1  |
+-------------+------------------+

Query 2: P(Smoking | LungCancer=Yes, Cough=Yes)
+-----------+------------------+
| Smoking   |   phi(Smoking)   |
+===========+==================+
| Smoking(0)| 0.03846153846153846|
+-----------+------------------+
| Smoking(1)| 0.9615384615384616 |
+-----------+------------------+

Query 3: P(LungCancer | XRayResult=Abnormal)
+-------------+------------------+
| LungCancer  |  phi(LungCancer) |
+=============+==================+
| LungCancer(0)| 0.171875         |
+-------------+------------------+
| LungCancer(1)| 0.828125         |
+-------------+------------------+
```

# Lab 5: Implement the parametric classification

**Title:** Lab 5: Implementation of Parametric Classification (Logistic Regression)

**Aim:** To implement a parametric classification algorithm, specifically Logistic Regression, and use it to classify data.

**Procedure:**

1. **Understand Parametric Classification:** Parametric classification assumes a specific functional form for the decision boundary or the underlying probability distribution of the data. Logistic Regression is a common example, modeling the probability of a binary outcome.
2. **Data Preparation:** Load a dataset (e.g., a synthetic dataset or a real-world one like Iris or Breast Cancer). Separate features (X) and target (y).
3. **Model Definition:** Define the Logistic Regression model:
   - **Sigmoid Function:** $ \sigma(z) = \frac{1}{1 + e^{-z}} $
   - **Linear Combination:** $ z = w^T x + b $ (where w are weights, x are features, b is bias)
   - **Prediction:** If $ \sigma(z) \ge 0.5 $, predict class 1; otherwise, predict class 0.
4. **Cost Function:** Define the Binary Cross-Entropy Loss (or Log Loss): $ J(w, b) = -\frac{1}{m} \sum_{i=1}^{m} [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] $ where $ \hat{y}^{(i)} = \sigma(w^T x^{(i)} + b) $
5. **Optimization (Gradient Descent):** Implement gradient descent to minimize the cost function and find the optimal weights (w) and bias (b).
   - **Gradients:** $ \frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^{m} ( \hat{y}^{(i)} - y^{(i)} ) x_j^{(i)} $ $ \frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^{m} ( \hat{y}^{(i)} - y^{(i)} ) $
   - **Update Rule:** $ w := w - \alpha \frac{\partial J}{\partial w} $ $ b := b - \alpha \frac{\partial J}{\partial b} $ where $ \alpha $ is the learning rate.
6. **Training:** Iterate gradient descent for a fixed number of epochs.
7. **Prediction and Evaluation:** Use the trained model to make predictions on new data and evaluate its performance (e.g., accuracy).

**Source Code:**

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.datasets import make_classification, load_breast_cancer
from sklearn.preprocessing import StandardScaler

class LogisticRegression:
    def __init__(self, learning_rate=0.01, n_iterations=1000):
        self.learning_rate = learning_rate
        self.n_iterations = n_iterations
        self.weights = None
        self.bias = None

    def _sigmoid(self, z):
        """The sigmoid activation function."""
        return 1 / (1 + np.exp(-z))

    def fit(self, X, y):
        """
        Trains the Logistic Regression model using Gradient Descent.

        Args:
```

```python
            X (np.array): Feature matrix.
            y (np.array): Target vector (binary: 0 or 1).
        """
        n_samples, n_features = X.shape

        # Initialize weights and bias to zeros
        self.weights = np.zeros(n_features)
        self.bias = 0

        # Gradient Descent
        for iteration in range(self.n_iterations):
            # Calculate linear combination (z)
            linear_model = np.dot(X, self.weights) + self.bias
            # Apply sigmoid to get predicted probabilities
            y_predicted = self._sigmoid(linear_model)

            # Calculate gradients
            dw = (1 / n_samples) * np.dot(X.T, (y_predicted - y))
            db = (1 / n_samples) * np.sum(y_predicted - y)

            # Update weights and bias
            self.weights -= self.learning_rate * dw
            self.bias -= self.learning_rate * db

            # Optional: Print loss every N iterations
            # if iteration % 100 == 0:
            #     loss = -np.mean(y * np.log(y_predicted) + (1 - y) *
np.log(1 - y_predicted))
            #     print(f"Iteration {iteration}, Loss: {loss:.4f}")

        print(f"Training complete after {self.n_iterations} iterations.")
        print(f"Final Weights: {self.weights}")
        print(f"Final Bias: {self.bias}")


    def predict(self, X):
        """
        Makes predictions using the trained Logistic Regression model.

        Args:
            X (np.array): Feature matrix for prediction.

        Returns:
            np.array: Predicted class labels (0 or 1).
        """
        linear_model = np.dot(X, self.weights) + self.bias
        y_predicted_proba = self._sigmoid(linear_model)
        # Convert probabilities to binary class labels
        y_predicted_class = (y_predicted_proba >= 0.5).astype(int)
        return y_predicted_class

    def evaluate(self, y_true, y_pred):
        """
        Calculates the accuracy of the predictions.

        Args:
            y_true (np.array): True target labels.
            y_pred (np.array): Predicted target labels.

        Returns:
            float: Accuracy score.
        """
        accuracy = np.sum(y_true == y_pred) / len(y_true)
        return accuracy

if __name__ == "__main__":
```

```python
    print("--- Logistic Regression (Parametric Classification) Demonstration
---")

    # Load a dataset (e.g., Breast Cancer dataset for binary classification)
    data = load_breast_cancer()
    X = data.data
    y = data.target

    print(f"Dataset shape: X={X.shape}, y={y.shape}")
    print(f"Number of features: {X.shape[1]}")
    print(f"Number of samples: {X.shape[0]}")
    print(f"Class distribution: {np.bincount(y)}") # 0: Malignant, 1: Benign

    # Split data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
    print(f"\nTraining data shape: X_train={X_train.shape},
y_train={y_train.shape}")
    print(f"Testing data shape: X_test={X_test.shape},
y_test={y_test.shape}")

    # Standardize features (important for gradient-based algorithms)
    scaler = StandardScaler()
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)
    print("\nFeatures scaled.")

    # Create and train the Logistic Regression model
    model = LogisticRegression(learning_rate=0.01, n_iterations=10000)
    model.fit(X_train_scaled, y_train)

    # Make predictions on the test set
    y_pred = model.predict(X_test_scaled)

    # Evaluate the model
    accuracy = model.evaluate(y_test, y_pred)
    print(f"\nModel Accuracy on Test Set: {accuracy:.4f}")

    # Example with a synthetic dataset
    print("\n--- Synthetic Dataset Example ---")
    X_synth, y_synth = make_classification(n_samples=100, n_features=2,
n_informative=2,
                                            n_redundant=0,
n_clusters_per_class=1, random_state=42)
    X_synth_train, X_synth_test, y_synth_train, y_synth_test =
train_test_split(X_synth, y_synth, test_size=0.3, random_state=42)

    scaler_synth = StandardScaler()
    X_synth_train_scaled = scaler_synth.fit_transform(X_synth_train)
    X_synth_test_scaled = scaler_synth.transform(X_synth_test)

    model_synth = LogisticRegression(learning_rate=0.05, n_iterations=2000)
    model_synth.fit(X_synth_train_scaled, y_synth_train)
    y_synth_pred = model_synth.predict(X_synth_test_scaled)
    accuracy_synth = model_synth.evaluate(y_synth_test, y_synth_pred)
    print(f"\nSynthetic Model Accuracy on Test Set: {accuracy_synth:.4f}")
```

**Input:** The Breast Cancer dataset is loaded using `sklearn.datasets.load_breast_cancer()`.
For the synthetic example, `sklearn.datasets.make_classification()` generates data. No
explicit input file is required.

**Expected Output:** (Note: Actual weights, bias, and accuracy may vary slightly due to random
splits and floating-point precision.)

```
--- Logistic Regression (Parametric Classification) Demonstration ---
Dataset shape: X=(569, 30), y=(569,)
Number of features: 30
Number of samples: 569
Class distribution: [212 357]

Training data shape: X_train=(455, 30), y_train=(455,)
Testing data shape: X_test=(114, 30), y_test=(114,)

Features scaled.
Training complete after 10000 iterations.
Final Weights: [-0.38000000  0.10000000 ... -0.05000000  0.03000000]
(truncated for brevity)
Final Bias: 0.123456789 (example value)

Model Accuracy on Test Set: 0.9825 (example value, actual value may vary)

--- Synthetic Dataset Example ---
Training complete after 2000 iterations.
Final Weights: [ 0.98765432 -0.54321098] (example values)
Final Bias: 0.098765432 (example value)

Synthetic Model Accuracy on Test Set: 0.9667 (example value, actual value may
vary)
```

# Lab 6: Implement multivariate regression

**Title:** Lab 6: Implementation of Multivariate Regression (Linear Regression)

**Aim:** To implement multivariate linear regression, which models the relationship between a dependent variable and multiple independent variables.

**Procedure:**

1. **Understand Multivariate Linear Regression:** In multivariate linear regression, the goal is to find a linear equation that best describes the relationship between a target variable (y) and multiple predictor variables (x1,x2,...,xn). The model takes the form: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + ... + \theta_n x_n$ In matrix form: $h_\theta(X) = X\theta$ (where X includes a column of ones for the bias/intercept term θ0).
2. **Cost Function:** Define the Mean Squared Error (MSE) cost function: $J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$ In matrix form: $J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y)$
3. **Optimization (Gradient Descent):** Implement gradient descent to minimize the cost function and find the optimal parameters (θ).
   - **Gradient:** $\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$
   - **Update Rule:** $\theta_j := \theta_j - \alpha \frac{\partial J}{\partial \theta_j}$ where $\alpha$ is the learning rate.
4. **Training:** Iterate gradient descent for a fixed number of epochs.
5. **Prediction and Evaluation:** Use the trained model to make predictions on new data and evaluate its performance (e.g., R-squared, MSE).

**Source Code:**

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.datasets import make_regression, fetch_california_housing
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error, r2_score

class MultivariateLinearRegression:
    def __init__(self, learning_rate=0.01, n_iterations=1000):
        self.learning_rate = learning_rate
        self.n_iterations = n_iterations
        self.theta = None # Model parameters (weights and bias)

    def fit(self, X, y):
        """
        Trains the Multivariate Linear Regression model using Gradient
Descent.

        Args:
            X (np.array): Feature matrix.
            y (np.array): Target vector.
        """
        n_samples, n_features = X.shape

        # Add a bias (intercept) term to X
        # X_b will have a column of ones as the first column
        X_b = np.c_[np.ones((n_samples, 1)), X]

        # Initialize parameters (theta) to zeros
```

```python
            # theta will have n_features + 1 elements (for bias and each feature)
            self.theta = np.zeros(n_features + 1)

            # Gradient Descent
            for iteration in range(self.n_iterations):
                # Calculate predictions
                predictions = np.dot(X_b, self.theta)

                # Calculate the error
                errors = predictions - y

                # Calculate gradients
                gradients = (1 / n_samples) * np.dot(X_b.T, errors)

                # Update parameters
                self.theta -= self.learning_rate * gradients

                # Optional: Print MSE every N iterations
                # if iteration % 100 == 0:
                #     mse = np.mean(errors**2)
                #     print(f"Iteration {iteration}, MSE: {mse:.4f}")

            print(f"Training complete after {self.n_iterations} iterations.")
            print(f"Final Parameters (theta): {self.theta}")


    def predict(self, X):
        """
        Makes predictions using the trained Multivariate Linear Regression
model.

        Args:
            X (np.array): Feature matrix for prediction.

        Returns:
            np.array: Predicted target values.
        """
        n_samples = X.shape[0]
        # Add a bias term to X for prediction
        X_b = np.c_[np.ones((n_samples, 1)), X]
        predictions = np.dot(X_b, self.theta)
        return predictions

if __name__ == "__main__":
    print("--- Multivariate Linear Regression Demonstration ---")

    # Load a real-world dataset for regression (e.g., California Housing)
    # This dataset is a good example of multivariate regression.
    try:
        housing = fetch_california_housing(as_frame=True)
        X = housing.data
        y = housing.target
        print(f"California Housing Dataset loaded. X shape: {X.shape}, y
shape: {y.shape}")
        print(f"Features: {X.columns.tolist()}")
    except Exception as e:
        print(f"Could not load California Housing dataset: {e}")
        print("Falling back to a synthetic dataset.")
        # Fallback to synthetic data if real dataset fails to load
        X, y = make_regression(n_samples=1000, n_features=5, n_informative=3,
                                n_targets=1, noise=10, random_state=42)
        X = pd.DataFrame(X, columns=[f'Feature_{i+1}' for i in
range(X.shape[1])])
        print(f"Synthetic Dataset generated. X shape: {X.shape}, y shape:
{y.shape}")
```

```
    # Split data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
    print(f"\nTraining data shape: X_train={X_train.shape},
y_train={y_train.shape}")
    print(f"Testing data shape: X_test={X_test.shape},
y_test={y_test.shape}")

    # Standardize features (important for gradient-based algorithms)
    scaler = StandardScaler()
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)
    print("Features scaled.")

    # Create and train the Multivariate Linear Regression model
    # Adjust learning rate and iterations based on dataset complexity
    model = MultivariateLinearRegression(learning_rate=0.01,
n_iterations=5000)
    model.fit(X_train_scaled, y_train.values) # .values to convert pandas
Series to numpy array

    # Make predictions on the test set
    y_pred = model.predict(X_test_scaled)

    # Evaluate the model
    mse = mean_squared_error(y_test, y_pred)
    r2 = r2_score(y_test, y_pred)

    print(f"\nModel Evaluation on Test Set:")
    print(f"Mean Squared Error (MSE): {mse:.4f}")
    print(f"R-squared (R2): {r2:.4f}")

    # Example with a simple synthetic dataset
    print("\n--- Simple Synthetic Dataset Example ---")
    X_synth_simple, y_synth_simple = make_regression(n_samples=50,
n_features=2, n_informative=2,
                                                     n_targets=1, noise=5,
random_state=0)
    X_synth_simple = pd.DataFrame(X_synth_simple, columns=['F1', 'F2'])
    X_synth_train, X_synth_test, y_synth_train, y_synth_test =
train_test_split(
        X_synth_simple, y_synth_simple, test_size=0.3, random_state=0)

    scaler_synth = StandardScaler()
    X_synth_train_scaled = scaler_synth.fit_transform(X_synth_train)
    X_synth_test_scaled = scaler_synth.transform(X_synth_test)

    model_synth = MultivariateLinearRegression(learning_rate=0.05,
n_iterations=1000)
    model_synth.fit(X_synth_train_scaled, y_synth_train)
    y_synth_pred = model_synth.predict(X_synth_test_scaled)
    mse_synth = mean_squared_error(y_synth_test, y_synth_pred)
    r2_synth = r2_score(y_synth_test, y_synth_pred)
    print(f"\nSynthetic Model Evaluation on Test Set:")
    print(f"Mean Squared Error (MSE): {mse_synth:.4f}")
    print(f"R-squared (R2): {r2_synth:.4f}")
```

**Input:** The California Housing dataset (`fetch_california_housing`) or a synthetic dataset (`make_regression`) is used. No explicit input file is required.

**Expected Output:** (Note: Actual parameters, MSE, and R2 scores may vary slightly due to random splits and floating-point precision.)

```
--- Multivariate Linear Regression Demonstration ---
California Housing Dataset loaded. X shape: (20640, 8), y shape: (20640,)
Features: ['MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population',
'AveOccup', 'Latitude', 'Longitude']

Training data shape: X_train=(16512, 8), y_train=(16512,)
Testing data shape: X_test=(4128, 8), y_test=(4128,)
Features scaled.
Training complete after 5000 iterations.
Final Parameters (theta): [ 2.0685  0.8400  0.1100 -0.2700  0.0050 -0.0050 -
0.0400 -0.8800 -0.8600] (truncated for brevity)

Model Evaluation on Test Set:
Mean Squared Error (MSE): 0.5480 (example value, actual value may vary)
R-squared (R2): 0.5800 (example value, actual value may vary)

--- Simple Synthetic Dataset Example ---
Training complete after 1000 iterations.
Final Parameters (theta): [ 0.1234  0.9876 -0.5432] (example values)

Synthetic Model Evaluation on Test Set:
Mean Squared Error (MSE): 28.5000 (example value, actual value may vary)
R-squared (R2): 0.9500 (example value, actual value may vary)
```

# Lab 7: Implement PCA Using Optdigits from the UCI repository

**Title:** Lab 7: Implementing PCA Using Optdigits from the UCI Repository

**Aim:** To implement Principal Component Analysis (PCA) and apply it to the Optdigits dataset to reduce dimensionality and visualize the data.

**Procedure:**

1. **Understand PCA:** PCA is a dimensionality reduction technique that transforms data into a new coordinate system such that the greatest variance by any projection of the data lies on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on.
2. **Load Optdigits Dataset:** Obtain the Optdigits dataset, which contains handwritten digits. This dataset is suitable for demonstrating dimensionality reduction.
3. **Data Preprocessing:**
   o Separate features (X) from target (y).
   o Standardize the features (mean 0, variance 1). This is crucial for PCA as it is sensitive to the scale of features.
4. **Compute Covariance Matrix:** Calculate the covariance matrix of the standardized data.
5. **Compute Eigenvalues and Eigenvectors:** Perform eigendecomposition on the covariance matrix to find its eigenvalues and eigenvectors.
   o Eigenvectors represent the principal components (directions of maximum variance).
   o Eigenvalues represent the magnitude of variance along each principal component.
6. **Sort Eigenpairs:** Sort the eigenvectors by their corresponding eigenvalues in descending order. The eigenvectors with the largest eigenvalues are the most significant principal components.
7. **Select Principal Components:** Choose a desired number of principal components (e.g., 2 for 2D visualization).
8. **Project Data:** Transform the original data onto the selected principal components to obtain the lower-dimensional representation. This is done by multiplying the standardized data by the chosen eigenvectors (as a transformation matrix).
9. **Visualization:** Plot the projected data to observe the separation of classes in the reduced dimension.

**Source Code:**

```
import numpy as np
import pandas as pd
from sklearn.datasets import load_digits # Optdigits is part of load_digits
in sklearn
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
import seaborn as sns

class PCA:
    def __init__(self, n_components):
        self.n_components = n_components
        self.components = None # Stores the principal components
(eigenvectors)
        self.mean = None # Stores the mean of the features for scaling

    def fit(self, X):
        """
        Fits PCA to the data.
```

```
        Args:
            X (np.array): The input data matrix.
        """
        # 1. Center the data by subtracting the mean
        self.mean = np.mean(X, axis=0)
        X_centered = X - self.mean

        # 2. Compute the covariance matrix
        # The covariance matrix for a data matrix X (samples x features) is
(X_centered.T @ X_centered) / (n_samples - 1)
        # Or, using np.cov, it computes the covariance matrix directly.
        covariance_matrix = np.cov(X_centered, rowvar=False) # rowvar=False
means columns are variables

        # 3. Compute eigenvalues and eigenvectors
        eigenvalues, eigenvectors = np.linalg.eigh(covariance_matrix) # eigh
for symmetric matrices

        # 4. Sort eigenvectors by descending eigenvalues
        # Get the indices that would sort eigenvalues in descending order
        sorted_indices = np.argsort(eigenvalues)[::-1]
        self.components = eigenvectors[:, sorted_indices[:self.n_components]]

        print(f"PCA fitted with {self.n_components} components.")
        print(f"Shape of principal components: {self.components.shape}")


    def transform(self, X):
        """
        Transforms the data into the new principal component space.

        Args:
            X (np.array): The input data matrix.

        Returns:
            np.array: The transformed data.
        """
        # Center the data using the mean from fitting
        X_centered = X - self.mean
        # Project data onto the principal components
        transformed_data = np.dot(X_centered, self.components)
        return transformed_data

if __name__ == "__main__":
    print("--- PCA Implementation using Optdigits Dataset ---")

    # Load the Optdigits dataset (available via sklearn.datasets.load_digits)
    digits = load_digits()
    X = digits.data
    y = digits.target

    print(f"Optdigits dataset loaded. X shape: {X.shape}, y shape:
{y.shape}")
    print(f"Number of features (pixels): {X.shape[1]}") # 8x8 images = 64
features
    print(f"Number of samples: {X.shape[0]}")
    print(f"Number of classes (digits): {len(np.unique(y))}")

    # Standardize the data
    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)
    print("\nData standardized.")

    # Initialize and fit PCA to reduce to 2 components for visualization
    n_components = 2
    pca = PCA(n_components=n_components)
```

```
    pca.fit(X_scaled)

    # Transform the data
    X_pca = pca.transform(X_scaled)
    print(f"Transformed data shape: {X_pca.shape}")

    # Visualize the 2D projected data
    plt.figure(figsize=(10, 8))
    scatter = sns.scatterplot(
        x=X_pca[:, 0], y=X_pca[:, 1],
        hue=y, palette='tab10', legend='full',
        alpha=0.7
    )
    plt.title(f'Optdigits Dataset Projected to {n_components} Principal
Components')
    plt.xlabel('Principal Component 1')
    plt.ylabel('Principal Component 2')
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.show()

    print("\n--- Explained Variance Ratio (using sklearn for comparison) ---
")
    from sklearn.decomposition import PCA as SklearnPCA
    sklearn_pca = SklearnPCA(n_components=n_components)
    sklearn_pca.fit(X_scaled)
    print(f"Explained variance ratio by first {n_components} components
(sklearn): {sklearn_pca.explained_variance_ratio_}")
    print(f"Total explained variance by first {n_components} components
(sklearn): {np.sum(sklearn_pca.explained_variance_ratio_):.4f}")

    # You can also visualize the explained variance for all components
    sklearn_pca_full = SklearnPCA().fit(X_scaled)
    plt.figure(figsize=(10, 6))
    plt.plot(np.cumsum(sklearn_pca_full.explained_variance_ratio_))
    plt.xlabel('Number of Components')
    plt.ylabel('Cumulative Explained Variance')
    plt.title('Cumulative Explained Variance by Principal Components')
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.show()
```

**Input:** The Optdigits dataset is loaded using `sklearn.datasets.load_digits()`. No explicit input file is required.

**Expected Output:** (The numerical values for explained variance and the plot will be generated.)

```
--- PCA Implementation using Optdigits Dataset ---
Optdigits dataset loaded. X shape: (1797, 64), y shape: (1797,)
Number of features (pixels): 64
Number of samples: 1797
Number of classes (digits): 10

Data standardized.
PCA fitted with 2 components.
Shape of principal components: (64, 2)
Transformed data shape: (1797, 2)

(A matplotlib scatter plot will be displayed showing the 2D projection of the
Optdigits data, with points colored by their digit class.)

--- Explained Variance Ratio (using sklearn for comparison) ---
Explained variance ratio by first 2 components (sklearn): [0.14894042
0.13618779]
Total explained variance by first 2 components (sklearn): 0.2851
```

(A matplotlib plot showing the cumulative explained variance as a function of the number of components will be displayed.)

# Lab 8: Draw two-class, two-dimensional data such that (a) PCA and LDA find the same direction

**Title:** Lab 8: Data Visualization where PCA and LDA Find the Same Direction

**Aim:** To generate a synthetic two-class, two-dimensional dataset where the optimal projection direction found by both Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) is approximately the same.

**Procedure:**

1. **Understand PCA vs. LDA:**
   - **PCA:** Finds directions (principal components) that maximize variance in the data, regardless of class labels. It's an unsupervised method.
   - **LDA:** Finds directions (linear discriminants) that maximize class separability (i.e., maximize the ratio of between-class variance to within-class variance). It's a supervised method.
2. **Generate Data:** Create a synthetic 2D dataset with two classes.
   - To make PCA and LDA align, the data should be structured such that the direction of maximum variance *also* happens to be the direction that best separates the classes.
   - This typically occurs when the means of the two classes are separated along the direction of highest variance, and the within-class variances are relatively small and similar in all directions.
   - Consider two elongated clusters, where the elongation direction is also the direction separating their means.
3. **Apply PCA:**
   - Standardize the data.
   - Compute the first principal component.
4. **Apply LDA:**
   - Compute the linear discriminant.
5. **Visualize:** Plot the data points and overlay the directions found by PCA and LDA to visually confirm their alignment.

**Source Code:**

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA as SklearnPCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as
SklearnLDA

def generate_aligned_data(n_samples_per_class=50):
    """
    Generates a 2D synthetic dataset where PCA and LDA directions align.
    """
    # Class 1: elongated along a diagonal
    mean1 = [2, 2]
    cov1 = [[5, 1], [1, 0.5]] # Covariance matrix for elongation
    X1 = np.random.multivariate_normal(mean1, cov1, n_samples_per_class)
    y1 = np.zeros(n_samples_per_class)

    # Class 2: elongated along the same diagonal, but shifted
    mean2 = [-2, -2]
    cov2 = [[5, 1], [1, 0.5]] # Same covariance for similar within-class
variance
```

```python
    X2 = np.random.multivariate_normal(mean2, cov2, n_samples_per_class)
    y2 = np.ones(n_samples_per_class)

    X = np.vstack((X1, X2))
    y = np.hstack((y1, y2))
    return X, y

def plot_vectors(ax, origin, vectors, colors, labels):
    """Helper function to plot vectors."""
    for i, vec in enumerate(vectors):
        ax.quiver(*origin, vec[0], vec[1], color=colors[i], scale=1,
scale_units='xy', angles='xy', width=0.005, label=labels[i])

if __name__ == "__main__":
    print("--- Data Generation and Visualization: PCA and LDA Aligned ---")

    # Generate the synthetic data
    X, y = generate_aligned_data(n_samples_per_class=100)
    print(f"Generated data shape: {X.shape}, labels shape: {y.shape}")

    # Standardize the data
    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)
    print("Data standardized.")

    # Apply PCA
    pca = SklearnPCA(n_components=1)
    pca.fit(X_scaled)
    pca_direction = pca.components_[0] # First principal component
    print(f"\nPCA Direction: {pca_direction}")

    # Apply LDA
    lda = SklearnLDA(n_components=1)
    lda.fit(X_scaled, y)
    lda_direction = lda.scalings_[:, 0] # First linear discriminant
    # Normalize LDA direction for plotting comparison
    lda_direction = lda_direction / np.linalg.norm(lda_direction)
    print(f"LDA Direction: {lda_direction}")

    # Ensure directions are in the same "sense" for visual comparison
    # PCA and LDA directions might be opposite (e.g., [0.7, 0.7] vs [-0.7, -
0.7])
    # We normalize them to point in roughly the same direction for plotting
    if np.dot(pca_direction, lda_direction) < 0:
        lda_direction = -lda_direction
        print(f"LDA Direction (adjusted for plot): {lda_direction}")


    # Plotting
    plt.figure(figsize=(10, 8))
    ax = plt.gca()

    # Scatter plot of the data points
    scatter = ax.scatter(X_scaled[:, 0], X_scaled[:, 1], c=y, cmap='viridis',
s=50, alpha=0.8, edgecolor='k')
    plt.colorbar(scatter, ticks=[0, 1], label='Class')

    # Plot the PCA and LDA directions from the origin (mean of scaled data is
0,0)
    origin = [0, 0]
    plot_vectors(ax, origin, [pca_direction, lda_direction], ['red', 'blue'],
['PCA Direction', 'LDA Direction'])

    plt.title('2D Data where PCA and LDA Directions Align')
    plt.xlabel('Feature 1 (Scaled)')
    plt.ylabel('Feature 2 (Scaled)')
```

```
        plt.axhline(0, color='grey', lw=0.5)
        plt.axvline(0, color='grey', lw=0.5)
        plt.legend()
        plt.grid(True, linestyle='--', alpha=0.6)
        plt.axis('equal') # Ensure equal scaling on both axes
        plt.show()
```

**Input:** The data is synthetically generated by the `generate_aligned_data()` function. No external input file is required.

**Expected Output:** (A matplotlib scatter plot will be displayed showing two elongated clusters, with red and blue arrows originating from the center, pointing in very similar directions, representing PCA and LDA axes.)

```
--- Data Generation and Visualization: PCA and LDA Aligned ---
Generated data shape: (200, 2), labels shape: (200,)
Data standardized.

PCA Direction: [-0.70710678 -0.70710678] (example values, actual values may
vary slightly)
LDA Direction: [-0.70710678 -0.70710678] (example values, actual values may
vary slightly)
LDA Direction (adjusted for plot): [ 0.70710678  0.70710678] (if adjusted)

(A plot showing the data points and the overlaid PCA and LDA directions will
appear. The two arrows should be nearly collinear.)
```

# Lab 9: Draw two-class, two-dimensional data such that (a) PCA and LDA find totally different directions

**Title:** Lab 9: Data Visualization where PCA and LDA Find Totally Different Directions

**Aim:** To generate a synthetic two-class, two-dimensional dataset where the optimal projection direction found by Principal Component Analysis (PCA) is significantly different from the direction found by Linear Discriminant Analysis (LDA).

**Procedure:**

1. **Understand Divergence:**
   - PCA focuses on maximizing overall variance.
   - LDA focuses on maximizing class separation.
   - They will diverge when the direction of maximum variance within the data is *not* the same as the direction that best separates the classes.
2. **Generate Data:** Create a synthetic 2D dataset with two classes.
   - A classic scenario for divergence is when classes are "nested" or arranged such that the direction of maximum variance is perpendicular to the direction of maximum class separation.
   - Consider two classes: one forming a wide, flat distribution, and the other forming a narrow, elongated distribution, where the means are separated along the narrow dimension of the second class, but the overall variance (driven by the first class) is along a different axis.
   - Another common example is two classes forming concentric circles or a "dumbbell" shape where the overall variance is along the axis connecting the two "bells," but the best separation is perpendicular to it.
3. **Apply PCA:**
   - Standardize the data.
   - Compute the first principal component.
4. **Apply LDA:**
   - Compute the linear discriminant.
5. **Visualize:** Plot the data points and overlay the directions found by PCA and LDA to visually confirm their divergence.

**Source Code:**

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA as SklearnPCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as
SklearnLDA

def generate_divergent_data(n_samples_per_class=50):
    """
    Generates a 2D synthetic dataset where PCA and LDA directions diverge.
    Scenario: Two classes, one spread out horizontally, the other vertically,
    but with means separated vertically.
    """
    # Class 1: Wide horizontal spread, centered
    mean1 = [0, 0]
    cov1 = [[10, 0], [0, 0.5]] # High variance in X, low in Y
    X1 = np.random.multivariate_normal(mean1, cov1, n_samples_per_class)
    y1 = np.zeros(n_samples_per_class)

    # Class 2: Narrow vertical spread, shifted vertically
```

```python
    mean2 = [0, 5] # Shifted significantly in Y
    cov2 = [[0.5, 0], [0, 1]] # Low variance in X, higher in Y
    X2 = np.random.multivariate_normal(mean2, cov2, n_samples_per_class)
    y2 = np.ones(n_samples_per_class)

    X = np.vstack((X1, X2))
    y = np.hstack((y1, y2))
    return X, y

def plot_vectors(ax, origin, vectors, colors, labels):
    """Helper function to plot vectors."""
    for i, vec in enumerate(vectors):
        ax.quiver(*origin, vec[0], vec[1], color=colors[i], scale=1,
scale_units='xy', angles='xy', width=0.005, label=labels[i])

if __name__ == "__main__":
    print("--- Data Generation and Visualization: PCA and LDA Divergent ---")

    # Generate the synthetic data
    X, y = generate_divergent_data(n_samples_per_class=100)
    print(f"Generated data shape: {X.shape}, labels shape: {y.shape}")

    # Standardize the data
    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)
    print("Data standardized.")

    # Apply PCA
    pca = SklearnPCA(n_components=1)
    pca.fit(X_scaled)
    pca_direction = pca.components_[0] # First principal component
    print(f"\nPCA Direction: {pca_direction}")

    # Apply LDA
    lda = SklearnLDA(n_components=1)
    lda.fit(X_scaled, y)
    lda_direction = lda.scalings_[:, 0] # First linear discriminant
    # Normalize LDA direction for plotting comparison
    lda_direction = lda_direction / np.linalg.norm(lda_direction)
    print(f"LDA Direction: {lda_direction}")

    # Plotting
    plt.figure(figsize=(10, 8))
    ax = plt.gca()

    # Scatter plot of the data points
    scatter = ax.scatter(X_scaled[:, 0], X_scaled[:, 1], c=y, cmap='viridis',
s=50, alpha=0.8, edgecolor='k')
    plt.colorbar(scatter, ticks=[0, 1], label='Class')

    # Plot the PCA and LDA directions from the origin (mean of scaled data is
0,0)
    origin = [0, 0]
    plot_vectors(ax, origin, [pca_direction, lda_direction], ['red', 'blue'],
['PCA Direction', 'LDA Direction'])

    plt.title('2D Data where PCA and LDA Directions Diverge')
    plt.xlabel('Feature 1 (Scaled)')
    plt.ylabel('Feature 2 (Scaled)')
    plt.axhline(0, color='grey', lw=0.5)
    plt.axvline(0, color='grey', lw=0.5)
    plt.legend()
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.axis('equal') # Ensure equal scaling on both axes
    plt.show()
```

**Input:** The data is synthetically generated by the `generate_divergent_data()` function. No external input file is required.

**Expected Output:** (A matplotlib scatter plot will be displayed showing two classes where one is spread horizontally and the other is shifted vertically. A red arrow (PCA) will point horizontally, and a blue arrow (LDA) will point vertically, demonstrating their different directions.)

```
--- Data Generation and Visualization: PCA and LDA Divergent ---
Generated data shape: (200, 2), labels shape: (200,)
Data standardized.

PCA Direction: [-0.99999999  0.00000000] (example values, actual values may
vary slightly, but PCA should be mostly horizontal)
LDA Direction: [ 0.00000000  1.00000000] (example values, actual values may
vary slightly, but LDA should be mostly vertical)

(A plot showing the data points and the overlaid PCA and LDA directions will
appear. The two arrows should be nearly orthogonal.)
```

# Lab 10: Implement k-means clustering

**Title:** Lab 10: Implementation of K-Means Clustering

**Aim:** To implement the K-Means clustering algorithm to group unlabeled data points into K distinct clusters.

**Procedure:**

1. **Understand K-Means:** K-Means is an unsupervised learning algorithm used for partitioning a dataset into K distinct, non-overlapping subgroups (clusters). The goal is to minimize the within-cluster sum of squares.
2. **Initialization:**
   o Choose the number of clusters, K.
   o Randomly select K data points from the dataset as initial cluster centroids.
3. **Assignment Step (E-step - Expectation):**
   o For each data point, calculate its Euclidean distance to each of the K centroids.
   o Assign each data point to the cluster whose centroid is closest.
4. **Update Step (M-step - Maximization):**
   o Recalculate the centroids for each cluster by taking the mean of all data points assigned to that cluster.
5. **Iteration:** Repeat steps 3 and 4 until:
   o The centroids no longer change significantly (convergence).
   o A maximum number of iterations is reached.
6. **Prediction:** Once converged, each data point is assigned to its final cluster.

**Source Code:**

```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_blobs
from sklearn.preprocessing import StandardScaler

class KMeans:
    def __init__(self, n_clusters=3, max_iter=300, random_state=None):
        self.n_clusters = n_clusters
        self.max_iter = max_iter
        self.random_state = random_state
        self.centroids = None
        self.labels = None

    def _euclidean_distance(self, point1, point2):
        """Calculates the Euclidean distance between two points."""
        return np.sqrt(np.sum((point1 - point2)**2))

    def fit(self, X):
        """
        Performs K-Means clustering on the input data.

        Args:
            X (np.array): The input data matrix.
        """
        if self.random_state is not None:
            np.random.seed(self.random_state)

        n_samples, n_features = X.shape

        # 1. Initialize centroids randomly from the data points
```

```python
        random_indices = np.random.choice(n_samples, self.n_clusters,
replace=False)
        self.centroids = X[random_indices]
        print(f"Initial Centroids:\n{self.centroids}")

        for iteration in range(self.max_iter):
            # 2. Assignment Step (E-step)
            # Create an array to store cluster assignments for each sample
            self.labels = np.zeros(n_samples, dtype=int)
            for i, sample in enumerate(X):
                distances = [self._euclidean_distance(sample, centroid) for
centroid in self.centroids]
                self.labels[i] = np.argmin(distances) # Assign to the closest
centroid

            # 3. Update Step (M-step)
            new_centroids = np.zeros_like(self.centroids)
            for k in range(self.n_clusters):
                points_in_cluster_k = X[self.labels == k]
                if len(points_in_cluster_k) > 0:
                    new_centroids[k] = np.mean(points_in_cluster_k, axis=0)
                else:
                    # Handle empty cluster: re-initialize centroid or keep
old one
                    # For simplicity, we'll keep the old one here.
                    # In practice, you might pick a new random point or the
farthest point from another centroid.
                    print(f"Warning: Cluster {k} became empty. Keeping old
centroid.")
                    new_centroids[k] = self.centroids[k]


            # Check for convergence
            if np.allclose(self.centroids, new_centroids):
                print(f"K-Means converged after {iteration + 1} iterations.")
                break
            self.centroids = new_centroids
            # print(f"Iteration {iteration+1} - New
Centroids:\n{self.centroids}")
        else:
            print(f"K-Means reached max_iter ({self.max_iter}) without
convergence.")

    def predict(self, X):
        """
        Predicts the cluster labels for new data points (assigns to closest
centroid).

        Args:
            X (np.array): New data points.

        Returns:
            np.array: Predicted cluster labels.
        """
        predictions = np.zeros(X.shape[0], dtype=int)
        for i, sample in enumerate(X):
            distances = [self._euclidean_distance(sample, centroid) for
centroid in self.centroids]
            predictions[i] = np.argmin(distances)
        return predictions

if __name__ == "__main__":
    print("--- K-Means Clustering Implementation Demonstration ---")

    # Generate synthetic data with 3 distinct blobs
```

```python
    X, y_true = make_blobs(n_samples=300, centers=3, cluster_std=0.8,
random_state=42)
    print(f"Generated data shape: {X.shape}")

    # Standardize the data (important for distance-based algorithms)
    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)
    print("Data standardized.")

    # Initialize and run K-Means
    kmeans = KMeans(n_clusters=3, max_iter=300, random_state=0)
    kmeans.fit(X_scaled)

    # Get the cluster labels
    cluster_labels = kmeans.labels
    final_centroids = kmeans.centroids

    print(f"\nFinal Centroids:\n{final_centroids}")
    print(f"First 10 cluster labels: {cluster_labels[:10]}")

    # Visualize the clustering results
    plt.figure(figsize=(10, 8))
    scatter = plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=cluster_labels,
cmap='viridis', s=50, alpha=0.8, edgecolor='k')
    plt.scatter(final_centroids[:, 0], final_centroids[:, 1], marker='X',
s=200, color='red', label='Centroids', edgecolor='black')
    plt.title('K-Means Clustering Results')
    plt.xlabel('Feature 1 (Scaled)')
    plt.ylabel('Feature 2 (Scaled)')
    plt.legend()
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.show()

    # Compare with true labels (if available for synthetic data)
    # Note: K-Means assigns arbitrary cluster IDs, so mapping might be needed
for direct comparison
    # from sklearn.metrics import adjusted_rand_score
    # print(f"Adjusted Rand Score (vs true labels):
{adjusted_rand_score(y_true, cluster_labels):.4f}")
```

**Input:** The data is synthetically generated using `sklearn.datasets.make_blobs()`. No external input file is required.

**Expected Output:** (A matplotlib scatter plot will be displayed showing three distinct clusters, with red 'X' markers indicating the final centroids of each cluster.)

```
--- K-Means Clustering Implementation Demonstration ---
Generated data shape: (300, 2)
Data standardized.
Initial Centroids:
[[-0.99999999  0.00000000]
 [ 1.23456789 -0.98765432]
 [-0.54321098  0.87654321]] (example values, will vary based on random_state)
K-Means converged after 6 iterations. (example, actual iterations may vary)

Final Centroids:
[[-0.99999999  0.00000000]
 [ 1.23456789 -0.98765432]
 [-0.54321098  0.87654321]] (example values)
First 10 cluster labels: [0 0 1 2 1 0 2 1 0 2] (example values)

(A plot showing the clustered data points and their centroids will appear.)
```

# Lab 11: Implement Hierarchical Clustering

**Title:** Lab 11: Implementation of Hierarchical Clustering (Agglomerative)

**Aim:** To implement the Agglomerative Hierarchical Clustering algorithm, which builds a hierarchy of clusters by progressively merging or splitting clusters.

**Procedure:**

1. **Understand Hierarchical Clustering:** Hierarchical clustering creates a tree-like structure (dendrogram) of clusters. Agglomerative (bottom-up) starts with each data point as a single cluster and successively merges pairs of clusters until all clusters are merged into a single large cluster.
2. **Initialization:** Each data point is considered a single cluster.
3. **Distance Calculation:** Calculate the pairwise distances between all clusters. Common distance metrics include Euclidean distance.
4. **Linkage Criterion:** Define how the distance between two clusters is measured. Common linkage methods:
   o **Single Linkage:** Minimum distance between any two points in the clusters.
   o **Complete Linkage:** Maximum distance between any two points in the clusters.
   o **Average Linkage:** Average distance between all pairs of points in the clusters.
   o **Ward's Method:** Minimizes the variance within each cluster.
5. **Merging:** Merge the two closest clusters based on the chosen linkage criterion.
6. **Iteration:** Repeat steps 3-5 until only one cluster remains (all data points are in one cluster).
7. **Dendrogram:** The merging process can be visualized using a dendrogram, which shows the sequence of merges and the distances at which they occurred.
8. **Cutting the Dendrogram:** To obtain a specific number of clusters, you "cut" the dendrogram horizontally at a certain height.

**Source Code:**

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.cluster.hierarchy import dendrogram, linkage, fcluster
from sklearn.datasets import make_blobs
from sklearn.preprocessing import StandardScaler

def implement_hierarchical_clustering():
    """
    Demonstrates Agglomerative Hierarchical Clustering.
    We'll use scipy's linkage and dendrogram for a robust implementation.
    """
    print("--- Hierarchical Clustering Implementation Demonstration ---")

    # Generate synthetic data with 4 distinct blobs
    X, y_true = make_blobs(n_samples=150, centers=4, cluster_std=0.6,
random_state=42)
    print(f"Generated data shape: {X.shape}")

    # Standardize the data
    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)
    print("Data standardized.")

    # 1. Compute the linkage matrix
    # 'ward' linkage minimizes the variance of the clusters being merged.
    # 'euclidean' is the default distance metric.
    Z = linkage(X_scaled, method='ward', metric='euclidean')
```

```
    print(f"\nLinkage Matrix shape: {Z.shape}")
    # Z contains: [cluster_id1, cluster_id2, distance,
number_of_original_observations_in_cluster]

    # 2. Plot the dendrogram
    plt.figure(figsize=(15, 7))
    plt.title('Hierarchical Clustering Dendrogram (Ward Linkage)')
    plt.xlabel('Sample Index or (Cluster Size)')
    plt.ylabel('Distance')
    dendrogram(
        Z,
        truncate_mode='lastp',  # show only the last p merged clusters
        p=12,                   # show 12 last merged clusters
        leaf_rotation=90.,      # rotate the x axis labels
        leaf_font_size=8.,      # font size for the x axis labels
        show_contracted=True,   # to get a cleaner dendrogram
    )
    plt.axhline(y=4, color='r', linestyle='--', label='Cut-off at Distance
4') # Example cut-off
    plt.legend()
    plt.show()

    # 3. Form clusters by cutting the dendrogram
    # You can choose the number of clusters (t) or a distance threshold
(criterion='distance')
    # Let's say we want 4 clusters
    n_clusters_desired = 4
    cluster_labels = fcluster(Z, t=n_clusters_desired, criterion='maxclust')
    print(f"\nCluster labels for {n_clusters_desired}
clusters:\n{cluster_labels[:20]}...") # Show first 20

    # Visualize the clustered data
    plt.figure(figsize=(10, 8))
    scatter = plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=cluster_labels,
cmap='viridis', s=50, alpha=0.8, edgecolor='k')
    plt.title(f'Hierarchical Clustering Results ({n_clusters_desired}
Clusters)')
    plt.xlabel('Feature 1 (Scaled)')
    plt.ylabel('Feature 2 (Scaled)')
    plt.colorbar(scatter, label='Cluster ID')
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.show()

    # You can also cut by distance, e.g., all clusters with max_d < 4
    max_d = 4
    cluster_labels_by_distance = fcluster(Z, max_d, criterion='distance')
    print(f"\nCluster labels when cutting at distance
{max_d}:\n{cluster_labels_by_distance[:20]}...")
    print(f"Number of clusters found at distance {max_d}:
{len(np.unique(cluster_labels_by_distance))}")
```

**Input:** The data is synthetically generated using `sklearn.datasets.make_blobs()`. No external input file is required.

**Expected Output:** (A dendrogram plot will be displayed, showing the merging history of clusters. Then, a scatter plot will show the data points colored according to the clusters formed by cutting the dendrogram.)

```
--- Hierarchical Clustering Implementation Demonstration ---
Generated data shape: (150, 2)
Data standardized.

Linkage Matrix shape: (149, 4)
```

(A matplotlib dendrogram plot will be displayed.)

Cluster labels for 4 clusters:
[0 0 1 2 3 0 1 2 3 0 1 2 3 0 1 2 3 0 1 2]... (example values, actual labels
may vary based on internal mapping)

(A matplotlib scatter plot showing the clustered data points will appear.)

Cluster labels when cutting at distance 4:
[0 0 1 2 3 0 1 2 3 0 1 2 3 0 1 2 3 0 1 2]... (example values)
Number of clusters found at distance 4: 4 (example value, depends on data and
cut-off)

# Lab 12: Generalize kernel smoother to multivariate data

**Title:** Lab 12: Generalizing Kernel Smoother to Multivariate Data

**Aim:** To implement a multivariate kernel smoother (e.g., Nadaraya-Watson kernel regression) for predicting a continuous target variable based on multiple input features.

**Procedure:**

1. **Understand Kernel Smoothing:** Kernel smoothing is a non-parametric regression technique. For a given query point, it estimates the target value as a weighted average of the target values of nearby training points. The weights are determined by a kernel function, which gives higher weights to closer points.
2. **Multivariate Extension:** For multivariate data, the distance calculation in the kernel function needs to be extended to multiple dimensions (e.g., Euclidean distance in N-dimensional space).
3. **Kernel Function:** Choose a kernel function. A common choice is the Gaussian (Radial Basis Function) kernel: $K(x, x_i) = \exp\left(-\frac{||x - x_i||^2}{2h^2}\right)$ where x is the query point, xi is a training point, ||.|| denotes Euclidean distance, and h is the bandwidth parameter.
4. **Nadaraya-Watson Estimator:** The prediction for a query point x is given by: $\hat{y}(x) = \frac{\sum_{i=1}^{m} K(x, x_i) y_i}{\sum_{i=1}^{m} K(x, x_i)}$ where m is the number of training points, yi is the target value for training point xi.
5. **Bandwidth (h) Selection:** The bandwidth parameter h is crucial. A small h leads to a very localized fit (high variance, low bias), while a large h leads to a smoother fit (low variance, high bias). It can be chosen via cross-validation.
6. **Implementation:** Implement the kernel function and the Nadaraya-Watson estimator.
7. **Prediction and Evaluation:** Apply the smoother to a dataset and evaluate its performance (e.g., MSE).

**Source Code:**

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error
from sklearn.datasets import make_regression # For generating synthetic
multivariate data

class MultivariateKernelSmoother:
    def __init__(self, bandwidth=1.0):
        self.bandwidth = bandwidth
        self.X_train = None
        self.y_train = None

    def _gaussian_kernel(self, distance):
        """Gaussian kernel function."""
        return np.exp(-0.5 * (distance / self.bandwidth)**2)

    def fit(self, X_train, y_train):
        """
        Stores the training data.

        Args:
            X_train (np.array): Training features.
            y_train (np.array): Training target values.
        """
```

```python
        self.X_train = X_train
        self.y_train = y_train
        print(f"Kernel Smoother initialized with bandwidth:
{self.bandwidth}")
        print(f"Training data shape: {self.X_train.shape},
{self.y_train.shape}")

    def predict(self, X_test):
        """
        Predicts target values for new data points using Nadaraya-Watson.

        Args:
            X_test (np.array): Test features.

        Returns:
            np.array: Predicted target values.
        """
        if self.X_train is None or self.y_train is None:
            raise ValueError("Model must be fitted before making
predictions.")

        n_test_samples = X_test.shape[0]
        y_pred = np.zeros(n_test_samples)

        for i in range(n_test_samples):
            query_point = X_test[i]

            # Calculate Euclidean distances from query_point to all training
points
            distances = np.linalg.norm(self.X_train - query_point, axis=1)

            # Apply kernel function to get weights
            weights = self._gaussian_kernel(distances)

            # Handle cases where all weights are zero (e.g., query point too
far from all training points)
            sum_of_weights = np.sum(weights)
            if sum_of_weights == 0:
                # If no points are close, return the mean of training targets
or a default value
                y_pred[i] = np.mean(self.y_train)
                # print(f"Warning: Sum of weights is zero for point {i}.
Returning mean of training targets.")
            else:
                # Calculate weighted average
                y_pred[i] = np.sum(weights * self.y_train) / sum_of_weights
        return y_pred

if __name__ == "__main__":
    print("--- Multivariate Kernel Smoother Demonstration ---")

    # Generate synthetic multivariate regression data
    # n_features > 1 for multivariate
    X, y = make_regression(n_samples=300, n_features=3, n_informative=2,
                           n_targets=1, noise=15, random_state=42)
    print(f"Generated data shape: X={X.shape}, y={y.shape}")

    # Split data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
    print(f"\nTraining data shape: X_train={X_train.shape},
y_train={y_train.shape}")
    print(f"Testing data shape: X_test={X_test.shape},
y_test={y_test.shape}")

    # Standardize features (important for distance-based methods)
```

```
    scaler = StandardScaler()
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)
    print("Features scaled.")

    # Initialize and train the Kernel Smoother
    # Bandwidth 'h' is crucial. Experiment with different values.
    # A smaller h means more local smoothing, larger h means more global
smoothing.
    bandwidth = 0.5 # Example bandwidth
    kernel_smoother = MultivariateKernelSmoother(bandwidth=bandwidth)
    kernel_smoother.fit(X_train_scaled, y_train)

    # Make predictions on the test set
    y_pred = kernel_smoother.predict(X_test_scaled)

    # Evaluate the model
    mse = mean_squared_error(y_test, y_pred)
    print(f"\nModel Evaluation on Test Set (Bandwidth={bandwidth}):")
    print(f"Mean Squared Error (MSE): {mse:.4f}")

    # --- Visualizing for a 1D slice (if features allow) or just show actual
vs predicted ---
    # For multivariate data, direct 2D/3D plotting of the function is hard.
    # Instead, we can plot actual vs. predicted values.
    plt.figure(figsize=(10, 6))
    plt.scatter(y_test, y_pred, alpha=0.7, edgecolor='k')
    plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], 'r--',
lw=2, label='Perfect Prediction')
    plt.xlabel('Actual Values')
    plt.ylabel('Predicted Values')
    plt.title('Actual vs. Predicted Values (Multivariate Kernel Smoother)')
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.legend()
    plt.show()

    # Experiment with different bandwidths
    print("\n--- Experimenting with different bandwidths ---")
    bandwidths_to_try = [0.1, 0.5, 1.0, 2.0]
    for h in bandwidths_to_try:
        smoother = MultivariateKernelSmoother(bandwidth=h)
        smoother.fit(X_train_scaled, y_train)
        y_p = smoother.predict(X_test_scaled)
        mse_h = mean_squared_error(y_test, y_p)
        print(f"Bandwidth={h:.1f}, MSE: {mse_h:.4f}")
```

**Input:** The data is synthetically generated using `sklearn.datasets.make_regression()`. No explicit input file is required.

**Expected Output:** (A scatter plot of actual vs. predicted values will be displayed. Numerical MSE values for different bandwidths will be printed.)

```
--- Multivariate Kernel Smoother Demonstration ---
Generated data shape: X=(300, 3), y=(300,)

Training data shape: X_train=(240, 3), y_train=(240,)
Testing data shape: X_test=(60, 3), y_test=(60,)
Features scaled.
Kernel Smoother initialized with bandwidth: 0.5
Training data shape: (240, 3), (240,)

Model Evaluation on Test Set (Bandwidth=0.5):
Mean Squared Error (MSE): 225.0000 (example value, actual value may vary)
```

(A matplotlib scatter plot showing Actual vs. Predicted Values will appear.)

--- Experimenting with different bandwidths ---
Kernel Smoother initialized with bandwidth: 0.1
Training data shape: (240, 3), (240,)
Bandwidth=0.1, MSE: 350.0000 (example value)
Kernel Smoother initialized with bandwidth: 0.5
Training data shape: (240, 3), (240,)
Bandwidth=0.5, MSE: 225.0000 (example value)
Kernel Smoother initialized with bandwidth: 1.0
Training data shape: (240, 3), (240,)
Bandwidth=1.0, MSE: 180.0000 (example value)
Kernel Smoother initialized with bandwidth: 2.0
Training data shape: (240, 3), (240,)
Bandwidth=2.0, MSE: 200.0000 (example value)

# Lab 13: Implement a tree induction algorithm with backtracking.

**Title:** Lab 13: Implementing a Tree Induction Algorithm with Backtracking (ID3-like)

**Aim:** To implement a decision tree induction algorithm that constructs a decision tree from a dataset, incorporating a conceptual understanding of backtracking (though in standard ID3, backtracking is not explicit as it's a greedy approach). This lab will focus on building a basic ID3-like algorithm.

**Procedure:**

1. **Understand Decision Trees:** Decision trees are non-parametric supervised learning methods used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features.
2. **ID3 Algorithm Overview:**
    - **Choose Best Attribute:** Select the attribute that best splits the data. This is typically done using Information Gain or Gini Impurity. Information Gain is based on entropy.
    - **Create Node:** Create a decision node for the chosen attribute.
    - **Create Branches:** For each possible value of the attribute, create a branch.
    - **Partition Data:** Partition the dataset based on the attribute's values.
    - **Recursion:** Recursively call the algorithm for each branch with the reduced dataset and remaining attributes.
    - **Stopping Conditions:** Stop recursion when:
        - All examples in a partition belong to the same class (create a leaf node with that class).
        - No more attributes left to split on (create a leaf node with the majority class).
        - No examples left in a partition (create a leaf node with the majority class of the parent's data).
3. **Information Gain Calculation:**
    - **Entropy:** $H(S) = -\sum_{c \in C} p(c) \log_2 p(c)$
    - **Information Gain:** $Gain(S, A) = H(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} H(S_v)$
4. **Backtracking (Conceptual):** While standard ID3 is greedy and doesn't explicitly backtrack, the "backtracking" element in tree induction often refers to pruning (post-pruning) or exploring alternative splits if a greedy choice leads to a poor sub-tree. For this lab, we will focus on the core ID3 construction. If a split leads to an empty subset or a subset with no further distinguishing attributes, the algorithm effectively "backtracks" by creating a leaf node based on the majority class of the parent node's data.

**Source Code:**

```
import numpy as np
import pandas as pd
from collections import Counter

class DecisionTreeID3:
    def __init__(self):
        self.tree = None
        self.feature_names = None

    def _entropy(self, y):
        """Calculates the entropy of a target variable."""
        class_counts = Counter(y)
```

```python
        total_samples = len(y)
        entropy = 0.0
        for count in class_counts.values():
            probability = count / total_samples
            entropy -= probability * np.log2(probability)
        return entropy

    def _information_gain(self, X, y, feature_index):
        """Calculates the information gain for a given feature."""
        total_entropy = self._entropy(y)
        weighted_avg_entropy = 0.0
        feature_values = np.unique(X[:, feature_index])

        for value in feature_values:
            subset_indices = np.where(X[:, feature_index] == value)
            y_subset = y[subset_indices]
            if len(y_subset) > 0: # Avoid division by zero for empty subsets
                weighted_avg_entropy += (len(y_subset) / len(y)) *
self._entropy(y_subset)
        return total_entropy - weighted_avg_entropy

    def _find_best_split(self, X, y, available_features):
        """Finds the feature that yields the highest information gain."""
        best_gain = -1
        best_feature_index = -1

        for feature_index in available_features:
            gain = self._information_gain(X, y, feature_index)
            if gain > best_gain:
                best_gain = gain
                best_feature_index = feature_index
        return best_feature_index, best_gain

    def _build_tree(self, X, y, available_features):
        """Recursively builds the decision tree."""
        # Stopping condition 1: All examples belong to the same class
        if len(np.unique(y)) == 1:
            return {'leaf': y[0]}

        # Stopping condition 2: No more features to split on, or no examples
        if len(available_features) == 0 or len(y) == 0:
            # Return majority class of the current subset
            return {'leaf': Counter(y).most_common(1)[0][0]}

        # Find the best feature to split on
        best_feature_index, gain = self._find_best_split(X, y,
available_features)

        # Stopping condition 3: No gain from any feature (e.g., all features
are identical)
        if best_feature_index == -1 or gain == 0:
            return {'leaf': Counter(y).most_common(1)[0][0]}


        node = {self.feature_names[best_feature_index]: {}}
        remaining_features = [f for f in available_features if f !=
best_feature_index]

        # Create branches for each value of the best feature
        for value in np.unique(X[:, best_feature_index]):
            subset_indices = np.where(X[:, best_feature_index] == value)
            X_subset = X[subset_indices]
            y_subset = y[subset_indices]

            if len(y_subset) == 0:
                # This is where "backtracking" might conceptually occur:
```

```python
                    # if a subset is empty, we can't split further.
                    # We assign the majority class of the parent node's data.
                    node[self.feature_names[best_feature_index]][value] =
{'leaf': Counter(y).most_common(1)[0][0]}
                else:
                    node[self.feature_names[best_feature_index]][value] = \
                        self._build_tree(X_subset, y_subset, remaining_features)
        return node

    def fit(self, X, y, feature_names):
        """
        Fits the decision tree model.

        Args:
            X (np.array): Feature matrix.
            y (np.array): Target vector.
            feature_names (list): List of names for each feature.
        """
        self.feature_names = feature_names
        available_features = list(range(X.shape[1]))
        self.tree = self._build_tree(X, y, available_features)
        print("Decision Tree built successfully.")

    def _predict_single(self, instance, tree):
        """Helper to predict a single instance."""
        if 'leaf' in tree:
            return tree['leaf']

        # Find the current decision node's feature name
        current_feature_name = list(tree.keys())[0]
        feature_index = self.feature_names.index(current_feature_name)
        feature_value = instance[feature_index]

        # Traverse the tree
        if feature_value in tree[current_feature_name]:
            return self._predict_single(instance,
tree[current_feature_name][feature_value])
        else:
            # Handle unseen feature values: return majority class of the
current node's children
            # or a default. For simplicity, we'll return the majority class
of the parent node's data
            # (which is not directly available here, so we'll just return a
common class).
            # In a real scenario, this would involve more sophisticated
handling.
            # For this basic implementation, we'll return the most common
class from the original training data.
            # A more robust approach would involve finding the closest known
branch or a default.
            # For demonstration, let's just return the most common class from
the entire training set.
            # This is a simplification for unseen values, in practice, you'd
need a strategy.
            # For simplicity, we'll assume the model has seen all possible
values in training.
            # If an unseen value is encountered, we'll return a placeholder.
            # For a proper ID3, this shouldn't happen if the test data
attributes are from the same domain.
            print(f"Warning: Unseen feature value '{feature_value}' for
feature '{current_feature_name}'.")
            # Fallback: return the majority class of the training data
            return Counter(self.y_train_for_fallback).most_common(1)[0][0]


    def predict(self, X_test):
```

```python
        """
        Makes predictions for a dataset.

        Args:
            X_test (np.array): Test feature matrix.

        Returns:
            list: Predicted class labels.
        """
        predictions = []
        # Store original y_train for fallback in _predict_single
        # This is a hack for the simple _predict_single, in a real model,
        # the tree structure itself would implicitly handle this.
        self.y_train_for_fallback = X_test # This is just to pass something,
not actually used for prediction logic.
                                           # The assumption for ID3 is
discrete features and known values.
        for instance in X_test:
            predictions.append(self._predict_single(instance, self.tree))
        return predictions

    def print_tree(self, tree=None, indent=""):
        """Prints the decision tree structure."""
        if tree is None:
            tree = self.tree

        if 'leaf' in tree:
            print(f"{indent}Leaf: {tree['leaf']}")
            return

        feature_name = list(tree.keys())[0]
        print(f"{indent}Decision on: {feature_name}")
        for value, subtree in tree[feature_name].items():
            print(f"{indent}  Value '{value}':")
            self.print_tree(subtree, indent + "    ")

if __name__ == "__main__":
    print("--- Decision Tree (ID3-like) Implementation Demonstration ---")

    # Sample Dataset: Play Tennis (from Machine Learning by Tom Mitchell)
    data = {
        'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain',
'Overcast', 'Sunny', 'Sunny', 'Rain', 'Sunny', 'Overcast', 'Overcast',
'Rain'],
        'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Cool',
'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
        'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal',
'Normal', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal', 'High'],
        'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong',
'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Strong'],
        'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No',
'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
    }
    df = pd.DataFrame(data)

    # Convert categorical features to numerical (for internal processing,
though ID3 handles categories directly)
    # For a pure ID3, we work directly with categorical values.
    # Let's keep them as strings for direct ID3 implementation.
    X = df[['Outlook', 'Temperature', 'Humidity', 'Wind']].values # Convert
to numpy array
    y = df['PlayTennis'].values
    feature_names = ['Outlook', 'Temperature', 'Humidity', 'Wind']

    print(f"Dataset shape: X={X.shape}, y={y.shape}")
    print(f"Features: {feature_names}")
```

```
    print(f"Target classes: {np.unique(y)}")

    # Create and train the Decision Tree
    tree_classifier = DecisionTreeID3()
    tree_classifier.fit(X, y, feature_names)

    print("\n--- Learned Decision Tree Structure ---")
    tree_classifier.print_tree()

    # Make predictions on some new instances
    print("\n--- Making Predictions ---")
    test_instances = np.array([
        ['Sunny', 'Cool', 'High', 'Strong'], # Should be 'No'
        ['Rain', 'Mild', 'Normal', 'Weak'],  # Should be 'Yes'
        ['Overcast', 'Hot', 'Normal', 'Weak'] # Should be 'Yes'
    ])
    test_feature_names = ['Outlook', 'Temperature', 'Humidity', 'Wind'] # For
consistency

    # To handle the fallback in predict_single, we temporarily assign
y_train_for_fallback
    # This is not ideal for general use, but for this specific lab's
simplified _predict_single
    tree_classifier.y_train_for_fallback = y # Assign the full training y for
fallback

    predictions = tree_classifier.predict(test_instances)

    for i, instance in enumerate(test_instances):
        print(f"Instance: {instance}, Predicted: {predictions[i]}")

    # Evaluate on training data (for demonstration purposes, not true
evaluation)
    train_predictions = tree_classifier.predict(X)
    accuracy = np.sum(train_predictions == y) / len(y)
    print(f"\nAccuracy on training data: {accuracy:.2f}")
```

**Input:** The `data` dictionary used to create the pandas DataFrame `df`.

```
{
    'Outlook': ['Sunny', 'Sunny', 'Overcast', 'Rain', 'Rain', 'Rain',
'Overcast', 'Sunny', 'Sunny', 'Rain', 'Sunny', 'Overcast', 'Overcast',
'Rain'],
    'Temperature': ['Hot', 'Hot', 'Hot', 'Mild', 'Cool', 'Cool', 'Cool',
'Cool', 'Mild', 'Cool', 'Mild', 'Mild', 'Mild', 'Hot', 'Mild'],
    'Humidity': ['High', 'High', 'High', 'High', 'Normal', 'Normal',
'Normal', 'Normal', 'High', 'Normal', 'Normal', 'High', 'Normal', 'High'],
    'Wind': ['Weak', 'Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong',
'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak', 'Strong'],
    'PlayTennis': ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes',
'Yes', 'Yes', 'Yes', 'Yes', 'No']
}
```

And the new instances for prediction:

```
test_instances = np.array([
    ['Sunny', 'Cool', 'High', 'Strong'],
    ['Rain', 'Mild', 'Normal', 'Weak'],
    ['Overcast', 'Hot', 'Normal', 'Weak']
])
```

**Expected Output:** (The exact tree structure and predictions will depend on the information gain calculations and the specific data. The output below is an example based on the typical ID3 tree for this dataset.)

```
--- Decision Tree (ID3-like) Implementation Demonstration ---
Dataset shape: X=(14, 4), y=(14,)
Features: ['Outlook', 'Temperature', 'Humidity', 'Wind']
Target classes: ['No' 'Yes']
Decision Tree built successfully.

--- Learned Decision Tree Structure ---
Decision on: Outlook
  Value 'Overcast':
    Leaf: Yes
  Value 'Rain':
    Decision on: Wind
      Value 'Strong':
        Leaf: No
      Value 'Weak':
        Leaf: Yes
  Value 'Sunny':
    Decision on: Humidity
      Value 'High':
        Leaf: No
      Value 'Normal':
        Leaf: Yes

--- Making Predictions ---
Instance: ['Sunny' 'Cool' 'High' 'Strong'], Predicted: No
Instance: ['Rain' 'Mild' 'Normal' 'Weak'], Predicted: Yes
Instance: ['Overcast' 'Hot' 'Normal' 'Weak'], Predicted: Yes

Accuracy on training data: 1.00
```

# Lab 14: Implement a rule induction algorithm for regression

**Title:** Lab 14: Implementing a Rule Induction Algorithm for Regression (Simple Rule-Based Regression)

**Aim:** To implement a basic rule induction algorithm for regression tasks. This involves creating rules that partition the feature space and assign a predicted continuous value to each partition.

**Procedure:**

1. **Understand Rule Induction for Regression:** Instead of predicting a class label, rule induction for regression predicts a continuous value. Each rule will be of the form "IF condition THEN predict value". The predicted value for a rule is typically the mean or median of the target values of the training instances covered by that rule.
2. **Algorithm Strategy (Greedy Approach):**
   o Start with a single rule covering all data, predicting the mean of the target.
   o Iteratively refine rules:
      ▪ Find the best split (attribute and threshold) that reduces the error (e.g., Mean Squared Error) most significantly.
      ▪ Split the data based on this rule.
      ▪ Create new rules for the partitioned data.
      ▪ This is similar to how a regression tree is built, but the output is a set of rules.
3. **Splitting Criteria:** For numerical features, consider splitting based on a threshold (e.g., `feature <= threshold` and `feature > threshold`). For categorical features, split based on individual categories. The "best" split is one that minimizes the sum of squared errors in the resulting subsets.
4. **Prediction:** For a new instance, traverse the rules. If multiple rules apply, a strategy like taking the average of their predictions or using the prediction of the most specific rule (covering the fewest training examples) can be used. For simplicity, we'll build a tree-like structure and convert it to rules.

**Source Code:**

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.datasets import make_regression # For synthetic data

class RegressionRuleInduction:
    def __init__(self, max_depth=3, min_samples_leaf=5):
        self.max_depth = max_depth
        self.min_samples_leaf = min_samples_leaf
        self.tree = None
        self.feature_names = None

    def _calculate_prediction(self, y):
        """Calculates the prediction for a leaf node (mean of target
values)."""
        if len(y) == 0:
            return 0.0 # Or some other default/fallback
        return np.mean(y)

    def _calculate_mse(self, y_true, y_pred):
        """Calculates Mean Squared Error."""
        return np.mean((y_true - y_pred)**2)
```

```python
    def _find_best_split(self, X, y, current_features):
        """
        Finds the best split (feature and threshold) that minimizes MSE.
        Assumes numerical features for simplicity.
        """
        best_mse = float('inf')
        best_split = None # (feature_index, threshold)

        n_samples, n_features = X.shape

        for feature_index in current_features:
            unique_values = np.unique(X[:, feature_index])
            # Consider potential thresholds between unique values
            if len(unique_values) > 1:
                thresholds = (unique_values[:-1] + unique_values[1:]) / 2.0
            else:
                thresholds = unique_values # Only one value, no split
possible

            for threshold in thresholds:
                left_indices = np.where(X[:, feature_index] <= threshold)
                right_indices = np.where(X[:, feature_index] > threshold)

                y_left = y[left_indices]
                y_right = y[right_indices]

                # Ensure minimum samples in each leaf
                if len(y_left) < self.min_samples_leaf or len(y_right) <
self.min_samples_leaf:
                    continue

                # Calculate MSE for potential split
                pred_left = self._calculate_prediction(y_left)
                pred_right = self._calculate_prediction(y_right)

                mse_left = self._calculate_mse(y_left, pred_left)
                mse_right = self._calculate_mse(y_right, pred_right)

                # Weighted average MSE for the split
                weighted_mse = (len(y_left) / n_samples) * mse_left + \
                               (len(y_right) / n_samples) * mse_right

                if weighted_mse < best_mse:
                    best_mse = weighted_mse
                    best_split = (feature_index, threshold)
        return best_split, best_mse

    def _build_tree(self, X, y, current_depth, current_features):
        """Recursively builds the regression tree."""
        # Stopping conditions
        if current_depth >= self.max_depth or \
           len(y) < self.min_samples_leaf or \
           len(np.unique(y)) == 1: # All targets are the same
            return {'leaf': self._calculate_prediction(y)}

        best_split, best_mse = self._find_best_split(X, y, current_features)

        if best_split is None: # No good split found
            return {'leaf': self._calculate_prediction(y)}

        feature_index, threshold = best_split
        node = {
            'feature': self.feature_names[feature_index],
            'threshold': threshold,
            'left': None,
            'right': None
```

```python
            }

        left_indices = np.where(X[:, feature_index] <= threshold)
        right_indices = np.where(X[:, feature_index] > threshold)

        node['left'] = self._build_tree(X[left_indices], y[left_indices],
                                        current_depth + 1, current_features)
        node['right'] = self._build_tree(X[right_indices], y[right_indices],
                                         current_depth + 1, current_features)
        return node

    def fit(self, X, y, feature_names):
        """
        Fits the regression rule induction model.

        Args:
            X (np.array): Feature matrix.
            y (np.array): Target vector.
            feature_names (list): List of names for each feature.
        """
        self.feature_names = feature_names
        available_features = list(range(X.shape[1]))
        self.tree = self._build_tree(X, y, 0, available_features)
        print("Regression Rule Induction model (Decision Tree based) built
successfully.")

    def _predict_single(self, instance, tree):
        """Helper to predict a single instance."""
        if 'leaf' in tree:
            return tree['leaf']

        feature_name = tree['feature']
        threshold = tree['threshold']
        feature_index = self.feature_names.index(feature_name)

        if instance[feature_index] <= threshold:
            return self._predict_single(instance, tree['left'])
        else:
            return self._predict_single(instance, tree['right'])

    def predict(self, X_test):
        """
        Makes predictions for a dataset.

        Args:
            X_test (np.array): Test feature matrix.

        Returns:
            np.array: Predicted target values.
        """
        predictions = np.array([self._predict_single(instance, self.tree) for
instance in X_test])
        return predictions

    def print_tree_as_rules(self, tree=None, rule_prefix="IF ", indent=""):
        """Prints the regression tree as a set of IF-THEN rules."""
        if tree is None:
            tree = self.tree

        if 'leaf' in tree:
            print(f"{indent}{rule_prefix} THEN Predict = {tree['leaf']:.4f}")
            return

        feature_name = tree['feature']
        threshold = tree['threshold']
```

```python
            # Left branch (<= threshold)
            left_rule_prefix = f"{rule_prefix}{feature_name} <= {threshold:.2f}"
            self.print_tree_as_rules(tree['left'], rule_prefix=left_rule_prefix +
" AND ", indent=indent)

            # Right branch (> threshold)
            right_rule_prefix = f"{rule_prefix}{feature_name} > {threshold:.2f}"
            self.print_tree_as_rules(tree['right'], rule_prefix=right_rule_prefix
+ " AND ", indent=indent)


if __name__ == "__main__":
    print("--- Regression Rule Induction (Decision Tree based) Demonstration
---")

    # Generate synthetic regression data
    X, y = make_regression(n_samples=200, n_features=2, n_informative=2,
                           n_targets=1, noise=10, random_state=42)
    feature_names = ['Feature_1', 'Feature_2']

    print(f"Generated data shape: X={X.shape}, y={y.shape}")

    # Split data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
    print(f"\nTraining data shape: X_train={X_train.shape},
y_train={y_train.shape}")
    print(f"Testing data shape: X_test={X_test.shape},
y_test={y_test.shape}")

    # Create and train the Regression Rule Induction model
    # max_depth controls complexity, min_samples_leaf controls purity of
leaves
    model = RegressionRuleInduction(max_depth=4, min_samples_leaf=10)
    model.fit(X_train, y_train, feature_names)

    print("\n--- Learned Rules (from Decision Tree structure) ---")
    model.print_tree_as_rules()

    # Make predictions on the test set
    y_pred = model.predict(X_test)

    # Evaluate the model
    mse = mean_squared_error(y_test, y_pred)
    print(f"\nModel Evaluation on Test Set:")
    print(f"Mean Squared Error (MSE): {mse:.4f}")

    # Visualize actual vs. predicted values for a simple 1D slice or just
scatter plot
    plt.figure(figsize=(10, 6))
    plt.scatter(y_test, y_pred, alpha=0.7, edgecolor='k')
    plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], 'r--',
lw=2, label='Perfect Prediction')
    plt.xlabel('Actual Values')
    plt.ylabel('Predicted Values')
    plt.title('Actual vs. Predicted Values (Regression Rule Induction)')
    plt.grid(True, linestyle='--', alpha=0.6)
    plt.legend()
    plt.show()
```

**Input:** The data is synthetically generated using `sklearn.datasets.make_regression()`. No explicit input file is required.

**Expected Output:** (The exact rules and MSE will vary based on the random data generation and max_depth/min_samples_leaf parameters. A scatter plot of actual vs. predicted values will be displayed.)

```
--- Regression Rule Induction (Decision Tree based) Demonstration ---
Generated data shape: X=(200, 2), y=(200,)

Training data shape: X_train=(160, 2), y_train=(160,)
Testing data shape: X_test=(40, 2), y_test=(40,)
Regression Rule Induction model (Decision Tree based) built successfully.

--- Learned Rules (from Decision Tree structure) ---
IF Feature_1 <= -0.12 AND Feature_2 <= -0.25 AND THEN Predict = -35.0000
IF Feature_1 <= -0.12 AND Feature_2 > -0.25 AND THEN Predict = -10.0000
IF Feature_1 > -0.12 AND Feature_1 <= 0.50 AND Feature_2 <= 0.10 AND THEN
Predict = 20.0000
IF Feature_1 > -0.12 AND Feature_1 <= 0.50 AND Feature_2 > 0.10 AND THEN
Predict = 45.0000
IF Feature_1 > 0.50 AND Feature_2 <= 0.80 AND THEN Predict = 70.0000
IF Feature_1 > 0.50 AND Feature_2 > 0.80 AND THEN Predict = 95.0000
(Example rules, actual rules will vary)

Model Evaluation on Test Set:
Mean Squared Error (MSE): 120.0000 (example value, actual value may vary)

(A matplotlib scatter plot showing Actual vs. Predicted Values will appear.)
```

# Lab 15: Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets

**Title:** Lab 15: Building an Artificial Neural Network with Backpropagation

**Aim:** To implement a simple Artificial Neural Network (ANN) and train it using the Backpropagation algorithm for a classification task.

**Procedure:**

1. **Understand ANN and Backpropagation:**
   - **ANN:** A network of interconnected nodes (neurons) organized in layers (input, hidden, output). Each connection has a weight, and each neuron has a bias.
   - **Forward Pass:** Input data propagates through the network, activating neurons and producing an output.
   - **Backpropagation:** An algorithm used to train ANNs by calculating the gradient of the loss function with respect to the weights and biases. It works by propagating the error backward from the output layer to the input layer, adjusting weights to minimize the error.
2. **Network Architecture:** Define the number of input features, hidden layers (and neurons per layer), and output neurons.
3. **Activation Functions:** Choose activation functions for hidden layers (e.g., ReLU, Sigmoid) and the output layer (e.g., Sigmoid for binary classification, Softmax for multi-class).
4. **Weight and Bias Initialization:** Initialize weights randomly (e.g., using Xavier/Glorot or He initialization) and biases to zeros or small constants.
5. **Loss Function:** Define the loss function (e.g., Binary Cross-Entropy for binary classification).
6. **Training Loop:**
   - **Forward Pass:** Compute predictions.
   - **Calculate Loss:** Compute the loss between predictions and true labels.
   - **Backward Pass (Backpropagation):** Calculate gradients of the loss with respect to weights and biases using the chain rule.
   - **Weight Update:** Adjust weights and biases using an optimizer (e.g., Gradient Descent or Adam).
7. **Prediction:** Use the trained network to make predictions.
8. **Evaluation:** Evaluate the model's performance (e.g., accuracy).

**Source Code:**

```
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.datasets import make_moons, load_iris
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.metrics import accuracy_score

class NeuralNetwork:
    def __init__(self, input_size, hidden_size, output_size,
learning_rate=0.01, n_iterations=10000):
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.output_size = output_size
        self.learning_rate = learning_rate
        self.n_iterations = n_iterations
```

```python
        # Initialize weights and biases
        # Weights from input to hidden layer (W1)
        # Using Xavier/Glorot initialization for better training stability
        self.W1 = np.random.randn(self.input_size, self.hidden_size) *
np.sqrt(2 / (input_size + hidden_size))
        self.b1 = np.zeros((1, self.hidden_size))

        # Weights from hidden to output layer (W2)
        self.W2 = np.random.randn(self.hidden_size, self.output_size) *
np.sqrt(2 / (hidden_size + output_size))
        self.b2 = np.zeros((1, self.output_size))

        print(f"Neural Network initialized: Input={input_size},
Hidden={hidden_size}, Output={output_size}")
        print(f"Learning Rate: {learning_rate}, Iterations: {n_iterations}")

    def _sigmoid(self, x):
        """Sigmoid activation function."""
        return 1 / (1 + np.exp(-x))

    def _sigmoid_derivative(self, x):
        """Derivative of the sigmoid function."""
        return x * (1 - x)

    def _relu(self, x):
        """ReLU activation function."""
        return np.maximum(0, x)

    def _relu_derivative(self, x):
        """Derivative of the ReLU function."""
        return (x > 0).astype(float)

    def _softmax(self, x):
        """Softmax activation function for multi-class output."""
        exp_x = np.exp(x - np.max(x, axis=1, keepdims=True)) # For numerical
stability
        return exp_x / np.sum(exp_x, axis=1, keepdims=True)

    def _binary_cross_entropy_loss(self, y_true, y_pred):
        """Binary Cross-Entropy Loss for binary classification."""
        # Add a small epsilon to avoid log(0)
        epsilon = 1e-10
        return -np.mean(y_true * np.log(y_pred + epsilon) + (1 - y_true) *
np.log(1 - y_pred + epsilon))

    def _forward_pass(self, X):
        """
        Performs the forward pass through the network.
        Returns activations at each layer.
        """
        # Hidden layer
        self.z1 = np.dot(X, self.W1) + self.b1
        self.a1 = self._relu(self.z1) # Using ReLU for hidden layer

        # Output layer
        self.z2 = np.dot(self.a1, self.W2) + self.b2
        # Use sigmoid for binary classification, softmax for multi-class
        if self.output_size == 1:
            self.a2 = self._sigmoid(self.z2)
        else:
            self.a2 = self._softmax(self.z2)
        return self.a2

    def _backward_pass(self, X, y_true, y_pred):
        """
```

```
        Performs the backward pass (backpropagation) to compute gradients.
        """
        n_samples = X.shape[0]

        # Output layer error
        # For binary cross-entropy with sigmoid: dL/dz2 = y_pred - y_true
        # For multi-class cross-entropy with softmax: dL/dz2 = y_pred -
y_true (if y_true is one-hot)
        delta2 = y_pred - y_true

        # Gradients for W2 and b2
        dW2 = (1 / n_samples) * np.dot(self.a1.T, delta2)
        db2 = (1 / n_samples) * np.sum(delta2, axis=0, keepdims=True)

        # Hidden layer error
        # delta1 = (delta2 @ W2.T) * relu_derivative(z1)
        delta1 = (np.dot(delta2, self.W2.T)) * self._relu_derivative(self.a1)

        # Gradients for W1 and b1
        dW1 = (1 / n_samples) * np.dot(X.T, delta1)
        db1 = (1 / n_samples) * np.sum(delta1, axis=0, keepdims=True)

        return dW1, db1, dW2, db2

    def fit(self, X, y):
        """
        Trains the neural network using backpropagation.

        Args:
            X (np.array): Input features.
            y (np.array): Target labels (should be one-hot encoded for multi-
class).
        """
        print("\nStarting training...")
        for iteration in range(self.n_iterations):
            # Forward pass
            y_pred = self._forward_pass(X)

            # Calculate loss
            if self.output_size == 1: # Binary classification
                loss = self._binary_cross_entropy_loss(y, y_pred)
            else: # Multi-class classification (assuming y is one-hot)
                # Cross-entropy loss for softmax
                epsilon = 1e-10
                loss = -np.mean(np.sum(y * np.log(y_pred + epsilon), axis=1))


            # Backward pass (compute gradients)
            dW1, db1, dW2, db2 = self._backward_pass(X, y, y_pred)

            # Update weights and biases
            self.W1 -= self.learning_rate * dW1
            self.b1 -= self.learning_rate * db1
            self.W2 -= self.learning_rate * dW2
            self.b2 -= self.learning_rate * db2

            if iteration % (self.n_iterations // 10) == 0:
                print(f"Iteration {iteration}/{self.n_iterations}, Loss:
{loss:.4f}")
        print(f"Training complete. Final Loss: {loss:.4f}")


    def predict(self, X):
        """
        Makes predictions using the trained network.
```

```python
        Args:
            X (np.array): Input features for prediction.

        Returns:
            np.array: Predicted class labels.
        """
        probabilities = self._forward_pass(X)
        if self.output_size == 1: # Binary classification
            return (probabilities >= 0.5).astype(int)
        else: # Multi-class classification
            return np.argmax(probabilities, axis=1)

if __name__ == "__main__":
    print("--- Artificial Neural Network with Backpropagation Demonstration -
--")

    # --- Example 1: Binary Classification (Make Moons Dataset) ---
    print("\n--- Binary Classification: Make Moons Dataset ---")
    X_moon, y_moon = make_moons(n_samples=200, noise=0.15, random_state=42)
    y_moon = y_moon.reshape(-1, 1) # Reshape y for binary output layer (1
neuron)

    # Split and scale data
    X_train_moon, X_test_moon, y_train_moon, y_test_moon = train_test_split(
        X_moon, y_moon, test_size=0.2, random_state=42
    )
    scaler_moon = StandardScaler()
    X_train_moon_scaled = scaler_moon.fit_transform(X_train_moon)
    X_test_moon_scaled = scaler_moon.transform(X_test_moon)
    print(f"Moons Dataset shapes: X_train={X_train_moon_scaled.shape},
y_train={y_train_moon.shape}")

    # Create and train NN for binary classification
    nn_binary = NeuralNetwork(input_size=2, hidden_size=10, output_size=1,
                              learning_rate=0.1, n_iterations=10000)
    nn_binary.fit(X_train_moon_scaled, y_train_moon)

    # Evaluate
    y_pred_moon = nn_binary.predict(X_test_moon_scaled)
    accuracy_moon = accuracy_score(y_test_moon, y_pred_moon)
    print(f"\nBinary Classification Accuracy on Moons Test Set:
{accuracy_moon:.4f}")

    # Plot decision boundary (for 2D data)
    plt.figure(figsize=(10, 7))
    h = .02 # step size in the mesh
    x_min, x_max = X_moon[:, 0].min() - 1, X_moon[:, 0].max() + 1
    y_min, y_max = X_moon[:, 1].min() - 1, X_moon[:, 1].max() + 1
    xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
                         np.arange(y_min, y_max, h))
    # Predict on meshgrid points
    Z = nn_binary.predict(scaler_moon.transform(np.c_[xx.ravel(),
yy.ravel()]))
    Z = Z.reshape(xx.shape)
    plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral, alpha=0.8)
    plt.scatter(X_moon[:, 0], X_moon[:, 1], c=y_moon.flatten(),
cmap=plt.cm.Spectral, edgecolor='k')
    plt.title('Neural Network Decision Boundary (Moons Dataset)')
    plt.xlabel('Feature 1')
    plt.ylabel('Feature 2')
    plt.show()


    # --- Example 2: Multi-Class Classification (Iris Dataset) ---
    print("\n--- Multi-Class Classification: Iris Dataset ---")
    iris = load_iris()
```

```
    X_iris = iris.data
    y_iris = iris.target

    # One-hot encode target labels for multi-class classification
    encoder = OneHotEncoder(sparse_output=False)
    y_iris_onehot = encoder.fit_transform(y_iris.reshape(-1, 1))

    # Split and scale data
    X_train_iris, X_test_iris, y_train_iris_onehot, y_test_iris_onehot =
train_test_split(
        X_iris, y_iris_onehot, test_size=0.2, random_state=42
    )
    scaler_iris = StandardScaler()
    X_train_iris_scaled = scaler_iris.fit_transform(X_train_iris)
    X_test_iris_scaled = scaler_iris.transform(X_test_iris)
    print(f"Iris Dataset shapes: X_train={X_train_iris_scaled.shape},
y_train={y_train_iris_onehot.shape}")

    # Create and train NN for multi-class classification
    nn_multi = NeuralNetwork(input_size=X_iris.shape[1], hidden_size=10,
output_size=len(iris.target_names),
                            learning_rate=0.05, n_iterations=15000)
    nn_multi.fit(X_train_iris_scaled, y_train_iris_onehot)

    # Evaluate
    y_pred_iris = nn_multi.predict(X_test_iris_scaled)
    # Convert one-hot encoded true labels back to single labels for
accuracy_score
    y_test_iris_labels = np.argmax(y_test_iris_onehot, axis=1)
    accuracy_iris = accuracy_score(y_test_iris_labels, y_pred_iris)
    print(f"\nMulti-Class Classification Accuracy on Iris Test Set:
{accuracy_iris:.4f}")
```

**Input:** The "Make Moons" dataset (`sklearn.datasets.make_moons`) and the "Iris" dataset (`sklearn.datasets.load_iris`) are used. No explicit input file is required.

**Expected Output:** (Numerical training loss and final accuracy will be printed. A scatter plot with the decision boundary for the Moons dataset will be displayed.)

```
--- Artificial Neural Network with Backpropagation Demonstration ---

--- Binary Classification: Make Moons Dataset ---
Neural Network initialized: Input=2, Hidden=10, Output=1
Learning Rate: 0.1, Iterations: 10000
Moons Dataset shapes: X_train=(160, 2), y_train=(160, 1)

Starting training...
Iteration 0/10000, Loss: 0.6931
Iteration 1000/10000, Loss: 0.3500
Iteration 2000/10000, Loss: 0.2000
Iteration 3000/10000, Loss: 0.1000
Iteration 4000/10000, Loss: 0.0500
Iteration 5000/10000, Loss: 0.0250
Iteration 6000/10000, Loss: 0.0150
Iteration 7000/10000, Loss: 0.0090
Iteration 8000/10000, Loss: 0.0060
Iteration 9000/10000, Loss: 0.0040
Training complete. Final Loss: 0.0030 (example values, actual values may
vary)

Binary Classification Accuracy on Moons Test Set: 0.9750 (example value,
actual value may vary)
```

(A matplotlib scatter plot showing the Moons dataset with the learned
decision boundary will appear.)

--- Multi-Class Classification: Iris Dataset ---
Neural Network initialized: Input=4, Hidden=10, Output=3
Learning Rate: 0.05, Iterations: 15000
Iris Dataset shapes: X_train=(120, 4), y_train=(120, 3)

Starting training...
Iteration 0/15000, Loss: 1.0986
Iteration 1500/15000, Loss: 0.1500
Iteration 3000/15000, Loss: 0.0500
Iteration 4500/15000, Loss: 0.0200
Iteration 6000/15000, Loss: 0.0100
Iteration 7500/15000, Loss: 0.0050
Iteration 9000/15000, Loss: 0.0030
Iteration 10500/15000, Loss: 0.0020
Iteration 12000/15000, Loss: 0.0010
Iteration 13500/15000, Loss: 0.0008
Training complete. Final Loss: 0.0005 (example values, actual values may
vary)

Multi-Class Classification Accuracy on Iris Test Set: 0.9667 (example value,
actual value may vary)