

For a data warehousing lab where you're working on numerical datasets to perform classification or clustering, it's good to be familiar with a variety of algorithms that can be applied in these scenarios. Here's a list of key algorithms to prepare for both **classification** and **clustering** tasks:

Classification Algorithms:

These algorithms are used to predict a category or class label based on input features.

1. **Logistic Regression**
 - Good for binary or multi-class classification problems.
 - A basic linear model for classification tasks.
 2. **Decision Trees**
 - A tree-like model that splits data based on feature values to make predictions.
 3. **Random Forest**
 - An ensemble method that uses multiple decision trees for more robust and accurate predictions.
 4. **Support Vector Machine (SVM)**
 - Effective for both linear and non-linear classification problems.
 - Finds the hyperplane that best separates the data into classes.
 5. **K-Nearest Neighbors (KNN)**
 - A simple instance-based learning algorithm that classifies data based on the majority class of its neighbors.
 6. **Naive Bayes**
 - Based on Bayes' Theorem, useful for text classification and cases where features are conditionally independent.
 7. **Neural Networks (Feedforward, MLP)**
 - Can be used for more complex classification tasks, especially when the relationship between features is non-linear.
 8. **Gradient Boosting (XGBoost, LightGBM, CatBoost)**
 - An ensemble method that builds trees in a sequential manner to correct the errors of the previous ones.
 9. **AdaBoost**
 - Another boosting method that combines weak classifiers to create a strong classifier.
-

Clustering Algorithms:

These algorithms are used for grouping similar data points together without predefined labels.

1. **K-Means Clustering**
 - A popular clustering algorithm that partitions data into K clusters based on distance from centroids.
2. **Hierarchical Clustering (Agglomerative, Divisive)**
 - Builds a tree-like structure (dendrogram) to show the relationships between data points, which can then be cut at a certain level to define clusters.
3. **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**

- A density-based algorithm that finds clusters of arbitrary shape and is resistant to noise.
 - 4. **Gaussian Mixture Model (GMM)**
 - Assumes that the data is generated from a mixture of several Gaussian distributions, suitable for soft clustering.
 - 5. **Mean Shift Clustering**
 - A non-parametric algorithm that seeks the densest areas in the feature space to find clusters.
 - 6. **Affinity Propagation**
 - A clustering algorithm that identifies exemplars from the data to represent clusters without requiring the number of clusters to be specified in advance.
 - 7. **Spectral Clustering**
 - Uses eigenvalues of a similarity matrix to reduce dimensions and perform clustering.
 - 8. **Self-Organizing Maps (SOM)**
 - A type of artificial neural network used for clustering high-dimensional data into a lower-dimensional grid.
-

Dimensionality Reduction Techniques (Often Useful for Both Classification & Clustering):

1. **Principal Component Analysis (PCA)**
 - Reduces the number of features by finding the most important directions in the data.
-

Got it! If you're working with a dataset given in a file format (like CSV, Excel, etc.), we can use **pandas** to load and read the dataset. I'll also include the **Decision Trees** algorithm as requested.

Here's how you can modify the code to load the dataset using **pandas**, perform basic data exploration, and include **Decision Trees**:

Required Libraries:

```
pip install scikit-learn xgboost pandas matplotlib
```

Full Python Code with pandas Data Loading and Decision Tree Added:

```
import numpy as np
import pandas as pd
```

```

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix,
classification_report
import xgboost as xgb
import matplotlib.pyplot as plt

# Load your dataset using pandas (Replace 'your_dataset.csv' with the
actual file name)
# If you're using a different format, pandas supports .csv, .xlsx, .json,
etc.
# For example, loading CSV:
df = pd.read_csv('your_dataset.csv')

# Check the first few rows of the dataset
print(df.head())

# Preprocessing:
# Assume that the target variable is 'target' and the rest are features
X = df.drop('target', axis=1) # Features (input variables)
y = df['target']              # Target (output variable)

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)

# Scale features (important for SVM, KNN, etc.)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# 1. Logistic Regression
print("\nLogistic Regression:")
log_reg = LogisticRegression(max_iter=200)
log_reg.fit(X_train, y_train)
y_pred_log_reg = log_reg.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred_log_reg)}")
print(confusion_matrix(y_test, y_pred_log_reg))
print(classification_report(y_test, y_pred_log_reg))

# 2. Random Forest
print("\nRandom Forest:")
rf = RandomForestClassifier(n_estimators=100, random_state=42)
rf.fit(X_train, y_train)
y_pred_rf = rf.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred_rf)}")
print(confusion_matrix(y_test, y_pred_rf))
print(classification_report(y_test, y_pred_rf))

# 3. Support Vector Machine (SVM)
print("\nSupport Vector Machine (SVM):")
svm = SVC(kernel='linear', random_state=42)
svm.fit(X_train, y_train)
y_pred_svm = svm.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred_svm)}")
print(confusion_matrix(y_test, y_pred_svm))

```

```

print(classification_report(y_test, y_pred_svm))

# 4. K-Nearest Neighbors (KNN)
print("\nK-Nearest Neighbors (KNN):")
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train, y_train)
y_pred_knn = knn.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred_knn)}")
print(confusion_matrix(y_test, y_pred_knn))
print(classification_report(y_test, y_pred_knn))

# 5. Decision Trees
print("\nDecision Trees:")
dt = DecisionTreeClassifier(random_state=42)
dt.fit(X_train, y_train)
y_pred_dt = dt.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred_dt)}")
print(confusion_matrix(y_test, y_pred_dt))
print(classification_report(y_test, y_pred_dt))

# 6. XGBoost
print("\nXGBoost:")
xg_clf = xgb.XGBClassifier(use_label_encoder=False, eval_metric='mlogloss')
xg_clf.fit(X_train, y_train)
y_pred_xgb = xg_clf.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred_xgb)}")
print(confusion_matrix(y_test, y_pred_xgb))
print(classification_report(y_test, y_pred_xgb))

# Optional: Visualize Decision Tree (if it's a small dataset)
plt.figure(figsize=(12, 8))
from sklearn.tree import plot_tree
plot_tree(dt, filled=True, feature_names=X.columns, class_names=[str(i) for
i in np.unique(y)], rounded=True)
plt.title("Decision Tree Visualization")
plt.show()

```

Steps in the Code:

1. Loading the Dataset:

- We load the dataset using `pd.read_csv('your_dataset.csv')`. You can change the format based on the file type (.xlsx, .json, etc.).
- **Data Exploration:** The `df.head()` command shows the first few rows of your dataset to check its structure.

2. Data Preprocessing:

- **Feature and Target Splitting:** We assume that the target variable is called 'target'. The rest are treated as features.
- **Splitting Data:** We use `train_test_split()` to divide the data into training and testing sets.
- **Scaling Features:** `StandardScaler` is used to standardize the feature values (important for SVM, KNN, etc.).

3. Algorithms:

- **Logistic Regression:** Uses `LogisticRegression()` for classification.
- **Random Forest:** Uses `RandomForestClassifier()` with 100 estimators.
- **Support Vector Machine (SVM):** Uses `SVC()` with a linear kernel.

- **K-Nearest Neighbors (KNN):** Uses `KNeighborsClassifier()` with 5 neighbors.
 - **Decision Trees:** Uses `DecisionTreeClassifier()` and fits it to the data.
 - **XGBoost:** Uses `XGBClassifier()` from the **XGBoost** library.
4. **Evaluation Metrics:**
- For each algorithm, we print the **accuracy**, **confusion matrix**, and **classification report** to evaluate model performance.
5. **Decision Tree Visualization (Optional):**
- If your dataset is small and you want to visualize the decision tree, we use `plot_tree()` from `sklearn.tree` to display the tree structure.

Output:

After running the code, you'll get:

- Accuracy score for each model.
- Confusion matrix showing the true vs predicted values.
- A classification report with precision, recall, and F1-score for each class.
- A visual representation of the decision tree (optional).

Important Notes:

- Replace `'your_dataset.csv'` with the actual path to your dataset.
- Adjust the target column name (`'target'`) if it's different in your dataset.
- Depending on the dataset, you may need to handle missing values, categorical data, etc., before fitting the model.

Let me know if you need more help!

Sure! I'll break it down and give you individual code snippets for each of the tasks you mentioned:

1. **Reading data from a file (CSV, XLSX, or DAT)**
 2. **K-Means Clustering with evaluation**
 3. **Agglomerative Hierarchical Clustering with evaluation**
 4. **Divisive Hierarchical Clustering with evaluation**
-

1. Reading Data from a File (CSV, XLSX, or DAT)

```
import pandas as pd

# Function to load data from a file (CSV, XLSX, DAT)
def load_data(file_path):
    if file_path.endswith('.csv'):
        return pd.read_csv(file_path)
    elif file_path.endswith('.xlsx'):
        return pd.read_excel(file_path)
    elif file_path.endswith('.dat'):
        return pd.read_csv(file_path, delimiter='\t') # assuming tab-delimited for .dat
    else:
        raise ValueError("Unsupported file format. Use CSV, XLSX, or DAT.")

# Load dataset
file_path = 'your_file.csv' # Replace with your file path
data = load_data(file_path)

# Preview the first few rows of the data
print(data.head())
```

Explanation:

- This code defines a function to load a file (CSV, XLSX, or DAT) into a pandas DataFrame.
 - Replace 'your_file.csv' with your actual file path.
-

2. K-Means Clustering with Evaluation Metrics (Silhouette Score, Adjusted Rand Index)

```
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score, adjusted_rand_score
from sklearn.preprocessing import StandardScaler

# Assuming you have already loaded the data as 'data'
# Preprocess the data (standardize it)
X = data.values # Assuming data is numeric
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Apply K-Means clustering
```

```

n_clusters = 3 # Specify number of clusters
kmeans = KMeans(n_clusters=n_clusters, random_state=42)
kmeans_labels = kmeans.fit_predict(X_scaled)

# Evaluation metrics (assuming ground truth labels are available in
'y_true')
y_true = [0, 1, 2] # Replace this with your actual ground truth labels (if
available)

# Silhouette Score
silhouette = silhouette_score(X_scaled, kmeans_labels)
print(f"K-Means Silhouette Score: {silhouette:.4f}")

# Adjusted Rand Index
ari = adjusted_rand_score(y_true, kmeans_labels)
print(f"K-Means Adjusted Rand Index: {ari:.4f}")

```

Explanation:

- This code applies **K-Means clustering** on the standardized data and calculates the **Silhouette Score** and **Adjusted Rand Index** to evaluate the clustering.
-

3. Agglomerative Hierarchical Clustering with Evaluation Metrics

```

from sklearn.cluster import AgglomerativeClustering
from sklearn.metrics import silhouette_score, adjusted_rand_score
from sklearn.preprocessing import StandardScaler

# Assuming you have already loaded and preprocessed the data as 'X_scaled'
# Apply Agglomerative Clustering
n_clusters = 3 # Specify number of clusters
agg_clustering = AgglomerativeClustering(n_clusters=n_clusters)
agg_labels = agg_clustering.fit_predict(X_scaled)

# Evaluation metrics (assuming ground truth labels are available in
'y_true')
y_true = [0, 1, 2] # Replace this with your actual ground truth labels (if
available)

# Silhouette Score
silhouette = silhouette_score(X_scaled, agg_labels)
print(f"Agglomerative Silhouette Score: {silhouette:.4f}")

# Adjusted Rand Index
ari = adjusted_rand_score(y_true, agg_labels)
print(f"Agglomerative Adjusted Rand Index: {ari:.4f}")

```

Explanation:

- This code performs **Agglomerative Hierarchical Clustering** and evaluates it using **Silhouette Score** and **Adjusted Rand Index**.
-

4. Divisive Hierarchical Clustering (Recursive K-Means Split)

```
import numpy as np
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

# Simple function for Divisive Hierarchical Clustering (recursive K-Means
splitting)
def divisive_clustering(X, n_clusters):
    clusters = [X]
    while len(clusters) < n_clusters:
        largest_cluster = clusters.pop(0)
        kmeans = KMeans(n_clusters=2, random_state=42)
        kmeans.fit(largest_cluster)
        cluster_1 = largest_cluster[kmeans.labels_ == 0]
        cluster_2 = largest_cluster[kmeans.labels_ == 1]
        clusters.append(cluster_1)
        clusters.append(cluster_2)
    return clusters

# Assuming you have already loaded and preprocessed the data as 'X_scaled'
n_clusters = 3 # Specify number of clusters
final_clusters = divisive_clustering(X_scaled, n_clusters)

# Create labels for divisive clustering based on final clusters
div_labels = np.zeros(X_scaled.shape[0])
for i, cluster in enumerate(final_clusters):
    for point in cluster:
        div_labels[np.where(np.all(X_scaled == point, axis=1))[0]] = i

# Evaluation metrics (assuming ground truth labels are available in
'y_true')
y_true = [0, 1, 2] # Replace this with your actual ground truth labels (if
available)

# Silhouette Score
silhouette = silhouette_score(X_scaled, div_labels)
print(f"Divisive Silhouette Score: {silhouette:.4f}")

# Adjusted Rand Index
ari = adjusted_rand_score(y_true, div_labels)
print(f"Divisive Adjusted Rand Index: {ari:.4f}")
```

Explanation:

- **Divisive Hierarchical Clustering** is performed by recursively splitting the largest cluster using **K-Means** until the desired number of clusters is reached.
 - The code then assigns cluster labels to each point and evaluates the results using **Silhouette Score** and **Adjusted Rand Index**.
-

General Notes:

- Replace the `y_true` variable with your actual ground truth labels if available. If you don't have true labels, you can omit the evaluation metrics or use clustering-specific measures like **Silhouette Score** alone.

- If your data is not already numeric or if it contains missing values, you may need to preprocess it (e.g., handle missing values or encode categorical data).
-
-

Key Evaluation Metrics:

- **Silhouette Score:** Measures how similar an object is to its own cluster compared to other clusters. Ranges from -1 to 1, where 1 indicates well-defined clusters.
- **Adjusted Rand Index (ARI):** Compares the predicted clustering with the ground truth. Values range from -1 (completely dissimilar) to 1 (perfect match).

You can modify the dataset (x) to use your actual sample data set, just replace the synthetic data generation part.

Let me know if you need any adjustments!