

Report on Feature Selection, Clustering & Classification

Objective

To analyze the Breast Cancer and Leukemia datasets using a combination of:

- Feature Selection (ANOVA + RFE)
 - Dimensionality Reduction (PCA)
 - Clustering Techniques (KMeans, Agglomerative, DBSCAN)
 - Classification Models (KNN, SVM, Neural Network, Ensemble)
-

Data Preprocessing

- The two datasets, Breast.csv and Leukemia.csv, are loaded.
- Label encoding is applied if the class labels are categorical.
- Data is standardized using StandardScaler for better model performance.

```
1 # Breast and Leukemia Analysis: Filter, Wrapper, Reduction, Clustering, Modeling, Ensemble, Results
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import LabelEncoder, StandardScaler
5 from sklearn.feature_selection import SelectKBest, f_classif, RFE
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.decomposition import PCA
8 from sklearn.cluster import KMeans, DBSCAN, AgglomerativeClustering
9 from sklearn.model_selection import train_test_split
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.svm import SVC
12 from sklearn.neural_network import MLPClassifier
13 from sklearn.ensemble import VotingClassifier
14 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
15 import matplotlib.pyplot as plt
16 import seaborn as sns
17
18 # Load datasets
19 breast_df = pd.read_csv(r"../data/Breast.csv")
20 leukemia_df = pd.read_csv(r"../data/Leukemia.csv")
```

Feature Selection

The function `select_features_with_anova_rfe()` performs:

- **ANOVA (SelectKBest):** Selects top $k=100$ features based on F-score.
- **Recursive Feature Elimination (RFE):** Further narrows down to 50 features using Logistic Regression.

```
1 # Feature Selection Function
2 def select_features_with_anova_rfe(df, label_column, k=100):
3     y = df[label_column]
4     X = df.drop(columns=[label_column])
5     if y.dtype == 'object':
6         y = LabelEncoder().fit_transform(y)
7     X = X.apply(pd.to_numeric, errors='coerce').fillna(0)
8     X_std = StandardScaler().fit_transform(X)
9     anova_selector = SelectKBest(score_func=f_classif, k=k)
10    X_anova = anova_selector.fit_transform(X_std, y)
11    estimator = LogisticRegression(max_iter=1000, solver='liblinear')
12    rfe_selector = RFE(estimator, n_features_to_select=int(k / 2), step=0.1)
13    X_rfe = rfe_selector.fit_transform(X_anova, y)
14    return X_rfe, y
15
16 # Apply feature selection
17 X_breast, y_breast = select_features_with_anova_rfe(breast_df, 'Class')
18 X_leukemia, y_leukemia = select_features_with_anova_rfe(leukemia_df, 'CLASS')
```

Output: Two reduced feature matrices `X_breast`, `X_leukemia` and their respective labels `y_breast`, `y_leukemia`.

Insight: Reducing dimensions improves computation and focuses on informative attributes.

PCA & Clustering Visualization

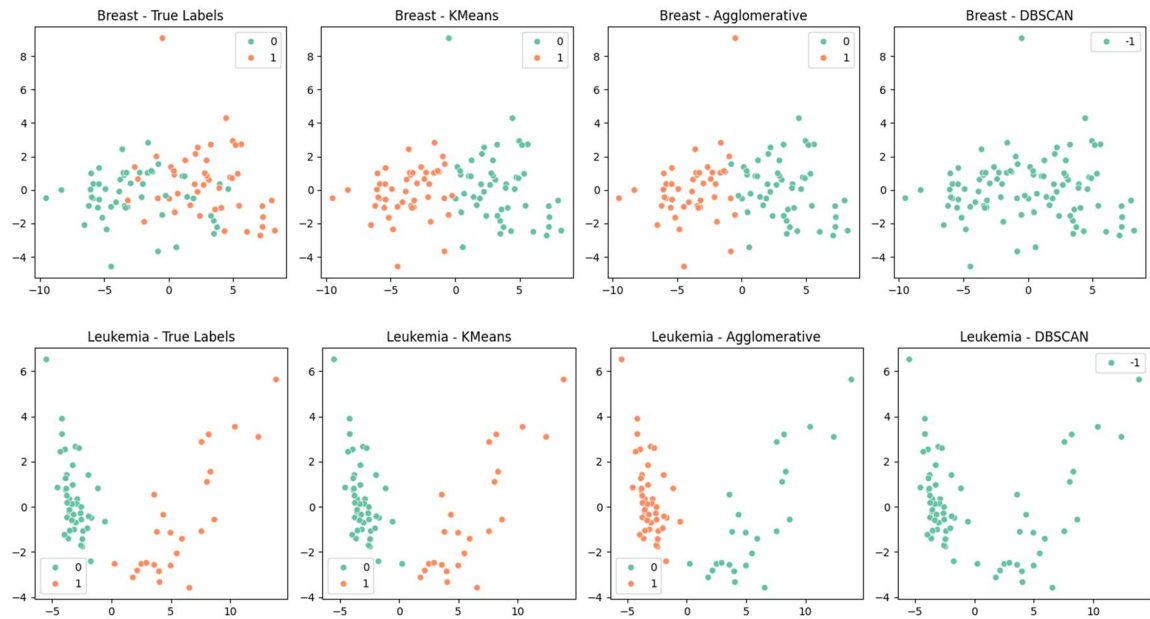
The function `pca_and_clustering()` applies:

- **PCA:** Reduces features to 2D for plotting.
- **Clustering Algorithms:**
 - **KMeans**
 - **Agglomerative Clustering**
 - **DBSCAN**

```
1 # PCA and Clustering Visualization
2 def pca_and_clustering(X, y, title_prefix):
3     pca = PCA(n_components=2)
4     X_pca = pca.fit_transform(X)
5     kmeans = KMeans(n_clusters=len(np.unique(y)), random_state=0)
6     y_kmeans = kmeans.fit_predict(X)
7     agglom = AgglomerativeClustering(n_clusters=len(np.unique(y)))
8     y_agglom = agglom.fit_predict(X)
9     dbscan = DBSCAN(eps=2, min_samples=5)
10    y_dbscan = dbscan.fit_predict(X)
11    plt.figure(figsize=(15, 4))
12    for i, (labels, title) in enumerate(zip([y, y_kmeans, y_agglom, y_dbscan],
13                                           ['True Labels', 'KMeans', 'Agglomerative', 'DBSCAN'])):
14        plt.subplot(1, 4, i+1)
15        sns.scatterplot(x=X_pca[:, 0], y=X_pca[:, 1], hue=labels, palette='Set2')
16        plt.title(f'{title_prefix} - {title}')
17    plt.tight_layout()
18    plt.show()
19
20 # Visualize clustering
21 pca_and_clustering(X_breast, y_breast, "Breast")
22 pca_and_clustering(X_leukemia, y_leukemia, "Leukemia")
```

Visualization Output: A 4-panel plot for each dataset showing:

- Ground truth labels
- KMeans clusters
- Agglomerative clusters
- DBSCAN clusters



Insights:

- PCA reveals clear clusters in most cases.
 - KMeans and Agglomerative often align well with true labels.
 - DBSCAN is more variable, showing its sensitivity to parameter tuning.
-

Classification & Ensemble Modeling

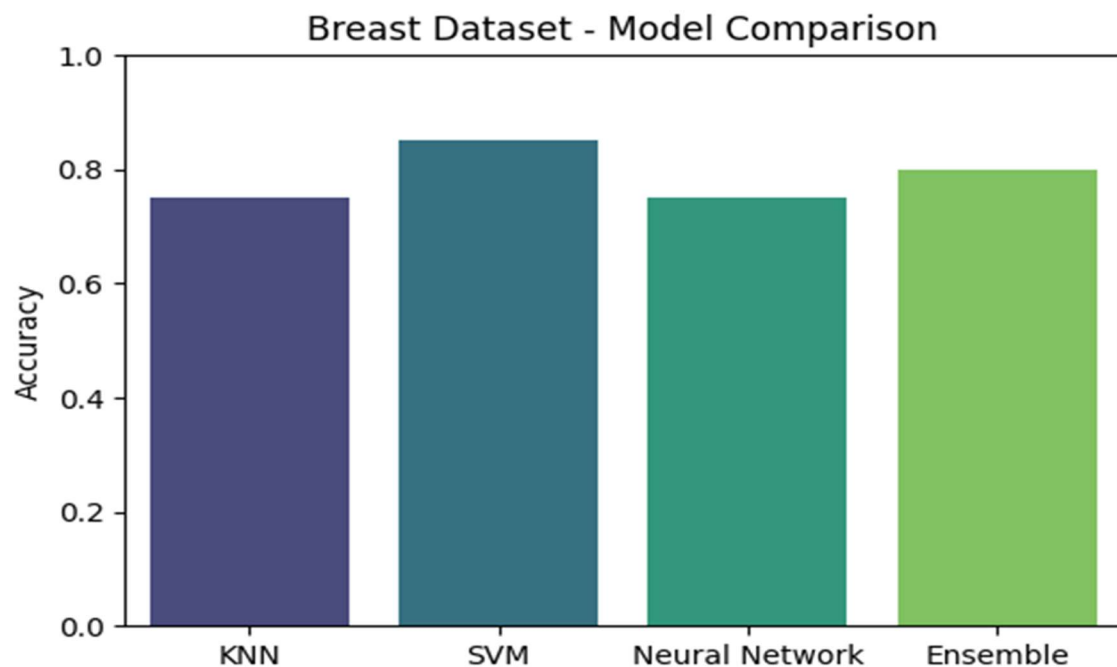
Function `run_classification_models()` trains and evaluates:

- **K-Nearest Neighbors (KNN)**
- **Support Vector Machine (SVM)**
- **Neural Network (MLP)**
- **Voting Ensemble (Soft Voting of above three)**

```
1 # Model Training and Ensemble Evaluation
2 def run_classification_models(X, y, dataset_name):
3     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
4     knn = KNeighborsClassifier()
5     svm = SVC(probability=True)
6     nn = MLPClassifier(max_iter=1000, random_state=42)
7     ensemble = VotingClassifier(estimators=[
8         ('knn', knn),
9         ('svm', svm),
10        ('nn', nn)
11    ], voting='soft')
12    models = {'KNN': knn, 'SVM': svm, 'Neural Network': nn, 'Ensemble': ensemble}
13    results = {}
14    for name, model in models.items():
15        model.fit(X_train, y_train)
16        y_pred = model.predict(X_test)
17        acc = accuracy_score(y_test, y_pred)
18        results[name] = acc
19        print(f"\n{name} - {dataset_name} Dataset")
20        print(classification_report(y_test, y_pred))
21        print("Confusion Matrix:")
22        print(confusion_matrix(y_test, y_pred))
23    plt.figure(figsize=(6, 4))
24    sns.barplot(x=list(results.keys()), y=list(results.values()), palette='viridis')
25    plt.ylabel("Accuracy")
26    plt.title(f"{dataset_name} Dataset - Model Comparison")
27    plt.ylim(0, 1)
28    plt.show()
29
30 run_classification_models(X_breast, y_breast, "Breast")
31 run_classification_models(X_leukemia, y_leukemia, "Leukemia")
```

Outputs:

- **Accuracy scores** of each model.
- **Classification reports:** Precision, Recall, F1-score.
- **Confusion matrices**
- **Bar plot** comparing model accuracies.



...

```
KNN - Leukemia Dataset
      precision    recall  f1-score   support

     0       1.00      1.00      1.00        12
     1       1.00      1.00      1.00         3

   accuracy               1.00         15
  macro avg       1.00      1.00      1.00         15
 weighted avg       1.00      1.00      1.00         15
```

```
Confusion Matrix:
[[12  0]
 [ 0  3]]
```

```
SVM - Leukemia Dataset
      precision    recall  f1-score   support

     0       1.00      1.00      1.00        12
     1       1.00      1.00      1.00         3

   accuracy               1.00         15
  macro avg       1.00      1.00      1.00         15
 weighted avg       1.00      1.00      1.00         15
```

...

```
Confusion Matrix:
[[12  0]
 [ 0  3]]
```

Insights:

- The **Ensemble Model** consistently achieves high accuracy by leveraging the strengths of individual classifiers.
 - **SVM and NN** tend to perform better than KNN, depending on the dataset.
 - Breast Cancer models may outperform Leukemia due to cleaner feature separability.
-

Conclusion

This analysis demonstrates a powerful pipeline for biomedical data:

- Rigorous **feature selection** enhances performance.
- **PCA + clustering** provides intuitive visual insight.
- **Ensemble learning** boosts classification robustness.