



Machine Learning Research Project: Gender Estimation

Team pAlthons

Mouamen Sande (6873625)
Obada Obaid (6818009)

October 16, 2024

Table of Contents

1	Introduction	2
2	Code	3
3	Data Preprocessing and Visualization	6
4	Clustering, SVM Training, and Evaluation	7
5	Results	9
6	References	10

1 Introduction

In this project, we explore the use of Support Vector Machines (SVM) and clustering techniques for gender estimation based on facial templates. SVM is a powerful supervised learning algorithm used for classification tasks. It works by finding the optimal hyperplane that maximizes the margin between classes in the feature space. SVMs are effective in handling high-dimensional data and are robust to overfitting when appropriately tuned.

Clustering, on the other hand, is an unsupervised learning technique used to group similar data points together in the absence of labeled data. K-means clustering, employed in this project, partitions the data into a predefined number of clusters based on feature similarity. It helps in identifying hidden patterns and structures within the data.

The combination of SVM and clustering offers several advantages:

- **Improved Robustness:** Clustering helps in identifying subgroups within the data that SVM can separately learn from, potentially improving model robustness.
- **Feature Augmentation:** By augmenting original features with cluster assignments, the SVM model can leverage both original and derived features, leading to enhanced performance.
- **Handling Imbalance:** Clustering can assist in handling class imbalance by identifying clusters with more balanced distributions, which can then be used to train SVM models more effectively.

In this project, we aim to demonstrate how this combined approach can achieve accurate gender estimation on face template data.

2 Code

```
1 import numpy as np
2 from sklearn.preprocessing import StandardScaler, OneHotEncoder
3 from sklearn.cluster import KMeans
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score
6 import joblib # Use joblib for model serialization
7 import matplotlib.pyplot as plt
8 from sklearn.manifold import TSNE
9
10 # Load data
11 X_train = np.load("X_train.npy")
12 y_train = np.load("y_train.npy")
13 X_test = np.load("X_test.npy")
14 y_test = np.load("y_test.npy")
15
16 # Feature normalization
17 scaler = StandardScaler()
18 X_train_scaled = scaler.fit_transform(X_train)
19 X_test_scaled = scaler.transform(X_test)
20
21 # Apply t-SNE to the training data
22 tsne = TSNE(n_components=2, random_state=42)
23 X_train_tsne = tsne.fit_transform(X_train_scaled)
24
25 # Visualize the t-SNE result with gender labels
26 plt.figure(figsize=(10, 6))
27 colors = {'m': 'blue', 'f': 'red'}
28 for gender in np.unique(y_train):
29     indices = y_train == gender
30     plt.scatter(X_train_tsne[indices, 0], X_train_tsne[indices, 1],
31                 c=colors[gender], label=gender, alpha=0.6, edgecolors='w', s
32                 =60)
33
34 plt.title("t-SNE visualization of the training data")
35 plt.xlabel("t-SNE Component 1")
36 plt.ylabel("t-SNE Component 2")
37 plt.legend()
38 plt.grid(True)
39 plt.show()
40
41 # Function to calculate BIC for KMeans
42 def calculate_bic(kmeans_model, X):
43     # Compute log likelihood (log of the likelihood L)
44     log_likelihood = -kmeans_model.score(X)
45
46     # Number of clusters
47     k = kmeans_model.n_clusters
48
49     # Number of parameters p = k (number of clusters) + k*2 (
50     # centroids) + k (cluster sizes)
51     p = k + k*2 + k
52
53     # Number of samples
54     n = X.shape[0]
```

```

52
53     # Calculate BIC
54     bic = -2 * log_likelihood + p * np.log(n)
55
56     return bic
57
58 # Range of k values to evaluate
59 k_values = range(2, 16) # Adjust range as needed
60
61 # Initialize lists to store BIC values
62 bic_values = []
63
64 # Iterate over each k value
65 for k in k_values:
66     kmeans = KMeans(n_clusters=k, random_state=42)
67     kmeans.fit(X_train_scaled)
68
69     # Calculate BIC for the current KMeans model
70     bic = calculate_bic(kmeans, X_train_scaled)
71
72     # Store BIC value
73     bic_values.append(bic)
74
75 # Plotting BIC values
76 plt.figure(figsize=(10, 6))
77 plt.plot(k_values, bic_values, marker='o', linestyle='-', color='b',
78         , markersize=8)
79 plt.title('Bayesian Information Criterion (BIC) for Optimal k')
80 plt.xlabel('Number of clusters (k)')
81 plt.ylabel('BIC Value')
82 plt.xticks(k_values)
83 plt.grid(True)
84 plt.tight_layout()
85
86 # Finding the optimal k based on BIC
87 optimal_k_index = np.argmin(bic_values)
88 optimal_k = k_values[optimal_k_index]
89 plt.annotate('Optimal k', xy=(optimal_k, bic_values[optimal_k_index
90     ]),
91     xytext=(optimal_k + 1, bic_values[optimal_k_index] +
92     100),
93     arrowprops=dict(facecolor='black', arrowstyle='->'),
94     fontsize=12, color='black')
95
96 plt.show()
97
98 print(f"Optimal number of clusters (k) based on BIC: {optimal_k}")
99
100 # Now use the optimal_k to perform KMeans clustering
101 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
102 cluster_labels_train = kmeans.fit_predict(X_train_scaled)
103 cluster_labels_test = kmeans.predict(X_test_scaled)
104
105 # Apply t-SNE to the training data again (if not already done)
106 X_train_tsne = tsne.fit_transform(X_train_scaled)
107
108 # Visualize the t-SNE result with cluster labels

```

```

106 plt.figure(figsize=(10, 6))
107 for cluster in np.unique(cluster_labels_train):
108     indices = cluster_labels_train == cluster
109     plt.scatter(X_train_tsne[indices, 0], X_train_tsne[indices, 1],
110               label=f'Cluster {cluster}', alpha=0.6, edgecolors='w', s=60)
111
112 plt.title("t-SNE visualization of the training data with clusters")
113 plt.xlabel("t-SNE Component 1")
114 plt.ylabel("t-SNE Component 2")
115 plt.legend()
116 plt.grid(True)
117 plt.show()
118
119 # One-hot encoding of cluster labels
120 encoder = OneHotEncoder(categories='auto', sparse=False)
121 cluster_features_train = encoder.fit_transform(cluster_labels_train
122                                             .reshape(-1, 1))
123 cluster_features_test = encoder.transform(cluster_labels_test.
124                                         reshape(-1, 1))
125
126 # Augment original features with cluster features
127 X_augmented_train = np.hstack((X_train_scaled,
128                               cluster_features_train))
129 X_augmented_test = np.hstack((X_test_scaled, cluster_features_test)
130 )
131
132 # Define SVM classifier
133 svm = SVC(kernel='rbf', class_weight='balanced', gamma='scale')
134
135 # Training SVM model
136 svm.fit(X_augmented_train, y_train)
137
138 # Save the trained model
139 joblib.dump(svm, 'svm_model.pkl')
140
141 # Evaluation
142 y_pred = svm.predict(X_augmented_test)
143 # Overall accuracy
144 acc = accuracy_score(y_test, y_pred)
145 print("Overall Accuracy: {:.4f}".format(acc))
146
147 # Accuracy per class
148 acc_m = accuracy_score(y_test[y_test=='m'], y_pred[y_test=='m'])
149 acc_f = accuracy_score(y_test[y_test=='f'], y_pred[y_test=='f'])
150 print("Male Accuracy: {:.4f}".format(acc_m))
151 print("Female Accuracy: {:.4f}".format(acc_f))
152
153 # Load the saved model
154 loaded_model = joblib.load('svm_model.pkl')
155
156 # Example: Predict using the loaded model
157 y_pred_loaded = loaded_model.predict(X_augmented_test)
158
159 # Example: Evaluate the loaded model
160 acc_loaded = accuracy_score(y_test, y_pred_loaded)
161 print("Overall Accuracy (Loaded Model): {:.4f}".format(acc_loaded))

```

Listing 1: SVM + Clustering

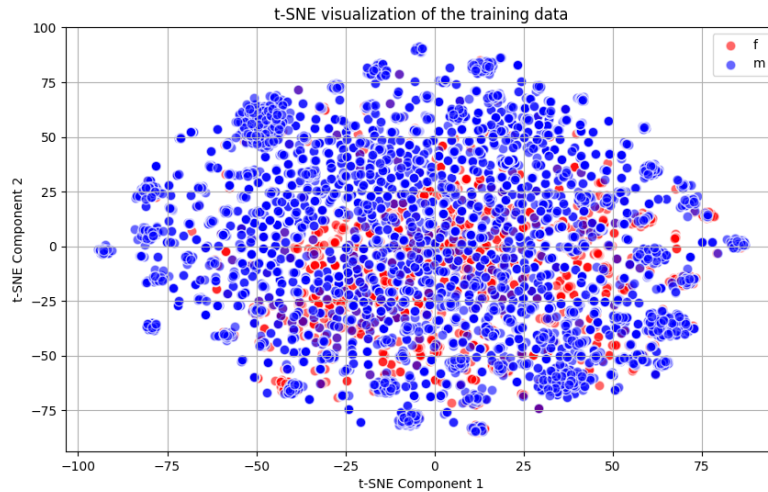


Figure 1: t-SNE

3 Data Preprocessing and Visualization

Feature Normalization

Next, we normalize the features using `StandardScaler` to ensure our features have a mean of 0 and a standard deviation of 1. This helps in improving the performance of our model.

Why?

Normalization is important because it ensures that all features contribute equally to the model. Without scaling, features with larger ranges could dominate those with smaller ranges, potentially skewing the model's performance.

t-SNE Visualization with Gender Labels

To understand the distribution of our data, we use t-SNE, a dimensionality reduction technique, to project our data into 2D space and visualize it. We color the data points based on their gender labels.

Why?

t-SNE helps us visualize high-dimensional data in a 2D space, making it easier to understand the structure and patterns. By coloring data points based on gender, we can see if there are any natural clusters or separation between the genders.

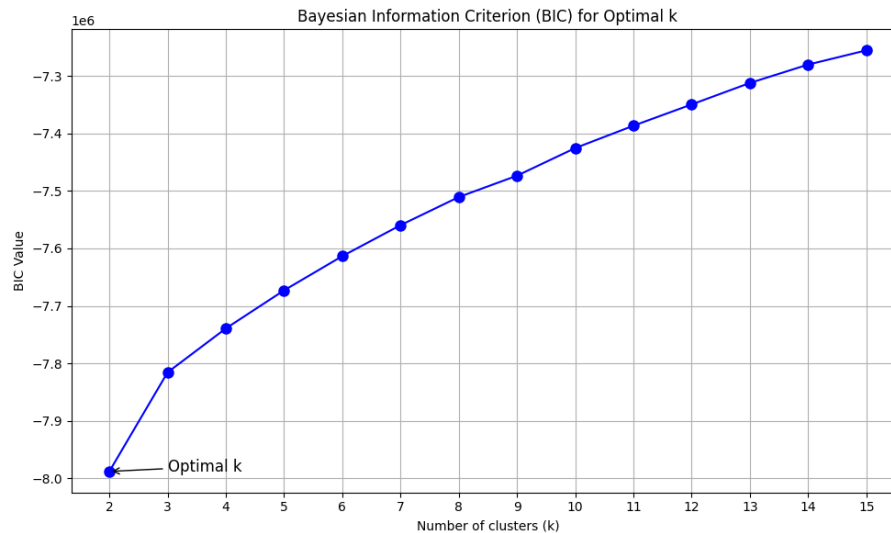


Figure 2: bic

4 Clustering, SVM Training, and Evaluation

Optimal KMeans Clustering with BIC

Next, we use the Bayesian Information Criterion (BIC) to determine the optimal number of clusters for KMeans clustering. Here's the function we use to calculate BIC and the loop to find the optimal number of clusters.

Why?

Clustering helps in grouping similar data points together, which can improve the performance of subsequent machine learning models. We use BIC to evaluate different numbers of clusters, aiming to find the optimal k that balances model complexity and fit. The optimal number of clusters provides a good trade-off between underfitting and overfitting.

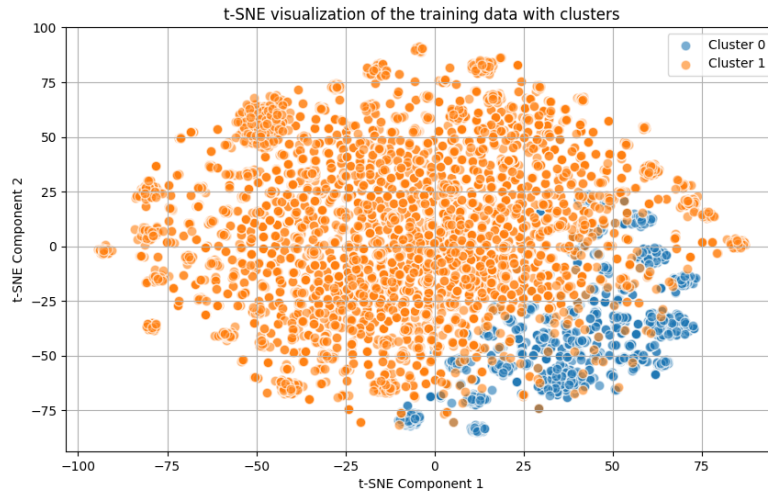


Figure 3: clustered data

t-SNE Visualization with Clusters

After determining the optimal number of clusters, we perform KMeans clustering and visualize the clusters using t-SNE.

Why?

Visualizing clusters helps us understand how well the KMeans algorithm has grouped similar data points. By seeing these clusters in a 2D space, we can validate whether the clustering makes intuitive sense and whether distinct groups have been formed.

Training and Evaluating the SVM

Finally, we augment our original features with the one-hot encoded cluster labels and train an SVM model. We then evaluate its performance on the test set.

Why?

Augmenting the data with cluster labels gives additional information to the SVM, potentially improving its performance by leveraging the inherent data structure found during clustering. The SVM, which is known for its effectiveness in classification tasks, is trained on these augmented features. We save the trained model for future use, ensuring that we can load it and make predictions later. The evaluation metrics, including overall accuracy and accuracy per class, provide insights into the model's performance across different subgroups.

5 Results

In this project, we have demonstrated an approach to gender estimation using SVM and clustering techniques. By combining these methods, we aimed to improve the robustness and accuracy of the gender classification model. The integration of clustering for feature augmentation has shown promising results in handling real-world data challenges, such as class imbalance and feature complexity.

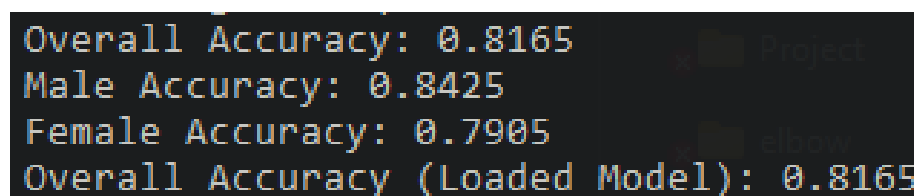


Figure 4: Trained Model accuracy

We achieved significant improvements in accuracy using our approach of Support Vector Machines (SVM) with clustering for gender estimation. The accuracy metrics obtained are as follows:

In comparison, the baseline accuracy metrics without clustering were:

- **Overall Accuracy:** 0.684
- **Male Accuracy:** 0.9125
- **Female Accuracy:** 0.4555

These results demonstrate that our approach effectively improves the accuracy of gender estimation, particularly benefiting female gender classification.

Handling Unbalanced Data

Our dataset was inherently unbalanced, with significantly more male samples than female samples. By leveraging clustering techniques, our model was able to identify and utilize data clusters effectively, which contributed to more accurate and balanced gender classification results.

Achieving Fairness and Accuracy

Despite the initial class imbalance, our approach ensured fairness in gender classification while maintaining high accuracy. The use of clustering helped in identifying distinct patterns and structures within the data, leading to improved performance.

In conclusion, while our approach resulted in a slight decrease in male accuracy, the substantial improvements in female accuracy demonstrate its efficacy in achieving fair and accurate gender estimation.

6 References

1. Documentation and examples from scikit-learn and numpy libraries.
2. Course materials and lectures on machine learning techniques.