

REPORT FOR THE STUDY**Classifying Bacteria Species with Dimension Reduction****Wonjun Lee¹ | Junseong Bang^{2,3}**¹Intelligence & Manufacturing Convergence Laboratory, ETRI, Daejun, Republic of Korea²Police Science & Public Safety ICT Research Center, ETRI, Daejun, Republic of Korea³Department of Computer Software, UST, Daejun, Republic of Korea**Abstract**

Bacteria are microscopic, single-celled organisms. But There are only a few types of bacteria that cause the disease and There are only a few types of bacteria that cause the disease. The goal is to predict bacteria species based on repeated lossy measurements of DNA snippets. Using this dataset, after apply dimension reduction method(PCA, LDA), by Random Forest Classifier. It compares and analyzes to apply Naive, PCA, and LDA.

KEYWORDS:

Bacterias, PDA, LDA

1 | INTRODUCTION

Bacteria are microscopic, single-celled organisms. There are thousands of different kinds of bacteria. Many bacteria are present in various body parts, including human and animal skin, airway, mouth, digestive tract, reproductive organs, and urethra, without causing harm. Such bacteria are called bacteriophages or microbial communities.

There are at least as many bacteria in the bacteriophage as cells in the body. A lot of bivalve fungi actually help people.

There are only a few types of bacteria that cause the disease. These bacteria are called pathogens. Sometimes bacteriophage fungi cause disease under certain conditions.

Bacteria can produce harmful substances (toxins) or invade tissues and cause diseases. Some bacteria can trigger inflammation that can affect the heart, nervous system, kidneys, or gastrointestinal tract.

Some bacteria increase the risk of cancer. There are only a few types of bacteria that cause the disease. These include anthrax, botulinum poisoning, plague, and bacteria that cause Yato disease.

It is important to classify bacteria that cause diseases in advance and bacteria that can be used as weapons.

For Tabular Playground Series - Feb 2022, the problem deals with classifying 10 different bacteria species using data from a genomic analysis technique that has some data compression and data loss. In this technique, 10-mer snippets of

DNA are sampled and analyzed to give the histogram of base count.

The goal is to predict bacteria species based on repeated lossy measurements of DNA snippets.

2 | DATASET

This is Tabular Playground Series - Feb 2022 dataset, obtained from Kaggle.

the problem deals with classifying 10 different bacteria species using data from a genomic analysis technique that has some data compression and data loss.

Each row of data contains a spectrum of histograms generated by repeated measurements of a sample, each row containing the output of all 286 histogram possibilities, which then has a bias spectrum (of totally random ATGC) subtracted from the results.

The data (both train and test) also contains simulated measurement errors (of varying rates) for many of the samples, which makes the problem more challenging.

3 | LINEAR DISCRIMINANT ANALYSIS (LDA)

Linear Discriminant Analysis is a dimensionality-reduction technique in the pre-processing step for pattern-classification and machine learning applications. The goal is to project

a dataset onto a lower-dimensional space with good class-separability in order to avoid overfitting and also reduce computational costs.

4 | PRINCIPAL COMPONENT ANALYSIS VS. LINEAR DISCRIMINANT ANALYSIS

Both Linear Discriminant Analysis (LDA) and Principal Component Analysis (PCA) are linear transformation techniques that are commonly used for dimensionality reduction. PCA can be described as an “unsupervised” algorithm, since it “ignores” class labels and its goal is to find the directions (the so-called principal components) that maximize the variance in a dataset. In contrast to PCA, LDA is “supervised” and computes the directions (“linear discriminants”) that will represent the axes that maximize the separation between multiple classes.

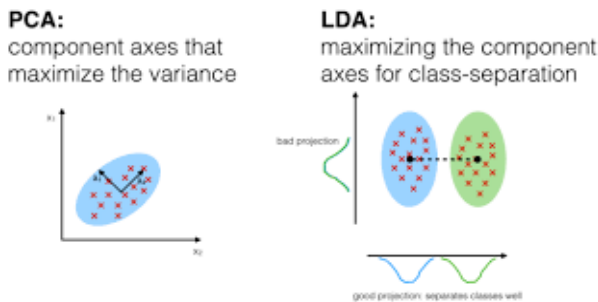


FIGURE 1 Comparison of PCA and LDA

5 | CODE

Import Library

```
1 #General Libraries
2 import numpy as np
3 import pandas as pd
4
5 # Machine Learning Libraries
6 from sklearn.preprocessing import scale
7 from sklearn.model_selection import
8     train_test_split
9 from sklearn.metrics import confusion_matrix,
10     accuracy_score, confusion_matrix,
11     plot_confusion_matrix
12 from sklearn.metrics import roc_auc_score,
13     roc_curve
14 import sklearn.metrics as metrics
15 from sklearn.preprocessing import StandardScaler
16
17 # Data Visualization Libraries
```

```
14 import seaborn as sns
15 import matplotlib.pyplot as plt
16 %matplotlib inline
17
18 from warnings import filterwarnings
```

Load Dataset, Split Data

```
1 test = pd.read_csv("./test.csv", index_col="
2     row_id")
3 train = pd.read_csv("./train.csv", index_col="
4     row_id")
5
6 TARGET = 'target'
7 FEATURES = [col for col in train.columns if col
8     not in [TARGET]]
9
10 print(f'Number of observations in TRAIN:{len(
11     train)}')
12 print(f'Number of observations in TEST:{len(test
13     )}')
14 train.target.value_counts()
15
16 print(train[FEATURES].duplicated().sum())
17 print(test[FEATURES].duplicated().sum()) #
18
19 train.drop_duplicates(keep='first', inplace=True
20     ) #
21
22 X = train.drop("target", axis=1).astype(np.
23     float32)
24
25 from sklearn.preprocessing import LabelEncoder
26
27 target_encoder = LabelEncoder()
28 y = pd.Series(target_encoder.fit_transform(train
29     ["target"]))
30
31 X_train, X_test, y_train, y_test =
32     train_test_split(train[FEATURES].values, y,
33         random_state=42)
```

Standard Scale

```
1 from sklearn.preprocessing import StandardScaler
2
3 std_scale = StandardScaler()
4 std_scale.fit(X_train)
5 X_train_std = std_scale.transform(X_train)
6 X_test_std = std_scale.transform(X_test)
```

LDA

```
1 from sklearn.discriminant_analysis import
2     LinearDiscriminantAnalysis
3 lda = LinearDiscriminantAnalysis()
4 lda.fit(X_train_std, y_train)
5 X_train_lda = lda.transform(X_train_std)
6 X_test_lda = lda.transform(X_test_std)
```

LDA Visualization

```

1 from sklearn.discriminant_analysis import
  LinearDiscriminantAnalysis as LDA
2
3 train_sub = train.sample(10000, random_state=
  42)
4 lda_data = LDA(n_components=2).fit_transform(
  train_sub.drop(columns='target'), train_sub.
  target)
5 plt.figure(figsize=(10,10))
6 sns.scatterplot(x = lda_data[:, 0], y = lda_data
 [:, 1], hue = 'target', data=train_sub)

```

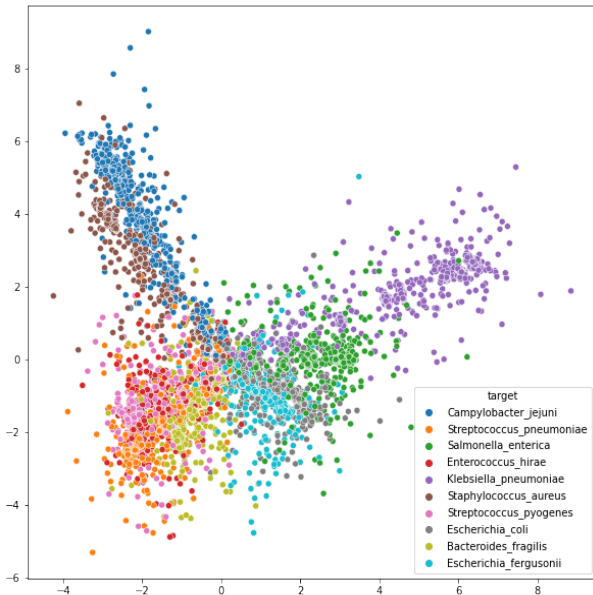


FIGURE 2 LDA Visualization

Random Forest Classifier

```

1 from sklearn.ensemble import
  RandomForestClassifier
2 # learning..
3 clf_rf_lda = RandomForestClassifier(max_depth=2,
  random_state=42)
4 clf_rf_lda.fit(X_train_lda, y_train)
5
6 # prediction
7 pred_rf_lda = clf_rf_lda.predict(X_test_lda)
8
9 from sklearn.metrics import accuracy_score
10 accuracy_lda = accuracy_score(y_test,
  pred_rf_lda)
11 print('# Prediction with LDA: ')
12 print(accuracy_lda)

```

6 | RESULT

This is the result of applying Random Forest Classifier to data applied with Naive which is version that have not been processed, PCA, and LDA.

	Naive	PCA	LDA
Testing Score	0.74	0.48	0.82

TABLE 1 Result of Naive, PCA, LDA

7 | CONCLUSION

In order to analyze classifying 10 different bacteria specie, the dimension was reduced to PCA and LDA, and then a model was created by applying a Random Forest Classifier .

this dataset, Random Forest Classifier showed the highest performance after applying LDA, and PCA showed the lowest performance.

Future research plans aim to increase accuracy and performance by using more diverse dimensional reduction techniques.