
Report of Final Project
Linear Discriminative Analysis using Power Iteration

Heinrich Dinkel

ID: 1140339107

E-mail: heinrich.dinkel@sjtu.edu.cn

1 Introduction

My topic I did my final Project about is related to my field of research, which is Speaker Recognition. In this field of research we try to find suitable model, which can predict if a given speech segment (called utterance) was spoken by a specific speaker. Generally speaking, we have usually a given set of $\mathbf{S} = \{S_1, S_2, \dots, S_n\}$ speakers of which each produce an arbitrary amount (q) speech segments, which we have as our basic data. We try to find the most probable hypothesis that a newly, unheard-of of speech segment, belongs to any of the given speaker set \mathbf{S} . To accomplish this task we assume having two different models, which try to identify a newly heard utterance u :

$$\begin{aligned} H_0 &= u \text{ is spoken by } S \\ H_1 &= u \text{ is spoken not by } S \end{aligned}$$

H_0, H_1 represent here the models for the correct speakers (H_0) and the wrong ones (or impostors H_1). So we need only to define a threshold t which indicates if the utterance u is either belonging to H_0 or H_1 . Then we can rewrite:

$$\frac{P(u|H_0)}{P(u|H_1)} = \begin{cases} \leq t, & \text{accept } H_0 \\ > t, & \text{accept } H_1 \end{cases}$$

The important part is to find robust and versatile models for H_0 and H_1 respectively.

In the field of speech recognition, the input consists of raw waveforms, which will be transformed into so called features. Features are basically an extracted vectorized form of the original signal. Their purpose is to reduce the dimensionality complexity for further processing, since the overall amount of raw data is simply too arduous to handle.

2 Motivation

The motivation for this Project comes from the real world application of feature extraction and speech Processing. Nowadays the amount of data for training any kinds of models is extremely large (hundreds of Gigabytes), so that one needs to extract only the meaningful information out of the data. This process is called feature extraction and results in a vectorized representation of the data. Using these feature vectors, we can train a model, which goal it is to classify a certain behaviour. This final model can then be post-processed by LDA to increase the distances between the classes. LDA is a general post-processing tool in speaker recognition. In this Report we show how LDA can also be used by using the power iteration instead of the more general QR-Decomposition.

3 Linear Discriminative Analysis

For speaker recognition we use LDA to improve our scores. It is widely used after the computation of the models H_0 and H_1 to project the resulting high-dimensional vectors into a lower dimensional space, where the separation between the involving classes is maximized.

To achieve a dimensionality reduction, we use

Assume having c classes in an arbitrary space. We calculate the proportion in between the between-class covariance and the within-class covariance, having the current feature vector \mathbf{x}_i :

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &= -\frac{\mathbf{w}^T \Sigma_{\text{between}} \mathbf{w}}{\mathbf{w}^T \Sigma_{\text{within}} \mathbf{w}} \\ \Sigma_{\text{within}} &= \sum_i^c \Sigma_i \\ \Sigma_i &= \sum_i^c (\mathbf{x} - \boldsymbol{\mu}_i) (\mathbf{x} - \boldsymbol{\mu}_i)^T \\ \Sigma_{\text{between}} &= N_i \sum_i^c (\boldsymbol{\mu}_i - \boldsymbol{\mu}) (\boldsymbol{\mu}_i - \boldsymbol{\mu})^T\end{aligned}$$

where $\boldsymbol{\mu}_i$ is the mean of the current class, $\boldsymbol{\mu}$ the overall mean and N_i is the sample mean of the i -th class.

As soon as the covariances are calculated, we need to focus on how to compute the weight vectors \mathbf{w} . We solve this generalized eigenvalue problem by calculating the eigenvalues of:

$$\sigma(\Sigma_{\text{within}}^{-1} \Sigma_{\text{between}}) = \{\lambda_1, \lambda_2, \dots, \lambda_n\}$$

With these eigenvalues we compute the corresponding eigenvectors. The clue of the whole algorithm is then to use the most important $m, m < n$ eigenvalues, to construct the overall weight matrix \mathbf{W} , which then projects the input features into the new m dimensional feature space.

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}$$

A large eigenvalue means that the corresponding eigenvector has also large values, which means that this eigenvector encompass most of the variance of the given class. This means we can simply sort the eigenvalues in descending order and pick the largest m eigenvalues out. Doing so allows us to reduce the dimensionality of the given data, by removing the less important parts.

Here we want to point out that there exist many different algorithms for computing the eigenvalues, whereas we have focused in this Report in the Power iteration.

3.1 Power Iteration

In mathematics, the power iteration is an eigenvalue algorithm: given a matrix \mathbf{A} , the algorithm will produce a number λ (the eigenvalue) and a nonzero vector \mathbf{v} (the eigenvector), such that $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$.

The power iteration algorithm starts with a vector \mathbf{b}_0 , which may be an approximation to the dominant eigenvector or a random vector. The method is described by the iteration $b_{k+1} = \frac{\mathbf{A}\mathbf{b}_k}{\|\mathbf{A}\mathbf{b}_k\|}$. So, at every iteration, the vector \mathbf{b}_k is multiplied by the matrix \mathbf{A} and normalized.

Under the assumptions:

- \mathbf{A} has an eigenvalue that is strictly greater in magnitude than its other eigenvalues
- The starting vector \mathbf{b}_0 has a nonzero component in the direction of an eigenvector associated with the dominant eigenvalue. then:
- A subsequence of (b_k) converges to an eigenvector associated with the dominant eigenvalue

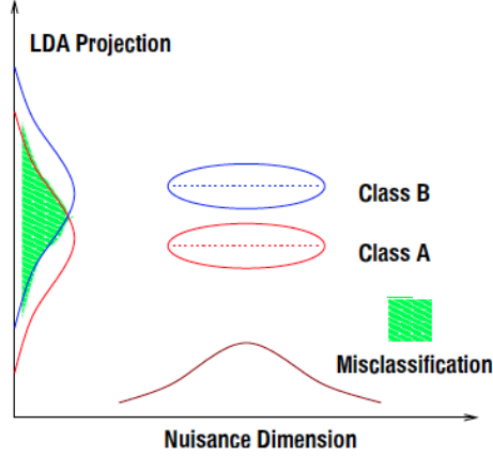


Figure 1: Using LDA to increase inner class distance in case of constant covariance

Note that the sequence (b_k) does not necessarily converge. It can be shown that: $b_k = e^{i\phi_k} v_1 + r_k$ where: v_1 is an eigenvector associated with the dominant eigenvalue, and $\|r_k\| \rightarrow 0$. The presence of the term $e^{i\phi_k}$ implies that (b_k) does not converge unless $e^{i\phi_k} = 1$. Under the two assumptions listed above, the sequence (μ_k) defined by: $\mu_k = \frac{b_k^* A b_k}{b_k^* b_k}$ converges to the dominant eigenvalue.

It needs to be said that b_k converges to a multiple of the eigenvector v_1 , at rate $\left| \frac{\lambda_2}{\lambda_1} \right|$, where λ_2 is the second dominant eigenvalue. Thus, the method converges slowly if there is an eigenvalue close in magnitude to the dominant eigenvalue.

4 Experiments

Generally speaking, when using real speech data, the dimensionality of the feature vectors is too large to efficiently display in any plot. Therefore the following experiments are artificial and only show how LDA works. We sampled n Gaussians with feature size of 2, which will be reduced to 1 dimensional features using LDA.

Moreover we have overall 1000 samples for each Gaussian, meaning that we have overall 3000 data points, characterized by a Gaussian distribution. Real data in speaker verification consists of 39 dimensional vectors. The amount of data is generally measured in hours, whereas we can expect that per 100 MB data, we achieve about 1 hour of raw wave speech.

We used Python2.7 with the help of numpy to do the calculations. Generating the artificial data, we can simply use:

```

1 import numpy as np
2
3 means = np.random.randint(-10,10,size=(dims))
4 variances = np.random.sample(size=(dims,dims))
5 #We sample 3 classes with 1000 datapoints
6 datapoints = np.random.multivariate_normal(means,variances,size=(3,1000)).T

```

In the chosen programming language, the power iteration is implemented as follows:

```

1 def powerIteration(A):
2     b = [random() for i in range(len(A))]
3     tmp = [0] * len(A)
4
5     for iteration in range(1000):
6
7         for i in range(0, len(A)):
8             tmp[i] = 0
9             for j in range(0, len(A)):
10                 tmp[i] += A[i][j] * b[j]
11
12         normSq = 0
13         for k in range(0, len(A)):
14             normSq += tmp[k] * tmp[k]
15         norm = sqrt(normSq)
16
17         for i in range(len(A)):
18             b[i] = tmp[i] / norm
19
20     return b

```

The LDA process was implemented in the following fashion, given the power iteration:

```

1 def lda(datapoints, tarDim = 2):
2     '''
3     datapoints is the given data, with shape [classsize samplesize featuredim]
4     tarDim is the amount of eigenvalues which will be kept, meaning the dimension, which will be
5     transferred to
6     '''
7     cl, samplesize, featuresize = datapoints.shape
8     mean_vectors = []
9     for i in range(cl):
10         mean_vectors.append(np.mean(datapoints[i], axis=0))
11     #     print('Mean Vector class %s: %s\samplesize' % (i, mean_vectors[i]))
12     #     Calculate Within Scatter matrix:
13     mean_vectors = np.array(mean_vectors)
14     s_w = np.zeros((featuresize, featuresize))
15
16     for i in range(cl):
17         curmat = datapoints[i]
18         for row in curmat:
19             #         Get the current row, transform to col just for dot product purpose
20             col = row.reshape(featuresize, 1)
21             meancol = mean_vectors[i].reshape(featuresize, 1)
22             s_w += (col - meancol).dot((col - meancol).T)
23     #     print('Within class covaricance Scatter : ', s_w[i])
24     overall_mean = np.mean(mean_vectors, axis=0)
25     s_b = np.zeros((featuresize, featuresize))
26     for i in range(cl):
27         samp_size = len(datapoints[i])
28         mean_row = mean_vectors[i].reshape(featuresize, 1)
29         s_b += samp_size * (mean_row - overall_mean).dot((mean_row - overall_mean).T)
30

```

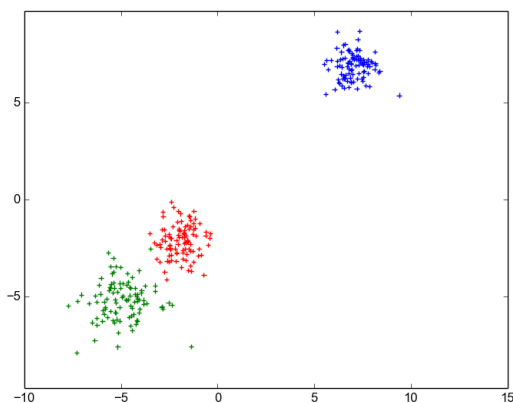


Figure 2: 3 Class Gaussian Randomly drawn out of 1000 samples

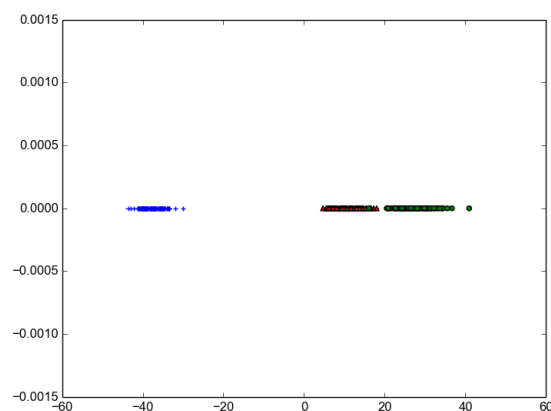


Figure 3: Projection of 3 Class Gaussian into 1 dimensional space

```

31 eig_vals,eig_vecs = powerIteration((s_w).dot(s_b))
32 eig_pairs = [(np.abs(eig_vals[i]), eig_vecs[:,i]) for i in range(len(eig_vals))]
33 eig_pairs = sorted(eig_pairs,key=lambda k: k[0] ,reverse=True)
34 w= np.hstack(eig_pairs[i][1].reshape(featuresize,1) for i in range(tarDim))
35 datapoints = datapoints.reshape(cl*samplesize,featuresize)
36 x_lda = w.T.dot(datapoints.T)
37 x_lda = x_lda.reshape(cl,samplesize)
38
39 log.debug("XLDA Dim {0}".format(x_lda.shape))
40 return x_lda

```

5 Results

As we can see, the newly projected space maximizes the distance in between the classes effectively. In this reduced space, a classification can be easier performed, leading to an overall better classification performance. Moreover the computational cost is minor, so that LDA should be used whenever possible.

A subset of the generated used numbers can be seen here:

```
#Source data:
[[ 7.90729539  8.17202412  7.47519462]
 [ 9.55144994  8.37470264  7.88132993]
 [ 8.80252236  7.41339669  6.96571076]
 [ 9.27757307  6.47732587  9.34183583]]
# LDA Projection:
[ 12.81458458 13.66795127 12.18943162 13.03469284]
```

In the figures shown, we can see how LDA projects the initial 2 dimensional space into a single dimension. The power iteration helped us to find this single dimension quickly, since it converged within the first 30 iterations.