# Text Mining Final-Project by Thomas Wagner, Alexander Allen, MingYi Wang

## Table of Contents

May 14 2015

# Input Data

```
clear all
close all

% load 10% sample data and the test data
load DDISample.mat

% convert to normal format instead of sparse
Classp_train=full(Classp_train);
Classm_train=full(Classm_train);
Classp_test=full(Classp_test);
Classm_test=full(Classm_test);


% Set random number to an initial seed
[r,c]=size(Classm_train);
s=RandStream('mt19937ar','Seed',550);
%generate a permutation of the data
p=randperm(s,r);
Classm_train=Classm_train(p,:);




Classm_train = Classm_train(1:380,:);
```

# PCA on data

```
Train_total = [Classp_train; Classm_train];

[mp,np] = size(Classp_train);    % size for Classp
[mm,nm] = size(Classm_train);    % size for Classm
[m,n] = size(Train_total);       % size for total
```

```matlab
train_mean = (1/m)*(ones(1,m)*Train_total);

Train_total2 = Train_total - ones(m,1)*train_mean;

[eigenvectors, scores, eigenvalues] = pca(Train_total);

trimmed_scores = scores(:,1:300);
classp_scores = trimmed_scores(1:mp,:);
classm_scores = trimmed_scores(mp+1:m,:);
```

# Fisher

```matlab
meanp=mean(classp_scores);
meanm=mean(classm_scores);

psize=size(classp_scores,1)
nsize=size(classm_scores,1)
Bp=classp_scores-ones(psize,1)*meanp;
Bn=classm_scores-ones(nsize,1)*meanm;

Sw=Bp'*Bp+Bn'*Bn;
wfisher = Sw\(meanp-meanm)';
wfisher=wfisher/norm(wfisher);

tfisher=(meanp+meanm)./2*wfisher
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

% Analyze training data  results of the Fisher Linear Discriminant

FisherPosErrorTrain = sum(classp_scores*wfisher <= tfisher);
FisherNegErrorTrain = sum(classm_scores*wfisher >= tfisher);

FisherTrainError= ((FisherPosErrorTrain + FisherNegErrorTrain)/(size(trimmed_score

% Histogram of Fisher Training Results
HistClass(classp_scores,classm_scores,wfisher,tfisher,...
    'Fisher Method Training Results',FisherTrainError);
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%


        psize =

            380


        nsize =

            380


        tfisher =

            -3.8519e-17
```
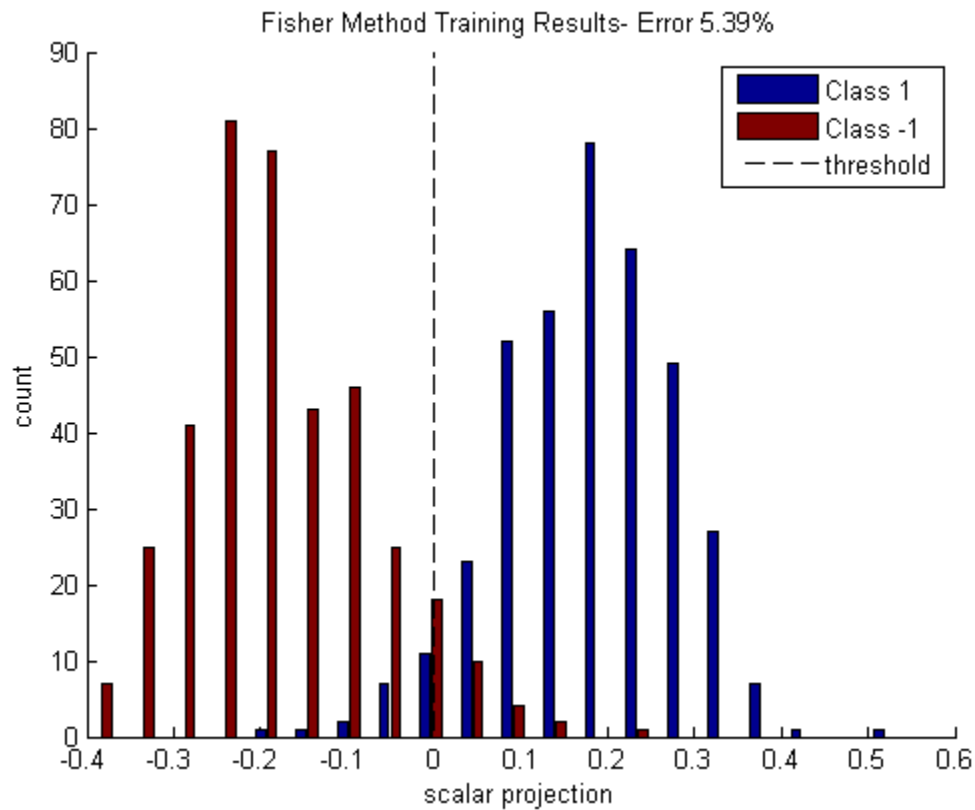
*FisherTrainError =*

*0.0539*



## Top 30 words

```
A = Train_total' * trimmed_scores * wfisher;
absA = abs(A);

words = cell(30,1);
word_values = zeros(30,1);

for i = 1:30
    [M,I] = max(absA);
    words(i,1) = featurenames(I);
    word_values(i,1) = A(I,1);
    absA(I,1) = 0;
end
```

## Compute Test Scores

```
Test_total = [Classp_test; Classm_test];
```

```
[mp_test,np_test] = size(Classp_test);     % size for Classp
[mm_test,nm_test] = size(Classm_test);      % size for Classm
[m_test,n_test] = size(Test_total);         % size for total

Test_total2 = Test_total - ones(m_test,1)*train_mean;
Classp_test2 = Test_total2(1:mp_test,:);
Classm_test2 = Test_total2((mp_test+1):end,:);

Classm_test_scores = Classm_test2 * eigenvectors;
Classp_test_scores = Classp_test2 * eigenvectors;

scores_test_total = [Classp_test_scores; Classm_test_scores];


trimmed_scores_test = scores_test_total(:,1:300);
classp_test_scores = trimmed_scores_test(1:mp_test,:);
classm_test_scores = trimmed_scores_test(mp_test+1:m_test,:);
```

# Fisher on Test

```
FisherPosErrorTest = sum(classp_test_scores*wfisher <= tfisher);
FisherNegErrorTest = sum(classm_test_scores*wfisher >= tfisher);

FisherTestError= ((FisherPosErrorTest + FisherNegErrorTest)/(size(trimmed_scores_t

% Histogram of Fisher Testing Results
HistClass(classp_test_scores,classm_test_scores,wfisher,tfisher,...
    'Fisher Method Testing Results',FisherTestError);


%RESULTS using 380 sentences from each class 5.39% training, 23.87% testing


      FisherTestError =

          0.2387
```
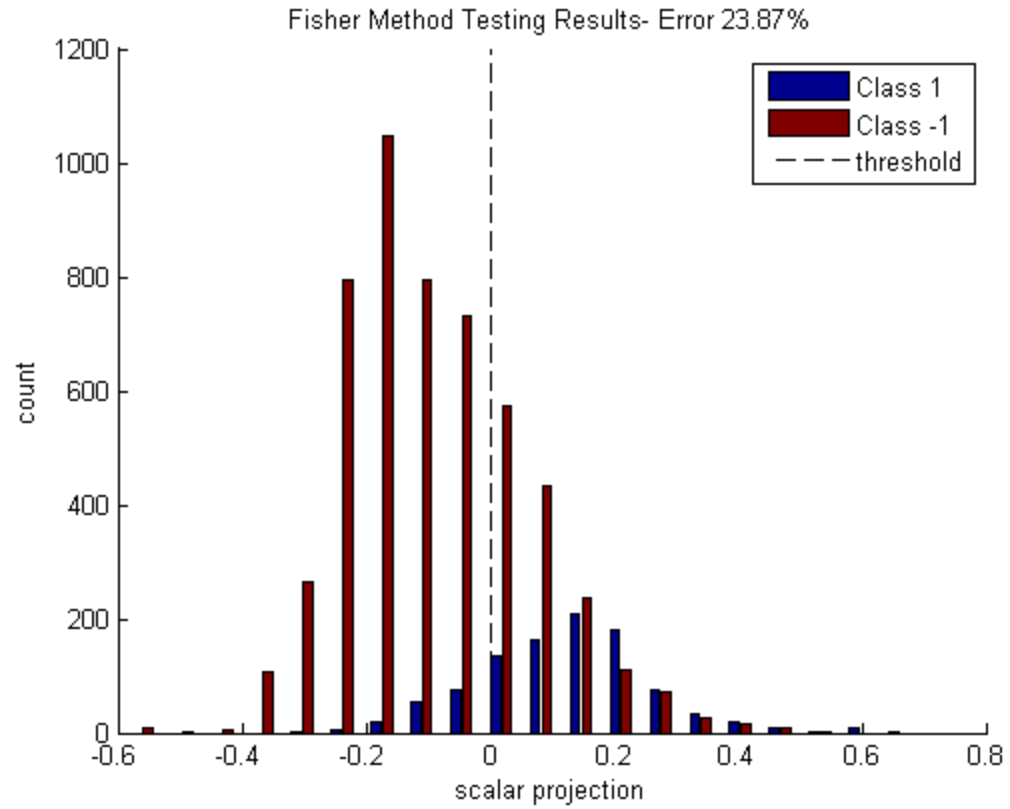
Fisher Method Testing Results- Error 23.87%

*Published with MATLAB® R2014a*