

# Assignment 2

---

Nabil Chowdhury

ID: 260622155

COMP 551

February 12, 2018

## 1 Generating DS1

The code explains this step clearly. The data is shuffled and then output to DS1\_test and DS1\_train. The test set contains 1200 examples (600 positive, 600 negative), and the train set contains 2800 examples (1400 positive, 1400 negative).

## 2 Probabilistic LDA Model

Using the data generated in Question 1, the parameters of the probabilistic LDA model were found and used to calculate the coefficients  $w_0$  and  $\mathbf{w}$ :

$$w_0 = -27.22691891296201$$

$$\mathbf{w} = \begin{bmatrix} -14.35081404 \\ 8.51556309 \\ 5.75326526 \\ 3.21580021 \\ 9.56280808 \\ 4.26881422 \\ -17.11861584 \\ 23.83916969 \\ 29.0419213 \\ -9.05879699 \\ 13.18260588 \\ 12.4169004 \\ -15.51799165 \\ -12.95808369 \\ 5.72465615 \\ -12.90280764 \\ -29.42539436 \\ 6.57706464 \\ 0.6836093 \\ 4.99936544 \end{bmatrix} \quad (1)$$

These coefficients were used to obtain the following metrics from the test set:

Metric	Score
Accuracy	0.9525
Precision	0.95
Recall	0.95477
F-Measure	0.95238

### 3 K Nearest Neighbors

Odd k values from 1 to 99 inclusive were tested. KNN performs significantly worse than the linear model for all values of k tested. Here is a plot of K vs the four metrics:



The best k for accuracy and recall was 77. The best k for precision and F-score was 17. It is clear from the plot that certain values of K perform better than others, but all are much worse when compared to the linear model.

The raw scores for all k is shown below, with k=17 and k=77 rows bolded:

K value	Accuracy	Precision	Recall	F-Measure
1	0.52917	0.51333	0.53012	0.52159
3	0.53833	0.53833	0.53833	0.53833
5	0.52000	0.53667	0.51935	0.52787
7	0.54083	0.55667	0.53958	0.54799
9	0.55083	0.55500	0.55041	0.55270
11	0.55167	0.57000	0.54984	0.55974
13	0.56417	0.56667	0.56385	0.56525
15	0.56500	0.56000	0.56566	0.56281
<b>17</b>	0.56417	<b>0.57833</b>	0.56240	<b>0.57025</b>
19	0.56167	0.56833	0.56086	0.56457
21	0.56333	0.57000	0.56250	0.56623
23	0.55333	0.55833	0.55281	0.55556
25	0.54250	0.53833	0.54286	0.54059
27	0.54667	0.53667	0.54762	0.54209
29	0.55667	0.54667	0.55782	0.55219
31	0.56333	0.55667	0.56419	0.56040
33	0.54917	0.53667	0.55043	0.54346
35	0.55750	0.54167	0.55938	0.55038
37	0.55583	0.53667	0.55806	0.54715
39	0.55083	0.52667	0.55342	0.53971
41	0.55000	0.53167	0.55190	0.54160
43	0.54583	0.52000	0.54833	0.53379
45	0.54750	0.52167	0.55009	0.53550
47	0.54583	0.51667	0.54867	0.53219
49	0.55667	0.52333	0.56071	0.54138
51	0.55167	0.51500	0.55576	0.53460
53	0.55667	0.53000	0.55986	0.54452
55	0.55917	0.53167	0.56261	0.54670
57	0.55583	0.53000	0.55888	0.54405
59	0.55583	0.52833	0.55908	0.54327
61	0.56333	0.53667	0.56690	0.55137
63	0.57500	0.54667	0.57951	0.56261
65	0.57500	0.53833	0.58094	0.55882
67	0.57500	0.54333	0.58007	0.56110
69	0.57750	0.54667	0.58259	0.56406
71	0.57833	0.54667	0.58363	0.56454
73	0.57250	0.53500	0.57838	0.55584
75	0.57000	0.53333	0.57554	0.55363
<b>77</b>	<b>0.57917</b>	0.54333	<b>0.58528</b>	0.56353
79	0.56250	0.52833	0.56708	0.54702
81	0.56250	0.52667	0.56732	0.54624
83	0.56083	0.52167	0.56600	0.54293
85	0.56333	0.52667	0.56835	0.54671
87	0.56750	0.53333	0.57245	0.55220
89	0.57250	0.53167	0.57895	0.55430
91	0.56750	0.52667	0.57350	0.54909
93	0.57167	0.52667	0.57875	0.55148
95	0.57083	0.52333	0.57827	0.54943
97	0.56250	0.51000	0.56983	0.53826
99	0.56417	0.51000	0.57196	0.53921

## 4 Generating DS2

The code explains this step clearly. The data is shuffled and then output to DS2.test and DS2.train. The test set contains 1200 examples (600 positive, 600 negative), and the train set contains 2800 examples (1400 positive, 1400 negative).

## 5 LDA and KNN on DS2

### 5.1 LDA

Using the data generated in Question 4, the parameters of the probabilistic LDA model were found and used to calculate the coefficients  $w_0$  and  $\mathbf{w}$ :

$$w_0 = 0.06778521423199102$$
$$\mathbf{w} = \begin{bmatrix} 0.09214283 \\ 0.00292239 \\ -0.01834667 \\ 0.02224508 \\ 0.01717231 \\ -0.02159144 \\ 0.09611934 \\ -0.11333148 \\ 0.02888671 \\ -0.02903803 \\ -0.04746149 \\ -0.00753071 \\ -0.09138599 \\ 0.00557748 \\ 0.03638741 \\ -0.03275205 \\ -0.05173444 \\ -0.00352907 \\ 0.01535743 \\ 0.03403121 \end{bmatrix} \quad (2)$$

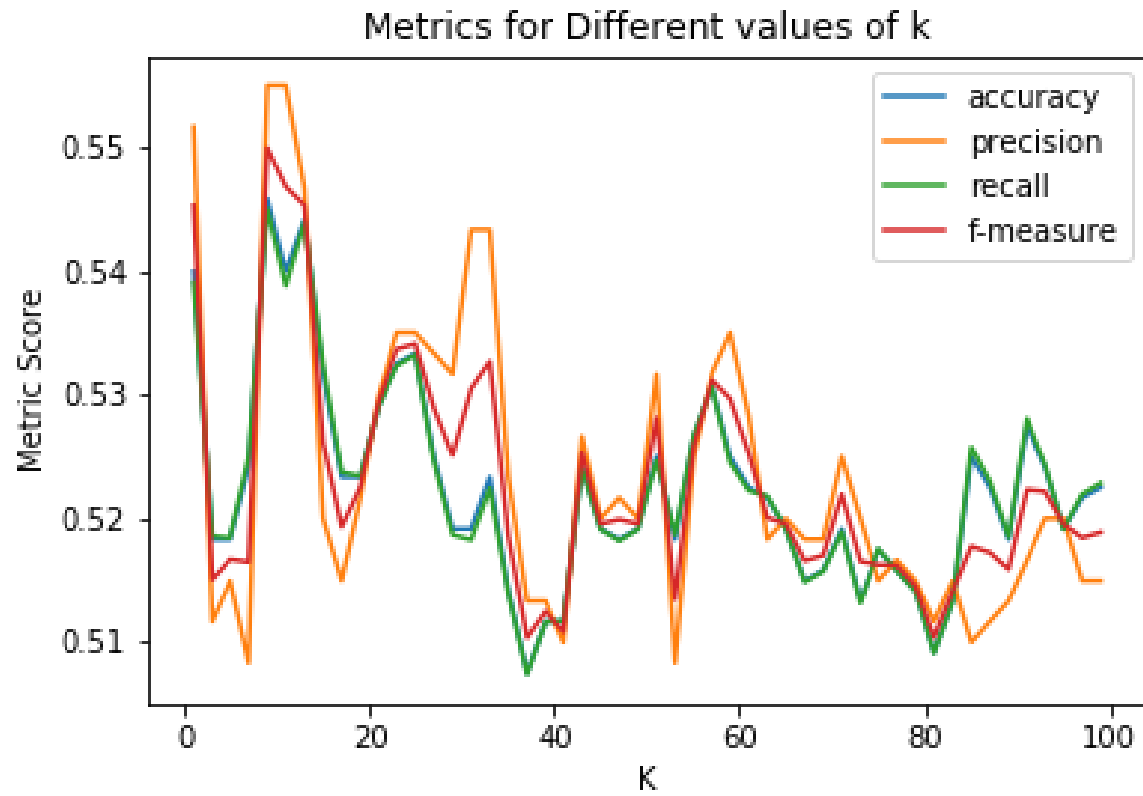
These coefficients were used to obtain the following metrics from the test set:

Metric	Score
Accuracy	0.53416
Precision	0.53166
Recall	0.53433
F-Measure	0.53299

We notice that when the data is generate by multiple Gaussian distributions, the performance of LDA takes a major hit. Let's look at KNN's performance before further comparison.

## 5.2 KNN

Here is a plot of K vs the four metrics:



The optimal K value was 9. The full performance metric for all Ks tested is shown below:

K value	Accuracy	Precision	Recall	F-Measure
1	0.54000	0.55167	0.53909	0.54530
3	0.51833	0.51167	0.51858	0.51510
5	0.51833	0.51500	0.51846	0.51672
7	0.52417	0.50833	0.52496	0.51651
9	<b>0.54583</b>	<b>0.55500</b>	<b>0.54501</b>	<b>0.54996</b>
11	0.54000	0.55500	0.53883	0.54680
13	0.54417	0.54667	0.54395	0.54530
15	0.53167	0.52000	0.53242	0.52614
17	0.52333	0.51500	0.52373	0.51933
19	0.52333	0.52167	0.52341	0.52254
21	0.52917	0.53000	0.52912	0.52956
23	0.53250	0.53500	0.53234	0.53367
25	0.53333	0.53500	0.53322	0.53411
27	0.52500	0.53333	0.52459	0.52893
29	0.51917	0.53167	0.51870	0.52510
31	0.51917	0.54333	0.51828	0.53051
33	0.52333	0.54333	0.52244	0.53268
35	0.51417	0.52333	0.51391	0.51858
37	0.50750	0.51333	0.50741	0.51036
39	0.51167	0.51333	0.51163	0.51248
41	0.51167	0.51000	0.51171	0.51085
43	0.52417	0.52667	0.52405	0.52535
45	0.51917	0.52000	0.51913	0.51957
47	0.51833	0.52167	0.51821	0.51993
49	0.51917	0.52000	0.51913	0.51957
51	0.52500	0.53167	0.52467	0.52815
53	0.51833	0.50833	0.51871	0.51347
55	0.52667	0.52500	0.52676	0.52588
57	0.53083	0.53167	0.53078	0.53122
59	0.52500	0.53500	0.52451	0.52970
61	0.52250	0.52833	0.52224	0.52527
63	0.52167	0.51833	0.52181	0.52007
65	0.51917	0.52000	0.51913	0.51957
67	0.51500	0.51833	0.51490	0.51661
69	0.51583	0.51833	0.51575	0.51704
71	0.51917	0.52500	0.51895	0.52196
73	0.51333	0.52000	0.51316	0.51656
75	0.51750	0.51500	0.51759	0.51629
77	0.51583	0.51667	0.51581	0.51624
79	0.51417	0.51500	0.51414	0.51457
81	0.50917	0.51167	0.50912	0.51039
83	0.51333	0.51500	0.51329	0.51414
85	0.52500	0.51000	0.52577	0.51777
87	0.52250	0.51167	0.52300	0.51727
89	0.51833	0.51333	0.51852	0.51591
91	0.52750	0.51667	0.52811	0.52233
93	0.52417	0.52000	0.52437	0.52218
95	0.51917	0.52000	0.51913	0.51957
97	0.52167	0.51500	0.52196	0.51846
99	0.52250	0.51500	0.52284	0.51889

## 6 Comparing performance of LDA and KNN on DS1 and DS2

We can see that KNN actually beats LDA for DS2 (only slightly)! In fact, the performances are quite similar to each other. When comparing the results to those of DS1, We see that KNN performs slightly worse for DS2 than DS1 (this is most likely due to randomness), whereas LDA performs much worse for DS2 than DS1. This experiment provides evidence that KNN is less affected when the dataset is not generated from a single Gaussian distribution. LDA performs worse because the in-class variances are bigger for DS2 as compared to DS1 due to multiple different Gaussian distributions being used instead of just one.

## References

[1] Only class notes were used for this assignment for all questions.