

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/2539230>

# On Discriminative vs. Generative Classifiers: A comparison of logistic regression and naive Bayes

Article in *Advances in neural information processing systems* · April 2002

Source: CiteSeer

---

CITATIONS

893

---

READS

523

2 authors, including:



[Michael Jordan](#)

University of California, Berkeley

652 PUBLICATIONS 88,894 CITATIONS

SEE PROFILE

# Comment on “On Discriminative vs. Generative Classifiers: A Comparison of Logistic Regression and Naive Bayes”

Jing-Hao Xue<sup>a,\*</sup>, D. Michael Titterington<sup>a</sup>

<sup>a</sup>*Department of Statistics, University of Glasgow, Glasgow G12 8QQ, UK*

---

## Abstract

Comparison of generative and discriminative classifiers is an ever-lasting topic. Based on their theoretical and empirical comparisons between the naïve Bayes classifier and linear logistic regression, Ng and Jordan (2001) claimed that there existed two distinct regimes of performance between the generative and discriminative classifiers with regard to the training-set size. However, our empirical and simulation studies, as presented in this paper, suggest that it is not so reliable to claim such an existence of the two distinct regimes. In addition, for real world datasets, so far there is no theoretically correct, general criterion for choosing between the discriminative and the generative approaches to classification of an observation  $\mathbf{x}$  into a class  $y$ ; the choice depends on the relative confidence you have in the correctness of the specification of either  $p(y|\mathbf{x})$  or  $p(\mathbf{x}, y)$ . This can be to some extent a demonstration of why Efron (1975) and O’Neill (1980) prefer LDA but other empirical studies may prefer linear logistic regression instead. Furthermore, we suggest that pairing of either LDA assuming a common diagonal covariance matrix (LDA- $\Lambda$ ) or the naïve Bayes classifier and linear logistic regression may not be perfect, and hence it may not be reliable for any claim that was derived from the comparison between LDA- $\Lambda$

or the naïve Bayes classifier and linear logistic regression to be generalised to all the generative and discriminative classifiers.

*Key words:* Asymptotic relative efficiency; Discriminative classifiers; Generative classifiers; Logistic regression; Normal-based Discriminant Analysis; Naïve Bayes classifier

---

## 1 Introduction

Generative classifiers, also termed the sampling paradigm (Dawid, 1976), such as normal-based discriminant analysis and the naïve Bayes classifier, model the joint distribution  $p(\mathbf{x}, y)$  of the measured features  $\mathbf{x}$  and the class labels  $y$  factorised in the form  $p(\mathbf{x}|y)p(y)$ , and learn the model parameters through maximisation of the likelihood with respect to  $p(\mathbf{x}|y)p(y)$ . Discriminative classifiers, also termed the diagnostic paradigm, such as logistic regression, model the conditional distribution  $p(y|\mathbf{x})$  of the class labels given the features, and learn the model parameters through maximising the conditional likelihood based on  $p(y|\mathbf{x})$ .

Comparison of generative and discriminative classifiers is an ever-lasting topic (Efron, 1975; O'Neill, 1980; Titterington et al., 1981; Rubinstein and Hastie, 1997; Ng and Jordan, 2001). Ng and Jordan (2001) presented some theoretical and empirical comparisons between linear logistic regression and the naïve Bayes classifier. The naïve Bayes classifier is a generative classifier, which assumes

---

\* Corresponding author. Tel.: +44 141 330 2474; fax: +44 141 330 4814.

*Email addresses:* `jinghao@stats.gla.ac.uk` (Jing-Hao Xue),  
`mike@stats.gla.ac.uk` (D. Michael Titterington).

1 statistically independent features  $\mathbf{x}$  within classes  $y$  and thus diagonal co-  
 2 variance matrices within classes; it is equivalent to normal-based linear (for  
 3 a common diagonal covariance matrix) or quadratic (for unequal within-class  
 4 covariance matrices) discriminant analysis, when  $\mathbf{x}$  is assumed normally dis-  
 5 tributed for each class. The results in Ng and Jordan (2001) suggested that,  
 6 between the two classifiers, there were two distinct regimes of discriminant  
 7 performance with respect to the training-set size. More precisely, they pro-  
 8 posed that the discriminative classifier had lower asymptotic error rate while  
 9 the generative classifier may approach its (higher) asymptotic error rate much  
 10 faster. In other words, the discriminative classifier performs better with larger  
 11 training sets while the generative classifier does better with smaller training  
 12 sets.

13 The setting for the theoretical proof and empirical evidence in Ng and Jordan  
 14 (2001) includes a binary class label  $y$ , *e.g.*,  $y \in \{1, 2\}$ , a  $p$ -dimensional feature  
 15 vector  $\mathbf{x}$  and the assumption of conditional independence amongst  $\mathbf{x}|y$ , the  
 16 features within a class.

17 In the case of discrete features, each feature  $x_i, i = 1, \dots, p$ , independent of  
 18 other features within  $\mathbf{x}$ , is assumed within a class to be a binomial variable  
 19 such that its value  $x_i \in \{0, 1\}$  within each class. We observe, however, this may  
 20 not guarantee the discriminant function  $\lambda(\alpha) = \log\{p(y = 1|\mathbf{x})/p(y = 2|\mathbf{x})\}$ ,  
 21 where  $\alpha$  is a parameter vector, to be linear; therefore, the naïve Bayes classifier  
 22 may not be a partner of linear logistic regression as a generative-discriminative  
 23 pair.

24 In the case of continuous features,  $\mathbf{x}|y$  is assumed to follow Gaussian distribu-  
 25 tions with equal covariance matrices across the two classes, *i.e.*,  $\Sigma_1 = \Sigma_2$  and,

1 in view of the conditional independence assumption, both covariance matrices  
2 are equal to a diagonal matrix  $\Lambda$ . All of the observed values of the features  
3 are rescaled so that  $x_i \in [0, 1]$ .

4 Based on such a setting, Ng and Jordan (2001) compared two so-called generative-  
5 discriminative pairs: one is for the continuous case, comparing LDA assuming  
6 a common diagonal covariance matrix  $\Lambda$  (denoted by LDA- $\Lambda$  hereafter) vs.  
7 linear logistic regression, and the other is for the discrete case, comparing the  
8 naïve Bayes classifier vs. linear logistic regression.

9 The conditional independence amongst the features within a class is a nec-  
10 essary condition for the naïve Bayes classifier and LDA- $\Lambda$ , but it is not a  
11 necessary condition for linear logistic regression. Therefore, the generative-  
12 discriminative pair of LDA with a common full covariance matrix  $\Sigma$  (denoted  
13 by LDA- $\Sigma$  hereafter) vs. linear logistic regression also merits investigation. In  
14 addition, a comparison of quadratic normal discriminant analysis (QDA) with  
15 unequal diagonal matrices  $\Lambda_1$  and  $\Lambda_2$  (denoted by QDA- $\Lambda_g$  hereafter) and un-  
16 equal full covariance matrices  $\Sigma_1$  and  $\Sigma_2$  (denoted by QDA- $\Sigma_g$  hereafter) with  
17 quadratic logistic regression may provide an interesting extension of the work  
18 of Ng and Jordan (2001).

19 Ng and Jordan (2001) reported experimental results on 15 real-world datasets,  
20 8 with only continuous and binary features and 7 with only discrete features,  
21 from the UCI machine learning repository (Newman et al., 1998); this reposi-  
22 tory stores more than 100 datasets contributed and widely used by the machine  
23 learning community, as a benchmark for empirical studies of machine learning  
24 approaches. As pointed out in that paper, there were a few cases (2 out of  
25 8 continuous cases and 4 out of 7 discrete cases) that did not support the

1 better asymptotic performance of the discriminative classifier, primarily be-  
2 cause of the lack of large enough training sets. However, it is known that the  
3 performance of a classifier varies to some extent with the features selected.

4 In this context, we first replicate experiments on these 15 datasets, with and  
5 without stepwise variable selection being performed on the full linear logistic  
6 regression model using all the observations of each dataset. In the stepwise  
7 variable selection process, the decision to include or exclude a variable is based  
8 on the calculation of the Akaike information criterion (AIC). Furthermore, in  
9 the 8 continuous cases, both LDA- $\Lambda$  and LDA- $\Sigma$  are compared with linear  
10 logistic regression. Then we will extend the comparison to between QDA and  
11 quadratic logistic regression for the 8 continuous UCI datasets and finally to  
12 simulated continuous datasets.

13 The implementations in R (<http://www.r-project.org/>) of LDA and QDA are  
14 rewritten from a Matlab function *cda* for classical linear and quadratic dis-  
15 criminant analysis (Verboven and Hubert, 2005). Logistic regression is imple-  
16 mented by an R function *glm* from a standard package **stats** in R, and the  
17 naïve Bayes classifier is implemented by an R function *naiveBayes* from a  
18 contributed package **e1071** for R.

19 In addition, similarly to what was done by Ng and Jordan (2001), for each  
20 sampled training-set size  $m$ , we perform 1000 random splits of each dataset  
21 into a training set of size  $m$  and a test set of size  $N - m$ , where  $N$  is the number  
22 of observations in the whole dataset, and report the average of the misclas-  
23 sification error rates over these 1000 test sets. The training set is required to  
24 have at least 1 sample for each of the two classes, and, for discrete datasets, to  
25 have all the levels of the features presented by the training samples, otherwise

1 the prediction for the test set may be asked to predict on some new levels for  
2 which no information has been provided in the training process.

3 Meanwhile, we observe that, in order to have all the coefficients of predictor  
4 variables in the model estimated in our implementation of logistic regression  
5 by *glm*, the number  $m$  of training samples should be larger than the number  
6  $\tilde{p}$  of predictor variables, where  $\tilde{p} = p$  for the continuous cases if all  $p$  features  
7 are used for the linear model. More attention should be paid to the discrete  
8 cases with multinomial features in the model, where more dummy variables  
9 have to be used as the predictor variables, with the consequence that  $\tilde{p}$  could  
10 be much larger than  $p$ , *e.g.*,  $\tilde{p} = 3p$  for the linear model if all the features have  
11 4 levels. In other words, although we may report misclassification error rates  
12 for logistic regression with small  $m$ , it is not reliable for us to base any general  
13 claim on those of  $m$  smaller than  $\tilde{p}$ , the actual number of predictor variables  
14 used by the logistic regression model.

## 15 **2 Linear Discrimination On Continuous Datasets**

16 For the continuous datasets, as was done by Ng and Jordan (2001), all the  
17 multinomial features are removed so that only continuous and binary features  
18  $x_i$  are kept and their values  $x_i$  are rescaled into  $[0, 1]$ . Any observation with  
19 missing features is removed from the datasets, as is any feature with only a  
20 single value for all the observations.

21 In addition, before carrying out the classification, we perform the Shapiro-  
22 Wilk test for within-class normality for each feature  $x_i|y$  and Levene's test for  
23 homogeneity of variance across the two classes. Levene's test is less sensitive

to deviations from normality than is the Bartlett test, another test for homogeneity of variance. For the following datasets, the significance level is set at 0.05, and we observe that null hypotheses of normality and homogeneity of variance are mostly rejected by the tests at that significance level.

Dataset	$N_0$	$N$	$p$	$p_{AIC}$	$p_{SW}$	$p_L$	$\mathbf{1}_{\{2R-\Lambda\}}$	$\mathbf{1}_{\{2R-\Sigma\}}$
Pima	768	768	8	7	8	5	1	0
Adult	32561	1000	6	6	6	4	1	1
Boston	506	506	13	10	13	12	1	1
Optdigits 0-1	1125	1125	52	5	52	45	1	1
Optdigits 2-3	1129	1129	57	9	57	37	1	0
Ionosphere	351	351	33	20	33	27	1	1
Liver disorders	345	345	6	6	6	1	1	1
Sonar	208	208	60	37	59	16	1	1

Table 1

Description of continuous datasets.

A brief description of the continuous datasets can be found in Table 1, which lists, for each dataset, the total number  $N_0$  of the observations, the number  $N$  of the observations that we use after the pre-processing mentioned above, the total number  $p$  of continuous or binary features, the number  $p_{AIC}$  of features selected by AIC, the number  $p_{SW}$  of features for which the null hypotheses were rejected by the Shapiro-Wilk test and the corresponding number  $p_L$  for Levene’s test, the indicator  $\mathbf{1}_{\{2R-\Lambda\}} \in \{1, 0\}$  of whether or not the two regimes are observed between LDA- $\Lambda$  and linear logistic regression and the indicator



1  $\mathbf{1}_{\{2R-\Sigma\}} \in \{1, 0\}$  with regard to LDA- $\Sigma$ . Note that, for some large datasets  
2 such as “Adult” (and “Sick” in Section 4), in order to reduce computational  
3 complexity without degrading the validity of the comparison between the clas-  
4 sifiers, we randomly sample observations with the class prior probability kept  
5 unchanged.

6 Our results are shown in Figure 1. Since with variable selection by AIC the  
7 results conform more to the claim of two regimes by Ng and Jordan (2001), we  
8 show such results if they are different from those without variable selection.  
9 Meanwhile, in the figures hereafter we use the same annotations of the vertical  
10 and horizontal axes and the same line type as those in Ng and Jordan (2001).  
11 All the observations from these figures are only valid for  $m > p$ , with the  
12 intercept in  $\lambda(\alpha)$  taken into account.

13 In general, our study of these continuous datasets suggests the following con-  
14 clusions.

- 15 (1) In the comparison of LDA- $\Lambda$  vs. linear logistic regression, the pattern of  
16 our results can be said to be similar to that of Ng and Jordan (2001).
- 17 (2) The performance of LDA- $\Sigma$  is worse than that of LDA- $\Lambda$  when the  
18 training-set size  $m$  is small, but better than that of the latter when  $m$  is  
19 large.
- 20 (3) The performance of LDA- $\Sigma$  is better than that of linear logistic regression  
21 when  $m$  is small, but is more or less comparable with that of the latter  
22 when  $m$  is large.
- 23 (4) Pre-processing with variable selection can reveal the distinction in per-  
24 formance of generative and discriminative classifiers with fewer training  
25 samples.

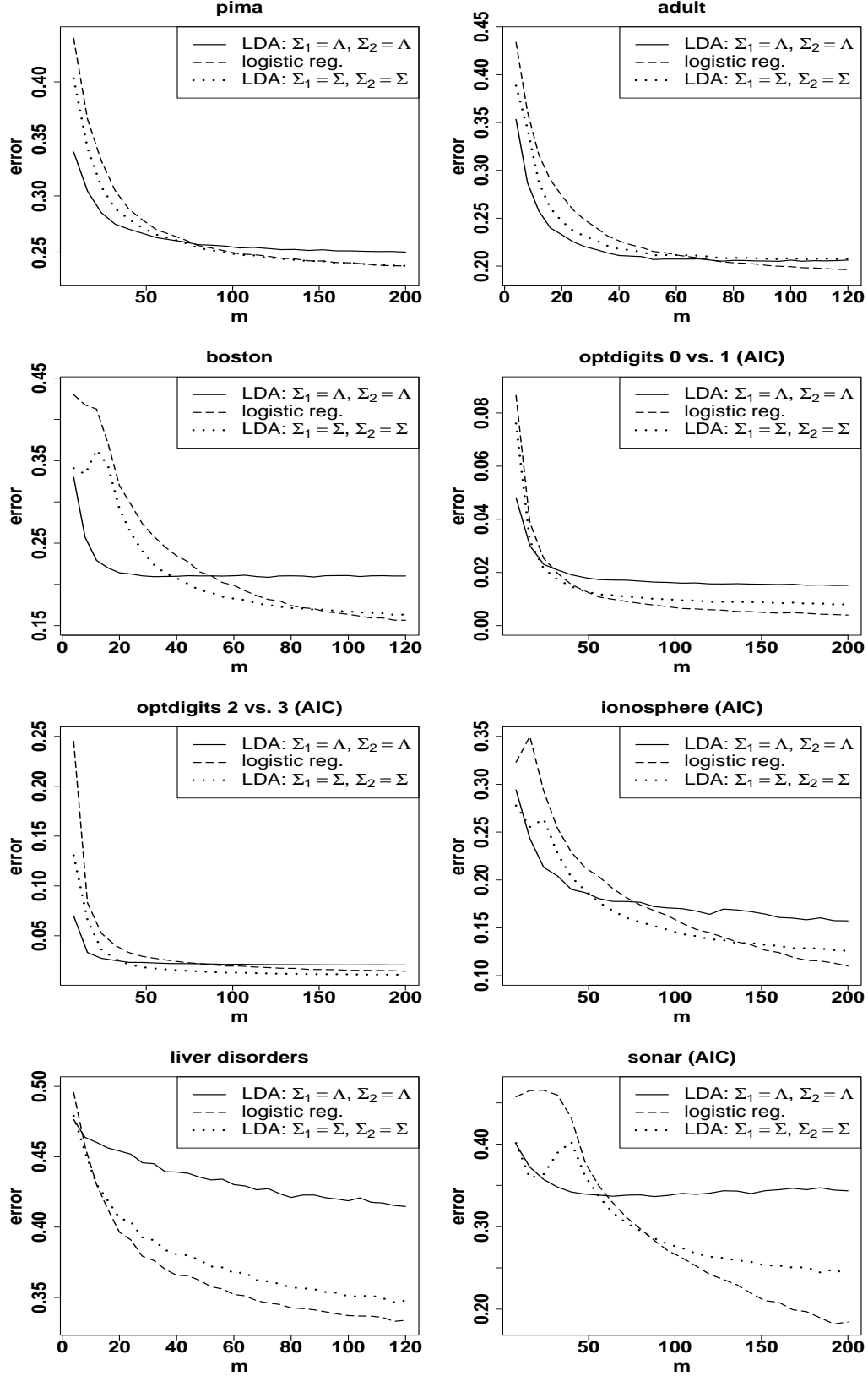


Figure 1. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 1000 random training/test set splits) on the continuous UCI datasets, with regard to linear discrimination.

(5) Therefore, considering LDA- $\Lambda$  vs. linear logistic regression, there is strong evidence to support the claim that the discriminative classifier has lower asymptotic error rate while the generative classifier may approach its (higher) asymptotic error rate much faster. However, considering LDA- $\Sigma$  vs. linear logistic regression, the evidence is not so strong, although the claim may still be made.

### 3 Quadratic Discrimination On Continuous Datasets

As a natural extension of the comparison between LDA- $\Lambda$  (with a common diagonal covariance matrix  $\Lambda$  across the two classes), LDA- $\Sigma$  (with a common full covariance matrix  $\Sigma$ ) and linear logistic regression that was presented in Section 2, this section presents the comparison between QDA- $\Lambda_g$  (with two unequal diagonal covariance matrices  $\Lambda_1$  and  $\Lambda_2$ ), QDA- $\Sigma_g$  (with two unequal full covariance matrices  $\Sigma_1$  and  $\Sigma_2$ ) and quadratic logistic regression.

Using the 8 continuous UCI datasets, all the settings are the same as those in Section 2 except for the following aspects.

First, considering that in the quadratic logistic regression model there are  $p(p-1)/2$  interaction terms between the features in a  $p$ -dimensional feature space, a large number of interactions when the dimensionality  $p$  is high, the model is constrained to contain only the intercept, the  $p$  features and their  $p$  squared terms, so as to make the estimation of the model more feasible and interpretable.

Secondly, for the same reason as explained at the end of Section 1, in the reported plots of misclassification error rate vs.  $m$  without variable selection,

1 only the results for  $m > 2p$  are reliable for comparison since there are  $2p$   
2 predictor variables in the quadratic logistic regression model.

3 Thirdly, the datasets are randomly split into training sets and test sets 100  
4 times rather than 1000 times for each sampled training-set size  $m$  because of  
5 the higher computational complexity of the quadratic models compared with  
6 that of the linear models.

7 In general, our study of these continuous datasets, as shown in Figure 2,  
8 suggests quite similar conclusions to those in Section 3, through substituting  
9 QDA- $\Lambda_g$  for LDA- $\Lambda$ , QDA- $\Sigma_g$  for LDA- $\Sigma$ , and quadratic logistic regression for  
10 linear logistic regression.

#### 11 4 Linear Discrimination On Discrete Datasets

12 For the discrete datasets, as was done by Ng and Jordan (2001), all the contin-  
13 uous features are removed and only the discrete features are used. The results  
14 are entitled ‘multinomial’ in following figures if a dataset includes multinomial  
15 features, and otherwise are entitled ‘binomial’. Meanwhile, any observation  
16 with missing features is removed from the datasets, as is any feature with  
17 only a single value for all the observations.

18 A brief description of the discrete datasets can be found in Table 2, which  
19 includes the indicator  $\mathbf{1}_{\{2R-NB\}} \in \{1, 0\}$  of whether or not the two regimes  
20 are observed between the naïve Bayes classifier and linear logistic regression.  
21 Our results are shown in Figure 3. All the observations from these figures  
22 are only valid for  $m > \tilde{p}$ , with dummy variables taken into account for the  
23 multinomial features.

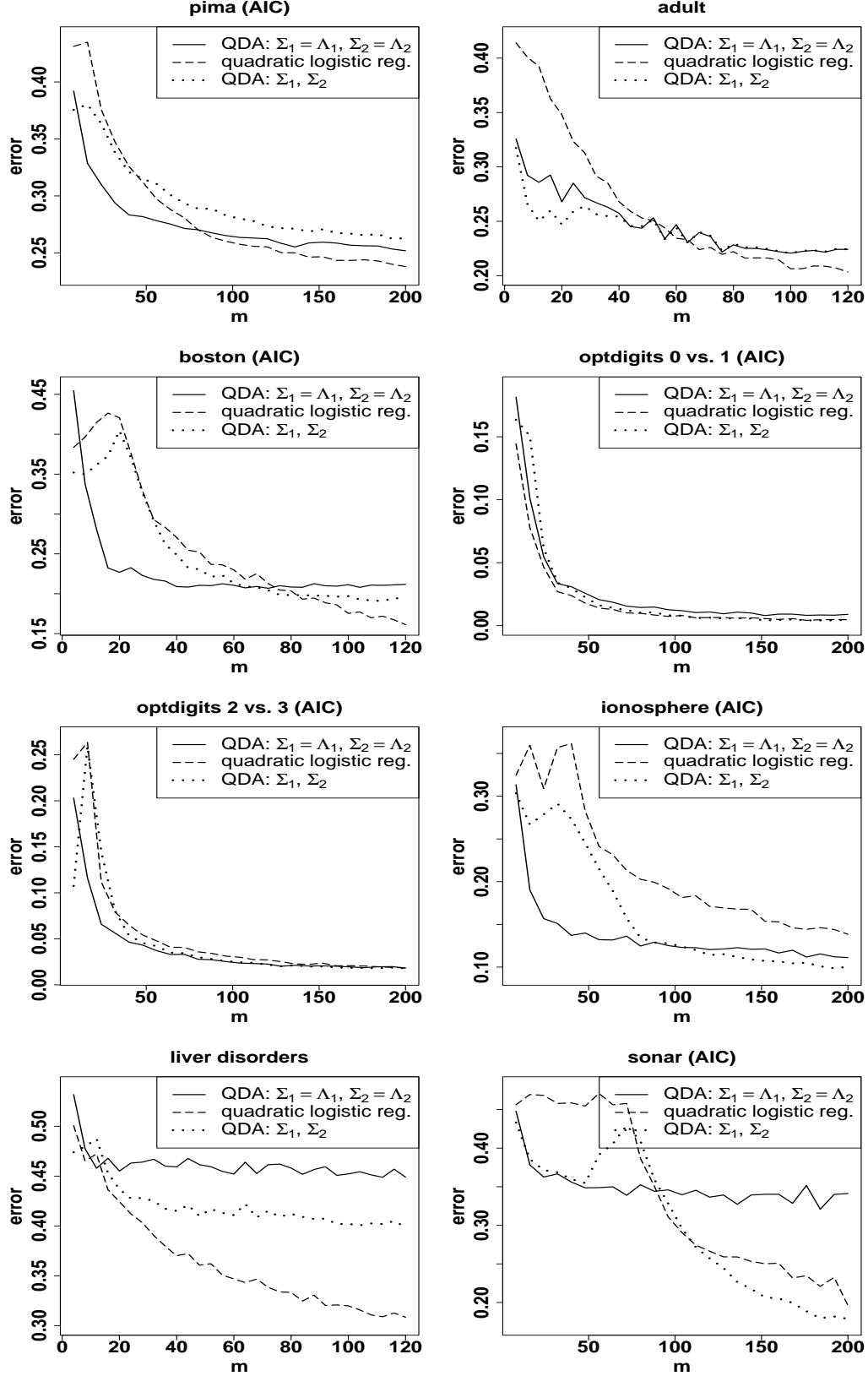


Figure 2. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 100 random training/test set splits) on the continuous UCI datasets, with regard to quadratic discrimination.

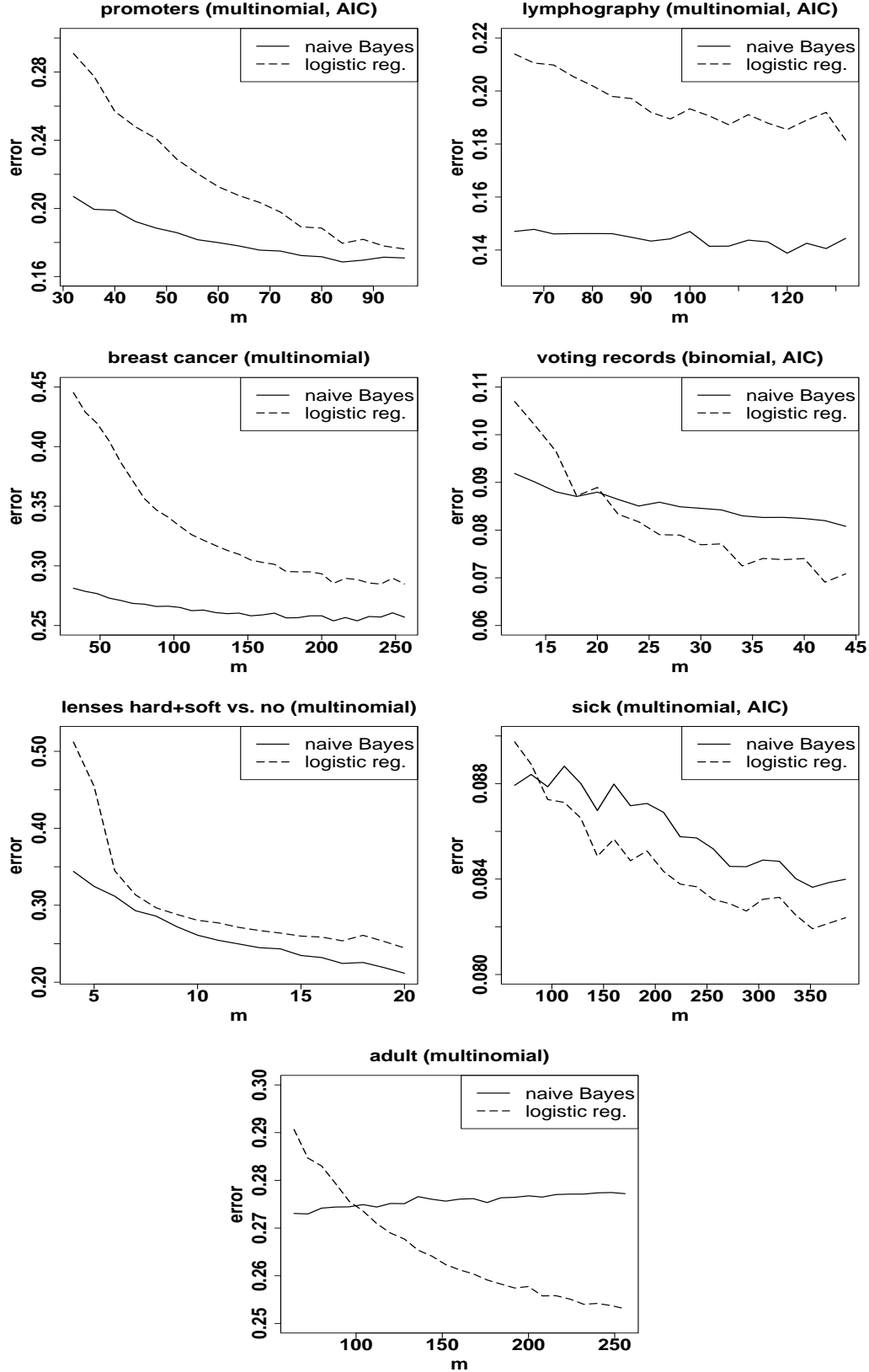


Figure 3. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 1000 random training/test set splits) on the discrete UCI datasets, with regard to linear discrimination.

Dataset	$N_0$	$N$	$p$	$p_{AIC}$	$\mathbf{1}_{\{2R-NB\}}$
Promoters	106	106	57	7	0
Lymphography	148	142	17	10	0
Breast cancer	286	277	9	4	0
Voting recorders	435	232	16	11	1
Lenses	24	24	4	1	0
Sick	2800	500	12	4	1
Adult	32561	1000	5	5	1

Table 2

Description of discrete datasets.

In general, our study of these discrete datasets suggests that, in the comparison of the naïve Bayes classifier vs. linear logistic regression, the pattern of our results can be said to be similar to that of Ng and Jordan (2001).

## 5 Linear Discrimination On Simulated Datasets

In this section, 16 simulated datasets are used to compare the performance of LDA- $\Lambda$ , LDA- $\Sigma$  and linear logistic regression. The samples are simulated from bivariate normal distributions, bivariate Student's  $t$ -distributions, bivariate log-normal distributions and mixtures of 2 bivariate normal distributions, with 4 datasets for each of these 4 types of distribution. Within each dataset there are 1000 simulated samples, which are divided equally into 2 classes. The simulations from the bivariate log-normal distributions and normal mixtures

are based on an R function *mvnrm* for simulating from a multivariate normal distribution from a contributed R package **MASS**, and the simulation from the bivariate Student’s *t*-distribution is implemented by an R function *rmvt* from a contributed R package **mvtnorm**. Differently from the UCI datasets, the simulated data are not rescaled into the range  $[0, 1]$  and no variable selection is used since the feature space is only of dimension two.

### 5.1 Normally Distributed Data

Four simulated datasets are randomly generated from two bivariate normal distributions,  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$ , where  $\mu_1 = (1, 0)^T$ ,  $\mu_2 = (-1, 0)^T$  and  $\Sigma_1$  and  $\Sigma_2$  are subject to four different types of constraint specified as having equal diagonal or full covariance matrices  $\Sigma_1 = \Sigma_2$  and having unequal diagonal or full covariance matrices  $\Sigma_1 \neq \Sigma_2$ .

Similarly to what was done for the UCI datasets, for each sampled training-set size  $m$ , we perform 1000 random splits of the 1000 samples of each simulated dataset into a training set of size  $m$  and a test set of size  $1000 - m$ , and report the average misclassification error rates over these 1000 test sets. The training set is required to have at least 1 sample from each of the two classes. In such a way, LDA- $\Lambda$  and LDA- $\Sigma$  are compared with linear logistic regression, in terms of misclassification error rate, with the following results shown in Figure 4.

The dataset for the top-left panel of Figure 4 has  $\Sigma_1 = \Sigma_2 = \Lambda$  with a diagonal matrix  $\Lambda = \text{Diag}(1, 1)$ , such that the data satisfy the assumptions underlying LDA- $\Lambda$ . The dataset for the top-right panel has  $\Sigma_1 = \Sigma_2 = \Sigma$



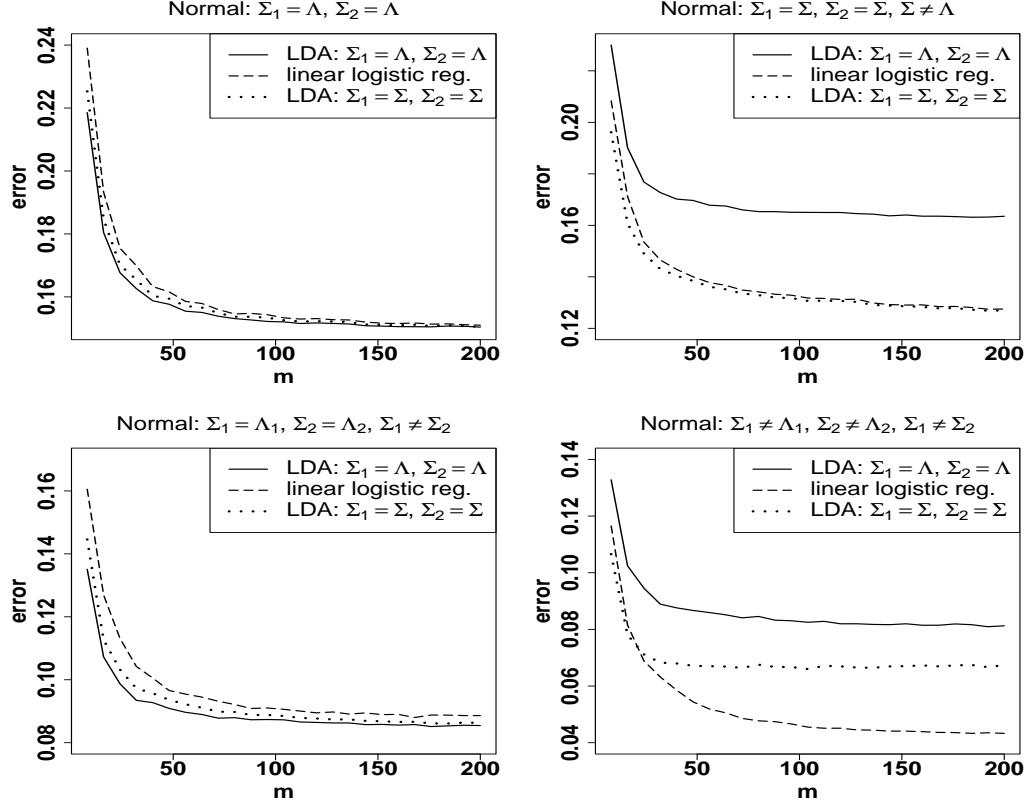


Figure 4. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 1000 random training/test set splits) on simulated bivariate normally distributed data for two classes.

- 1 with a full matrix  $\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$ , such that the data satisfy the assumptions
- 2 underlying LDA- $\Sigma$ . The dataset for the bottom-left panel has  $\Sigma_1 = \Lambda_1, \Sigma_2 =$
- 3  $\Lambda_2$  with diagonal matrices  $\Lambda_1 = \text{Diag}(1, 1)$  and  $\Lambda_2 = \text{Diag}(0.25, 0.75)$ , such
- 4 that the homogeneity of the covariance matrices is violated. The dataset for
- 5 the bottom-right panel has  $\Sigma_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$  and  $\Sigma_2 = \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1.75 \end{bmatrix}$ , such
- 6 that both the homogeneity of the covariance matrices and the conditional
- 7 independence (uncorrelatedness) of the features within a class are violated.

## 5.2 Student's t-Distributed Data

Four simulated datasets are randomly generated from two bivariate Student's  $t$ -distributions, both distributions with degrees of freedom  $\nu = 3$ . The values of class means  $\mu_1$  and  $\mu_2$ , the four types of constraint on  $\Sigma_1$  and  $\Sigma_2$ , and other settings of the experiments are all the same as those in Section 5.1.

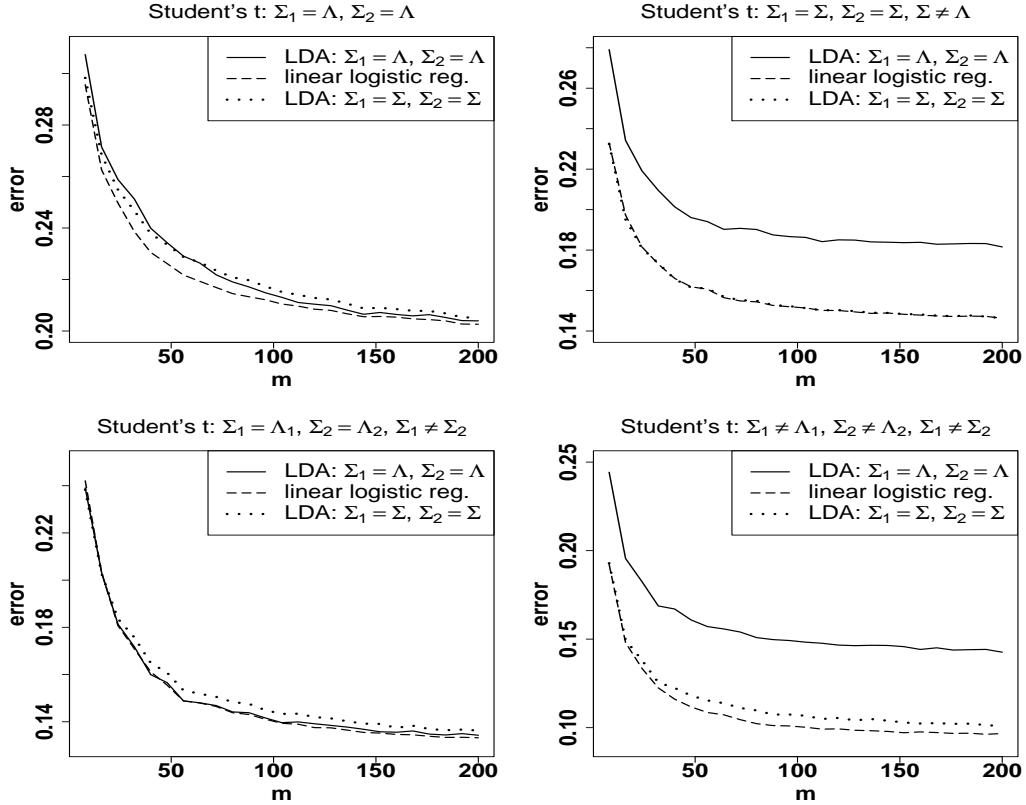


Figure 5. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 1000 random training/test set splits) on simulated bivariate Student's  $t$ -distributed data for two classes.

The results are shown in Figure 5, where for each panel the constraint with regard to  $\Sigma_1$  and  $\Sigma_2$  is the same as the corresponding one in Figure 4, except for a scalar multiplier  $\nu/(\nu - 2)$ .

### 1 5.3 Log-normally Distributed Data

2 Four simulated datasets are randomly generated from two bivariate log-normal  
 3 distributions, whose logarithms are normally distributed as  $\mathcal{N}(\mu_1, \Sigma_1)$  and  
 4  $\mathcal{N}(\mu_2, \Sigma_2)$ , respectively. The values of  $\mu_1$  and  $\mu_2$ , the four types of constraint  
 5 on  $\Sigma_1$  and  $\Sigma_2$ , and other settings of the experiments are all the same as those  
 6 in Section 5.1.

7 By definition, if a  $p$ -variate random vector  $\mathbf{x} \sim \mathcal{N}(\mu(\mathbf{x}), \Sigma(\mathbf{x}))$ , then a  $p$ -  
 8 variate vector  $\tilde{\mathbf{x}}$  of the exponentials of the components of  $\mathbf{x}$  follows a  $p$ -variate  
 9 log-normal distribution, *i.e.*,  $\tilde{\mathbf{x}} = \exp(\mathbf{x}) \sim \log \mathcal{N}(\mu(\tilde{\mathbf{x}}), \Sigma(\tilde{\mathbf{x}}))$ , where the  $i$ -th  
 10 element  $\mu^{(i)}(\tilde{\mathbf{x}})$  of the mean vector and the  $(i, j)$ -th element  $\Sigma^{(i,j)}(\tilde{\mathbf{x}})$  of the  
 11 covariance matrix,  $i, j = 1, \dots, p$ , are

$$\mu^{(i)}(\tilde{\mathbf{x}}) = e^{\mu^{(i)}(\mathbf{x}) + \frac{\Sigma^{(i,i)}(\mathbf{x})}{2}},$$

12

$$\Sigma^{(i,j)}(\tilde{\mathbf{x}}) = (e^{\Sigma^{(i,j)}(\mathbf{x})} - 1)e^{\mu^{(i)}(\mathbf{x}) + \mu^{(j)}(\mathbf{x}) + \frac{\Sigma^{(i,i)}(\mathbf{x}) + \Sigma^{(j,j)}(\mathbf{x})}{2}}.$$

13 It follows that, if the components of its logarithm  $\mathbf{x}$  are independent and nor-  
 14 mally distributed, the components of the log-normally distributed multivari-  
 15 ate random variable  $\tilde{\mathbf{x}}$  are uncorrelated. In other words, if  $\mathbf{x} \sim \mathcal{N}(\mu(\mathbf{x}), \Lambda(\mathbf{x}))$ ,  
 16 then  $\tilde{\mathbf{x}} = \exp(\mathbf{x}) \sim \log \mathcal{N}(\mu(\tilde{\mathbf{x}}), \Lambda(\tilde{\mathbf{x}}))$ . However, as shown by the equations  
 17 above,  $\Lambda(\tilde{\mathbf{x}})$  is determined by both  $\mu(\mathbf{x})$  and  $\Lambda(\mathbf{x})$ , so that  $\Sigma_1(\mathbf{x}) = \Sigma_2(\mathbf{x})$   
 18 may not mean  $\Sigma_1(\tilde{\mathbf{x}}) = \Sigma_2(\tilde{\mathbf{x}})$ . Therefore, considering in our cases  $\mu_1 \neq \mu_2$ ,  
 19 it can be expected that the pattern of performance of the classifiers for the  
 20 datasets with equal covariance matrices  $\Sigma_1 = \Sigma_2$  in the underlying normal  
 21 distributions could be similar to that for the datasets with unequal covariance  
 22 matrices  $\Sigma_1 \neq \Sigma_2$ , since in both cases the covariance matrices of the log-  
 23 normally distributed variables are in fact unequal. In this context, it makes

- 1 more sense to compare the classifiers in situations with diagonal and full co-  
2 variance matrices of the underlying normally distributed data, respectively,  
3 rather than those with equal and unequal covariance matrices.

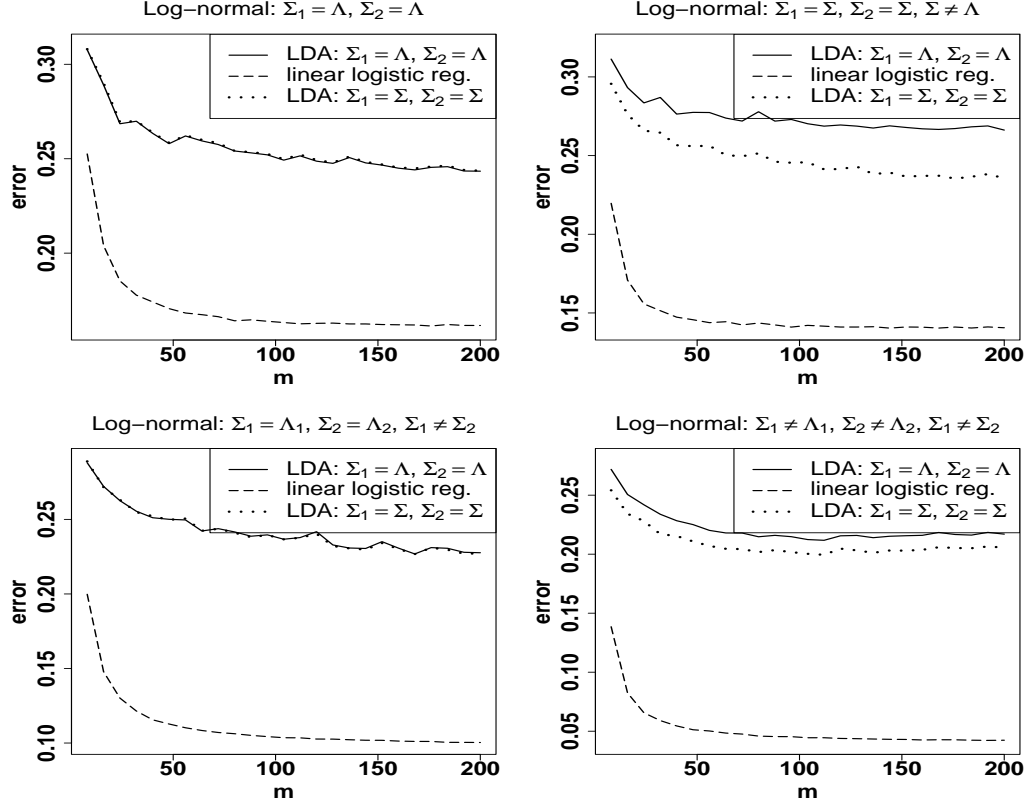


Figure 6. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 1000 random training/test set splits) on simulated bivariate log-normally distributed data for two classes.

- 4 The results are shown in Figure 6, where for each panel the constraint with  
5 regard to  $\Sigma_1$  and  $\Sigma_2$  is the same as the corresponding one in Figure 4.

#### 6 5.4 Normal Mixture Data

- 7 Compared with the normal distribution, the Student's  $t$ -distribution and the  
8 log-normal distribution used in Sections 5.1, 5.2 and 5.3 for the comparison of

1 the classifiers, the mixture of normal distributions is a better approximation  
 2 to real data in a variety of situations. In this section, 4 simulated datasets,  
 3 each consisting of 1000 samples, are randomly generated from two mixtures,  
 4 each of bivariate normal distributions, with 250 samples from each mixture  
 5 component. The two components,  $A$  and  $B$ , of the mixture for Class 1 are  
 6 normally distributed with distributions  $\mathcal{N}(\mu_{1A}, \Sigma_1)$  and  $\mathcal{N}(\mu_{1B}, \Sigma_1)$ , respec-  
 7 tively, where  $\mu_{1A} = (1, 0)^T$  and  $\mu_{1B} = (3, 0)^T$ ; and the two components,  $C$  and  
 8  $D$ , of the mixture for Class 2 are normally distributed with probability den-  
 9 sity functions  $\mathcal{N}(\mu_{2C}, \Sigma_2)$  and  $\mathcal{N}(\mu_{2D}, \Sigma_2)$ , respectively, where  $\mu_{2C} = (-1, 0)^T$   
 10 and  $\mu_{2D} = (-3, 0)^T$ . In such a way, when  $\Sigma_1$  and  $\Sigma_2$  are subject to the four  
 11 different types of constraint with regard to  $\Sigma_1$  and  $\Sigma_2$  as previously discussed,  
 12 the covariance matrices of the two mixtures will be subject to the same con-  
 13 straints. Other settings of the experiments are all the same as that in Section  
 14 5.1.

15 The results are shown in Figure 7, where for each panel the constraint with  
 16 regard to  $\Sigma_1$  and  $\Sigma_2$  is the same as the corresponding one in Figure 4.

## 17 5.5 Summary of Linear Discrimination on Simulated Datasets

18 In general, our study of these simulated continuous datasets suggests the fol-  
 19 lowing conclusions.

- 20 (1) When the data are consistent with the assumptions underlying LDA- $\Lambda$   
 21 or LDA- $\Sigma$ , both methods can perform the best among them and linear  
 22 logistic regression, throughout the range of the training-set size  $m$  in  
 23 our study; in these cases, there is no evidence to support the claim that

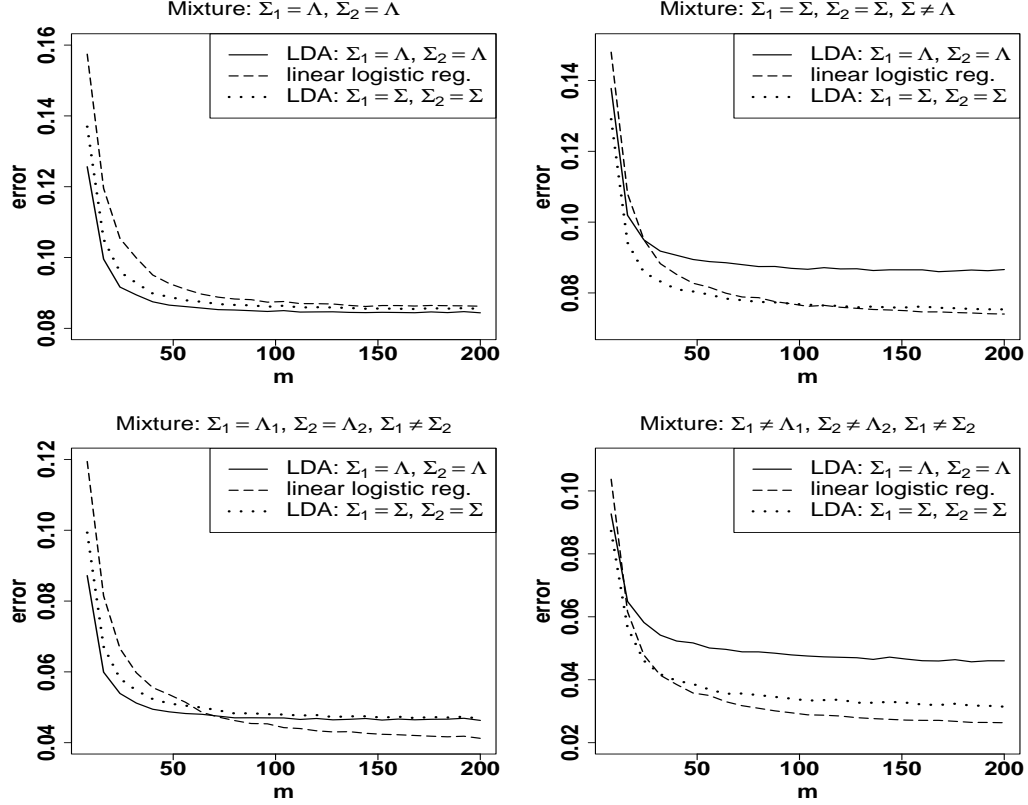


Figure 7. Plots of misclassification error rate vs. training-set size  $m$  (averaged over 1000 random training/test set splits) on simulated bivariate 2-component normal mixture data for two classes.

- 1 the discriminative classifier has lower asymptotic error rate while the
- 2 generative classifier may approach its (higher) asymptotic error rate much
- 3 faster.
- 4 (2) When the data violate the assumptions underlying the LDAs, linear lo-
- 5 gistic regression generally performs better than the LDAs, in particular
- 6 when  $m$  is large; in this case, there is strong evidence to support the
- 7 claim that the discriminative classifier has lower asymptotic error rate,
- 8 but there is no convincing evidence to support the claim that the gen-
- 9 erative classifier may approach its (higher) asymptotic error rate much
- 10 faster.

(3) When the covariance matrices are non-diagonal, LDA- $\Sigma$  performs remarkably better than LDA- $\Lambda$  and more remarkably when  $m$  is large; when the covariance matrices are diagonal, LDA- $\Lambda$  performs generally better than LDA- $\Sigma$  and more so when  $m$  is large.

## 6 Comments on Comparison of Discriminative and Generative Classifiers

Based on the theoretical analysis and empirical comparison between LDA- $\Lambda$  or the naïve Bayes classifiers and linear logistic regression, Ng and Jordan (2001) claim that there are two distinct regimes of performance with regard to the training-set size. Such a claim can be clarified further through commenting on the reliability of the two regimes and the parity between the compared classifiers.

### 6.1 On the Two Regimes of Performance regarding Training-Set Size

Suppose we have a training set  $\{(y_{tr}^{(i)}, \mathbf{x}_{tr}^{(i)})\}_{i=1}^m$  of  $m$  independent observations and a test set  $\{(y_{te}^{(i)}, \mathbf{x}_{te}^{(i)})\}_{i=1}^{N-m}$  of  $N-m$  independent observations, where  $\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)})^T$  is the  $i$ -th observed  $p$ -variate feature vector  $\mathbf{x}$ , and  $y^{(i)} \in \{1, 2\}$  is its observed univariate class label. Let us also assume that each observation  $\{(y^{(i)}, \mathbf{x}^{(i)})\}$  follows an identical distribution so that the testing based on the training results makes sense. In order to simplify the notation, let  $\underline{\mathbf{x}}_{tr}$  denote  $\{(\mathbf{x}_{tr}^{(i)})\}_{i=1}^m$ , and similarly define  $\underline{\mathbf{x}}_{te}$ ,  $\underline{y}_{tr}$  and  $\underline{y}_{te}$ . Meanwhile, a discriminant function  $\lambda(\alpha) = \log\{p(y=1|\mathbf{x})/p(y=2|\mathbf{x})\}$ , which is equivalent to a Bayes classifier  $\hat{y}(\mathbf{x}) = \operatorname{argmax}_y p(y|\mathbf{x})$ , is used for the 2-class classification.

1 Discriminative classifiers estimate the parameter  $\alpha$  of the discriminant func-  
2 tion  $\lambda(\alpha)$  through maximising a conditional probability  $\arg\max_{\alpha} p(\underline{y}_{tr}|\underline{\mathbf{x}}_{tr}, \alpha)$ ;  
3 such an estimation procedure can be regarded as a kind of maximum likeli-  
4 hood estimation with  $p(\underline{y}_{tr}|\underline{\mathbf{x}}_{tr}, \alpha)$  as the likelihood function. It is well known  
5 that, if the 0 – 1 loss function is used so that the misclassification error rate  
6 is the total risk, the Bayes classifiers will attain the minimum error rate. This  
7 implies that, under such a loss function, the discriminative classifiers are in  
8 fact using the same criterion to optimise the estimation of the parameter  $\alpha$   
9 and the performance of classification.

10 In this context, the following claims, supported by the simulation study in  
11 Section 5, can be proposed.

- 12 • If the same dataset is used to train and test, *i.e.*,  $\underline{\mathbf{x}}_{tr}$  as  $\underline{\mathbf{x}}_{te}$  and  $\underline{y}_{tr}$  as  $\underline{y}_{te}$ , then  
13 the discriminative classifiers should always provide the best performance,  
14 no matter how large the training-set size  $m$  is.
- 15 • If  $m$  is large enough to make  $(\underline{y}_{tr}, \underline{\mathbf{x}}_{tr})$  representative of all the observations  
16 including  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ , then the discriminative classifiers should also provide  
17 the best prediction performance on  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ , *i.e.*, with the best asymptotic  
18 performance.
- 19 • We note that all of the above claims are based on the premise that the  
20 modelling of  $p(y|\mathbf{x}, \alpha)$ , such as the linearity of  $\lambda(\alpha)$ , is correctly specified  
21 for all the observations, and thus the only work that remains is to estimate  
22 accurately the parameter  $\alpha$ .
- 23 • If  $m$  is not large enough to make  $(\underline{y}_{tr}, \underline{\mathbf{x}}_{tr})$  representative of all the observa-  
24 tions, and  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$  is not exactly the same as  $(\underline{y}_{tr}, \underline{\mathbf{x}}_{tr})$ , then the discrimina-  
25 tive classifiers may not necessarily provide the best prediction performance



1 on  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ , even though the modelling of  $p(y|\mathbf{x}, \alpha)$  may be correct.

2 Generative classifiers estimate the parameter  $\alpha$  of the discriminant function  
3  $\lambda(\alpha)$  through first maximising a joint probability  $\text{argmax}_{\theta} p(\underline{y}_{tr}, \underline{\mathbf{x}}_{tr}|\theta)$  to ob-  
4 tain a maximum likelihood estimate (MLE)  $\hat{\theta}$  of  $\theta$ , the parameter of the joint  
5 distribution of  $(y, \mathbf{x})$ , and then calculate  $\hat{\alpha}$  as a function  $\alpha(\theta)$  at  $\hat{\theta}$ . Under some  
6 regularity conditions, such as the existence of the first and second derivatives  
7 of the log-likelihood function and the inverse of the Fisher information matrix  
8  $I(\theta)$ , the MLE  $\hat{\theta}$  is asymptotically unbiased, efficient and normally distributed.  
9 Accordingly, by the delta method,  $\hat{\alpha}$  is also asymptotically normally distrib-  
10 uted, unbiased and efficient, given the existence of the first derivative of the  
11 function  $\alpha(\theta)$ .

12 Therefore, the following claims, supported by the simulation study in Section  
13 5, can be proposed.

- 14 • Asymptotically, the generative classifiers will provide the best prediction  
15 performance on  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ . However, this is dependent on the premise that  
16  $p(y, \mathbf{x}|\theta)$  is correctly specified for all the observations.
- 17 • If  $m$  is large enough to make  $(\underline{y}_{tr}, \underline{\mathbf{x}}_{tr})$  representative of all the observa-  
18 tions including  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ , then the generative classifiers should also provide  
19 the best prediction performance on  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ , *i.e.*, with the best asymptotic  
20 performance.
- 21 • We note that all of the above claims are based on the premise that that  
22  $p(y, \mathbf{x}|\theta)$  is correctly specified for all the observations.
- 23 • If  $m$  is not large enough to make  $(\underline{y}_{tr}, \underline{\mathbf{x}}_{tr})$  representative of all the obser-  
24 vations, then the generative classifiers may not necessarily provide the best  
25 prediction performance on  $(\underline{y}_{te}, \underline{\mathbf{x}}_{te})$ .

1 In summary, it is not so reliable to claim the existence of the two distinct  
2 regimes of performance between the generative and discriminative classifiers  
3 with regard to the training-set size  $m$ . For real world datasets such as those  
4 demonstrated in Section 2 and 4, there is no theoretically correct, general  
5 criterion for choosing between the discriminative and the generative classifiers;  
6 the choice depends on the relative confidence we have in the correctness of  
7 the specification of either  $p(y|\mathbf{x})$  or  $p(y, \mathbf{x})$ . This can be to some extent a  
8 demonstration of why Efron (1975) and O’Neill (1980) prefer LDA but other  
9 empirical studies may prefer linear logistic regression instead.

## 10 6.2 *On the Pairing of LDA- $\Lambda$ /Naïve Bayes and Linear Logistic Regression/GAM*

11 As mentioned in Section 1, first, the naïve Bayes classifier cannot guarantee the  
12 linear formulation of the discriminant function  $\lambda(\alpha) = \log\{p(y = 1|\mathbf{x})/p(y = 2|\mathbf{x})\}$ ,  
13 and, secondly, the conditional independence amongst the multiple features  
14 within a class is a necessary condition for the naïve Bayes classifier and LDA-  
15  $\Lambda$  with a diagonal covariance matrix  $\Lambda$  but not for linear logistic regression,  
16 although in the latter the discriminant function  $\lambda(\alpha)$  is modelled as a lin-  
17 ear combination of separate features. Therefore, the comparison between a  
18 generative-discriminative pair of LDA- $\Lambda$ /naïve Bayes classifier vs. linear lo-  
19 gistic regression should be interpreted with caution, in particular when the  
20 data do not support the assumption of conditional independence of  $\mathbf{x}|y$  that  
21 may shed unfavourable light on the simplified generative side, LDA- $\Lambda$  and the  
22 naïve Bayes classifier.

23 In this section, we will illustrate such pairing of two generative-discriminative  
24 pairs: one is LDA- $\Lambda$  vs. linear logistic regression (Ng and Jordan, 2001),

1 and the other is the naïve Bayes classifier vs. generalised additive model  
 2 (GAM) (Rubinstein and Hastie, 1997).

### 3 6.2.1 LDA- $\Lambda$ vs. Linear Logistic Regression

4 Consider a feature vector  $\mathbf{x} = (x_1, \dots, x_p)^T$  and a binary class label  $y = 1, 2$ .

5 Linear logistic regression, one of the discriminative classifiers that do not as-  
 6 sume any distribution  $p(\mathbf{x}|y)$  of the data, is modelled directly with a linear  
 7 discriminant function as

$$\lambda_{\text{dis}}(\alpha) = \log \frac{p(y=1|\mathbf{x})}{p(y=2|\mathbf{x})} = \log \frac{\pi_1}{\pi_2} + \log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=2)} = \beta_0 + \beta^T \mathbf{x} ,$$

8 where  $p(y=k) = \pi_k$ ,  $\alpha^T = (\beta_0, \beta^T)$  and  $\beta$  is a parameter vector of  $p$  elements.

9 By “linear”, we mean a scalar-valued function of a linear combination of the  
 10 features  $x_1, \dots, x_p$  of an observed feature vector  $\mathbf{x}$ .

11 In contrast, LDA- $\Lambda$ , one of the generative classifiers, assumes that the data  
 12 arise from two  $p$ -variate normal distributions with different means but the  
 13 same diagonal covariance matrix such that  $(\mathbf{x}|y=k;\theta) \sim \mathcal{N}(\mu_k, \Lambda)$ ,  $k=1,2$ ,  
 14 where  $\theta = (\mu_k, \Lambda)$ ; this implies an assumption of conditional independence  
 15 between any two features  $x_i|y$  and  $x_j|y$ ,  $i \neq j$ , within a class. The density  
 16 function of  $(\mathbf{x}|y=k;\theta)$  can be written as

$$p(\mathbf{x}|y=k;\theta) = \left\{ e^{\mu_k^T \Lambda^{-1} \mathbf{x}} \right\} \left\{ \frac{1}{\sqrt{(2\pi)^p |\Lambda|}} e^{-\frac{1}{2} \mu_k^T \Lambda^{-1} \mu_k} \right\} \left\{ e^{-\frac{1}{2} \mathbf{x}^T \Lambda^{-1} \mathbf{x}} \right\} ,$$

17 which leads to a linear discriminant function

$$\lambda_{\text{gen}}(\alpha) = \log \frac{p(y=1|\mathbf{x})}{p(y=2|\mathbf{x})} = \log \frac{\pi_1}{\pi_2} + \log \frac{A(\theta_1, \eta)}{A(\theta_2, \eta)} + (\theta_1 - \theta_2)^T \mathbf{x} ,$$

18 where  $\theta_k = \mu_k^T \Lambda^{-1}$ ,  $\eta = \Lambda^{-1}$  and  $A(\theta_k, \eta) = \frac{1}{\sqrt{(2\pi)^p |\Lambda|}} e^{-\frac{1}{2} \mu_k^T \Lambda^{-1} \mu_k}$ .

1 Similarly, by assuming that the data arise from two  $p$ -variate normal distri-  
 2 butions with different means but the same full covariance matrix such that  
 3  $(\mathbf{x}|y = k; \theta) \sim \mathcal{N}(\mu_k, \Sigma)$ ,  $k = 1, 2$ , we can obtain the same formula as  $\lambda_{\text{gen}}(\alpha)$   
 4 but with  $\theta_k = \mu_k^T \Sigma^{-1}$ ,  $\eta = \Sigma^{-1}$  and  $A(\theta_k, \eta) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} e^{-\frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k}$ , which  
 5 leads to the linear discriminant function of LDA- $\Sigma$ . Therefore, we could rewrite  
 6  $\theta$  as  $\theta = (\theta_k, \eta)$ , where  $\theta_k$  is a class-dependent parameter vector while  $\eta$  is a  
 7 common parameter vector across the classes.

8 It is clear that the assumption of conditional independence amongst the fea-  
 9 tures within a class is not a necessary condition for a generative classifier to  
 10 attain a linear  $\lambda_{\text{gen}}(\alpha)$ . In fact, as pointed out by O'Neill (1980), if the fea-  
 11 ture vector  $\mathbf{x}$  follows a multivariate exponential family distribution with the  
 12 density or probability mass function within a class being

$$p(\mathbf{x}|y = k, \theta_k) = e^{\theta_k^T \mathbf{x}} A(\theta_k, \eta) h(\mathbf{x}, \eta), k = 1, 2 ,$$

13 the generative classifiers will attain a linear  $\lambda_{\text{gen}}(\alpha)$ .

### 14 6.2.2 Naïve Bayes vs. Generalised Additive Model (GAM)

15 As with logistic regression, a GAM does not assume any distribution  $p(\mathbf{x}|y)$   
 16 for the data; it is modelled directly with a discriminant function as a sum of  
 17  $p$  functions  $f(x_i)$ ,  $i = 1, \dots, p$ , of the  $p$  features  $x_i$  separately (Rubinstein and  
 18 Hastie, 1997); that is

$$\lambda_{\text{dis}}(\alpha) = \log \frac{p(y = 1|\mathbf{x})}{p(y = 2|\mathbf{x})} = \log \frac{\pi_1}{\pi_2} + \sum_{i=1}^p f(x_i) .$$

19 Meanwhile, besides the assumption of the distribution of  $(\mathbf{x}|y)$ , a fundamental  
 20 assumption underlying the naïve Bayes classifier is the conditional indepen-

1 dence amongst the  $p$  features within a class, so that the joint probability is  
 2  $p(\mathbf{x}|y) = \prod_{i=1}^p p(x_i|y)$ . It follows that the discriminant function  $\lambda(\alpha)$  is

$$\lambda_{\text{gen}}(\alpha) = \log \frac{p(y=1|\mathbf{x})}{p(y=2|\mathbf{x})} = \log \frac{\pi_1}{\pi_2} + \sum_{i=1}^p \log \frac{p(x_i|y=1)}{p(x_i|y=2)}.$$

3 It is clear, as pointed out by Rubinstein and Hastie (1997), that the naïve  
 4 Bayes classifier is a specialised case of a GAM, with  $f(x_i) = \log\{p(x_i|y=1)/p(x_i|y=2)\}$ .  
 5 Furthermore, GAMs may not necessarily assume conditional independence.

6 One sufficient condition that leads to another specialised case of a GAM (we  
 7 call it Q-GAM) is that  $p(\mathbf{x}|y) = q(\mathbf{x}) \prod_{i=1}^p q(x_i|y)$ , where  $q(\mathbf{x})$  is common  
 8 across the classes but cannot be further factorised into a product of func-  
 9 tions of individual features as  $\prod_{i=1}^p q(x_i)$ . In such a case, the assumption of  
 10 conditional independence between  $x_i|y$  and  $x_j|y$ ,  $i \neq j$ , is invalid but we still  
 11 have  $f(x_i) = \log\{q(x_i|y=1)/q(x_i|y=2)\}$ , where  $q(x_i|y)$  is different from the  
 12 marginal probability  $p(x_i|y)$  that is used by the naïve Bayes classifier.

13 In summary, considering the parity between  $\lambda_{\text{gen}}(\alpha)$  and  $\lambda_{\text{dis}}(\alpha)$  and thus that,  
 14 between two pairs, LDA- $\Sigma$  vs. linear logistic regression and Q-GAM vs. GAM  
 15 in terms of classification, neither classifier assumes conditional independence  
 16 of  $\mathbf{x}|y$  amongst the features within a class, which is an elementary assumption  
 17 underlying LDA- $\Lambda$  and the naïve Bayes classifier. Therefore, it may not be  
 18 reliable for any claim that is derived from the comparison between LDA- $\Lambda$  or  
 19 the naïve Bayes classifier and linear logistic regression to be generalised to all  
 20 the generative and discriminative classifiers.

## 1 Acknowledgements

The authors thank Andrew Y. Ng for communication about the implementation of the empirical studies in this paper.

## 2 References

- 3 Dawid, A. P., 1976. Properties of diagnostic data distributions. *Biometrics*  
4 32 (3), 647–658.
- 5 Efron, B., 1975. The efficiency of logistic regression compared to normal dis-  
6 criminant analysis. *Journal of the American Statistical Association* 70 (352),  
7 892–898.
- 8 Newman, D. J., Hettich, S., Blake, C. L., Merz, C. J., 1998. UCI  
9 Repository of machine learning databases. University of Cal-  
10 ifornia, Irvine, Dept. of Information and Computer Sciences,  
11 <http://www.ics.uci.edu/~mlearn/MLRepository.html>.
- 12 Ng, A. Y., Jordan, M. I., 2001. On discriminative vs. generative classifiers: a  
13 comparison of logistic regression and naive Bayes. In: *NIPS*. pp. 841–848.
- 14 O’Neill, T. J., 1980. The general distribution of the error rate of a classification  
15 procedure with application to logistic regression discrimination. *Journal of*  
16 *the American Statistical Association* 75 (369), 154–160.
- 17 Rubinstein, Y. D., Hastie, T., 1997. Discriminative vs. informative learning.  
18 In: *KDD*. pp. 49–53.
- 19 Titterton, D. M., Murray, G. D., Murray, L. S., Spiegelhalter, D. J., Skene,  
20 A. M., Habbema, J. D. F., Gelpke, G. J., 1981. Comparison of discrimi-  
21 nation techniques applied to a complex data set of head injured patients

- 1 (with discussion). Journal of the Royal Statistical Society. Series A (Gen-  
2 eral) 144 (2), 145–175.
- 3 Verboven, S., Hubert, M., 2005. LIBRA: a MATLAB library for robust analy-  
4 sis. Chemometrics and Intelligent Laboratory Systems 75 (2), 127–136.