#### **Avoiding Bias from Feature Selection**

Longhai Li

Based on a paper coauthored with Jianguo Zhang and Radford Neal

longhai@math.usask.ca
http://math.usask.ca/~longhai

Department of Mathematics and Statistics University of Saskatchewan Saskatoon, Saskatchewan, Canada

CRiSM Workshop "Bayesian Analysis of High-dimensional Data" University of Warwick, United Kingdom

16 April 2008

## **Outline of the talk**

- Why do We Need to Select Features?
- Problems with Feature Selection the optimistic bias
- Our Solution for Avoiding Bias from Feature Selection
- Application to Naive Bayes Classification Model
  - Definition of a Naive Bayes Classification Model
  - Computation of Adjustment Factor for Avoiding Bias
  - Demonstration with Simulated and Real Datasets
- Conclusions and Discussions
- Future Work

• Goal: Given a feature vector x, we want to predict the associated response y, i.e. find a predictive function C from x to y:

 $\hat{y} = C(x).$ 

• Goal: Given a feature vector x, we want to predict the associated response y, i.e. find a predictive function C from x to y:

$$\hat{y} = C(x).$$

• How do we find  $\hat{y}$ ?

• Goal: Given a feature vector x, we want to predict the associated response y, i.e. find a predictive function C from x to y:

$$\hat{y} = C(x).$$

- How do we find  $\hat{y}$ ?
  - First find the predictive distribution of y given x:

P(y|x)

• Goal: Given a feature vector x, we want to predict the associated response y, i.e. find a predictive function C from x to y:

$$\hat{y} = C(x).$$

- How do we find  $\hat{y}$ ?
  - First find the predictive distribution of y given x:

P(y|x)

• Given x, predict y with  $\hat{y}$  that minimizes expected loss  $E(L(y, \hat{y})|x)$ . e.g., y is binary, and  $L(y = 0, \hat{y} = 1)/L(y = 1, \hat{y} = 0) = r$ ,

$$\hat{y} = \begin{cases} 1 & \text{if } P(y=1|x) \ge 1 - \frac{1}{1+r}, \\ 0 & \text{if } P(y=1|x) < 1 - \frac{1}{1+r} \end{cases}$$

• Goal: Given a feature vector x, we want to predict the associated response y, i.e. find a predictive function C from x to y:

$$\hat{y} = C(x).$$

- How do we find  $\hat{y}$ ?
  - First find the predictive distribution of y given x:

P(y|x)

• Given x, predict y with  $\hat{y}$  that minimizes expected loss  $E(L(y, \hat{y})|x)$ . e.g., y is binary, and  $L(y = 0, \hat{y} = 1)/L(y = 1, \hat{y} = 0) = r$ ,

$$\hat{y} = \begin{cases} 1 & \text{if } P(y=1|x) \ge 1 - \frac{1}{1+r}, \\ 0 & \text{if } P(y=1|x) < 1 - \frac{1}{1+r} \end{cases}$$

• Statistical Method: Estimate P(y|x) by learning from the available data  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ , collectively  $\{x^{\text{train}}, y^{\text{train}}\}$ , called training data.

• Goal: Given a feature vector x, we want to predict the associated response y, i.e. find a predictive function C from x to y:

$$\hat{y} = C(x).$$

- How do we find  $\hat{y}$ ?
  - First find the predictive distribution of y given x:

P(y|x)

• Given x, predict y with  $\hat{y}$  that minimizes expected loss  $E(L(y, \hat{y})|x)$ . e.g., y is binary, and  $L(y = 0, \hat{y} = 1)/L(y = 1, \hat{y} = 0) = r$ ,

$$\hat{y} = \begin{cases} 1 & \text{if } P(y=1|x) \ge 1 - \frac{1}{1+r}, \\ 0 & \text{if } P(y=1|x) < 1 - \frac{1}{1+r} \end{cases}$$

- Statistical Method: Estimate P(y|x) by learning from the available data  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ , collectively  $\{x^{\text{train}}, y^{\text{train}}\}$ , called training data.
- Example: Given gene expression data of a patient, classify the type of tumour.

• Examples of High Dimensional Features

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.
- Difficulties with High Dimensional Features:

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.
- Difficulties with High Dimensional Features:
  - In Time: Computation time grows with the number of features.

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.
- Difficulties with High Dimensional Features:
  - In Time: Computation time grows with the number of features.
  - In Money: Costly in measuring the features for future cases.

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.
- Difficulties with High Dimensional Features:
  - In Time: Computation time grows with the number of features.
  - In Money: Costly in measuring the features for future cases.
  - In Statistics:

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.
- Difficulties with High Dimensional Features:
  - In Time: Computation time grows with the number of features.
  - In Money: Costly in measuring the features for future cases.
  - In Statistics:
    - When the number of observations is smaller than the number of parameters, the likelihood function  $P(x^{\text{train}}, y^{\text{train}} | \theta)$  based on a simple model, such as a linear model, will favour "too good"  $\theta$ .

- Examples of High Dimensional Features
  - Gene Expression Data: Measure the expression levels of tens of thousands of genes.
  - Document Classification: Count the frequency (times of occurrence) of all words in a large dictionary, which may have, for example, 30,000 words.
- Difficulties with High Dimensional Features:
  - In Time: Computation time grows with the number of features.
  - In Money: Costly in measuring the features for future cases.
  - In Statistics:
    - When the number of observations is smaller than the number of parameters, the likelihood function  $P(x^{\text{train}}, y^{\text{train}} | \theta)$  based on a simple model, such as a linear model, will favour "too good"  $\theta$ .
    - Our models may fail when applied to high-dimensional features, e.g., we may not be able to model the dependency between features very well; our priors for high-dimensional parameters may be inappropriate.

### **Optimistic Bias from Feature Selection**

Feature Selection:

For previous reasons or others, one may like to select a subset of features to use based on some measure of the dependency (e.g. absolute correlation) between the features and y. There are also more sophisticated Bayesian/non-Bayesian methods by looking at how well the subset can predict y (e.g.,LASSO, Bayesian methods using a prior with mass at 0 for coefficients).

### **Optimistic Bias from Feature Selection**

Feature Selection:

For previous reasons or others, one may like to select a subset of features to use based on some measure of the dependency (e.g. absolute correlation) between the features and y. There are also more sophisticated Bayesian/non-Bayesian methods by looking at how well the subset can predict y (e.g.,LASSO, Bayesian methods using a prior with mass at 0 for coefficients).

Optimistic Bias:

However, such procedures will make the relationship between y and x appears stronger than it actually is, i.e., the response appears more predictable from the features. We will be overconfident with our prediction.

## **Optimistic Bias from Feature Selection**

Feature Selection:

For previous reasons or others, one may like to select a subset of features to use based on some measure of the dependency (e.g. absolute correlation) between the features and y. There are also more sophisticated Bayesian/non-Bayesian methods by looking at how well the subset can predict y (e.g.,LASSO, Bayesian methods using a prior with mass at 0 for coefficients).

Optimistic Bias:

However, such procedures will make the relationship between y and x appears stronger than it actually is, i.e., the response appears more predictable from the features. We will be overconfident with our prediction.

An Extreme Example:

All features  $x_1, \ldots, x_p$  and the response y are actually independent. We select only k features  $x_1^*, \ldots, x_k^*$  using the criterion of absolute correlation. These k features will exhibit strong relationship with y in training data, when  $p \gg k$ .

### What Causes Bias from Feature Selection?

Let's look at the absolute correlations of 10,000 binary features with a binary response from a data set simulated using a method to be discussed later.



# **Problems Resulting from Feature Selection Bias (I)**

The following graph displays a typical plot of predictive probabilities based on selected features versus the actual predictive probabilities based on all information available:



# **Problems Resulting from Feature Selection Bias (I)**

The following graph displays a typical plot of predictive probabilities based on selected features versus the actual predictive probabilities based on all information available:



Predictive probabilities based on selected features

The predictive probabilities are overconfident. For example, if we use  $\frac{1}{2}$  as cutoff for predicting the class labels, expected error rate,

$$\frac{1}{n} \sum_{i=1}^{n} \left( \hat{p}_{s}^{(i)} I(\hat{p}_{s}^{(i)} < \frac{1}{2}) + (1 - \hat{p}_{s}^{(i)}) I(\hat{p}_{s}^{(i)} \ge \frac{1}{2}) \right),$$

is smaller than actual error rate

$$\frac{1}{n} \sum_{i=1}^{n} I(y^{(i)} \neq \hat{y}_s^{(i)})$$

# **Problems Resulting from Feature Selection Bias (II)**

More seriously, when the loss by misclassifying 0 to 1 is different from the loss by misclassifying 1 to 0, the prediction based on selected features results in more loss in practice than that based on all information available, i.e.,



$$E(L(\hat{y}_{s}^{(*)}, y^{(*)})) \ge E(L(\hat{y}_{B}^{(*)}, y^{(*)}))$$

where,

$$\hat{y}_s^{(*)} = \begin{cases} 1 & \text{if } \hat{p}_s^{(*)} \ge \text{cutoff} \\ 0 & \text{if } \hat{p}_s^{(*)} < \text{cutoff} \end{cases},$$

and cutoff =  $1 - \frac{1}{1+r}$ , and  $r = \frac{L_{0 \to 1}}{L_{1 \to 0}}$ .  $\hat{y}_B^{(*)}$  is similar, but using the predictive probability  $\hat{p}_B^{(*)}$  based on all information available.

### **Our Method for Avoiding Features Selection Bias**

• All available information: Our predictions should condition not only on the retained features  $x_{1:k}^{\text{train}}$ , but also on the fact that the other p-k features have sample correlation with the response that is less than  $\gamma$  in absolute value:

 $|y^{\text{train}}, x_{1:k}^{\text{train}}, |\text{COR}(y^{\text{train}}, x_t^{\text{train}})| \leq \gamma \text{ for } t = k+1, \dots, p$ 

### **Our Method for Avoiding Features Selection Bias**

• All available information: Our predictions should condition not only on the retained features  $x_{1:k}^{\text{train}}$ , but also on the fact that the other p-k features have sample correlation with the response that is less than  $\gamma$  in absolute value:

 $|y^{\text{train}}, x_{1:k}^{\text{train}}, |\text{COR}(y^{\text{train}}, x_t^{\text{train}})| \leq \gamma \text{ for } t = k+1, \dots, p$ 

 Models: The response and the predictors are modeled jointly. Given the response *y* and a model parameter *α*, the features *x*<sub>1</sub>,..., *x<sub>p</sub>*, are modeled to be independent and has identical distribution:

$$P(x_1, \cdots, x_p | y, \alpha) = \prod_{t=1}^p \left[ P(x_t | y, \alpha) \right]$$

### **Our Method for Avoiding Features Selection Bias**

• All available information: Our predictions should condition not only on the retained features  $x_{1:k}^{\text{train}}$ , but also on the fact that the other p-k features have sample correlation with the response that is less than  $\gamma$  in absolute value:

 $|y^{\text{train}}, x_{1:k}^{\text{train}}, |\text{COR}(y^{\text{train}}, x_t^{\text{train}})| \leq \gamma \text{ for } t = k+1, \dots, p$ 

Models: The response and the predictors are modeled jointly. Given the response *y* and a model parameter *α*, the features *x*<sub>1</sub>,..., *x*<sub>p</sub>, are modeled to be independent and has identical distribution:

$$P(x_1, \cdots, x_p | y, \alpha) = \prod_{t=1}^p \left[ P(x_t | y, \alpha) \right]$$

• Adjustment factor: The likelihood function of  $\alpha$  based on  $y^{\text{train}}, x_{1:k}^{\text{train}}$  is multiplied by:

$$P(|\mathsf{COR}(y^{\mathsf{train}}, x_t^{\mathsf{train}})| \le \gamma \text{ for } t = k+1, \dots, p|\alpha, y^{\mathsf{train}})$$
$$= \left[P(|\mathsf{COR}(y^{\mathsf{train}}, x_t^{\mathsf{train}})| \le \gamma |\alpha, y^{\mathsf{train}})\right]^{p-k}$$

#### A Bayesian Naive Bayes Model for Binary Data



#### **Computation of the adjustment factor (I)**

 $COR(x_t^{train}, y^{train})$  can be written in terms of  $I_0 = \sum_{i=1}^n I(y^{(i)} = 0, x_t^{(i)} = 1)$  and  $I_1 = \sum_{i=1}^n I(y^{(i)} = 1, x_t^{(i)} = 1)$ :

$$\mathsf{COR}(x_t^{\mathsf{train}}, y^{\mathsf{train}}) = \frac{(0 - \overline{y}) I_0 + (1 - \overline{y}) I_1}{\sqrt{n\overline{y}(1 - \overline{y})} \sqrt{I_0 + I_1 - (I_0 + I_1)^2/n}}$$

 $I_0, I_1$  are visualized:



### **Computation of the adjustment factor (II)**

The following graph displays the values of  $|COR(x_t^{train}, y^{train})|$  in terms of  $I_0$  and  $I_1$  for a particular example, with the threshold of magnitude of correlation in selecting features is  $\gamma = 0.2$ :

	14	+1.00	+0.90	+0.81	+0.72	+0.62	+0.53	+0.42	+0.29	0.00
	13	+0.91	+0.80	+0.70	+0.60	+0.49	+0.38	+0.25	+0.09	-0.16
	12	+0.83	+0.72	+0.61	+0.50	+0.39	+0.27	+0.13	-0.03	-0.24
	11	+0.76	+0.64	++0.52	+0.41	+0.30	+0.17	+0.04	-0.11	-0.30
	10	+0.69	+0.57	+0.45	+0.33	+0.21	+0.09	-0.04	-0.18	-0.36
	9	+0.63	+0.50	+0.38	+0.26	+0.14	+0.02	-0.11	-0.25	-0.41
	8	+0.57	+0.44	+0.31	+0.19	+0.07	-0.05	-0.18	-0.31	-0.46
1	7	+0.52	+0.38	+0.24	+0.12	0.00	-0.12	-0.24	-0.37	-0.52
	6	+0.46	+0.31	+0.18	+0.05	-0.07	-0.19	-0.31	-0.44	-0.57
	5	+0.41	+0.25	+0.11	-0.02	-0.14	-0.26	-0.38	-0.50	-0.63
	4	+0.36	+0.18	+0.04	-0.09	-0.21	-0.33	-0.45	<u>_</u> 0.57	-0.69
	3	+0.30	+0.11	-0.04	-0.17	-0.30	-0.41	-0.52	-0.64	-0.76
	2	+0.24	+0.03	-0.13	-0.27	-0.39	-0.50	-0.61	-0.72	-0.83
	1	+0.16	-0.09	-0.25	-0.38	-0.49	-0.60	-0.70	-0.80	-0.91
	0	0.00	-0.29	-0.42	-0.53	-0.62	-0.72	-0.81	-0.90	-1.00
		0	1	2	3	4	5	6	7	8
						I <sub>0</sub>				

Ι

#### **Computation of the adjustment factor (III)**

- $P(I_0, I_1 \mid \alpha, y^{\text{train}})$  is symmetric for  $H_+$  and  $H_-$ , so  $P(|\mathsf{COR}(x_t^{\text{train}}, y^{\text{train}})| \leq \gamma \mid \alpha, y^{\text{train}}) = 1 - 2 \sum_{(I_0, I_1) \in H_+} P(I_0, I_1 \mid \alpha, y^{\text{train}})$
- Conditioning on  $\theta_t$ :

$$\sum_{(I_0,I_1)\in H_+} P(I_0, I_1 \mid \alpha, y^{\text{train}}) = \int_0^1 \sum_{(I_0,I_1)\in H_+} P(I_0, I_1 \mid \alpha, \theta_t, y^{\text{train}}) d\theta_t$$

•  $|\mathsf{COR}(x_t^{\text{train}}, y^{\text{train}})|$  is monotone with respect to  $I_0$  or  $I_1$ , so

$$\sum_{(I_0,I_1)\in H_+} P(I_0, I_1 \mid \alpha, \theta_t, y^{\text{train}}) = \sum_{I_1=b_0}^{n\overline{y}} \sum_{I_0=0}^{r_{I_1}} P(I_0, I_1 \mid \alpha, \theta_t, y^{\text{train}})$$

•  $I_0$  and  $I_1$  are independent given  $lpha, heta_t, y^{ ext{train}}$ , so

$$P(I_0, I_1 \mid \alpha, \theta_t, y^{\text{train}}) = P(I_1 \mid \alpha, \theta_t, y^{\text{train}})P(I_0 \mid \alpha, \theta_t, y^{\text{train}})$$

#### **Computation of the adjustment factor (IV)**

Finally, we can easily compute probability of  $I_0$  and  $I_1$ :

$$P(I_1 \mid \alpha, \theta_t, y^{\text{train}}) = \binom{N_1}{I_1} U(\alpha \theta_t, \alpha(1-\theta_t), I_1, N_1 - I_1)$$

and

$$P(I_0 \mid \alpha, \theta_t, y^{\text{train}}) = \binom{n - N_1}{I_0} U(\alpha \theta_t, \alpha (1 - \theta_t), I_0, n - N_1 - I_0)$$

## **A Simulation Experiment**

- Generating data:  $\alpha = 300$ , p = 10,000, 100 training cases, 2000 test cases
- Selecting features: 4 subsets with only 1, 10, 100 and 1000 features were selected, with smallest absolute value of correlation being 0.36,0.27,0.21, and 0.13.
- Prior:  $f_0 = f_1 = 1$ , a = 0.5, b = 5

#### **Plot of Predictive Probabilities**



#### **Calibration of Predictive Probabilities**

	100 features selected out of 10000						
		Corre	cted	Uncorrected			
Category	#	Pred	Actual	#	Pred	Actual	
0.0 - 0.1	155	0.067	0.077	717	0.017	0.199	
0.1 – 0.2	247	0.151	0.162	133	0.150	0.391	
0.2 – 0.3	220	0.247	0.286	70	0.251	0.429	
0.3 – 0.4	225	0.352	0.356	68	0.351	0.515	
0.4 – 0.5	237	0.450	0.494	58	0.451	0.500	
0.5 – 0.6	227	0.545	0.586	78	0.552	0.603	
0.6 – 0.7	202	0.650	0.728	77	0.654	0.532	
0.7 – 0.8	214	0.749	0.785	80	0.746	0.662	
0.8 – 0.9	182	0.847	0.857	98	0.852	0.633	
0.9 – 1.0	91	0.935	0.923	621	0.979	0.818	

#### **Expected versus Actual Error Rates**



### **Comparison of Average Loss**



#### **Comparison of Average Minus the Log Probabilities**

$$\mathsf{AMLP} = -\frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right)$$



## **Posterior Distribution of** $log(\alpha)$

Red: Uncorrected Green: Corrected Black: Based on all features

1 feature selected out of 10000

10 features selected out of 10000











10

9

# **Computation Times**

# of Features Selected	1	10	100	1000	10000
Uncorrected	12.66	28.80	203.60	2065.73	20627.71
Corrected	12.72	29.55	204.54	2076.21	

### **A Test with Gene Expression Data**

- Data set: 2000 genes, 62 cases (40 Cancerous vs 22 Normal tissues)
- Converted into binary data by thresholding at medians of features
- Split 2000 genes randomly into 10 subsets each with 200 genes.
- 5 Genes were selected out of 200 genes
- Used leave-one-out cross-validation to obtain predictive probabilities
- Prior the same as previous simulation experiment

#### **Plot of Predictive Probabilities**



#### **Expected versus Actual Error Rates**



#### **Comparison of Average Loss**



#### **Comparison of Average Minus the Log Probabilities**

$$\mathsf{AMLP} = -\frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right)$$



• We propose a correction method to avoid the bias from feature selection.

- We propose a correction method to avoid the bias from feature selection.
- We've applied the method to binary naive Bayes models. The simulation results show that it does avoid the bias from feature selection and it is faster than using all features. The results from microarray datasets show that the corrected method improves the predictive performance.

- We propose a correction method to avoid the bias from feature selection.
- We've applied the method to binary naive Bayes models. The simulation results show that it does avoid the bias from feature selection and it is faster than using all features. The results from microarray datasets show that the corrected method improves the predictive performance.
- We have devised a fast method for computing it in binary naive Bayes models. This method can be generalized to all discrete naive Bayes models. However, the computation of the adjustment factor is difficult generally.

- We propose a correction method to avoid the bias from feature selection.
- We've applied the method to binary naive Bayes models. The simulation results show that it does avoid the bias from feature selection and it is faster than using all features. The results from microarray datasets show that the corrected method improves the predictive performance.
- We have devised a fast method for computing it in binary naive Bayes models. This method can be generalized to all discrete naive Bayes models. However, the computation of the adjustment factor is difficult generally.
- When there is enough data, we better use different subsets to select features and fit the models, if we do not know how to correct for the feature selection bias.

- We propose a correction method to avoid the bias from feature selection.
- We've applied the method to binary naive Bayes models. The simulation results show that it does avoid the bias from feature selection and it is faster than using all features. The results from microarray datasets show that the corrected method improves the predictive performance.
- We have devised a fast method for computing it in binary naive Bayes models. This method can be generalized to all discrete naive Bayes models. However, the computation of the adjustment factor is difficult generally.
- When there is enough data, we better use different subsets to select features and fit the models, if we do not know how to correct for the feature selection bias.
- More generally, if we do not know how to correct for the feature selection bias, we better avoid doing feature selection.

### **Future Work**

• For real data sets, such as gene expression data, the independence assumption may not be true. We must apply the bias correction method to the independent factors governing the observed data. Some simple methods, such as principal component analysis, may be appropriate. Further work will be done to find a good way to do this.

## **Future Work**

- For real data sets, such as gene expression data, the independence assumption may not be true. We must apply the bias correction method to the independent factors governing the observed data. Some simple methods, such as principal component analysis, may be appropriate. Further work will be done to find a good way to do this.
- One could pretty easily extend the method to other simple models (e.g. Gaussian models for continuous data) and other selection criteria (e.g. *t* statistic).

## **Future Work**

- For real data sets, such as gene expression data, the independence assumption may not be true. We must apply the bias correction method to the independent factors governing the observed data. Some simple methods, such as principal component analysis, may be appropriate. Further work will be done to find a good way to do this.
- One could pretty easily extend the method to other simple models (e.g. Gaussian models for continuous data) and other selection criteria (e.g. *t* statistic).
- We've applied the correction method to binary mixture models and factor analysis models. Computation of the adjustment factor for these two models is more difficult but still feasible. Further work is necessary to improve efficiency of computing adjustment factor for such models, since we need to perform a large amount of this computation (e.g. in each iteration of Markov chain sampling). Some numerical methods, such as using importance sampling framework, are very promising.

## Thank You! Questions and Comments?

References:

[1] Li, L., Zhang, J., and Neal, R. M. (2008), A Method for Avoiding Bias from Feature Selection with Application to Naive Bayes Classification Models, *Bayesian Analysis*, volume 3, number 1, pp 171-196.

[2] Li, L., (2007), *Bayesian Classification and Regression with High Dimensional Features*, Ph.D. thesis, University of Toronto.