



Frequentist Approximation of the Bayesian Posterior Inclusion Probability by Stochastic Subsampling

V. Ly¹ and E. Fokoué^{2*}

¹*M & T Bank, Buffalo, NY, USA.*

²*School of Mathematical Sciences, Rochester Institute of Technology, 98 Lomb Memorial Drive, Rochester, NY 14623, USA.*

Authors' contributions

This work was carried out in collaboration between both authors. Author VL set up and methodically ran all the computer simulations, while author EF provided the theoretical insights along with the write up of the technical derivations. Both authors read and approved the final manuscript.

Article Information

DOI: 10.9734/BJMCS/2016/27023

Editor(s):

(1) Andrej V. Plotnikov, Department of Applied and Calculus Mathematics and CAD, Odessa State Academy of Civil Engineering and Architecture, Ukraine.

Reviewers:

(1) Thomas L. Toulas, Technological Educational Institute of Athens, Greece.

(2) Louis Asiedu, University of Ghana, Ghana.

Complete Peer review History: <http://www.sciencedomain.org/review-history/16212>

Received: 16th May 2016

Accepted: 28th August 2016

Published: 16th September 2016

Original Research Article

Abstract

This paper uses stochastic subsampling of the dataset to provide a frequentist approximation to what is known in the Bayesian framework as the posterior inclusion probability (PIP). The distinct merit of this contribution lies in the fact that it makes it easier for typically non-Bayesian-minded practitioners, of which there are many, to relate to the way the Bayesian paradigm allows the computation of the nicely interpretable variable importance. Despite its computationally intensive nature, due to the need to fitting a very large number of models, the proposed approach is readily applicable to both classification and regression tasks, and can be done in comparatively competitive computational times thanks to the availability of parallel computing facilities through cloud and cluster computing. Finally, the scheme proposed is very general and can therefore be easily adapted to all kinds of statistical prediction tasks. Application of the proposed method to some very famous benchmark datasets shows that it mimics the Bayesian counterpart quite well in the important context of variable selection.

*Corresponding author: E-mail: epfeqa@rit.edu;

Keywords: Generalized linear model; bayesian variable selection; posterior inclusion probability; stochastic subsampling; binary classification; prediction; logistic; logit; probit, Cauchit.

2010 Mathematics Subject Classification: 62H30, 62H25.

1 Introduction

Modern statistical machine learning is replete with thousands of studies where the main statistical task revolves around estimations and predictions based on the traditional generalized linear model (GLM) given by

$$g(\mathbb{E}[\mathbf{Y}|\mathbf{X}, \boldsymbol{\beta}, \alpha]) = \alpha \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta}, \quad (1)$$

where $g(\cdot)$ is the so-called link function, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$, $\mathbf{Y} = (y_1, \dots, y_n)^\top$, the design matrix \mathbf{X} is an $n \times p$ matrix, and $\mathbf{1}_n = (1, 1, \dots, 1)^\top$ is a $n \times 1$ dimensional vector of 1's. We shall refer to (1) as the *full model*. We assume that many of the β_j 's are essentially zero, so that the intrinsic rank of the design matrix \mathbf{X} is a number $q \in \mathbb{N}$ with $q \ll p$. Many data mining problems do exhibit such a characteristic of rank deficiency, mainly because variables are typically pick up as they are available, and therefore will turn out to be either noise variable (no relationship with the response) or redundant variables. Let \mathbf{I}_n denote the $n \times n$ identity matrix. A basic result in GLM analysis shows that when \mathbf{X} is rank deficient, and $g(\cdot)$ is the identity function, and the density of \mathbf{Y} is gaussian, ie, $[\mathbf{Y}|\mathbf{X}, \boldsymbol{\beta}, \alpha] \sim N_n(\alpha \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$, then the ordinary least squares estimator

$$\hat{\boldsymbol{\beta}}^{(\text{OLS})} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (2)$$

of $\boldsymbol{\beta}$ will tend to exhibit a high (inflated) variance, thereby corrupting all predictions and inferences with the computed model. It is therefore crucial to determine (if possible) the intrinsic model that generated the data, i.e. the model made up of only the most significant and non redundant variables. For many decades, both frequentist and Bayesian statisticians have contributed substantially to this topic of *variable selection*. In elementary statistical regression analysis courses, the method of choice for variable selection has been overwhelmingly frequentist with *stepwise regression heuristic* occupying a prominent place, and *best subsets selection* occasionally used whenever possible. While a heuristic like stepwise regression does provide a workable approach to variable selection, it is not a principled method, and does have the extra limitation of not providing any measure of variable importance. In recent years, both Bayesians and non-Bayesians have developed new methods for handling some of the most formidable variable selection tasks, many of which arose from the statistical learning and data mining community.

2 Bayesian Approach to Variable Selection

The vast majority of Bayesian contributions to variable selection of late have concentrated on the use of conjugate prior, with the typical choice of prior on $\boldsymbol{\beta}$ being a Gaussian prior of the form

$$[\boldsymbol{\beta}|\sigma^2, \mathbf{W}] \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbf{W}^{-1}), \quad (3)$$

where \mathbf{W} is the prior precision matrix. In other words, we have

$$p(\boldsymbol{\beta}|\sigma^2, \mathbf{W}) = \frac{1}{\sqrt{((2\pi)\sigma^2)^p |\mathbf{W}|}} \exp \left\{ -\frac{1}{2\sigma^2} \boldsymbol{\beta}^\top \mathbf{W} \boldsymbol{\beta} \right\}. \quad (4)$$

Of course, the use of a zero mean prior expresses the assumption of many insignificant coefficients. However, even more important is the use of a vector of indicator variables that ultimately provides a mechanism (device) for performing variable selection. One of the key building blocks of the

Bayesian variable selection machinery is the use of a vector of indicator variables. With the p original predictor variables, there are $2^p - 1$ non empty models corresponding each to a subset of the provided variables. We shall use a vector $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \dots, \gamma_p)^\top$ to denote the index of a given model, with each γ_j being an indicator of the variable's presence in the model under consideration, namely

$$\gamma_j = \begin{cases} 1 & \text{If variable } X_j \text{ appears in the model} \\ 0 & \text{Otherwise} \end{cases} \quad (5)$$

Clearly, $\boldsymbol{\gamma} = (1, \dots, 1)^\top$ corresponds to the *full model* M_f , while $\boldsymbol{\gamma} = (0, \dots, 0)^\top$ corresponds to the empty model also referred to as the *null model*, and given by

$$M_n : \quad g(\mathbb{E}[\mathbf{Y}|\alpha]) = \alpha \mathbf{1}_n. \quad (6)$$

Equipped with this index, $p_\gamma = \sum_{j=1}^p \gamma_j$ is the number of predictor variables in model M_γ , and $\boldsymbol{\beta}_\gamma$ is the subset of $\boldsymbol{\beta}$ made up of only the β_j 's picked up by $\boldsymbol{\gamma}$. Finally, \mathbf{X}_γ is the submatrix of \mathbf{X} whose columns are only those p_γ columns of \mathbf{X} picked up by $\boldsymbol{\gamma}$, so that \mathbf{X}_γ is really an $n \times p_\gamma$ matrix, and the corresponding model M_γ is given by

$$M_\gamma : \quad g(\mathbb{E}[\mathbf{Y}|\mathbf{X}_\gamma, \boldsymbol{\beta}_\gamma, \alpha]) = \alpha \mathbf{1}_n + \mathbf{X}_\gamma \boldsymbol{\beta}_\gamma. \quad (7)$$

For the normal linear model, we have $[\mathbf{Y}|\alpha, \boldsymbol{\beta}_\gamma, \sigma^2, M_\gamma] \sim N_n(\alpha \mathbf{1}_n + \mathbf{X}_\gamma \boldsymbol{\beta}_\gamma, \sigma^2 \mathbf{I}_n)$, which means that

$$p(\mathbf{Y}|\boldsymbol{\theta}_\gamma, M_\gamma) = \frac{1}{\sqrt{((2\pi)\sigma^2)^n}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}_\gamma \boldsymbol{\beta}_\gamma)^\top (\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}_\gamma \boldsymbol{\beta}_\gamma) \right\}$$

where $\boldsymbol{\theta}_\gamma = \{\alpha, \boldsymbol{\beta}_\gamma, \sigma^2\}$. When it comes to Bayesian variable selection, arguably the most crucial ingredient is the posterior density of a given model, given by

$$p(M_\gamma|\mathbf{y}) = \frac{p(\mathbf{y}|M_\gamma)p(M_\gamma)}{\sum_{\boldsymbol{\gamma} \in \Gamma} p(\mathbf{y}|M_\gamma)p(M_\gamma)}$$

where $\Gamma = \{0, 1\}^p$ and $p(\mathbf{y}|M_\gamma)$ is the marginal density of the data, also referred to as the marginal likelihood of model M_γ , and defined by

$$p(\mathbf{y}|M_\gamma) = \int_{\Theta} p(\mathbf{y}|\boldsymbol{\theta}_\gamma, M_\gamma) p(\boldsymbol{\theta}_\gamma|M_\gamma) d\boldsymbol{\theta}_\gamma.$$

In some special cases, it is possible to derive closed-form (analytical) expressions for $p(\mathbf{y}|M_\gamma)$, but in general, it must be approximated using a variety of schemes. The posterior probability $p(M_\gamma|\mathbf{y})$ of model M_γ , plays a central role in Bayesian learning.

$$p(\mathbf{z}|\mathbf{y}) = \sum_{\boldsymbol{\gamma} \in \Gamma} p(\mathbf{z}|M_\gamma, \mathbf{y}) p(M_\gamma|\mathbf{y})$$

and also

$$\mathbb{E}(\mathbf{z}|\mathbf{y}) = \sum_{\boldsymbol{\gamma} \in \Gamma} \mathbb{E}(\mathbf{z}|M_\gamma, \mathbf{y}) p(M_\gamma|\mathbf{y})$$

Among Bayesian statisticians, there are those who suggest that when it comes to model selection, one must choose the model with the highest posterior density model, i.e.,

$$\boldsymbol{\gamma}_{\text{HPM}} = \underset{\boldsymbol{\gamma} \in \Gamma}{\operatorname{argmax}} \{p(M_\gamma|\mathbf{y})\}$$

[1] have suggested selecting instead the so-called median probability model (MPM) given $\boldsymbol{\gamma}_{\text{MPM}}$, such that

$$[\gamma_j]_{\text{MPM}} = \begin{cases} 1 & \text{if } \pi_j \equiv \Pr[\gamma_j = 1|\mathbf{y}] \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

where the posterior inclusion probability is given by

$$\text{PIP}_j = \pi_j = \Pr[\gamma_j = 1|\mathbf{y}] = \sum_{\gamma \in \Gamma} 1(\gamma_j = 1)p(M_\gamma|\mathbf{y}), \quad (9)$$

In practice, the estimated posterior inclusion probability is given by

$$\widehat{\text{PIP}}_j = \hat{\pi}_j = \sum_{\gamma:\gamma_j=1} p(M_\gamma|\mathbf{y}), \quad (10)$$

which means that the importance of a variable is measured in terms of its relative frequency of appearance in models. In Equation (10), it is crucial to be able to compute the posterior density of a given model. Empirically, this means that it is able to generate, at least a representative subset (sample) of all the models, and then to compute estimates of the posterior density from them, [2].

3 Frequentist Approximation of the Bayesian PIP

It turns out that the theoretical derivation of the Bayesian Posterior Inclusion Probability defined in (8) is technically very complicated. In fact, the estimation indicated in (10) often requires sophisticated Markov Chain Monte Carlo (MCMC) simulations. Fortunately, from its definition, it turns out that the Bayesian PIP can be approximated quite well by frequentist stochastic subsampling, an alternative that is both viable, desirable and feasible thanks to the availability of computing resources. Besides, we demonstrate and argue later that the advantage of this approximation lies in the fact that any practitioner who understands GLM can readily implement the idea of this paper and gain great benefits. Our frequentist's alternative to the Bayesian PIP proceeds by building M different models based on random subsamples of the data. For each one of the M models, the corresponding indicator vector is estimated by simply identifying the significant variables based on their P-values or T-values, or (expensively) on heuristic like stepwise regression when the computing resources allow it. After the M models are built, one should consider a matrix of the form

$$\hat{\mathbf{\Gamma}} = \begin{bmatrix} \hat{\gamma}_{11} & \hat{\gamma}_{12} & \dots & \hat{\gamma}_{1j} & \dots & \hat{\gamma}_{1p} \\ \hat{\gamma}_{21} & \hat{\gamma}_{22} & \dots & \hat{\gamma}_{2j} & \dots & \hat{\gamma}_{2p} \\ \vdots & \vdots & \ddots & \ddots & \dots & \vdots \\ \hat{\gamma}_{m1} & \hat{\gamma}_{m2} & \dots & \hat{\gamma}_{mj} & \dots & \hat{\gamma}_{mp} \\ \vdots & \vdots & \ddots & \ddots & \dots & \vdots \\ \hat{\gamma}_{M1} & \hat{\gamma}_{M2} & \dots & \hat{\gamma}_{Mj} & \dots & \hat{\gamma}_{Mp} \end{bmatrix}, \quad (11)$$

where $\hat{\gamma}_{mj}$ is set to 1 if variable X_j appeared in model m or was significant in model m , and zero otherwise. In other words, $\hat{\gamma}_{mj} \in \{0, 1\}$. Once the matrix $\hat{\mathbf{\Gamma}} \in \{0, 1\}^{M \times p}$ is built, the posterior inclusion probability $\text{PIP}_j = \pi_j$ of variable X_j can be approximated using

$$\hat{\pi}_j = \frac{1}{M} \sum_{m=1}^M \hat{\gamma}_{mj}.$$

Finally, if $\hat{\gamma}^{\text{FREQ}} = (\hat{\gamma}_1^{\text{FREQ}}, \dots, \hat{\gamma}_p^{\text{FREQ}})$ denotes the indicator vector, our final approximating model, where

$$\hat{\gamma}_j^{\text{FREQ}} = 1 \text{ if } \hat{\pi}_j > 0.5.$$

4 Computational Demonstrations

4.1 Computational demonstration on pattern recognition

According to Microsoft, spam is a term used to classify unwanted email. Spam may contain viruses or other malicious programs that can harm a computer. Furthermore, spam may be used as scams to acquire vital personal information such as credit card accounts, bank accounts, social security numbers, etc. Spam filters have been developed as preventative measures to protect the end user from ever opening these emails. Text categorization is one of the several techniques used to create spam filters. A number of terms are identified as indicators of spam/non-spam from a training set of emails. In the simplest spam filters, the frequencies of these terms in an email are determined and used to flag emails as spam/non-spam, [3]. The purpose of this subsection is threefold:

1. Build generalized linear models with a variety of link functions in the binary (Bernoulli) family, and use the built model to accurately and precisely classify emails as spam or nonspam.
2. Present a frequentist alternative to Bayesian Posterior Inclusion Probability (PIP) for variable selection.
3. Compare logistic regression classification accuracy with accuracies of newer machine learning algorithms.

The data used for this pattern recognition analysis came from the ubiquitous Spambase dataset at the UCI Machine Learning Repository. This dataset was donated by [4] from Hewlett-Packard Labs. The collection of emails was provided from [4]'s email account at Hewlett-Packard. The dataset consisted of 4,601 observations with 57 explanatory variables and 1 response variable. 54 of the 57 explanatory variables measure the percentage of an email in which specific words or characters appear. The remaining 3 explanatory variables measure the average length of uninterrupted sequences of capital letters, length of longest uninterrupted sequence of capital letters, and the total number of capital letters in an email. The response variable was binary coded 1 for spam and 0 for non-spam. Out of the 4,601 observations, 2,788 emails were non-spam (60.60%) and 1,813 emails were categorized as spam (39.40%). The Spambase dataset webpage on the UCI Machine Learning Repository cited an average $\sim 7\%$ misclassification error. The goal of this project is to generate an alternative classification model using regression techniques.

The figure below shows the linear pairwise correlations between all of the explanatory variables. Based on the correlation plot, there is potential for multicollinearity to affect regression results. Due to the multicollinearity, variable selection will be required during the model building phase. Even with variable selection, there are still some questions that remain unanswered. What is the optimal model size to achieve model complexity-testing accuracy tradeoff? How much confidence should be placed on the variables identified by variable selection as significant variables? Furthermore, how does one characterize the importance of a variable's contribution to the model? We will provide some insight into these questions through a frequentist approach to variable selection.

The frequentist's alternative approach to PIP will now be applied to the Spambase dataset. The entire dataset consisted of 4,601 observations. For step 1, 70% of the observations (3,221 observations) were randomly sampled from the entire dataset to form the training set. The remaining 30% of the observations (1,380 observations) formed the test set. This process was repeated 100 times to form 100 replicates of 70/30 training/test split. Since the main goal will be binary classification of spam / non-spam, the logistic regression model was selected for step 2. A model was built for each of the 100 training replicates using logistic regression with Logit link function on all 57 explanatory variables. The 100 models were then applied on their corresponding test sets to calculate the out-of-sample accuracies. The left-most figure below is a comparative boxplot between the training and test set accuracies of the Logit link function on all 57 explanatory variables. The same process was

Table 1. Accuracies of training sets for different link functions

Model	Min	1st Qu	Median	Mean	3rd Qu	Max
Logit	83.64%	92.92%	93.14%	92.83%	92.83%	93.32%
Probit	81.50%	88.66%	92.50%	90.50%	92.90%	93.73%
Cauchit	69.95%	94.65%	94.88%	94.28%	95.16%	95.87%

Table 2. Accuracies of test sets for different link functions

Link function	Min	1st Qu	Median	Mean	3rd Qu	Max
Logit	80.87%	92.39%	92.97%	92.69%	93.48%	94.42%
Probit	81.38%	88.73%	92.17%	90.43%	92.90%	94.28%
Cauchit	69.86%	93.99%	94.42%	93.84%	94.78%	95.65%

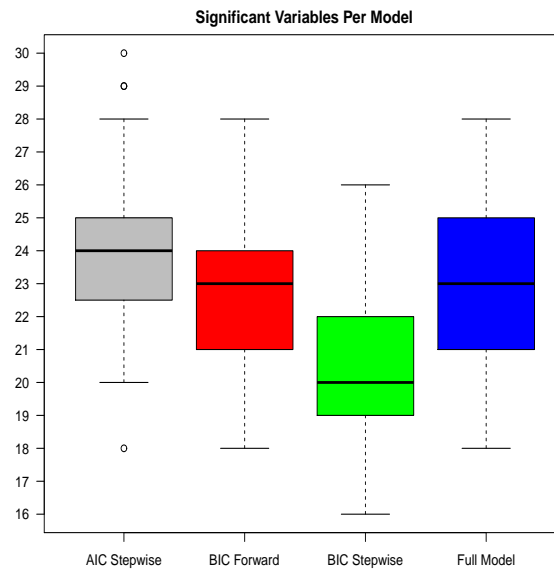


Fig. 3. Number of variables deemed significant in replications

The first example below shows 2 of the 100 replicates after performing AIC stepwise regression. In the replicate on the left, 18 variables were deemed significant, while 30 variables were deemed significant in the replicate on the right.

The second example below shows 2 of the 100 replicates after performing BIC stepwise regression. In the replicate on the left, 17 variables were deemed significant, while 26 variables were deemed significant in the replicate on the right.

In the next section, we will examine how to find a more robust method that helps truly identify significant variables.

The percentage in which each explanatory variable was deemed significant ($p\text{-value} \leq 0.05$) out of the 100 training replicates was calculated and plotted in the figure below (step 4). There

are several advantages to this approach. Firstly, we have identified, with high confidence, the 18 significant variables to comprise our core model. These variables are wf.remove, wf.hp, wf.re, wf.our, wf.free, wf.edu, cf.exclamation, cf.dollar, wf.business, wf.george, wf.project, cap.run.length.total, wf.your, wf.000, wf.internet, wf.receive, wf.over, and wf.money (note: wf is acronym for word frequency, cf is acronym for character frequency, and cap.run.length.total is the total number of capital letters in an email); for the remainder of this analysis, they will be referred to as the 18 core variables. There is high confidence that these 18 variables are significant variables due to their robustness. They were consistently identified as significant in $\geq 50\%$ of the replicates in all 4 of the variable selection methods despite the variability in the training sets' observations.

Coefficients:				(Intercept)					
	Estimate	Std. Error	z value	Pr(> z)		Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-4.059e+00	5.171e-01	-7.849	< 2e-15 ***	wf_address	-2.291e+00	3.563e-01	-6.430	1.28e-10 ***
wf_address	-3.365e-01	2.332e-01	-1.443	0.149060	wf_26	-2.196e-01	2.423e-01	-0.906	0.364794
wf_our	-3.696e-01	2.859e-01	-1.293	0.190107	wf_our	0.522e+00	1.059e+01	0.008	0.419887
wf_remove	1.302e+00	3.103e-01	3.873	0.000107 ***	wf_remove	1.707e+00	2.925e-01	5.833	5.39e-09 ***
wf_remove	3.870e+01	9.080e+00	4.262	2.03e-05 ***	wf_order	3.123e+00	7.097e-01	4.400	1.08e-05 ***
wf_remove	1.137e+00	4.932e-01	2.268	0.023039 *	wf_order	-2.934e+00	1.534e+00	-4.673	2.96e-06 ***
wf_order	9.449e-01	7.058e-01	1.339	0.180618	wf_free	1.503e+00	1.129e-01	4.805	1.53e-06 ***
wf_receive	1.366e-01	1.105e-01	1.570	0.116346	wf_bustress	4.444e+00	9.039e-01	4.924	8.64e-07 ***
wf_free	-2.765e+00	8.873e-01	-3.114	0.001847 **	wf_email	-6.050e-01	3.057e-01	-1.979	0.047845 *
wf_will	-1.860e-01	1.640e-01	-1.134	0.256687	wf_you	-2.697e-01	1.028e-01	-2.624	0.008683 **
wf_receive	2.075e+00	5.221e-01	3.974	0.000055 ***	wf_credit	5.886e+00	1.939e+00	3.046	0.002322 **
wf_bustress	3.442e+00	8.488e-01	4.058	4.99e-05 ***	wf_your	3.046e-01	1.121e-01	2.718	0.005652 **
wf_bustress	1.874e-01	7.099e-02	2.639	0.008301 **	wf_000	4.817e+00	1.815e+00	2.692	0.007104 **
wf_credit	2.113e+00	1.673e+00	1.262	0.206779	wf_money	4.007e-01	1.648e-01	2.431	0.015050 *
wf_your	2.082e-01	1.082e-01	1.924	0.054407	wf_hp	-2.093e+01	4.716e+00	-4.438	9.07e-06 ***
wf_000	1.407e+00	9.058e-01	1.533	0.129398	wf_hp1	-3.653e+00	3.608e+00	-1.005	0.314943
wf_hp	-1.803e+01	4.463e+00	-4.039	5.36e-05 ***	wf_george	-2.358e+01	4.864e+00	-4.848	1.25e-06 ***
wf_hp1	-6.307e+00	4.012e+00	-1.547	0.121745	wf_650	3.030e+00	8.277e-01	3.653	1.47e-05 ***
wf_george	-1.579e+02	2.515e+01	-6.280	3.39e-10 ***	wf_lab	-1.508e+01	1.528e+00	-9.867	0.326760
wf_lab	-1.926e+01	3.411e+01	-0.565	0.572241	wf_data	-1.894e+00	6.934e-01	-2.732	0.006301 **
wf_data	-1.597e+00	1.172e+00	-1.363	0.172840	wf_85	2.732e+00	9.478e-01	2.882	0.003947 **
wf_85	-2.956e+00	3.248e+00	-0.910	0.362674	wf_technology	-1.749e+00	8.462e-01	-2.054	0.041818 *
wf_technology	2.697e+00	8.061e-01	3.346	0.000821 ***	wf_direct	-1.793e+00	1.541e+00	-1.164	0.244463
wf_parts	-5.438e+00	4.424e+00	-1.229	0.219025	wf_cs	-1.229e+02	9.936e+01	-1.237	0.219024
wf_pm	-1.597e+00	1.114e+00	-1.423	0.151929	wf_meeting	-4.634e+00	1.556e+00	-2.978	0.002897 **
wf_cs	-1.504e+01	1.443e+01	-1.040	0.298309	wf_project	-6.814e+00	3.184e-01	-2.182	1.26e-07 ***
wf_meeting	-6.747e+00	4.402e+00	-1.533	0.125359	wf_re	-5.782e+00	1.421e+00	-4.070	4.71e-05 ***
wf_project	-2.965e+01	3.345e+01	-0.884	0.376408	wf_table	-6.934e+00	8.222e+00	-0.833	0.404711
wf_re	-1.377e+00	4.040e-01	-3.408	0.000654 ***	wf_confidence	-1.450e+01	8.436e+00	-1.719	0.085558 *
wf_re	-2.856e-01	5.274e-02	-5.340	0.28e-08 ***	cf_semicolon	-2.021e+00	6.809e-01	-2.968	0.002988 **
wf_table	-2.467e+01	4.616e+01	-0.534	0.593007	cf_parenthesis	-3.120e+00	1.576e+00	-1.952	0.020206 *
wf_confidence	-1.043e+01	8.982e+00	-1.161	0.245452	cf_exclamation	3.640e+00	5.245e-01	6.941	3.89e-12 ***
cf_semicolon	-2.009e+00	0.951e+00	-2.111	0.035008 *	cf_dollar	1.886e+00	2.202e-01	8.553	1.41e-13 ***
cf_parenthesis	-1.584e+00	6.276e-01	-2.523	0.011623 *	cf_pound	1.807e+01	3.418e+00	5.288	1.24e-07 ***
cf_exclamation	1.886e+00	2.202e-01	8.553	1.41e-13 ***	cap_run_length_avg	8.937e-01	1.404e-01	6.368	1.92e-10 ***
cf_dollar	1.807e+01	3.418e+00	5.288	1.24e-07 ***	cap_run_length_max	3.499e-03	9.004e-04	3.866	0.000107 ***
cf_pound	1.807e+01	3.418e+00	5.288	1.24e-07 ***	cap_run_length_total	3.394e-03	7.376e-04	4.601	4.20e-06 ***
cap_run_length_avg	8.937e-01	1.404e-01	6.368	1.92e-10 ***					
cap_run_length_max	3.499e-03	9.004e-04	3.866	0.000107 ***					
cap_run_length_total	3.394e-03	7.376e-04	4.601	4.20e-06 ***					

Fig. 4. 2 replicates after AIC stepwise regression

Coefficients:				Coefficients:					
	Estimate	Std. Error	z value	Pr(> z)		Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-4.197e+00	4.633e-01	-9.060	< 2e-16 ***	wf_our	4.628e+00	4.893e-01	-9.457	< 2e-16 ***
wf_our	4.206e-01	1.076e-01	3.909	9.26e-05 ***	wf_over	2.675e+00	4.090e-01	6.541	6.12e-11 ***
wf_remove	1.289e+01	2.499e+00	5.157	2.50e-07 ***	wf_remove	3.031e+00	7.866e-01	3.853	0.000117 ***
wf_remove	1.759e+00	6.568e-01	2.678	0.007402 **	wf_remove	9.283e+00	2.094e+00	4.433	9.29e-06 ***
wf_remove	-2.693e+00	7.370e-01	-3.654	0.000259 ***	wf_remove	1.732e+00	4.436e-01	3.903	9.41e-05 ***
wf_free	1.684e+00	3.626e-01	4.645	3.40e-06 ***	wf_order	3.114e+00	1.120e+00	2.781	0.005424 **
wf_bustress	5.114e+00	1.028e+00	4.976	6.49e-07 ***	wf_receive	-1.420e+00	5.867e-01	-2.420	0.015524 *
wf_you	2.397e-01	5.909e-02	4.056	5.00e-05 ***	wf_free	7.416e-01	2.289e-01	3.239	0.001198 **
wf_hp	-2.469e+01	5.373e+00	-4.595	4.33e-06 ***	wf_bustress	2.961e+00	6.893e-01	4.331	1.49e-05 ***
wf_hp1	-6.203e+00	5.081e+00	-1.221	0.222160	wf_credit	2.866e-01	5.756e-02	5.173	2.32e-07 ***
wf_george	-2.907e+01	3.860e+00	-7.530	5.07e-14 ***	wf_hp	3.949e+00	1.936e+00	2.040	0.041350 *
wf_lab	-6.730e+01	3.088e+01	-2.179	0.029219 *	wf_hp1	-1.349e+01	3.786e+00	-3.581	0.000342 ***
wf_pm	-1.106e+00	6.194e-01	-1.785	0.071837 *	wf_george	-9.716e+00	4.123e+00	-2.357	0.018431 **
wf_cs	-1.101e+02	8.226e+01	-1.338	0.180791	wf_george	-2.718e+00	6.011e+00	-0.453	4.35e-08 ***
wf_meeting	-8.776e+00	3.676e+00	-2.387	0.016964 *	wf_650	7.750e+00	1.214e+00	6.382	1.74e-10 ***
wf_project	-2.311e+00	1.248e+00	-1.852	0.064015 *	wf_tabs	-3.435e+00	1.058e+00	-3.246	0.001172 **
wf_re	-1.269e+00	4.065e-01	-3.123	0.001790 **	wf_85	-9.342e+00	2.287e+00	-4.084	4.42e-05 ***
wf_re	-3.547e+01	5.907e+00	-6.005	1.91e-09 ***	wf_pm	-3.565e+00	1.380e+00	-2.583	0.009792 **
wf_confidence	-6.943e+00	4.649e+00	-1.493	0.135340	wf_meeting	-2.209e+00	2.369e+00	-0.942	0.002348 **
cf_semicolon	-1.964e+00	1.061e+00	-1.852	0.064049 *	wf_project	-7.105e+00	2.447e+00	-2.903	0.003690 **
cf_exclamation	1.849e+00	3.236e-01	5.715	1.10e-08 ***	wf_re	-2.332e+00	4.058e-01	-5.746	9.12e-09 ***
cf_dollar	2.533e+01	3.987e+00	6.317	7.19e-11 ***	wf_re	-6.190e+00	1.414e+00	-4.381	6.90e-09 ***
cap_run_length_avg	9.439e-01	1.273e-01	7.417	1.20e-13 ***	wf_confidence	-5.539e+00	3.239e+00	-1.713	0.088641 *
cap_run_length_max	3.394e-03	7.376e-04	4.601	4.20e-06 ***	cf_semicolon	-1.267e+00	6.267e-01	-2.022	0.043180 *
cap_run_length_total	3.394e-03	7.376e-04	4.601	4.20e-06 ***	cf_exclamation	2.072e+00	3.489e-01	5.938	2.88e-09 ***
					cf_dollar	2.423e+01	3.340e+00	7.254	4.06e-13 ***
					cap_run_length_avg	8.939e-01	1.403e-01	6.368	1.92e-10 ***
					cap_run_length_max	1.079e-03	3.716e-04	2.903	0.003674 **
					cap_run_length_total	3.394e-03	7.376e-04	4.601	4.20e-06 ***

Fig. 5. 2 replicates after BIC stepwise regression

Secondly, this approach also affords flexibility to the end user in modeling. There were 7 variables (marked by red dotted vertical lines) which were deemed significant in $\geq 50\%$ of the replicates by at least 1 of the 4 variable selection methods but not by all 4. These variables are wf.technology, cap.run.length.avg, wf.meeting, wf.order, wf.your, wf.credit, and cap.run.length.max. Depending on the end user's threshold for model complexity-accuracy tradeoff, the user can experiment building models with any combination of these 7 variables in addition to the 18 core variables. Furthermore, each variable's importance can now be characterized by the percent of replicates in which they are deemed significant. For example, it may not be cost-effective for a business to measure all 18 core variables or may require too much computing power; as a result, the user may be constrained to using only 10 variables. How would the end user decide which 10 variables to use? Based on the

plot, the user should select the first 10 variables (wf.remove, wf.hp, wf.re, wf.our, wf.free, wf.edu, cf.exclamation, cf.dollar, wf.business, and wf.george) because they were statistically significant in > 90% of the replicates in all 4 variable selection techniques. One can view this approach as an alternative to Mallows's Cp.

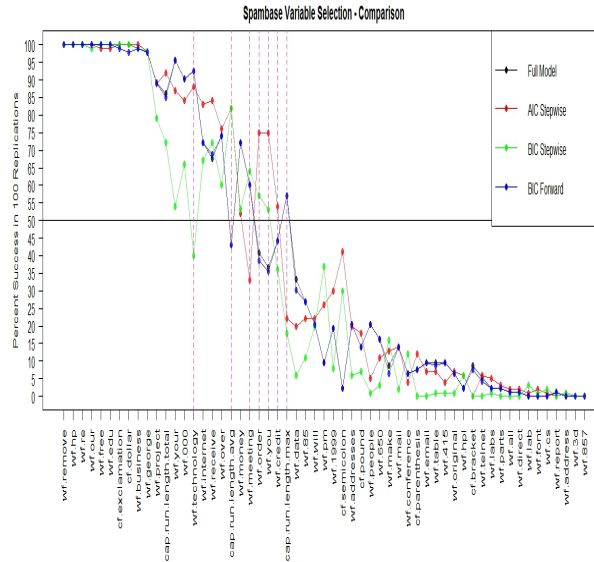


Fig. 6. Percentage in which each explanatory variable deemed significant

The remaining 32 variables (the variables to the right of the rightmost red dotted line) can be discarded because they never exceeded 50% in the replicate sets for any of the 4 selection methods. These 32 variables are more susceptible to the randomness of the observations in the training replicates. When one compares multiple replicates, these variables will not be consistently of deemed significance. The potential error in using only one variable selection method becomes evident. Suppose a modeler selected the model with the highest test set accuracy after performing only BIC stepwise regression (represented by green line in the plot above). There is ~ 20% chance the chosen model would not include wf.project (the percentage of an email in which the word "project" appears). However, when looking at variable selection from an ensemble point of view, there is high confidence that wf.project is a significant variable despite not being deemed significant in a replicate. Alternatively, there is a 30% chance that the chosen model would include cf.semicolon (the percentage of an email in which the character ";" appears). However, when looking at variable selection from an ensemble point of view, there is low confidence that cf.semicolon is a significant variable because it never reached $\geq 50\%$ in any of the variable selection methods.

Five families of models were built in increasing complexity (step 5). For the first family, a model was built for each training replicate using the 18 core variables. The accuracies were calculated for both training and test sets. This process was repeated 4 more times, in which different variables were added onto the 18 core variables (the variables used are listed below).

Let M_c denote the core model containing the 18 variables that always appear in every replication.

Table 3. Variable selection for various model complexities

Model Complexity	Variables Used
18 Variables	M_c
19 Variables	$M_{19} = M_c \cup \{\text{wf.technology}\}$
20 Variables	$M_{20} = M_{19} \cup \{\text{cap.run.length.avg}\}$
21 Variables	$M_{21} = M_{20} \cup \{\text{wf.meeting}\}$
23 Variables	$M_{23} = M_{21} \cup \{\text{wf.you, wf.credit}\}$

Table 4. Test set accuracies For various model complexities

Model Complexity	Min	1st Qu	Median	Mean	3rd Qu	Max
18	91.09%	93.04%	93.51%	93.50%	93.93%	94.93%
19	91.09%	93.04%	93.48%	93.46%	93.91%	94.86%
20	91.67%	93.32%	93.94%	93.78%	94.22%	95.29%
21	91.09%	93.62%	94.13%	94.05%	94.49%	95.65%
23	90.94%	93.84%	94.13%	94.12%	94.42%	93.51%

For nonlinear models like the ones that arise in the generalized linear model framework, it is often the case that the Fisher scoring algorithm used to estimate the parameters of the model does not converge. When that happens, the solution delivered is at best suboptimal, and may lead to misleading conclusions. For instance, the gist of the method proposed in this paper lies on scanning the variables and choosing the model made up of variables whose p -values are less than 0.05. With a suboptimal solution, it is unwise and misleading to consider that p -values are meaningful. For that reason, we systematically tract all the estimations throughout the totality of our random replications, and we provide an estimate of the percentage of times the estimates of the model are meaningful. It makes sense to us that only the cases where convergence is achieved should be used for inference, because -in a sense- that measures an aspect of the quality of the model space search. Below is a partial table of the percentages:

Table 5. Percentage of convergence/non-convergence for various model complexities

Model Complexity	% Converge	% Non-Converge
18	73	27
19	80	20
20	68	32
21	37	63
23	28	72

The table above shows the test set accuracies for all 5 models built. The last two columns in the table above list the percentage of 100 models which converged and the percentage of 100 models which did not converge. All 5 of the models have very respectable accuracies; furthermore, there is relatively little variation in test set accuracies throughout the 100 replicates. The end user now has several models to choose from. If the end user sought the most parsimonious model and is willing to accept a slight loss in accuracy, he/she can select the 18 variable model. If, on the other hand, the end user sought the highest prediction rate, he/she would select the 23 variable model. If the modeler sought the most computationally stable model, he/she would select the 19 variable model; this

model had 80% of its replicates converge. For comparison with current machine learning methods, we selected the 20 variable model. The 20 variable model achieved the best tradeoff between model accuracy and convergence rate; this model complexity is in agreement with the predicted optimal model size. The following two figures plot the ROC curves for all 5 models. In the first figure, all 5 ROC curves achieve "right angle" shape. The second figure provides a close-up of the upper-left corner of the first figure. The 20 variable model (green line) tracks well with the 23 variable model (purple dotted line); both of these have more area under the curve compared to the remaining 3 models.

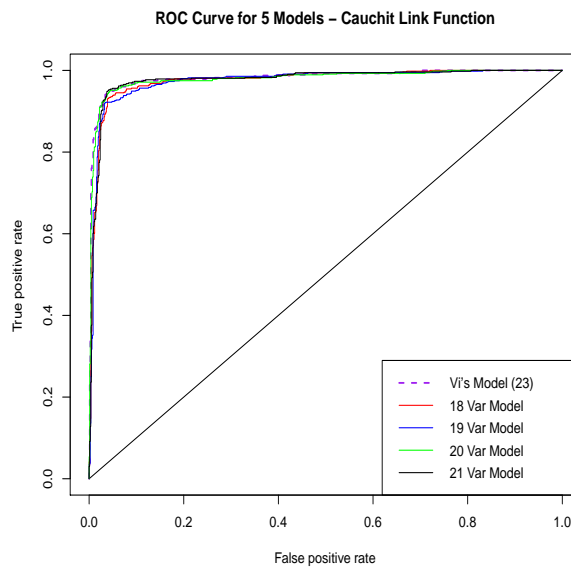


Fig. 7. ROC curves for 5 models

The 20 explanatory variables used to generate the model were `wf.remove`, `wf.hp`, `wf.re`, `wf.our`, `wf.free`, `wf.edu`, `cf.exclamation`, `cf.dollar`, `wf.business`, `wf.george`, `wf.project`, `cap.run.length.total`, `wf.your`, `wf.000`, `wf.internet`, `wf.receive`, `wf.over`, `wf.money`, `wf.technology`, and `cap.run.length.avg`. 6 out of the 20 explanatory variables are indicators of non-spam (as their frequency in an email increases, the more likely an email is not spam); these variables are `wf.hp`, `wf.re`, `wf.george`, `wf.edu`, `wf.project`, and `wf.technology`. 14 out of the 20 explanatory variables are indicators of spam (as their frequency in an email increases, the email is more likely to be spam); these variables are `wf.000`, `wf.money`, `cf.dollar`, `wf.free`, `wf.remove`, `wf.business`, `wf.your`, `cf.exclamation`, `cap.run.length.avg`, `cap.run.length.total`, `wf.internet`, `wf.over`, `wf.our`, and `wf.receive`.

When one looks at both classes of indicators, certain patterns arise which can provide further explanation. The variables `wf.hp`, `wf.re`, `wf.george`, `wf.project`, and `wf.technology` in the non-spam indicators, suggest a personal or professional relationship with the recipient. Since the dataset was donated by [4] at Hewlett Packard labs, it makes sense that emails containing "hp" and "George" indicate that the sender either knew the recipient and/or the email was work-related. The characters "re" are often used in emails as replies. Therefore, email replies are flagged as non-spam because the recipient is receiving a reply to an earlier email sent out by the recipient. Additionally, the words meeting and project are usually work-related terms and hence why they are also indicators of non-spam. The variables `wf.000`, `wf.money`, `cf.dollar`, `wf.free`, and `wf.business` are associated with money. This makes sense since most spam emails are attempts to get money from the recipient.

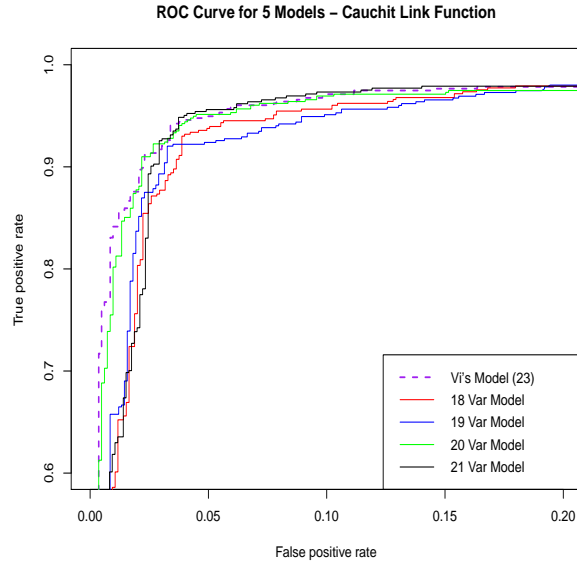


Fig. 8. Close up of ROC curves for 5 models

In [5], different machine learning techniques were applied to the Spambase dataset with the goal of optimizing correct classification rate. The following table lists the testing set accuracy for the different machine learning techniques used in the paper. Our 20 variable model, with a median test accuracy of 93.94%, was outperformed by Ensemble Decision Tree and Adaboost methods.

Table 6. Test set accuracies for various machine learning techniques

Classifier	Accuracy
Ensemble Decision Tree (Nb of trees = 25)	96.40%
Adaboost	95.00%
Stacking	93.80%
SVM	93.40%
Bagging	92.80%
Decision Tree	92.58%
Neural Network	90.80%
Naive Bayes	89.57%
Nearest Neighbor (k=5)	89.40%

Independently, [6] applied Generalized Additive Neural Networks (GANNs) to the Spambase dataset to classify email as spam/non-spam. The spambase dataset was broken in 67% training and 33% test. The AutoGANN method used by [6] attained a 95.8% accuracy. While this methods accuracy beat the 20 variable model (95.8% > 93.94%), it required higher model complexity with 41 explanatory variables used.

[7] used a neural network method called Multi-Layer Perceptron (MLP). In this paper, the Spambase dataset was broken down into training set with 4,025 observations (~ 87.48%) and test set with 576 observations (~ 12.59%). The author applied this technique on 3 different scenarios. In the

first scenario, [7] used MLP on the training set using all 57 exploratory variables. Due to the high dimensionality, the author noted that the MLP method had trouble converging and consequently, generated poor classification rates. A similar effect was observed in our approach. For our 5 models, the percentage of models that converged drastically decreased when model complexity increased above 20 variables. The following table lists the accuracies for different models built by the author using all 57 variables.

Table 7. Accuracies for MLP using 57 variables

MLP Architecture	Learning Rate	Momentum	Average Classification Rate
20 – 10 – 10 – 10 – 5	0.1	0.80	62.40%
20 – 10 – 10 – 10 – 5	0.1	0.95	63.20%
20 – 10 – 10 – 5	0.1	0.85	60.61%
15 – 15 – 15 – 5	0.1	0.85	60.59%

The author used the same MLP method after reducing dimensionality to 21 variables. The test accuracies listed below are very comparable to the test set accuracies for our 20 variable model.

Table 8. Accuracy for MLP using 21 variables

MLP Architecture	Learning Rate	Momentum	Average Classification Rate
20 – 10 – 10 – 10 – 5	0.1	0.80	93.50%
20 – 10 – 10 – 10 – 5	0.1	0.95	90.30%
20 – 10 – 10 – 10 – 7	0.1	0.80	93.80%

Lastly, in the third attempt, MLP was applied after reducing dimensionality to 9 variables. The test accuracies are listed below.

Table 9. Accuracy for MLP using 9 variables

MLP Architecture	Learning Rate	Momentum	Average Classification Rate
20 – 10 – 10 – 8	0.1	0.80	92.40%
20 – 10 – 10 – 5	0.1	0.80	91.70%
20 – 10 – 10 – 5	0.1	0.95	91.80%

In a fourth comparison, [8] utilized 9 different machine learning algorithms for their paper *Adaptive Approach for Spam Detection*. In their approach, the data was transformed into 1's and 0's. If a certain word appeared, that exploratory variable was a 1; if the certain word did not appear in an email, the exploratory variable was a 0. This was done for 55 of the 57 exploratory variables. The following table lists the performance of different algorithms after ten-fold cross validation. Our 20 variable model, with a median test accuracy of 93.94%, was out-performed by Random Committee and essentially tied with Random Forest technique.

There are several advantages to performing frequentist approach. First, this method provides a more robust variable selection by examining how often a variable is deemed significant by multiple traditional variable selection methods given random samples of observation data. Second, it also provides an approximation to the optimal model size. Third, it allows the modeler to characterize the importance of a variable to the model through the frequency in which a variable is deemed significant. Lastly, it affords the modeler flexibility in choosing certain variables to retain or discard depending on the modeler's threshold for model-complexity accuracy tradeoff. The main downside

Table 10. Accuracy for various machine learning algorithms

Algorithm	Accuracy
Bayes Network	88.56%
Logic Boost	89.70%
Random Tree	91.54%
JRip	92.32%
J48	92.34%
Multilayer Perceptron	93.28%
Kstar	93.56%
Random Forest	93.89%
Random Committee	94.28%

to this method is computational intensity. On a Windows 7 64-Bit Laptop with Intel i7 2.7GHz processor and 16GB RAM, this process required ~ 12 hours to just apply BIC stepwise regression to 100 replicates and another 12 hours to apply AIC stepwise regression to the 100 replicates. Running the full model and forward selection via BIC on the 100 replicates was markedly faster and completed within minutes. By incorporating more variable selection methods to the ensemble, the modeler will have a serious tradeoff in computing time. However, this dilemma may be alleviated through the use of parallel processing in which multiple tasks are dispersed over multiple workstations rather than running the tasks sequentially on one computer. With further advances in parallel processing and increases in computing power, the ensemble variable selection method's advantages will significantly dominate over its main weakness.

4.2 Computational demonstration on regression analysis

In the Spambase dataset, the frequentist approach was applied to classification. In the next example, the frequentist approach was applied to Multiple Linear Regression (MLR) on the Bodyfat dataset. The dataset, which was originally donated by [9], attempts to estimate body fat percentage by underwater weighing and various body circumference measurements for 252 men; this dataset may be found in the R package `mfp`. The dataset contained 2 response variables: `brozek` and `siri`. The `brozek` response variable calculated body fat percentage through the equation:

$$\text{brozek} = \frac{457}{\text{density}} - 414.2$$

The `siri` response variable calculated body fat percentage through the equation:

$$\text{siri} = \frac{495}{\text{density}} - 450$$

There were 14 explanatory variables. The first 3 variables are `density` (density determined from underwater weighing), `age`, and `weight`. The remaining 11 explanatory variables are body circumference measurements for `neck`, `chest`, `abdomen`, `hip`, `thigh`, `knee`, `ankle`, `biceps`, `forearm`, and `wrist`. The linear pairwise correlation plot below indicates a significant amount of multicollinearity and redundant variables.

This bodyfat data set can also be found at <http://lib.stat.cmu.edu/datasets/bodyfat>

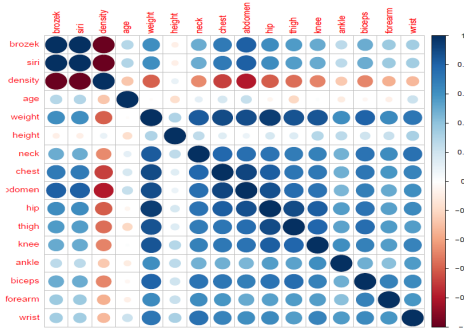


Fig. 9. Linear Correlation Plot for Bodyfat Dataset

The purpose of this section will be to compare the results from the frequentist approach for variable selection against results from Bayesian PIP. There were 5 cases identified in the dataset as erroneous observations/outliers, and as a result, were excluded during analyses. Cases 48, 76, and 96 were identified as having errors in body fat values. In case 42, the man weighed 200 lbs with a height of 3 ft. For case 182, the body fat percentage was rounded to 0 after having a negative body fat percentage. Additionally, the **density** variable was excluded from MLR because formulas to calculate **brozek** and **siri** response variables require **density** variable; consequently, the inclusion of **density** variable in MLR would dominate and bias the selection of other significant variables. The dataset was broken into 500 replicates of 70% training / 30% test sets (step 1 of frequentist approach). The linear model (MLR) was selected for the class of models (step 2). AIC stepwise regression, BIC stepwise regression BIC forward selection, and full model were applied to the 500 training sets (step 3). This section will concentrate solely on the **brozek** response variable. After scanning through the 500 replicates, the percentage in which each variable was deemed significant is listed below for each of the variable selection methods (step 4).

Table 11. Estimated percentage of inclusion of each variable

Variable	Full Model	Stepwise BIC	Forward BIC	Stepwise AIC
abdomen	1	1	1	1
wrist	0.8	0.698	0.8	0.842
weight	0.13	0.576	0.13	0.53
forearm	0.386	0.36	0.386	0.516
neck	0.21	0.222	0.21	0.304
hip	0.082	0.22	0.082	0.226
height	0.01	0.166	0.01	0.156
age	0.124	0.08	0.0124	0.2
biceps	0.036	0.08	0.036	0.1
thigh	0.07	0.038	0.07	0.208
chest	0.012	0.028	0.012	0.052
ankle	0.014	0.016	0.014	0.034
knee	0	0.002	0	0.006

The following section will compare the results from our frequentist approach against the results from Bayesian PIP. In order for a more objective comparison, only the stepwise BIC portion will be compared against the Bayesian PIP. The Bayesian PIP results were acquired using the Bayesian Model Selection (BMS) package in R and shown below.

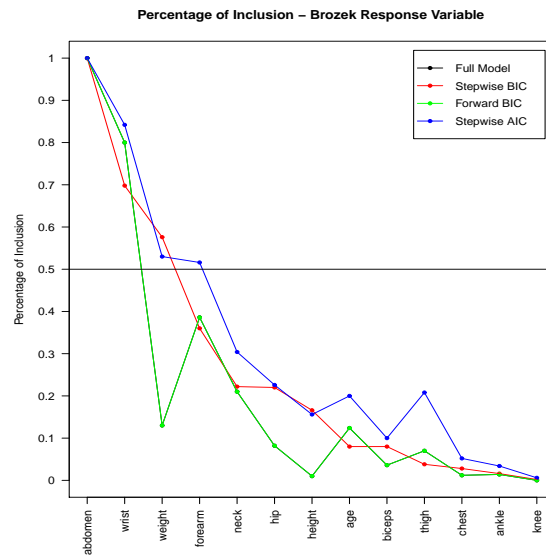


Fig. 10. Percentage of inclusion of explanatory variables for brozek response variable - frequentist approach

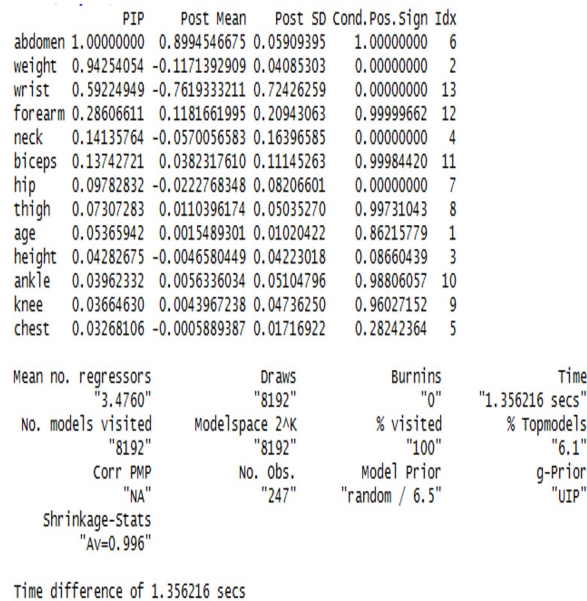


Fig. 11. R Output - Bayesian PIP of explanatory variables for brozek response variable

The frequentist approach (looking at stepwise BIC only) and the Bayesian PIP both identified abdomen as the most significant explanatory variable; in both methods, the abdomen variable was deemed significant in 100% of the 500 training sets. Additionally, in both methods, weight and wrist variables exceeded the median probability model (deemed significant in $\geq 50\%$ of 500 replicates)

and will be retained as significant variables. The main difference between the frequentist approach and PIP is evident in the **weight** variable. The **weight** variable was deemed significant in $\sim 94\%$ of the 500 models by Bayesian PIP but only $\sim 58\%$ by the frequentist method. The following plot shows the distribution of model size across the 500 replicates for the frequentist approach. The optimal model size should include 3 to 4 variables. The average model size across 500 replicates was 3.49.

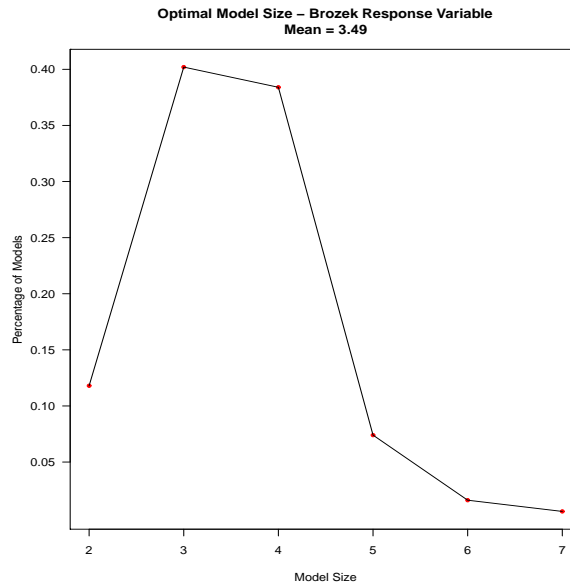


Fig. 12. Optimal model size - frequentist approach

The following plot shows the distribution of model size across 500 replicates for the Bayesian PIP approach. The PIP approach also indicated an optimal model size between 3 to 4 variables with an average model size of 3.48 across 500 replicates. The results between the two methods are very similar.

The frequentist approach (looking at stepwise BIC only) and the Bayesian PIP both identified **abdomen**, **weight** and **wrist** as significant variables. Building 500 linear models using only the 3 variables, we obtain the following results shown below. As expected, the out of sample prediction results for the frequentist and Bayesian PIP because we are applying the same variables to the training sets.

We will now build new sets of models to compare the frequentist and Bayesian PIP approaches. From a meta-analysis standpoint, **abdomen** and **wrist** variables were deemed significant in $\geq 50\%$ of the replicates for all 4 variable selection techniques in the frequentist approach. As a result, 500 linear models were built using only the **abdomen** and **wrist** explanatory variables. From the PIP standpoint, **abdomen**, **weight**, and **wrist** variables exceeded the median probability model. As a result, 500 linear models were built using **abdomen**, **weight**, and **wrist** explanatory variables. We could have justifiably included **weight** as an additional third variable in our frequentist approach because it was deemed significant by at least one variable selection technique; however, in doing so, we would get the same end results as the PIP linear models since both methods would now use the same variables. By comparing a two variable model (**abdomen** and **wrist**) against a three

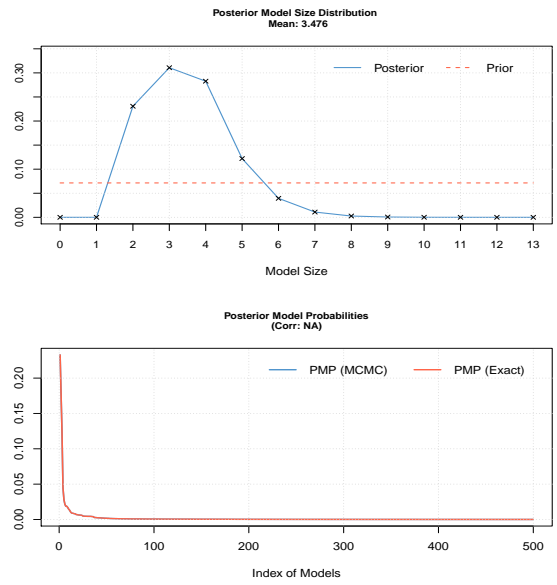


Fig. 13. Optimal model size - Bayesian PIP approach

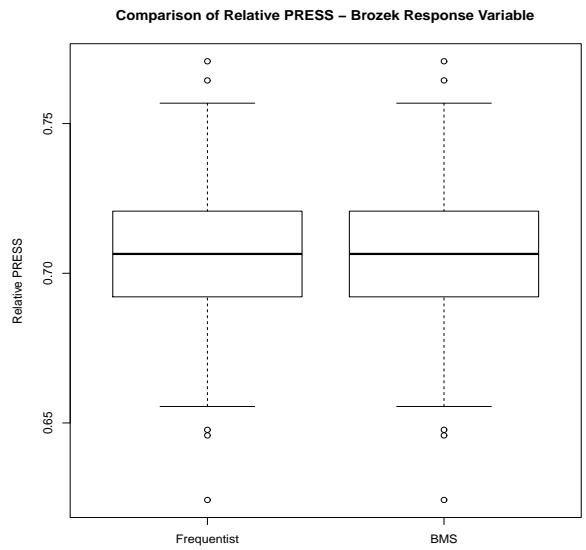


Fig. 14. Comparison of relative PRESS - both approaches use abdomen, weight and wrist Variables

variable model (abdomen, weight, and wrist), we hope to achieve a more distinction between the two model types. In the next section, we will compare the out of sample predictive performance of the frequentist and the PIP approaches. The following table and plot will compare the relative PRESS values between the two methods.

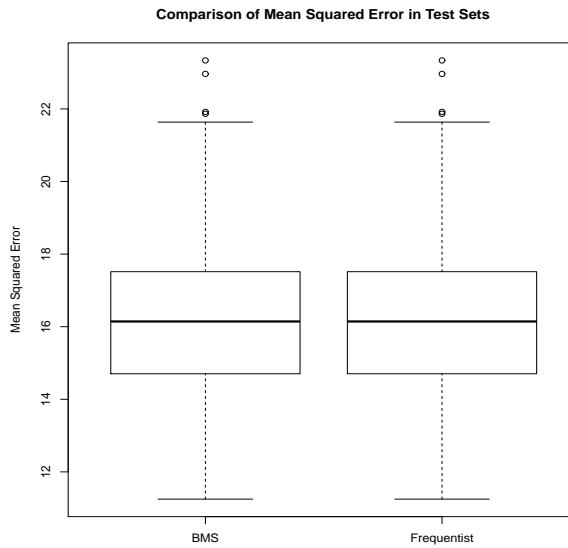


Fig. 15. Comparison of MSE - Both Approaches Use abdomen, weight and wrist Variables

Table 12. Comparison of PRESS between Frequentist Approximation and Bayesian PIP

Model	Min	1st Qu	Median	Mean	3rd Qu	Max
Frequentist	0.6025	0.6638	0.6801	0.6812	0.6983	0.7585
BMS PIP	0.6243	0.6921	0.7065	0.7061	0.7208	0.7708

The 500 linear models built for each of the approaches were applied on their corresponding test sets (step 5). The following table and plot will compare the Mean Squared Error (MSE) across the 500 test sets for both methods.

Table 13. Comparison of the MSE between Frequentist Approximation and Bayesian PIP

Model	Min	1st Qu	Median	Mean	3rd Qu	Max
Frequentist	11.31	15.84	17.52	17.68	19.27	25.31
BMS PIP	11.25	14.71	16.14	16.18	17.51	23.34

While *weight* achieved a PIP of 0.94, its addition into the three variable model did not provide a practical improvement. The results between the PIP and frequentist approaches are very comparable and provides validity of using the frequentist approach as an alternative to the Bayesian PIP.

The results when using *siri* as the response variable are almost identical to the results attained above using *brozek* as the response variable. As a result, the comparison between the frequentist and PIP approaches for *siri* response variable will not be provided.

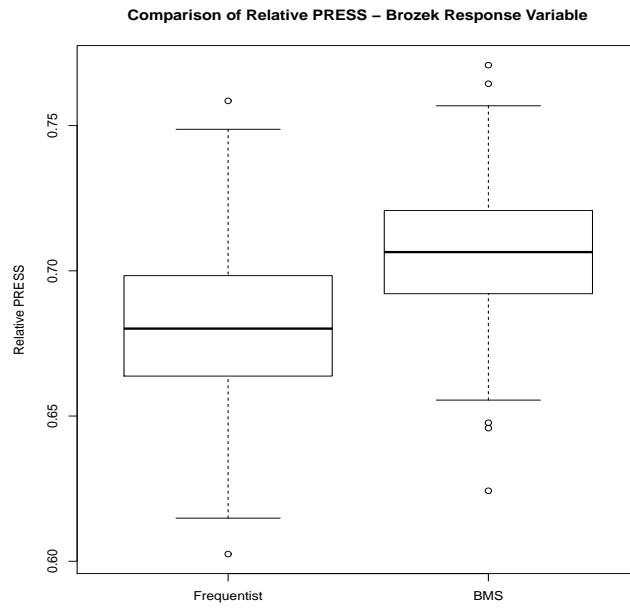


Fig. 16. Comparison of Relative PRESS

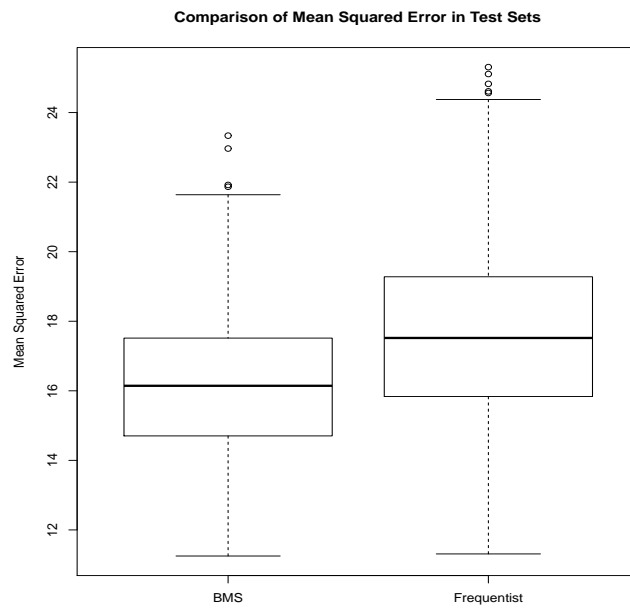


Fig. 17. Comparison of the MSE between Frequentist Approximation and Bayesian PIP

Applying the frequentist approach on linear models was significantly faster than on generalized linear models. Running stepwise regression via BIC and AIC, forward selection via BIC, and the full model on the 500 replicates for both `brozek` and `siri` response variables only took minutes.

5 Conclusion and Discussion

We have used a straightforward, quite general and easily interpretable subsampling scheme to provide a frequentist approximation of the celebrated Bayesian posterior inclusion probability. Despite the relatively higher computational burden arising in the use of the proposed method on high dimensional classification tasks, it is fair to say that the present method mimics the Bayesian framework quite well. All the scores, judging from the numerical values and the corresponding plots appear identical or at least very similar in shape and form. One would be particularly be excited to use this approach because it is easier to understand since it builds up on the widely used framework of variable selection by the stepwise regression heuristic. Even more importantly, it does not run into the some of the challenges of the Bayesian framework like the difficulty in computing the marginal density of the data. As we said earlier, the great challenge for this method is the heavy computational burden. However, with the availability of distributed and high performance parallel computing resources, this method becomes even more attractive for high dimensional data mining problems since one can perform the independent random split on different CPUs. Indeed, our future work will focus on substantially reducing the computing time by a careful use of the parallel computing resources.

Acknowledgement

Ernest Fokoué wishes to express his heartfelt gratitude and infinite thanks to Our Lady of Perpetual Help for Her ever-present support and guidance, especially for the uninterrupted flow of inspiration received through Her most powerful intercession.

Competing Interests

Authors have declared that no competing interests exist.

References

- [1] Barbieri MM, Berger JO. Optimal predictive model selection. *Annals of Statistics*. 2004;32:870-897.
- [2] Fokoué E. Estimation of Atom Prevalence for Optimal Prediction. *Contemporary Mathematics* (The American Mathematical Society). 2007;443:103-129.
- [3] Fumera G, Pillai I, Roli F. Spam filtering based on the analysis of text information embedded into images. *Journal of Machine Learning Research*. 2006;7:2699-2720.
- [4] Forman G, Hopkins M, Reeber E. UCI machine learning repository; 1999.
- [5] Kiran R, Atmosukarto I. Spam or not spam that is the question. Technical report, University of Washington; 2005.
- [6] Goosen JC, Du Toit JV. Spam detection with generalised additive neural networks. In *Southern Africa Telecommunication Networks and Applications Conference*; 2009.
- [7] Sivanadyan T. Spam? not any more! detecting spam emails using neural networks. Technical report, University of Wisconsin; 2003.

- [8] Sharma S, Arora A. Adaptive approach for spam detection. International Journal of Computer Science Issues. 2013;10(1):23-26.
- [9] Penrose KW, Nelson AG, Fisher AG. Generalized body composition prediction equation for men using simple measurement techniques. Medicine and Science in Sports and Exercise. 1985;17(2):189.

©2016 Ly and Fokoué; This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Peer-review history:

The peer review history for this paper can be accessed here (Please copy paste the total link in your browser address bar)

<http://sciencedomain.org/review-history/16212>