
On Bayesian Learning of Sparse Classifiers

Wen-Hua Ju

Avaya Labs Research
233 Mount Airy Road
Basking Ridge, NJ 07920
whju@avaya.com

David Madigan

Department of Statistics
Rutgers University
Piscataway, NJ 08855
madigan@stat.rutgers.edu

Steven L. Scott

University of Southern California
Bridge Hall 401-H
Los Angeles, CA 90089
sscott@marshall.usc.edu

Abstract

Figueiredo (2001) and Figueiredo and Jain (2001) described a particular sparseness-inducing Bayesian model for probit regression. For several standard datasets, they reported predictive performance for their model that was as good as, or better than, previously reported results. This paper explores several aspects of the Figueiredo and Jain model in an attempt to better understand its performance. We modify the Figueiredo and Jain approach in three ways. First, we introduce an alternative prior distribution. Second, we propose a fully Bayesian MCMC learning algorithm. Third, we replace their kernel based classifier with a linear classifier. We measure the impact of these modifications on three publicly available test data sets. Preliminary results indicate that while each change can produce a noticeable impact on test error rates, no one approach dominates the others in all cases.

1 Introduction

Inductive supervised learning infers a functional relation $y = f(\mathbf{x})$ from a set of training examples $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. In what follows the inputs are vectors $[x_{i_1}, \dots, x_{i_d}]^T$ in \mathfrak{R}^d , $y \in \{-1, 1\}$, and we refer to f as a classifier. We assume that a vector of parameters, β defines f and we write $f(\mathbf{x}; \beta)$. The learning procedures we consider output either a point estimate of β or a posterior distribution for β .

Our objective is to produce a classifier that makes accurate predictions for future input vectors. Typically this requires the learning procedure to control complexity and avoid over-fitting the training data. The Bayesian approach to complexity places a prior distribution on

β , typically resulting in estimates that shrink towards zero. Using so-called Laplacian (i.e., double exponential) priors can result in posterior modes for some components of β that are exactly zero. Such “sparse classifiers” can yield excellent predictive performance and are closely related to support vector machines (SVM) (see Girosi, 1998, and Zhang and Oles, 2001), relevance vector machines (RVM) (Tipping, 2001), and to the lasso (Tibshirani, 1995).

Figueiredo (2001) and Figueiredo and Jain (2001), hereafter FJ, considered classifiers of the form:

$$p(y = 1|\mathbf{x}) = \psi(\beta^T h(\mathbf{x})).$$

The logistic link function is a common choice for ψ , but following FJ, we instead adopt the probit model, $\psi(z) = \Phi(z)$, where Φ denotes the standard Gaussian distribution function. (see Chambers and Cox, 1967, for a comparison of the logistic and probit models.) We consider two forms for \mathbf{h} , namely $\mathbf{h}(\mathbf{x}) = [1, x_1, \dots, x_d]^T$ (i.e., the original input variables), and $\mathbf{h}(\mathbf{x}) = [1, K(\mathbf{x}, \mathbf{x}_1), \dots, K(\mathbf{x}, \mathbf{x}_n)]^T$, where $K(\mathbf{x}, \mathbf{y})$ is some symmetric kernel function not necessarily satisfying Mercer’s condition. SVMs and RVMs generally adopt the kernel representation.

FJ adopted a hierarchical prior for β , specifically, for $i = 1, \dots, d$:

$$p(\beta_i|\tau_i) = N(0, \tau_i)$$

and

$$p(\tau_i) \propto 1/\tau_i.$$

The prior on the τ ’s is the Jeffreys’ prior. Replacing it with an exponential prior

$$p(\tau_i|\gamma) = \frac{\gamma}{2} \exp(-\frac{\gamma}{2}\tau_i),$$

induces the following Laplacian prior on β

$$p(\beta_i|\gamma) = \frac{\sqrt{\gamma}}{2} \exp(-\sqrt{\gamma}|\beta_i|).$$

The maximum *a posteriori* (MAP) estimates under the Laplacian prior are the lasso estimates. Note that the Laplacian prior requires a choice for the hyperparameter γ , perhaps via cross-validation, whereas the FJ prior is more convenient insofar as it has no tunable hyperparameters.

In what follows we explore three questions:

1. FJ make “plug-in” predictions at the posterior mode of β . That is, plug-in predictions estimate the predictive density of a future observable z by $p(z; \hat{\beta})$ where $\hat{\beta}$ denotes the posterior mode. Alternatively one can adopt a fully Bayesian approach and integrate out β when making predictions. That is, fully-Bayesian predictions estimate the predictive density as $\int p(z; \beta)\pi(\beta|D)d\beta$ where D denotes the past data and π the posterior density of β . FJ find the posterior mode using an efficient EM algorithm. By contrast, integration over the posterior distribution of β requires a computationally expensive Monte Carlo algorithm. Does the fully Bayesian approach lead to better predictions?
2. FJ only report results using a kernel representation (specifically a Gaussian kernel). Does the kernel representation provide better predictive performance than the simpler non-kernel representation?
3. For the Laplacian prior, do there exist choices for γ that provide better predictive performance than that provided by the FJ prior.

2 The Learning Algorithms

Here we describe FJ’s EM algorithm for finding the posterior mode and a proposed MCMC algorithm for finding the posterior mean. Both algorithms make use of the latent variable representation for probit regression introduced by Albert and Chib (1993) and we describe this first. Define n independent latent variables z_1, \dots, z_n where z_i has distribution $N(\beta^T h(x_i), 1)$. Then define $y_i = 1$ if $z_i > 0$ and $y_i = -1$ if $z_i \leq 0$. It is straightforward to show that the y_i are then independent Bernoulli variables with $p(y_i = 1) = \Phi(\beta^T h(x_i))$, $i = 1, \dots, n$, and we recover the standard probit model. Figure 1 shows the corresponding graphical Markov model using the BUGS plate notation.

Observe that if the z_i are known, and β has a multivariate normal prior, then the posterior distribution for β is available in closed form using standard normal linear model results. The z_i are latent, but given the data y_i , z_i has truncated normal distribution. This

setup facilitates both the EM algorithm of FJ and our proposed MCMC algorithm in an obvious way.

For both the Laplacian and FJ priors, FJ derive an EM algorithm that estimates β by treating both τ and \mathbf{z} as missing. The complete data log-posterior is:

$$\begin{aligned} \log p(\beta|\mathbf{y}, \tau, \mathbf{z}) &\propto \log p(\mathbf{z}|\beta) + \log p(\beta|\tau) \\ &\propto -\|\mathbf{H}\beta - \mathbf{z}\|^2 - \beta^T \Gamma(\tau)\beta, \end{aligned}$$

where $\Gamma(\tau) = \text{diag}(\tau_1^{-1}, \dots, \tau_d^{-1})$ and \mathbf{H} is the design matrix with rows $\mathbf{h}(\mathbf{x}_1)^T, \dots, \mathbf{h}(\mathbf{x}_n)^T$. The E-step requires computing $E[\tau_i^{-1}|\hat{\beta}^{(t)}, \mathbf{y}, \gamma]$ and $E[z_i|\hat{\beta}^{(t)}, \mathbf{y}, \gamma]$. For the former, under the Laplacian prior, we have that:

$$p(\tau_i|\hat{\beta}^{(t)}, \mathbf{y}, \gamma) \propto p(\hat{\beta}^{(t)}|\tau_i)p(\tau_i|\gamma).$$

Then

$$\begin{aligned} \omega_i \equiv E[\tau_i^{-1}|\hat{\beta}^{(t)}, \mathbf{y}, \gamma] &= \frac{\int_0^\infty \frac{1}{\tau_i} p(\tau_i|\gamma)p(\hat{\beta}^{(t)}|\tau_i)d\tau_i}{\int_0^\infty p(\tau_i|\gamma)p(\hat{\beta}^{(t)}|\tau_i)d\tau_i} \\ &= \gamma|\hat{\beta}^{(t)}|^{-1}. \end{aligned}$$

For the latter, z_i is Gaussian with mean $\beta^T h(x_i)$, but left-truncated at zero if $y_i = 1$ and right-truncated at zero if $y_i = -1$. Thus:

$$\begin{aligned} v_i \equiv E[z_i|\hat{\beta}^{(t)}, \mathbf{y}, \gamma] &= \\ \begin{cases} \beta^T h(x_i) + \frac{N(\beta^T h(x_i))}{1 - \Phi(-\beta^T h(x_i))} & \text{if } y_i = 1 \\ \beta^T h(x_i) + \frac{N(\beta^T h(x_i))}{\Phi(-\beta^T h(x_i))} & \text{if } y_i = -1. \end{cases} \end{aligned}$$

So the E-step replaces Γ by its conditional expectation $\Omega = \text{diag}(\omega_1, \dots, \omega_d)$ and replaces \mathbf{z} by $\mathbf{v} = (v_1, \dots, v_n)^T$.

The M-step carries out a maximization with respect to β leading to:

$$\hat{\beta}^{(t+1)} = (\Omega + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{v}.$$

The EM procedure for the FJ prior requires minor modifications of the above algorithm. We refer the reader to FJ for details.

For fully Bayesian inference, we modify the Gibbs sampling procedure of Albert and Chib (1993). The basic idea is to draw a dependent sample from the target posterior distribution by drawing in turn from the conditional distribution of each of the unknowns given all the knowns and remaining unknowns. Specifically, the sampler draws in turn from:

$$z_i | z_{-i}, \mathbf{y}, \beta, \tau \sim TN_{y_i}(\beta^T h(x_i))$$

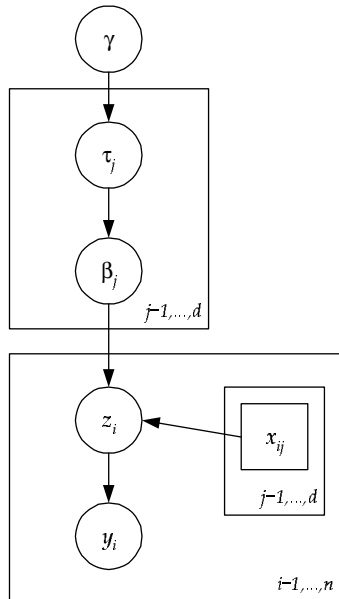


Figure 1: *Probit Model with Hierarchical Prior and Latent Variables.*

where “ TN ” is a truncated normal as described above, and

$$\beta|\mathbf{z}, \mathbf{y}, \tau \sim N((\mathbf{\Gamma} + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{z}, (\mathbf{\Gamma} + \mathbf{H}^T \mathbf{H})^{-1}).$$

The conditional density of τ is not available in closed form so we introduce a Metropolis-within-Gibbs step (see, for example, Gilks *et al.*, 1996). For each component $i = 1, \dots, d$, propose a candidate $\tau_i^{(t+1)}$ uniformly from a interval around $\tau_i^{(t)}$ of width h . For the Laplacian prior, accept the proposed move with probability:

$$\min \left\{ 1, \frac{N(\beta_i, 0, \tau_i^{(t+1)}) \exp(-\tau_i^{(t+1)} \gamma / 2)}{N(\beta_i, 0, \tau_i^{(t)}) \exp(-\tau_i^{(t)} \gamma / 2)} \right\}$$

otherwise set $\tau_i^{(t+1)} = \tau_i^{(t)}$. For the FJ prior, the acceptance probability is:

$$\min \left\{ 1, \frac{N(\beta_i, 0, \tau_i^{(t+1)}) 1 / \tau_i^{(t+1)}}{N(\beta_i, 0, \tau_i^{(t)}) 1 / \tau_i^{(t)}} \right\}.$$

In the experiments reported below, $h = 1$, and further tuning minimally impacted convergence behavior.

3 Classification Experiments

Following FJ and to facilitate comparison, we conducted experiments using three well-known benchmark datasets, namely Pima Indians Diabetes, *Leptograpus* Crabs, and Wisconsin Breast Cancer (WBC). The first two datasets are available divided into training and test sets. Pima Indians is at <http://www.stats.ox.ac.uk/pub/PRNN/>.

Crabs is at <http://www.inference.phy.cam.ac.uk/is/data/ripley-class/>. Pima has 200 training observations and 332 test observations. Crabs has 80 training observations and 120 test observations. The WBC data set comes from the UCI repository. Following FJ, we consider random splits of the WBC data into 300 training observations and 269 test observations. In each case we standardized the input variables. As with FJ, our kernel models used a Gaussian kernel with bandwidth 4 for Pima and Crabs and 12 for WBC. All MCMC runs are of length 10^6 with 10^3 burn-in. The EM algorithm terminated when $\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| / \|\hat{\beta}^{(t)}\| < 10^{-3}$. Table 1 shows the results with the FJ (i.e., Jeffreys) prior.

We were unable to reproduce the results of FJ (see lines 1 and 2 of Table 1). We created a second implementation in a different programming language to verify our results.

NOTE TO REVIEWERS: WE HAVE REQUESTED FIGEUEIREDO’S CODE AND HE HAS PROMISED TO SEND IT TO US.

3.1 Plug-in or Fully Bayesian Predictions?

Fully Bayesian inference via Markov chain Monte Carlo required at least 100 times the computing effort of the plug-in EM algorithm in these experiments. Generally, the EM algorithm provided parameter estimates in less than one minute on a 1GHz processor. By contrast, MCMC runs of 10^6 iterations took anything from 20 hours for the crabs data to 8 days for the WBC data with a kernel. For experiments not

	Method	Kernel	Bayes?	Pima	Crabs	WBC
1	Probit (F&J)	yes	plug-in	61	0	3.2%
2	Probit	yes	plug-in	72	2	4.1%
3	Probit	no	plug-in	70	3	5.0%
4	Probit	yes	full	72	3	2.8%
5	Probit	no	full	70	3	3.3%

Table 1: Number of test set errors (Pima and Crabs) or error rate (WBC). Line 1 reproduces the results reported in Figueiredo (2001). WBC results are averaged over 5 random splits into 300 training examples and 269 test examples. “Full Bayes” results used MCMC runs of length 10^6 . “Plug-in” used the EM algorithm described in the text.

involving a kernel, MCMC convergence was obtained with as few as 25,000 iterations so the computing times are somewhat excessive. For those experiments that used a kernel, however, mixing was considerably less rapid, and parameter estimates exhibited some instability even at 10^6 iterations. However, in each case the number of test set errors had stabilized. One run of 10^7 iterations for the Crabs data showed identical test set errors after each million iterations.

Despite this computing effort, comparison of the predictive performance of the fully Bayesian versus plug-in performance yields mixed results. For the Pima data, performance was identical for the plug-in and fully Bayesian predictions. For the Crabs data, the plug-in predictions for kernel-based model outperformed the fully Bayesian predictions. Without a kernel, the two approaches performed identically. For the WBC data, the fully Bayesian predictions outperformed their plug-in counterparts.

We note that in the context of text categorization with a Bayesian multinomial model, Rennie (2001) reported that the Bayesian plug-in approach predictively outperformed the fully Bayesian approach.

3.2 Kernel or No Kernel?

Using a Gaussian kernel instead of the original input variables gives mixed results. For MAP plug-in predictions, use of the kernel provides improved predictive performance for Crabs and WBC but poorer performance for Pima. For fully Bayesian prediction, the kernel provides improved predictive performance for WBC, poorer performance for Pima, and identical performance for Crabs.

We only experimented with a Gaussian kernel, and used the bandwidths suggested by FJ. Different kernels or cross-validated bandwidth selection may show the kernel classifiers in a better light.

3.3 Jeffreys Prior or Exponential Prior?

Table 2 shows our results for the Laplacian prior for different settings of the hyperparameter γ .

Again the results are mixed. For the Pima data and no kernel, $\gamma = 10$ provides the best result. For the Pima data with a kernel, $\gamma = 4$ is optimal amongst the values tried so far, and $\gamma = 10$ performs relatively poorly. In every case (kernel/no kernel; plug-in/full Bayes), there exist settings of γ that out-perform the Jeffreys prior.

For the Crabs data, γ less than one is optimal in every case, and no setting for γ outperforms the Jeffreys prior.

For the WBC data, the results are qualitatively similar to the results with the Pima data.

In general, kernel-based prediction shows more sensitivity to the choice of γ than non kernel-based prediction. This is to be expected since shrinkage plays a greater role in the kernel-based models.

3.4 Other Comparisons

Table 3 reproduces the results from Seeger (2000) along with our own implementation of logistic regression via maximum likelihood and the best dataset-specific results from the tables above.

Linear discriminant analysis and logistic regression show relatively poor predictive performance with the WBC data, but otherwise no clear winner emerges.

4 Discussion

The probit model with the Laplacian prior is the lasso version of probit regression (Tibshirani, 1996, §5). Although software for the lasso has been available for many years, we are not aware of any large-scale evaluation of its predictive performance.

Taking advantage of the Gaussian-exponential hierar-

	Method	Kernel	Bayes?	γ	Pima	Crabs	WBC
1	<i>Probit</i>	<i>no</i>	<i>full</i>	<i>Jeffreys</i>	70	3	3.3%
2	Probit	no	full	0.1	66	4	-
3	Probit	no	full	0.5	65	3	-
4	Probit	no	full	1	65	3	2.8%
5	Probit	no	full	2	65	4	-
6	Probit	no	full	10	66	7	-
7	<i>Probit</i>	<i>no</i>	<i>plug-in</i>	<i>Jeffreys</i>	70	3	4.0%
8	Probit	no	plug-in	0.001	66	4	4.9%
9	Probit	no	plug-in	0.01	66	4	4.7%
10	Probit	no	plug-in	0.1	66	4	4.3%
11	Probit	no	plug-in	0.5	66	3	3.9%
12	Probit	no	plug-in	1	65	4	3.6%
13	Probit	no	plug-in	2	66	4	3.9%
14	Probit	no	plug-in	10	63	17	3.4%
15	<i>Probit</i>	<i>yes</i>	<i>full</i>	<i>Jeffreys</i>	72	3	2.8%
16	Probit	yes	full	1	-	6	-
17	<i>Probit</i>	<i>yes</i>	<i>plug-in</i>	<i>Jeffreys</i>	72	2	4.1%
18	Probit	yes	plug-in	0.001	91	3	4.1%
19	Probit	yes	plug-in	0.01	85	2	3.0%
20	Probit	yes	plug-in	0.1	74	8	3.5%
21	Probit	yes	plug-in	0.5	67	10	3.7%
22	Probit	yes	plug-in	1	70	64	4.4%
23	Probit	yes	plug-in	2	68	60	5.8%
24	Probit	yes	plug-in	4	64	*	7.1%
25	Probit	yes	plug-in	10	109	*	10.5%

Table 2: Number of test set errors (Pima and Crabs) or error rate (WBC). WBC results are averaged over 5 random splits into 300 training examples and 269 test examples. MCMC runs of length 10^6 in each case. The two runs marked with a “*” produced no non-zero parameter estimates. Runs marked with a “-” are in progress.

chical representation of the Laplacian prior, FJ replace the exponential component with a parameter-free Jeffreys prior. The only motivation FJ provided was to avoid setting a hyperparameter, and indeed we know of no theoretical justification for the Jeffreys prior in this context. Our experiments show that while the Jeffreys prior performs adequately in some situations, judicious hyperparameter choice can sometimes improve predictive performance.

Fully Bayesian versus Plug-in

The primary motivation of this work was to empirically compare the performance of plug-in (or MAP) prediction versus the more fully Bayesian procedure of integrating out the parameters to make predictions. The finding that the plug-in approach is competitive with its computationally intensive counterpart is perhaps surprising and certainly comforting for large scale applications such as text categorization where an MCMC approach is impractical.

There exists a small literature comparing the theo-

retical predictive performance of plug-in versus fully-Bayesian prediction. For gamma and multinormal models, Aitchison (1975) showed that fully Bayesian is better than plug-in from the point of view of minimizing the Kullback-Leibler distance between the true and estimated predictive distributions.

Smith (1999) considers the following setup: Let $X = (X_1, \dots, X_n)$ denote an i.i.d. sample from an exponential density $\theta e^{-\theta x}$, $x, \theta > 0$ and let Z denote a future observation from the same density. Consider the predictive probability that $Z > z$ for some z . When θ is known, this probability is given by $\phi = e^{-\theta z}$. Smith shows that, for mean square error loss, the plug-in estimate of ϕ beats the fully-Bayesian estimate with Jeffreys prior if $\phi < 0.1801$. However for other reasonable loss functions, fully-Bayesian dominates plug-in. Smith also presents a more general argument and concludes that “one should not abandon Bayesian methods in favor of much more simple-minded approaches such as [plug-in prediction], but that one needs to give much more careful attention to the formulation of the

	Method	Kernel	Bayes?	Pima	Crabs	WBC
1	SVM	yes	plug-in	64	4	3.3%
2	Linear Discriminant	yes	plug-in	67	3	7.1%
3	Gaussian Process	yes	plug-in	67	3	3.0%
4	Logistic Regression	no	mle	66	5	7.7%
5	Best Probit Result	n/a	n/a	63	2	2.8%

Table 3: Number of test set errors (Pima and Crabs) or error rate (WBC). Lines 1 through 3 are from Seeger (2000).

Bayesian approach and in particular the choice of loss function.”

If instead of predicting ϕ , one looks at the squared loss for the next single observation, then plug-in and fully-Bayesian with Jeffreys are identical.

Possible Extensions

Variational methods can provide a more computationally efficient approach to fully Bayesian inference than MCMC. Bishop and Tipping (2000) developed such methods for RVMs.

Albert and Chib (1993) discussed a number of extensions to the basic probit model, such as replacing the Gaussian link with a mixture of Gaussians, one possibility being a t -distribution with random degrees of freedom. In fact the usual logistic link is close to a t -distribution with approximately 8 degrees of freedom. Another extension concerns a hierarchical model where β is distributed *a priori* $N(\mathbf{A}\beta_0, \sigma^2\mathbf{I})$ with β_0 having dimension less than d . This represents the prior belief that β lie on linear subspace $\mathbf{A}\beta_0$. A prior distribution for (β_0, σ^2) completes the hierarchy.

Bayesian model averaging (see, for example, Madigan and York, 1995), averaging over different sets of input features may also improve predictive performance. Here, a complete MCMC scheme requires a somewhat complex dimension-jumping sampler (see Green, 1995). However, an algorithm that only uses MCMC to move through model space and uses EM for parameter estimation is substantially simpler and may be just as effective predictively. Raftery *et al.* (1996) used a similar approach in the context of survival analysis.

Clearly our findings are preliminary and much further experimentation and theoretical exploration is required. Thusfar our findings are that “results are mixed.” This seems to indicate that plug-in methods with default priors, while requiring relatively modest computational effort, may provide competitive predictive performance in many situations.

Software for both algorithms described in this paper

available from the second author.

Acknowledgements

We thank David D. Lewis, Colin Mallows, Greg Ridgeway, and Werner Stuetzle for helpful discussions.

References

- Aitchison, J. (1975). Goodness of prediction fit. *Biometrika*, **62**, 547–554.
- Albert, J.H. and Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, **88**, 669–679.
- Bishop, C.M. and Tipping, M. (2000). Variational relevance vector machines. In: *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, 46–53.
- Carlin, B.P. and Louis, T.A. (2000). *Bayes and Empirical Bayes Methods for Data Analysis (second edition)* Chapman and Hall, New York.
- Chambers, F.A. and Cox, D.R. (1967). Discrimination between alternative binary response models. *Biometrika*, **54**, 573–578.
- Figueiredo, M.A.T. (2001). Adaptive sparseness using Jeffreys prior. *Neural Information Processing Systems*, Vancouver, December 2001.
- Figueiredo, M.A.T. and Jain, A.K. (2001). Bayesian learning of sparse classifiers. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Hawaii, December 2001.
- Gelman, A., Carlin, J.B., Stern, H.S., and Rubin, D.B. (1995). *Bayesian Data Analysis*. Chapman and Hall, London.
- Gilks, W.R., Richardson, S., and Spiegelhalter, D.J. (1996). *Markov chain Monte Carlo in Practice*. Chapman and Hall, London.
- Girosi, F. (1998). An equivalence between sparse approximation and support vector machines. *Neural Computation*, **10**, 1445–1480.

- Green, P.J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, **82**, 711-32.
- Madigan, D. and York, J. (1995). Bayesian graphical models for discrete data. *International Statistical Review*, **63**, 215-232.
- Raftery, A.F., Madigan, D., and Volinsky, C.T. (1996). Accounting for model uncertainty in survival analysis improves predictive performance. In: Bernardo, J. M., Berger, J. O., Dawid, A. P. and Smith A. F. M., (eds.), *Bayesian Statistics V*, Oxford University Press, 323-350.
- Rennie, J.D.M. (2001). Improving Multi-class Text Classification with Naive Bayes. *AI Technical Report, Massachusetts Institute of Technology AITR-2001-004*.
- Seeger, M. (2000). Bayesian model selection for support vector machines, Gaussian Processes and other kernel classifier. In: *Advances in Neural Information Processing*, MIT Press, 603-609.
- Smith, R.L. (1999). Bayesian and frequentist approaches to parametric predictive inference (with discussion). In: *Bayesian Statistics 6*, edited by J.M. Bernardo, J.O. Berger, A.P. Dawid and A.F.M. Smith. Oxford University Press, pp. 589-612.
- Tibshirani, R. (1995). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society (Series B)*, **58**, 267-288.
- Tipping, M.F. (2001). Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, **1**, 211-244.
- Zhang, T. and Oles, F. (2001). Text categorization based on regularized linear classifiers. *Information Retrieval*, **4**, 5-31.