# A Bayesian Beta Kernel Model for Binary Classification and Online Learning Problems

**Cameron A. MacKenzie[1]\*, Theodore B. Trafalis[2] and Kash Barker[2]**

[1]*Defense Resources Management Institute, Naval Postgraduate School, Monterey, CA 93943, USA*

[2]*School of Industrial and Systems Engineering, University of Oklahoma, Norman, OK 73019, USA*

**Abstract:**　Recent advances in data mining have integrated kernel functions with Bayesian probabilistic analysis of Gaussian distributions. These machine-learning approaches can incorporate prior information with new data to calculate probabilistic rather than deterministic values for unknown parameters. This article extensively analyzes a specific Bayesian kernel model that uses a kernel function to calculate a posterior beta distribution that is conjugate to the prior beta distribution. Numerical testing of the beta kernel model on several benchmark datasets reveals that this model's accuracy is comparable with those of the support vector machine (SVM), relevance vector machine, naive Bayes, and logistic regression, and the model runs more quickly than all the other algorithms except for logistic regression. When one class occurs much more frequently than the other class, the beta kernel model often outperforms other strategies to handle imbalanced datasets, including under-sampling, over-sampling, and the Synthetic Minority Over-Sampling Technique. If data arrive sequentially over time, the beta kernel model easily and quickly updates the probability distribution, and this model is more accurate than an incremental SVM algorithm for online learning. © 2014 Wiley Periodicals, Inc. Statistical Analysis and Data Mining, 2014

**Keywords:**　data mining; kernel; Bayesian; beta distribution; online learning

## 1. INTRODUCTION

Since advances in the mid-1990s, kernel-based approaches to machine learning and pattern recognition have revolutionized the field of data mining [1]. Kernel functions map input data to a higher dimensional space, called the feature space, where the dot product between the two vectors in the feature space is replaced by a kernel function. This approach enables algorithms designed to detect linear relationships in data, such as support vector machines (SVMs) and least-squares regression, to detect nonlinear relationships and patterns through their use on the feature space [2–4].

More recently, kernel functions have been integrated with Bayesian analysis to produce a new subset of machine-learning tools that produce probabilistic rather than deterministic solutions. Probabilistic outcomes can better express uncertainty in underlying data relative to deterministic outcomes. Most previous Bayesian kernel models, such as the relevance vector machine (RVM), have assumed Gaussian prior distributions over model parameters [4–7]. Different approaches, such as assuming the mean and variance of the Gaussian distribution are randomly chosen from other distributions, have been deployed to increase the flexibility and accuracy of the Gaussian kernel model [8–10]. These approaches carry additional computational complexities and generally require simulation algorithms such as Markov Chain Monte Carlo to solve for the optimal parameters.

At least two issues can pose challenges to kernel-based binary classifiers, including the Gaussian Bayesian models. First, imbalanced datasets, where one class appears much more frequently than another class, create difficulties for many machine-learning algorithms because the classifiers tend to classify almost all of the unknown data points in the class that occurs most frequently. Second, most classifiers are designed to process data in a single batch, and a classifier may have difficulty in incrementally updating if data arrive sequentially. Although Bayesian models can

　\* *Correspondence to:* Cameron A. MacKenzie (camacken@nps.edu)

typically use new information to update a probability distribution, the complexity of Gaussian kernel models limits their ability to update quickly.

A Bayesian kernel model using the beta rather than the normal distribution can potentially address both of these challenges. The beta kernel model appears to have been first presented by Montesano and Lopes [11] in order to predict a robot's ability to grasp objects. The authors used empirical probabilities calculated from previous experiments with the robot to update the beta distribution. Although the beta kernel model was developed specifically for this robotic application, we believe the model deserves a fuller exploration. This article analyzes the beta kernel model for a wider range of binary classification problems where empirical probabilities are unavailable, the data are heavily imbalanced, and data arrive incrementally. Additionally, the beta kernel model does not require solutions to optimization problems such as the SVM or RVM, which makes the beta kernel model extremely fast to calculate.

This article offers several unique contributions to analyze whether beta kernel models should become part of the machine-learning toolkit for binary classification problems. We explicitly relate the beta kernel model to the beta-binomial Bayesian model and discuss how to select parameters for the prior distribution. We generalize the beta kernel model to a Dirichlet kernel model that could be used for multiclass classification problems. This article also explores the similarities of the beta kernel and the well-known Parzen [12] window classifier and discusses how the beta kernel model can overcome some of the difficulties with the Parzen window. Our inclusion of weighting parameters with the likelihood function in the beta kernel model or beginning with a nonuniform prior distribution increases the predictive accuracy for imbalanced datasets. Finally, the posterior probabilities from the model can act as prior probabilities to incorporate additional information, making the model useful for online or incremental learning.

Section 2 focuses on binary classification problems. We review the existing Gaussian kernel models such as the RVM and the beta kernel model presented by Montesano and Lopes [11]. The model can be extended to imbalanced datasets through the addition of weighting parameters or a nonuniform prior. Section 3 extensively tests the beta kernel model, the RVM, and the SVM for standard binary classification problems, heavily imbalanced datasets, and online learning. Tests on imbalanced datasets include additional algorithms such as under-sampling and over-sampling in combination with the RVM and SVM.

## 2. BAYESIAN MODELS

Binary classification machine-learning tools seek to assign an unknown data point $y$ either to the negative class,

$y = -1$, or to the positive class, $y = 1$. The assignment is based on the input data $\mathbf{x}$, a vector with $d$ components, where each component is an attribute. Rather than assigning $y$ to either the positive or negative class, Bayesian binary classification models calculate the probability that $y$ belongs to each class given $\mathbf{x}$. Because most Bayesian kernel models assume Gaussian prior distributions, we first present the basic Gaussian kernel model and the popular RVM (a variation of the basic model) and next the beta kernel model.

### 2.1. Gaussian Kernel Model

Gaussian Bayesian kernel models assume a function $t$ maps the input data $\mathbf{x}$ to a target value that corresponds to an output $y$, where $y \in \{-1, +1\}$. The range of $t(\mathbf{x})$ is the set of all real numbers, and the logit function maps $t(\mathbf{x})$ to a probability that $y = 1$.

$$P(y = 1 | t(\mathbf{x})) = \frac{1}{1 + \exp(-t(\mathbf{x}))} \tag{1}$$

If the $m \times d$ data matrix $\mathbf{X}$ has $m$ rows (observed data points) each with $d$ attributes, the function $\mathbf{t(X)}$ can be thought of as a random vector of length $m$. The Gaussian model assumes that $\mathbf{t}$ follows a multivariate normal distribution where $E[\mathbf{t}] = \mathbf{0}$ and $\text{Cov}(\mathbf{t}) = \mathbf{K}$ [4]. The matrix $\mathbf{K}$ is positive definite where $K_{ij}$ is the kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$ between the $i$th and $j$th data points.

$$P(\mathbf{t}) = \frac{1}{\sqrt{(2\pi)^m}} (\det \mathbf{K})^{-1/2} \exp\left(-\frac{1}{2}\mathbf{t}^{\mathrm{T}}\mathbf{K}^{-1}\mathbf{t}\right) \tag{2}$$

Calculating the inverse of $\mathbf{K}$ is computationally expensive, and a new vector $\boldsymbol{\omega}$ of length $m$ is introduced such that $t(\mathbf{x}_i) = \sum_{j=1}^{m} k(\mathbf{x}_i, \mathbf{x}_j)\omega_j = \mathbf{k}(\mathbf{x}_i, \mathbf{X})\boldsymbol{\omega}$ [4]. The prior probability function for $\boldsymbol{\omega}$ is also a multivariate normal distribution but eliminates the need to take the inverse of $\mathbf{K}$.

$$P(\boldsymbol{\omega}) = \frac{1}{\sqrt{(2\pi)^m}} (\det \mathbf{K})^{-1/2} \exp\left(-\frac{1}{2}\boldsymbol{\omega}^{\mathrm{T}}\mathbf{K}\boldsymbol{\omega}\right) \tag{3}$$

Because $\frac{1}{\sqrt{(2\pi)^m}}(\det \mathbf{K})^{-1/2}$ does not depend on $\boldsymbol{\omega}$, the prior can be written without this constant. With the likelihood from Eq. (1) and the prior from Eq. (3), the posterior becomes a function of $\boldsymbol{\omega}$ and the observed output values $\mathbf{y}$.

$$\begin{aligned}
P(\boldsymbol{\omega}|\mathbf{y}) \propto \prod_{i=1}^{m} &\left[\frac{1}{1 + \exp(-\mathbf{k}(\mathbf{x}_i, \mathbf{X})\boldsymbol{\omega})}\right]^{0.5+0.5y_i} \\
&\times \left[\frac{1}{1 + \exp(\mathbf{k}(\mathbf{x}_i, \mathbf{X})\boldsymbol{\omega})}\right]^{0.5-0.5y_i} \\
&\times \exp\left(-\frac{1}{2}\boldsymbol{\omega}^{\mathrm{T}}\mathbf{K}\boldsymbol{\omega}\right) \tag{4}
\end{aligned}$$

Most solution methods seek to maximize the posterior or equivalently, to minimize the negative of the log of the posterior [4]. The Newton-Raphson method can be used to find $\omega$ that maximizes the posterior probability in Eq. (4). The posterior probability is identical to the objective function in non-Bayesian kernel logistic regression, which can be solved quickly using a truncated Newton method [13].

An alternative to the logit link function in Eq. (1) is the probit model, which generally assumes that the sign of $t(\mathbf{x})$ determines the class of $y$, except for Gaussian noise. If $y = \text{sgn}(t(\mathbf{x} + \xi))$ where $\xi \sim \mathcal{N}(0, \sigma)$, then $P(y = 1|t(\mathbf{x})) = \Phi(yt(\mathbf{x})/\sigma)$ where $\Phi(\cdot)$ is the standard normal cumulative distribution [4]. Figueiredo [8] uses the probit model to build a hierarchical Bayesian kernel model where most of the weights—which play a role similar to $\omega$ in Eq. (3)—are 0.

## 2.2. Relevance Vector Machine

Although several variations on the Gaussian Bayesian kernel model exist (see for e.g., refs [8–10]), the most popular extension of this model is perhaps the RVM. The RVM seeks to combine the sparsity aspects of the SVM with the probabilistic advantages of Bayesian methods [6,7]. The SVM assigns nonzero weights to only a small fraction of the total number of data points, and the RVM seeks to find a similarly sparse set.

The RVM has the same logit likelihood function as in Eq. (1); however, the prior on $\omega$ becomes a function of a hyperparameter $\mathbf{s}$, where each $s_j$ is the inverse of the variance of $\omega_j$ [6]. The conditional probability density function of $\omega$ given $\mathbf{s}$ is expressed in Eq. (5) where $\mathbf{S} = \text{diag}(\mathbf{s})$ [7].

$$P(\omega|\mathbf{s}) = \frac{1}{\sqrt{(2\pi)^m}} (\det \mathbf{S})^{1/2} \exp\left(-\frac{1}{2}\omega^{\mathsf{T}}\mathbf{S}\omega\right) \quad (5)$$

Traditionally, each $s_j \geq 0$ is assumed to be drawn from the same gamma distribution as given in Eq. (6) where $\Gamma(\cdot)$ is the gamma function.

$$P(s_j) = \frac{b^a}{\Gamma(a)} s_j^{a-1} e^{-bs_j} \quad (6)$$

The shape and scale parameters, $a$ and $b$, respectively, are usually set close to 0 to ensure a flat or noninformative prior over $\omega$. If $a \to 0$ and $b \to 0$, the gamma distribution becomes degenerate such that $s_i = 0$ in the limit, which implies infinite variance for $\omega_i$ [4].

The RVM algorithm [14] begins by selecting the data point that maximizes the posterior probability, and it sequentially adds additional data points until the posterior probability begins to converge to a value. Because $\text{E}[\omega_i] =$ 0, only a few $\omega_i$'s are nonzero. Testing the RVM on several databases reveals that approximately 10% of the $\omega_i$'s are nonzero and the RVM's error rate is comparable to that of the SVM [6,7]. The RVM has been applied to a wide variety of problems such as analyzing remote sensor data [15], forecasting stock indices [16], and estimating battery reliability [17].

The Gaussian distribution enables analytically tractable solutions, and the RVM delivers sparse solutions. As justification for a normal distribution over the prior, each target value $t_i$ is a linear-weighted sum of a multivariate random vector $\omega$ where $\omega \sim N(\mathbf{0}, \mathbf{K}^{-1})$ for the basic Gaussian kernel model and $\omega \sim N(\mathbf{0}, \mathbf{S}^{-1})$ in the RVM. Thus, each $t_i$ also follows a normal distribution.

## 2.3. Beta Kernel Model

Alternatively, the beta distribution can serve as a prior distribution for Bayesian models. The beta distribution is a natural model for binary outcomes because the distribution's two parameters can represent the number of times each outcome has occurred or is expected to occur [18], and the distribution provides the probability distribution of a parameter $\theta \in [0, 1]$ where $\theta = P(y = 1)$ and $1 - \theta = P(y = -1)$. We begin the development of the beta distribution as a prior distribution for classification problems by examining the model as applied to predict a robot's ability to grasp objects. Generalizing this model and adding weighting parameters enables us to apply this model to a broad spectrum of binary classification problems.

The beta distribution can model a Bernoulli process where $\alpha$ is the number of successes and $\beta$ is the number of failures. If the prior on $\theta$ follows a beta distribution with prior parameters $\alpha$ and $\beta$ and the Bernoulli process results in $r$ successes out of $m$ trials, which implies a binomial likelihood function, the posterior distribution over $\theta$ also follows a beta distribution.

$$\theta|r \sim \text{beta}(\alpha + r, \ \beta + m - r) \quad (7)$$

The beta distribution is known as a conjugate prior because the posterior is also a beta distribution.

Montesano and Lopes [11] and Mason and Lopes [19] exploit the beta-binomial relationship by adapting it to a data mining problem that slightly differs from the binary classification problem we have been discussing. In their problem, each $\mathbf{x}_j$ has $N_j$ trials, with $R_j$ positive classifications and $U_j$ negative classifications, where $R_j + U_j = N_j$. An empirical probability of a positive classification $y_j = 1$ exists for each data point, $\hat{\theta}_j = R_j/N_j$.

For a data point $\mathbf{x}_i$ whose empirical probability is unobserved or unknown, the kernel function, $k(\mathbf{x}_i, \mathbf{x}_j)$, serves as a measure of similarity between $\mathbf{x}_i$ and $\mathbf{x}_j$. In

this manner, the kernel function can be used to estimate the posterior distribution for $\theta_i$ based on the observed $R_j$ and $U_j$ for $m$ data points.

$$\theta_i | R_j, U_j \sim \text{beta} \left( \alpha + \sum_{j=1}^m k\left(\mathbf{x}_i, \mathbf{x}_j\right) R_j, \beta + \sum_{j=1}^m k\left(\mathbf{x}_i, \mathbf{x}_j\right) U_j \right) \tag{8}$$

The most likely estimate $\overline{\theta}_i = \text{E}\left[P(y_i = 1)\right]$ is the expected value of the beta posterior distribution.

$$\overline{\theta}_i = \frac{\alpha + \sum_{j=1}^m k\left(\mathbf{x}_i, \mathbf{x}_j\right) R_j}{\alpha + \sum_{j=1}^m k\left(\mathbf{x}_i, \mathbf{x}_j\right) R_j + \beta + \sum_{j=1}^m k\left(\mathbf{x}_i, \mathbf{x}_j\right) U_j} \tag{9}$$

An optimization algorithm selects parameters for the kernel function by minimizing the squared difference between the $\overline{\theta}_j$ as calculated by Eq. (9) and the empirical probability $\hat{\theta}_j$ for a training set [11]. This Bayesian kernel regression allows the authors to measure a robot's ability to grasp different objects (which have a set of features) and then use those results to estimate the probability that a robot will grasp a different object that has not been empirically measured.

Empirical probabilities are generally not available for binary classification problems because each data point $\mathbf{x}_j$ often has a single trial rather than multiple trials, i.e. $N_j = 1$ for all $j = 1, \ldots, m$. Another potential problem that plagues many binary classifiers is inaccurate classifications if the two classes do not contain the same number of data points [20–22]. Many machine-learning tools label too many points in the class that occurs most frequently. We resolve this problem of imbalanced datasets by including weighting parameters $m_-/m$ and $m_+/m$, where $m_-$ and $m_+$ are the number of negative and positive labels, respectively, in the training set.

$$\theta_i|\mathbf{y} \sim \text{beta} \left( \alpha + \frac{m_-}{m} \sum_{\{j|y_j=1\}} k(\mathbf{x}_i, \mathbf{x}_j), \right.$$
$$\left. \beta + \frac{m_+}{m} \sum_{\{j|y_j=-1\}} k(\mathbf{x}_i, \mathbf{x}_j) \right) \tag{10}$$

These weighting parameters balance the likelihood function so that the posterior probability moves to whichever class is closest to $\mathbf{x}_i$ in the feature space. These parameters follow typical cost penalties or weights for the SVM [23,24]. In a weighted SVM, the ratio of misclassification in the positive class to misclassification in the negative class weight often is $m_-/m_+$.

As a Bayesian formulation, Eq. (10) calculates a probability distribution over $\theta_i = P(y_i = 1)$ based on a set of known data points. The parameters of the beta distribution are updated for $\theta_i$ so that $\alpha_i = \alpha + (m_-/m) \sum_{\{j|y_j=1\}} k(\mathbf{x}_i, \mathbf{x}_j)$ and $\beta_i = \beta + (m_+/m) \sum_{\{j|y_j=-1\}} k(\mathbf{x}_i, \mathbf{x}_j)$. The expected value of the posterior distribution $\overline{\theta}_i = \alpha_i/(\alpha_i + \beta_i)$ can be used for making predictions.

This beta kernel model can be generalized to multiclass classification problems by deploying a Dirichlet distribution. The Dirichlet distribution describes the probability of $N$ discrete outcomes. Given a Dirichlet prior with prior parameters $\alpha_1, \alpha_2, \ldots, \alpha_N$, the weighted kernel approach in Eq. (11) derives the posterior probability for $\theta_i$, which also follows a Dirichlet distribution. The weight to update $\alpha_n$ is given by $m_{-n}/m$, the fraction of points not in the $n$th class where $n = 1, 2, \ldots, N$.

$$\theta_i|\mathbf{y} \sim \text{Dir} \left( \alpha_1 + \frac{m_{-1}}{m} \sum_{\{j|y_j=1\}} k(\mathbf{x}_i, \mathbf{x}_j), \ldots, \alpha_N \right.$$
$$\left. + \frac{m_{-N}}{m} \sum_{\{j|y_j=N\}} k(\mathbf{x}_i, \mathbf{x}_j) \right) \tag{11}$$

The expected probability that an unknown data point $y_i$ belongs to the $n$th class is given by $\text{E}\left[\theta_{i,n}\right] = \frac{\alpha_{i,n}}{\alpha_{i,0}}$ where $\alpha_{i,n} = \alpha_n + (m_{-n}/m) \sum_{\{j|y_j=n\}} k(\mathbf{x}_i, \mathbf{x}_j)$ and $\alpha_{i,0} = \sum_{l=1}^N \alpha_{i,l}$. The posterior marginal distribution for each $\theta_{i,n}$ follows a beta distribution where $\theta_{i,n} \sim \text{beta}\left(\alpha_{i,n}, \alpha_{i,0} - \alpha_{i,n}\right)$.

This article focuses on the specific case of the beta distribution and binary classification problems and examines numerical results for the beta kernel model. Future research can test the more general Dirichlet kernel model for multiclass classification problems.

We have previously been assuming that the prior $\alpha$ and $\beta$ are given, and neither Montesano and Lopes [11] nor Mason and Lopes [19] devote much discussion to selecting prior distributions. If nothing is known *a priori* about the probability of a positive or negative class, choosing a uniform prior where $\alpha = 1$ and $\beta = 1$ may be the best. However, many situations arise where the overall probability of an event may be estimated based on past experience or data. For example, if the data mining problem is to predict breast cancer in women age 40 to 49, we might know that a randomly selected female in her 40s has a 2% chance of being diagnosed with breast cancer. The prior distribution can reflect that knowledge through the prior mean $\frac{\alpha}{(\alpha+\beta)} = 0.02$.

Data from the training set can be used to generate a prior distribution. The parameters $\alpha$ and $\beta$ can be selected so that the mean equals the fraction of positively classified data points in the training set $\frac{m_+}{m}$ and so that the variance of the

prior $\frac{\alpha\beta}{\left[(\alpha+\beta)^2(\alpha+\beta+1)\right]}$ matches the variance in the training set. An empirical Bayes method estimates $\alpha$ and $\beta$ by maximizing the log-likelihood of the beta distribution. The Jeffreys prior where $\alpha = 0.5$ and $\beta = 0.5$ has modes at the two ends of the distribution 0 and 1, but if a large number of data points exist in the training set (e.g., $m > 20$), the posterior distribution from a Jeffreys prior will closely resemble the posterior from a uniform prior. (Carlin and Louis [25] present a helpful discussion on selecting priors for Bayesian analysis.)

The prior may also influence if the weighting parameters, $m_-/m$ and $m_+/m$, are used. Probabilistic data mining models often include a threshold probability that divides the positive and negative classes. With a uniform prior, if the expected value of the posterior distribution is greater than 0.5, the unknown point should be positively classified. The weighting parameters help ensure that this classification is due to the unknown data point's similarity with the known data as opposed to one class being more numerous than the other class. Including weights is not appropriate with a nonuniform prior distribution, however, because the weights would generate a posterior whose expectation is too close to the class with fewer data points. If a nonuniform prior is used, no weighting parameters are necessary. Instead, the expectation of the prior distribution can be used as the threshold probability so that if the expectation of the posterior is greater than the expectation of the prior, the point should be positively classified. Choosing between a uniform prior with weighting parameters and a nonuniform prior without weighting parameters can result in the same classification scheme, as the follow proposition illustrates.

**PROPOSITION 1.**

The following classification rules are equivalent.

1. Given a uniform prior where $\alpha = 1$ and $\beta = 1$ and the weights as depicted in Eq. (10), an unknown point $y_i$ should be positively classified if $\overline{\theta}_i > 0.5$ and negatively classified if $\overline{\theta}_i < 0.5$.

2. Given a nonuniform prior where $\frac{\alpha}{(\alpha+\beta)} = m_+/m$ and if weights are not deployed to update $\theta_i$, an unknown point $y_i$ should be positively classified if $\overline{\theta}_i > m_+/m$ and negatively classified if $\overline{\theta}_i < m_+/m$.

**Proof:** See Appendix.

Although a uniform prior and weights will generate different posterior probabilities than a nonuniform prior without weights, the final classification schemes will be the same if the expected value of the nonuniform prior equals $\frac{m_+}{m}$.

The beta kernel model closely resembles the Parzen [12] window classifier, which classifies an unknown data point based on whether it is closer in the feature space to the mean of the positive or negative class. An unknown data point $y_i$ is positively classified if $\frac{\sum_{\{j|y_j=1\}} k(\mathbf{x}_i,\mathbf{x}_j)}{m_+} > \frac{\sum_{\{j|y_j=-1\}} k(\mathbf{x}_i,\mathbf{x}_j)}{m_-}$. An alternative interpretation of the Parzen window is that the unknown data point is labeled according to whether the probability is greater for the positive or negative class [26]. The probability that $y_i$ is in the positive class is $\frac{\sum_{\{j|y_j=1\}} k(\mathbf{x}_i,\mathbf{x}_j)}{\sum_{j=1}^m k(\mathbf{x}_i,\mathbf{x}_j)}$. The Parzen window probability echoes the expectation of an unweighted posterior beta distribution.

A problem with using the Parzen window probability and the expectation of the beta posterior distribution occurs when a point $\mathbf{x}_i$ is far away in the feature space from both positively and negatively classified points. If the point is slightly closer to the positive class, the Parzen window and the expectation of the beta distribution will likely calculate a high probability that $y_i = 1$ [26]. In the Bayesian model, the posterior's variance can express uncertainty in the classifier. If $\mathbf{x}_i$ is dissimilar to the known points, both $\sum_{\{j|y_j=1\}} k(\mathbf{x}_i, \mathbf{x}_j)$ and $\sum_{\{j|y_j=-1\}} k(\mathbf{x}_i, \mathbf{x}_j)$ should be close to zero, which implies the variance of the posterior distribution $P(y_i|\mathbf{y})$ will be large, assuming $\alpha$ and $\beta$ are relatively small. Although the posterior distribution's expectation may be high, the large variance would indicate a significant likelihood that the unknown point should be negatively classified.

## 3. NUMERICAL RESULTS

### 3.1. Binary Classification for the Beta Kernel, RVM, and SVM

We test the beta kernel model on several datasets and compare the results to the RVM, the traditional soft-margin SVM [1,2], a weighted soft-margin SVM [27], the naive Bayes algorithm, and logistic regression. The SVM is a kernel-based linear classifier that uses a relatively small number of vectors to create a boundary between the classes in the feature space. The soft-margin SVM assigns a cost parameter for misclassifications. In the weighted SVM, we assign a different cost for the misclassification of each class: $\frac{Cm_-}{m}$ for the positive class and $\frac{Cm_+}{m}$ for the negative class where $C$ is a constant cost parameter to be optimized. We use LIBSVM 3.0 [28] for the SVM models and the code developed by Tipping [6] for the RVM. The classification rule for the beta kernel model follows Proposition 1. Matlab [29] provides a naive Bayes classification algorithm that uses a kernel smoothing density estimate based on a normal distribution. The logistic regression algorithm is a non-kernel generalized linear regression model with a logit

**Table 1.**　Binary classification datasets.

| Dataset | Percentage of positively labeled points | Number of attributes | Training set size | Tuning set size | Testing set size |
|---|---|---|---|---|---|
| Parkinson | 75 | 22 | 98 | 39 | 58 |
| Haberman's survival | 74 | 3 | 153 | 61 | 92 |
| Satellite | 53 | 36 | 509 | 203 | 305 |
| Arcene | 44 | 9961 | 100 | 40 | 60 |
| Spam | 39 | 57 | 230 | 92 | 138 |
| Colon cancer | 35 | 2000 | 31 | 12 | 19 |
| Adult | 24 | 13 | 799 | 320 | 480 |
| Transfusion | 24 | 4 | 374 | 150 | 224 |
| Breast cancer | 20 | 32 | 69 | 28 | 41 |
| Tornado | 7 | 83 | 541 | 216 | 325 |

link function that positively labels an unknown data point if the calculated probability exceeds the fraction of positively classified data points in the training set.

Table 1 shows characteristics of the 10 datasets used for comparing among the different classifiers. The Parkinson dataset contains biomedical voice measurements that correspond to individuals with Parkinson's disease and those without the disease [30]. Haberman's survival database consists of patients who survived or died after undergoing surgery for breast cancer [31]. The satellite dataset contains spectral values for pixels in order to classify land images as red or gray soil. The arcene data contain patients with ovarian or prostrate cancer and healthy patients, where the attributes are mass-spectrometry features [32]. The spam database is a collection of emails classified as either spam or not spam, and the attributes consist of frequency counts of words and characters. The adult data consist of census information from the U.S. 1994 census to predict whether or not an individual earns more or less than $50,000 [33]. The transfusion dataset contains blood donor attributes to predict whether or not an individual donated blood in a specific month [34]. The breast cancer data use image characteristics of breast mass to predict if women will have a recurrence of breast cancer within two years of treatment [35]. These eight datasets can be downloaded from the University of California-Irvine Machine Learning Repository [36]. The colon cancer data consist of gene expressions that either come from a tumor biopsy or a healthy biopsy, and this data derives from the Princeton University Gene Expression Project [37]. Finally, the tornado database contains weather characteristics corresponding to a tornado or no tornado as collected by the National Weather Center at the University of Oklahoma [38].

With one exception, the radial basis function (RBF) is used as the kernel function throughout this article, where $\sigma > 0$ is tuned to optimize each classifier.

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{||\mathbf{x}_i - \mathbf{x}_j||^2}{2\sigma^2}\right) \qquad (12)$$

The RBF is perhaps the most popular kernel function because the image of the function lies between zero and one and the kernel matrix has full rank (Schölkopf and Smola 2002).

The polynomial kernel has performed well in classifying text and spam data [39,40]. In addition to the RBF, we use the polynomial kernel on the spam dataset for the beta kernel, RVM, and SVM algorithms. The polynomial kernel has two parameters: the degree of the polynomial function $p > 0$ and a constant parameter $\kappa > 0$.

$$k(\mathbf{x}_i, \mathbf{x}_j) = \left(\mathbf{x}_i^{\mathsf{T}}\mathbf{x}_j + \kappa\right)^p \qquad (13)$$

Each of the 10 datasets is divided into a training, tuning, and testing set. The training set comprises 50% of each dataset, the tuning set 20%, and the testing set 30%. In each individual trial, $\sigma$ in the RBF or $p$ and $\kappa$ in the polynomial kernel (as well as the cost parameter $C$ in the SVM) is selected that achieves the highest accuracy score in the tuning set. (The accuracy score is described in the next paragraph.) The training and tuning set are combined to retrain the classifier using the optimal $\sigma$ or $p$ and $\kappa$ (and $C$) and test it on the testing set.

We repeat this procedure 200 times for each classifier, randomly selecting the training, tuning, and testing set for each trial. Table 2 displays the mean performance for the true positive (TP) rate, the true negative (TN) rate, and the accuracy score where $\mathrm{Acc} = \sqrt{\mathrm{TP} * \mathrm{TN}}$ is the geometric mean [41]. The geometric mean explicitly penalizes a classifying algorithm that performs badly in classifying one of the classes. The receiver operating characteristic (ROC) curve can judge the performance of probabilistic classifiers such as the beta kernel and RVM and can be extended to include deterministic classifiers such as the SVM. We also calculate the area under the ROC curve (AUC) for each of the six classifiers for each of the 200 runs.

The results show that the beta kernel model performs comparably well to the other data mining algorithms for these 10 datasets. The beta kernel model has the highest

**Table 2.**  Binary classification results.

| Dataset | Performance metric | Beta kernel | RVM | Traditional SVM | Weighted SVM | Naive Bayes | Logistic regression |
|---|---|---|---|---|---|---|---|
| Parkinson | Acc | **0.870** | 0.786 | 0.869 | 0.863 | 0.760 | 0.758 |
| | TP | 0.80 | 0.93 | 0.96 | 0.92 | 0.76 | 0.82 |
| | TN | 0.95 | 0.68 | 0.80 | 0.83 | 0.77 | 0.71 |
| | AUC | **0.946*** | 0.925 | 0.876 | 0.872 | 0.770 | 0.864 |
| Haberman's survival | Acc | 0.566 | 0.441 | 0.436 | 0.589 | 0.469 | **0.625**** |
| | TP | 0.81 | 0.92 | 0.86 | 0.73 | 0.91 | 0.78 |
| | TN | 0.40 | 0.22 | 0.24 | 0.50 | 0.27 | 0.51 |
| | AUC | 0.670 | 0.669 | 0.543 | 0.607 | 0.556 | **0.679** |
| Satellite | Acc | 0.987 | 0.985 | **0.988** | 0.988 | 0.971 | 0.978 |
| | TP | 0.98 | 0.98 | 0.99 | 0.98 | 0.96 | 0.98 |
| | TN | 0.99 | 0.99 | 0.99 | 0.99 | 0.98 | 0.98 |
| | AUC | **0.999** | 0.998 | 0.988 | 0.988 | 0.971 | 0.988 |
| Arcene | Acc | 0.783 | 0.753 | **0.842** | 0.840 | — | 0.453 |
| | TP | 0.93 | 0.73 | 0.86 | 0.86 | — | 0.50 |
| | TN | 0.66 | 0.79 | 0.83 | 0.82 | — | 0.50 |
| | AUC | 0.830 | **0.846** | 0.844 | 0.843 | — | 0.503 |
| Spam (RBF) | Acc | 0.815 | 0.874 | 0.883 | **0.885** | 0.875 | 0.860 |
| | TP | 0.83 | 0.82 | 0.83 | 0.84 | 0.81 | 0.85 |
| | TN | 0.83 | 0.93 | 0.94 | 0.93 | 0.95 | 0.87 |
| | AUC | 0.929 | **0.942**** | 0.888 | 0.891 | 0.863 | 0.867 |
| Spam (polynomial) | Acc | **0.880** | 0.839 | 0.865 | 0.879 | 0.677 | 0.855 |
| | TP | 0.88 | 0.82 | 0.79 | 0.83 | 0.98 | 0.84 |
| | TN | 0.88 | 0.87 | 0.95 | 0.94 | 0.49 | 0.88 |
| | AUC | **0.937** | 0.933 | 0.868 | 0.880 | 0.682 | 0.862 |
| Colon cancer | Acc | **0.800** | 0.691 | 0.799 | 0.771 | 0.490 | 0.598 |
| | TP | 0.81 | 0.64 | 0.75 | 0.75 | 0.37 | 0.56 |
| | TN | 0.80 | 0.80 | 0.88 | 0.85 | 0.93 | 0.68 |
| | AUC | **0.867**** | 0.808 | 0.814 | 0.801 | 0.646 | 0.641 |
| Adult | Acc | 0.782 | 0.726 | 0.720 | 0.799 | 0.798 | **0.807*** |
| | TP | 0.81 | 0.57 | 0.57 | 0.83 | 0.76 | 0.82 |
| | TN | 0.76 | 0.93 | 0.92 | 0.77 | 0.84 | 0.80 |
| | AUC | 0.867 | 0.890 | 0.742 | 0.800 | 0.799 | **0.896** |
| Transfusion | Acc | 0.649 | 0.532 | 0.525 | 0.664 | 0.537 | **0.689**** |
| | TP | 0.66 | 0.31 | 0.31 | 0.65 | 0.33 | 0.76 |
| | TN | 0.66 | 0.94 | 0.91 | 0.68 | 0.89 | 0.63 |
| | AUC | 0.734 | 0.741 | 0.607 | 0.669 | 0.607 | **0.750** |
| Breast cancer | Acc | **0.615*** | 0.335 | 0.470 | 0.571 | 0.559 | 0.503 |
| | TP | 0.65 | 0.17 | 0.29 | 0.47 | 0.44 | 0.39 |
| | TN | 0.61 | 0.95 | 0.89 | 0.76 | 0.79 | 0.71 |
| | AUC | **0.657*** | 0.628 | 0.578 | 0.608 | 0.610 | 0.579 |
| Tornado | Acc | **0.862** | 0.772 | 0.794 | 0.857 | 0.848 | 0.795 |
| | TP | 0.77 | 0.61 | 0.64 | 0.82 | 0.75 | 0.66 |
| | TN | 0.97 | 0.99 | 0.99 | 0.92 | 0.97 | 0.97 |
| | AUC | **0.959** | 0.949 | 0.816 | 0.869 | 0.858 | 0.845 |

*Notes:* Emdash (—) indicates that algorithm fails to terminate in a reasonable amount of time.
*The difference in accuracy is significant at the 0.1 level.
**The difference in accuracy is significant at the 0.01 level.
Bold indicates the best performance metric among the six classifiers.

average accuracy for five of the eleven data trials. (The spam dataset is tested twice, once with the RBF kernel and once with the polynomial kernel.) Because we are making multiple comparisons of mean accuracy levels where the sample sizes are equal (200 repetitions), we use Tukey's method to assess the statistical significance of the difference between the best performing classifier and the other classifiers. The beta kernel model's mean accuracy

is significantly different at the 0.1 level for the breast cancer dataset. Logistic regression has the highest average accuracy for three datasets (Haberman's survival, adult, and transfusion) and the difference in accuracy is significant at the 0.01 level for two datasets. However, logistic regression performs quite badly for the arcene and breast cancer data. The traditional SVM has the highest average accuracy for the satellite and arcene data, and the weighted SVM

**Table 3.**   Average run-time in seconds for binary classification models.

| Dataset | Beta kernel | RVM | Traditional SVM | Weighted SVM | Naive Bayes | Logistic regression |
|---|---|---|---|---|---|---|
| Parkinson | 0.12 | 8.67 | 1.78 | 1.89 | 8.76 | 0.03 |
| Haberman's survival | 0.26 | 3.06 | 2.79 | 2.49 | 1.34 | 0.06 |
| Satellite | 2.77 | 134.57 | 15.16 | 15.55 | 35.65 | 0.80 |
| Arcene | 6.30 | 46.23 | 488.25 | 489.94 | — | 0.26 |
| Spam | 0.66 | 73.18 | 11.08 | 11.61 | 25.39 | 0.26 |
| Colon cancer | 0.12 | 3.36 | 10.62 | 10.51 | 723.51 | 0.07 |
| Adult | 6.30 | 239.46 | 50.27 | 51.66 | 21.52 | 0.01 |
| Transfusion | 1.30 | 8.75 | 13.20 | 10.85 | 3.26 | 0.34 |
| Breast cancer | 0.08 | 5.36 | 1.54 | 1.56 | 11.82 | 0.02 |
| Tornado | 3.86 | 440.73 | 46.38 | 62.06 | 74.97 | 0.23 |

*Notes:* Emdash (—) indicates that algorithm fails to terminate in a reasonable amount of time

has the highest average accuracy for spam with the RBF kernel.

The beta kernel model performs quite well according to the AUC metric. The beta kernel has the largest average AUC for six datasets (Parkinson, satellite, spam with the polynomial kernel, colon cancer, breast cancer, and tornado). Two of those sets have differences in the means that are significant at the 0.1 level, and one is significant at the 0.01 level. Logistic regression has the largest AUC for three datasets, corresponding to the datasets for which it had the best average accuracy, and RVM has the highest average AUC for two datasets (arcene and spam with RBF kernel). Because the SVM does not calculate probabilities, it performs worse on this metric.

We record the average run-time from the 200 different runs for each algorithm in Table 3. Logistic regression is the fatest algorithm because it does not have any parameters to tune. The beta kernel's most complicated operation is calculating the kernel matrix, and the average run-time is a few seconds. The SVM algorithms are also very fast, but the SVM algorithm has two parameters to tune ($\sigma$ and $C$) which multiplies the number of times the SVM optimization problem is solved. Although the RVM usually relies on fewer vectors than the SVM to classify unknown data points, the RVM algorithm as developed by [14] sequentially selects vectors $\mathbf{x}_j$ that are used to classify unknown data points. This sequential process can take a long time, as with the satellite, adult, and tornado datasets. The naive Bayes algorithm takes a very long time when the data have a large number of attributes, and the arcene dataset has so many attributes that the naive Bayes algorithm does not terminate in a reasonable amount of time.

### 3.2. Imbalanced Datasets

Imbalanced datasets can cause problems for data mining and statistical learning because algorithms like the SVM

tend to classify almost all of data points in the class that occurs most frequently, the majority class [42–45]. Under-sampling the majority class or over-sampling the minority class so that both classes have an equal number of data points when training the algorithm can help the machine-learning tool more accurately identify data points in the minority class [46,47]. The Synthetic Minority Over-sampling Technique (SMOTE) creates new data points through linear combinations of two input data points belonging to the minority class [48,49] so that the minority class has an equivalent number of points as the majority class. These sampling techniques can be used in combination with machine-learning algorithms such as the SVM and RVM [50].

Although none of the datasets used previously have the same number of data points in both classes, we create much more imbalanced datasets from these 10 datasets. We randomly select only 5% in the minority class for the eight databases. For example, the Parkinson data have 75% positively labeled and 25% negatively labeled points, and the new imbalanced Parkinson data have 95% positive and 5% negative. Conversely, 95% of the data are negatively classified and 5% are positively classified in the breast cancer dataset. Because the colon cancer data only contain 62 data points, too few negatively classified data points exist to have the 19:1 ratio between the negative and positive class in a training, tuning, and testing set. Therefore, the imbalanced colon cancer data contain 10% positive and 90% negative.

The beta kernel model with the weighting parameters is compared to several other methods designed to address these heavily imbalanced datasets. In addition to the weighted SVM (which was previously deployed), these other methods are under-sampling the majority class, over-sampling the minority class, and SMOTE. Each of these three sampling techniques are applied separately to the RVM and SVM. We test the eight different classifiers for the 10 heavily imbalanced datasets, with the same rules as

**Table 4.** Binary classification results with 5% in the minority class.

| Dataset | Performance metric | Beta kernel | Weighted SVM | Under-sampling | | Over-sampling | | SMOTE | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | RVM | SVM | RVM | SVM | RVM | SVM |
| Parkinson | Acc | **0.743**** | 0.586 | 0.529 | 0.578 | 0.513 | 0.578 | 0.604 | 0.449 |
| | TP | 0.81 | 0.85 | 0.69 | 0.77 | 0.95 | 0.92 | 0.92 | 0.97 |
| | TN | 0.78 | 0.58 | 0.63 | 0.63 | 0.44 | 0.52 | 0.54 | 0.38 |
| Haberman's survival | Acc | 0.473 | 0.378 | **0.475** | 0.467 | 0.329 | 0.339 | 0.420 | 0.290 |
| | TP | 0.72 | 0.76 | 0.59 | 0.62 | 0.84 | 0.82 | 0.84 | 0.89 |
| | TN | 0.41 | 0.35 | 0.52 | 0.49 | 0.27 | 0.26 | 0.32 | 0.22 |
| Satellite | Acc | **0.975*** | 0.945 | 0.909 | 0.946 | 0.931 | 0.943 | 0.937 | 0.940 |
| | TP | 0.98 | 0.99 | 0.96 | 0.98 | 0.99 | 1.00 | 0.99 | 0.99 |
| | TN | 0.97 | 0.91 | 0.88 | 0.92 | 0.92 | 0.90 | 0.89 | 0.91 |
| Arcene | Acc | **0.502**** | 0.115 | 0.263 | 0.000 | 0.118 | 0.115 | 0.119 | 0.115 |
| | TP | 0.64 | 0.12 | 0.42 | 1.00 | 0.12 | 0.12 | 0.12 | 0.12 |
| | TN | 0.59 | 1.00 | 0.63 | 0.01 | 1.00 | 1.00 | 1.00 | 1.00 |
| Spam (RBF) | Acc | 0.528 | 0.576 | 0.673 | **0.691** | 0.500 | 0.510 | 0.656 | 0.660 |
| | TP | 0.41 | 0.49 | 0.70 | 0.65 | 0.38 | 0.38 | 0.56 | 0.58 |
| | TN | 0.95 | 0.91 | 0.74 | 0.83 | 0.95 | 0.97 | 0.90 | 0.89 |
| Spam (polynomial) | Acc | **0.755**** | 0.454 | 0.640 | 0.350 | 0.592 | 0.445 | 0.648 | 0.584 |
| | TP | 0.78 | 0.34 | 0.66 | 0.29 | 0.53 | 0.31 | 0.63 | 0.44 |
| | TN | 0.77 | 0.98 | 0.74 | 0.94 | 0.82 | 0.98 | 0.79 | 0.96 |
| Colon cancer[a] | Acc | **0.529** | 0.292 | 0.273 | 0.454 | 0.198 | 0.220 | 0.189 | 0.197 |
| | TP | 0.55 | 0.37 | 0.73 | 0.60 | 0.20 | 0.23 | 0.20 | 0.20 |
| | TN | 0.84 | 0.83 | 0.44 | 0.64 | 0.99 | 0.99 | 0.98 | 0.99 |
| Adult | Acc | 0.760 | **0.765** | 0.758 | 0.758 | — | 0.743 | — | 0.690 |
| | TP | 0.74 | 0.75 | 0.77 | 0.79 | — | 0.68 | — | 0.55 |
| | TN | 0.80 | 0.79 | 0.74 | 0.75 | — | 0.83 | — | 0.89 |
| Transfusion | Acc | **0.634** | 0.612 | 0.585 | 0.591 | 0.581 | 0.608 | 0.354 | 0.468 |
| | TP | 0.68 | 0.63 | 0.62 | 0.64 | 0.54 | 0.58 | 0.22 | 0.32 |
| | TN | 0.62 | 0.66 | 0.62 | 0.59 | 0.70 | 0.69 | 0.89 | 0.84 |
| Breast cancer | Acc | 0.295 | 0.174 | 0.262 | **0.328** | 0.114 | 0.203 | 0.123 | 0.111 |
| | TP | 0.54 | 0.42 | 0.44 | 0.52 | 0.12 | 0.27 | 0.14 | 0.12 |
| | TN | 0.35 | 0.51 | 0.57 | 0.43 | 0.97 | 0.73 | 0.81 | 0.95 |
| Tornado | Acc | 0.800 | 0.761 | **0.806** | 0.800 | 0.744 | 0.778 | 0.742 | 0.772 |
| | TP | 0.67 | 0.66 | 0.82 | 0.79 | 0.64 | 0.64 | 0.67 | 0.70 |
| | TN | 0.98 | 0.95 | 0.81 | 0.86 | 0.92 | 0.97 | 0.88 | 0.94 |

*Notes:* Emdash (—) indicates that algorithm fails to terminate in a reasonable amount of time.
[a]Because colon cancer is a small dataset, the imbalanced version has 10% positive.
**The difference in accuracy with the other models is significant at the 0.01 level.
*The difference in accuracy with the other models is significant at the 0.1 level.
Bold indicates the best accuracy among the eight classifiers.

before for training, tuning, and testing. Each algorithm is tested 200 times for each dataset.

With the imbalanced data, the beta kernel model performs the best in 6 of the 11 data trials (Table 4), and the differences in means are statistically significant at the 0.01 level for three datasets (Parkinson, arcene, and spam with the polynomial kernel). The beta kernel model's accuracy for four other datasets compares well to the best accuracy. The beta kernel's average accuracy is within 0.01 of the best mean accuracy for Haberman's survival, adult, and tornado and within 0.04 for breast cancer. The beta kernel performs poorly in comparison to the other classifiers on the spam dataset with the RBF kernel, but the beta kernel classifier performs much better on the same data using the polynomial kernel. Overall, under-sampling performs

better than over-sampling and has the additional benefit of being faster than over-sampling or SMOTE although the beta kernel is the fastest. The mean accuracies for under-sampling with the RVM and under-sampling with the SVM are each the best for two datasets, but under-sampling with SVM performs extremely poorly for the arcene data.

### 3.3. Comparison of Beta Kernel and SVM

Comparing the performance of the beta kernel model with the SVM can shed light on situations where the beta kernel model may be more effective than the SVM. We use a Gaussian mixture model as the basis for this comparison. Both classes are drawn from multivariate normal distributions with unit covariance matrices. The
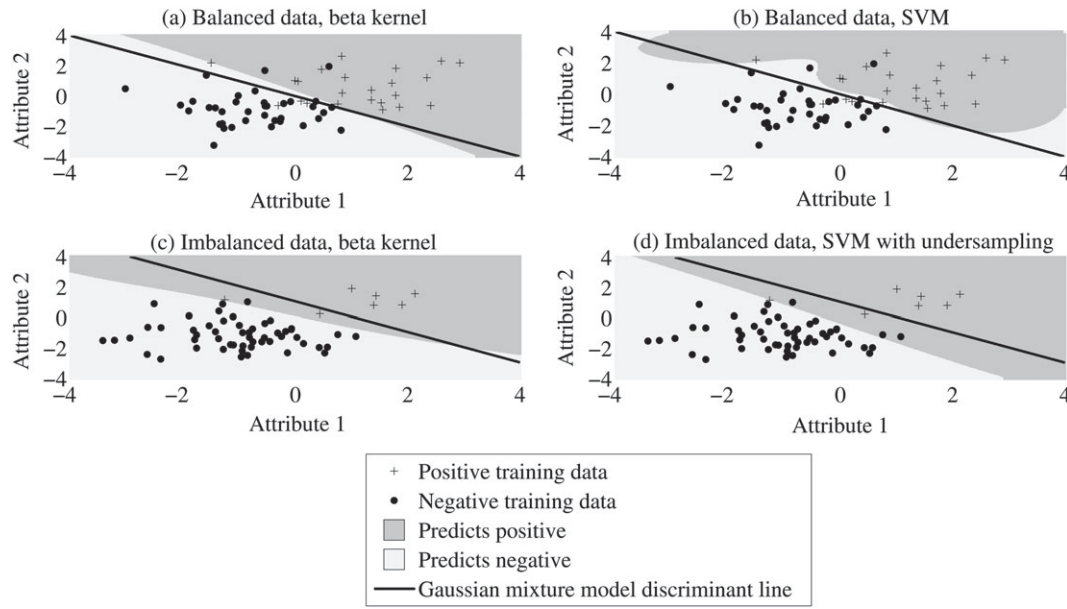
Fig. 1 Beta kernel and SVM classification of the Gaussian mixture model based on (a, b) 50% chance of positive and 50% chance of negative and (c, d) 10% chance of positive and 90% chance of negative.

attribute means corresponding to the positive class are $(1.0, 1.0)$ and the attribute means corresponding to the negative class are $(-1.0, -1.0)$. The beta kernel model and the SVM are compared for two training sets each with 60 total data points: (i) a balanced dataset with a 50% chance of sampling from the positive class and a 50% chance of sampling from the negative class and (ii) an imbalanced dataset with a 10% chance of sampling from the positive class and a 90% chance of sampling from the negative class. An unweighted SVM is used for the balanced training set, and under-sampling is used with the SVM for the second training set because this combination performed well for the imbalanced datasets in Table 4.

Figure 1 depicts how the beta kernel model and SVM classify data ranging between $-4.0$ and $4.0$ for each of the two attributes, given the two training datasets. The beta kernel model establishes a linear dividing line for both the balanced and imbalanced datasets. The dividing line predicted by the beta kernel model in the balanced dataset corresponds closely to the discriminant line in the Gaussian mixture model. A point on the discriminant line (the black line in Fig. 1) indicates a 50% chance the point is positively labeled and 50% chance it is negatively labeled according to the Gaussian mixture model.

The SVM establishes a nonlinear division for the balanced dataset in Fig. 1b. Because of the variation in the training data, the SVM creates areas around the positive training points for which it positively classifies unknown points and also creates areas around the negative training points for which it negatively classifies unknown points.

Because the SVM under samples the negative class for the imbalanced dataset, the SVM establishes a linear dividing line for the particular instance depicted in Fig. 1d, but a different sample may lead to a nonlinear division.

This demonstration with the Gaussian mixture model suggests that if the training data contain significant random variation, the beta kernel model may do a better job of predicting the underlying classification pattern than the SVM. The beta kernel model classifies unknown data points based on difference measures as determined by the RBF kernel. Consequently, the beta kernel model is less impacted by a point in the training set whose classification is a low-probability (i.e., less than 50%) event. If a point in the training set is labeled because of physical reasons rather than random variations—which may signify that other points in the vicinity should also be labeled similarly—the SVM algorithm does a better job than the beta kernel model of uncovering that discrepancy.

### 3.4. Online Learning

One of the principle advantages of Bayesian methods is the ability to rapidly incorporate new information into the analysis. Under Bayesian analysis, the prior distribution is updated with data, which produces a posterior distribution, and that posterior distribution serves as the prior distribution in the next iteration. The Gaussian Bayesian kernel methods discussed in Section 2 are not easily adaptable to this online learning pattern because the solution algorithms generally rely on maximizing the posterior distribution. Two Gaussian

Bayesian data mining tools for online learning or incremental updating include a Bayesian online perceptron [51,52] and a sparse representation for a Gaussian process model [53]. The Bayesian online perceptron approximates a posterior distribution by incrementally updating mean weights and covariances for the input data $\mathbf{x}_j$, but this method does not employ kernel functions. The sparse representation for a Gaussian process model uses the kernel function as the covariance matrix as depicted in Eq. (2) and determines whether to include a new input data vector $\mathbf{x}_j$ based on the degree to which the new data point changes the mean of the posterior distribution. Studying these two Bayesian classifiers for text classification reveals that they perform well in comparison with the SVM for relatively balanced classes, but SVM performs the best for imbalanced datasets [54].

The sequential RVM relies on a Bayesian Kalman filter and a least-squares approach to iteratively revise the RVM weights and hyperparameters for regression and time series forecasting problems [55]. Online or incremental learning has been investigated more fully for nonprobabilistic classifiers and includes algorithms such as LASVM [56], ALMA$_p$ [57], and NORMA [58]. Jin *et al*. [59] develop an algorithm that combines online learning and kernel learning which simultaneously trains an SVM classifier and assigns weights for multiple kernels.

The beta kernel model provides an easy tool for online learning because the model does not require solving an optimization problem. The posterior distribution at the end of one iteration acts as the prior distribution in the next iteration. Because the model relies on the kernel function to calculate posterior probabilities, the attributes for a data point $\mathbf{x}_i$ should be known although the outcome $y_i$ is not. The outcome remains uncertain as new data with known outcomes are collected. For example, an oil and gas company may have geological characteristics about several different areas. As it drills for oil in certain places and discovers whether or not those places contain oil, it can update its probability about whether or not a specific area contains oil based on the geological similarity between the known and unknown areas.

We explore the application of the beta kernel model to online learning by (i) demonstrating how $\alpha$, $\beta$, and $\theta$ change as new data arrive and (ii) comparing the beta kernel's accuracy to two other online learning tools. We depict how the beta kernel model's parameters change using the twonorm dataset as downloaded from the Delve project [60] at the University of Toronto. Unlike the previous example, this experiment keeps all 20 attributes for the two-norm data. Although the original dataset has an equal number of positively and negatively labeled outcomes, we purposely imbalance the data so that only 25% of the outcomes are positive.

**Table 5.** Updated parameters for beta kernel model with twonorm data.

| Iteration | Data point 1 | | | Data point 2 | | |
| | $\alpha$ | $\beta$ | $\overline{\theta}_1$ | $\alpha$ | $\beta$ | $\overline{\theta}_2$ |
|---|---|---|---|---|---|---|
| Prior | 1 | 1 | 0.50 | 1 | 1 | 0.50 |
| 1 | 1.21 | 1.35 | 0.47 | 1.04 | 2.32 | 0.31 |
| 2 | 2.05 | 1.54 | 0.57 | 1.28 | 3.35 | 0.28 |
| 5 | 2.18 | 3.01 | 0.42 | 1.31 | 6.66 | 0.16 |
| 10 | 4.92 | 4.97 | 0.50 | 1.70 | 10.47 | 0.14 |
| 20 | 8.29 | 8.40 | 0.50 | 2.59 | 19.54 | 0.12 |
| 30 | 13.50 | 11.71 | 0.54 | 3.59 | 27.73 | 0.11 |

We select two data points for which we assume the attributes but not the outcomes are known. At each iteration, a unique set of 10 data points whose outcomes are known is used to update $\alpha$ and $\beta$ for the each of the two unknown data points where $\alpha_i = \alpha + \sum_{\{j|y_j=1\}} k(\mathbf{x}_i, \mathbf{x}_j)$ and $\beta_i = \beta + \sum_{\{j|y_j=-1\}} k(\mathbf{x}_i, \mathbf{x}_j)$. Table 5 depicts the updated $\alpha$ and $\beta$ and the expected posterior probability $\overline{\theta}_i = \mathrm{E}[P(y_i = 1)]$. Figure 2 displays the beta distribution's probability density function as $\alpha$ and $\beta$ are updated for each of these two data points.

As the classifier receives more information, the first data point is much more likely to result in a positive outcome than the second data point. The expected probability for the first data point is close to 0.5 during all the iterations. Even after 30 iterations, the beta distribution's density function (the dark solid line in Fig. 2a) is still wide enough that the posterior probability could be between 0.25 and 0.75. The first data point's expected probability is 0.54. Much uncertainty exists over whether this data point is positively or negatively labeled; however, the posterior probability is much greater than 0.25, the fraction of positively labeled data points in the dataset.
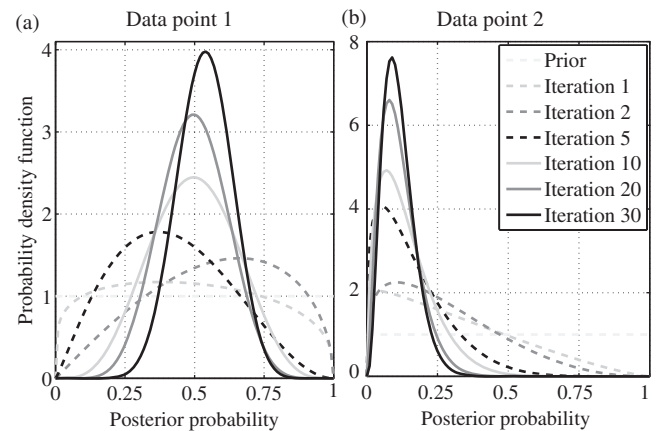


Fig. 2  Posterior probability distributions for two different data points in the twonorm dataset. (a) Data point 1, (b) data point 2.

Updating the parameters for the second data point significantly reduces the uncertainty of this data point's outcome. After only five iterations or 50 data points, the expected probability is 0.16. After 30 iterations, the expected probability is only 0.11, and most of the beta distribution's density function is less than 0.25.

### 3.5.  Online Learning Comparison

We compare the ability of three algorithms to correctly classify data incrementally. The weighted beta kernel model begins with a uniform prior where $\alpha = \beta = 1$, and these parameters and the weights for imbalanced data are updated with the arrival of each new set of data. Because the sequential RVM is designed for regression as opposed to classification problems, we slightly modify the standard RVM algorithm [61] so that it learns incrementally. The set of vectors with nonzero weights after one iteration is assumed to have nonzero weights with the arrival of a new batch of data, and only data in the new batch can be added to the set of vectors with nonzero weights. Once a data point is excluded from this set, it is discarded. The third algorithm is LASVM, an online SVM algorithm [56].

We use the 10 databases that were described in Section 3.1 but randomly divide each dataset into batches, where each batch contains either 5 or 10 data points, depending on size of the dataset. Five datasets have 50 different batches, three datasets have 30 batches, and the colon cancer and breast cancer data only have 10 batches because these latter sets contain a fewer data. At each iteration, each algorithm incorporates a new batch to classify unknown points in the testing set. Any data point not in a batch is in the testing set. The experiment is repeated 200 times for each dataset.

The median optimal values that were generated by the initial binary classification numerical experiments were used for $\sigma$ in the RBF kernel and $C$ in the SVM for this online learning experiment. Table 6 depicts the average accuracy (the square root of the product of the TP and TN rates) for different iterations. Because online learning typically does not begin with a large training set that can be used to generate optimal parameters for kernel functions, algorithms have been developed that linearly combine multiple kernels for a single classifier. As the classifier learns on data, the weights for each possible kernel are also updated to improve accuracy [59,62]. Future research could apply these concepts to the beta kernel model in an online-learning environment.

After 30 iterations, the three algorithms generally return accuracies consistent with the mean accuracies depicted in Table 2. After 1 or 5 iterations, the beta kernel outperforms the other two algorithms except on the spam data, but LASVM's accuracy exceeds that of the beta kernel by the 30th or 50th iteration for Parkinson, arcene, colon

**Table 6.**  Online learning results.

| Dataset | Data points in each batch | Iteration | Beta kernel | RVM | LASVM |
|---|---|---|---|---|---|
| Parkinson | 5 | 1 | **0.60** | 0.25 | 0.35 |
| | | 5 | **0.75** | 0.50 | 0.59 |
| | | 10 | **0.81** | 0.63 | 0.72 |
| | | 20 | **0.86** | 0.75 | 0.83 |
| | | 30 | 0.87 | 0.81 | **0.88** |
| Haberman's survival | 5 | 1 | **0.47** | 0.20 | 0.34 |
| | | 5 | **0.53** | 0.28 | 0.46 |
| | | 10 | **0.54** | 0.32 | 0.41 |
| | | 20 | **0.56** | 0.38 | 0.50 |
| | | 30 | **0.57** | 0.39 | 0.48 |
| | | 50 | **0.58** | 0.43 | 0.50 |
| Satellite | 10 | 1 | **0.93** | 0.81 | 0.91 |
| | | 5 | **0.98** | 0.96 | 0.97 |
| | | 10 | **0.98** | 0.97 | 0.98 |
| | | 20 | **0.98** | 0.98 | 0.98 |
| | | 30 | **0.99** | 0.98 | 0.99 |
| | | 50 | **0.99** | 0.98 | 0.99 |
| Arcene | 5 | 1 | **0.51** | 0.12 | 0.38 |
| | | 5 | **0.68** | 0.53 | 0.64 |
| | | 10 | **0.74** | 0.60 | 0.73 |
| | | 20 | 0.78 | 0.68 | **0.82** |
| | | 30 | 0.79 | 0.74 | **0.87** |
| Spam (RBF) | 10 | 1 | 0.45 | 0.43 | **0.63** |
| | | 5 | 0.66 | 0.76 | **0.81** |
| | | 10 | 0.71 | 0.79 | **0.84** |
| | | 20 | 0.76 | 0.79 | **0.86** |
| | | 30 | 0.79 | 0.79 | **0.88** |
| Colon cancer | 5 | 1 | **0.58** | 0.12 | 0.25 |
| | | 5 | **0.75** | 0.61 | 0.64 |
| | | 10 | 0.81 | 0.73 | **0.85** |
| Adult | 10 | 1 | **0.54** | 0.37 | 0.51 |
| | | 5 | **0.70** | 0.56 | 0.60 |
| | | 10 | **0.74** | 0.61 | 0.64 |
| | | 20 | **0.76** | 0.65 | 0.66 |
| | | 30 | **0.76** | 0.66 | 0.66 |
| | | 50 | **0.77** | 0.68 | 0.67 |
| Transfusion | 10 | 1 | **0.49** | 0.33 | 0.43 |
| | | 5 | **0.57** | 0.37 | 0.51 |
| | | 10 | **0.61** | 0.41 | 0.52 |
| | | 20 | **0.64** | 0.47 | 0.51 |
| | | 30 | **0.64** | 0.50 | 0.52 |
| | | 50 | **0.66** | 0.54 | 0.51 |
| Breast cancer | 5 | 1 | **0.50** | 0.08 | 0.42 |
| | | 5 | **0.57** | 0.25 | 0.47 |
| | | 10 | **0.59** | 0.31 | 0.50 |
| Tornado | 10 | 1 | **0.63** | 0.13 | 0.46 |
| | | 5 | **0.76** | 0.41 | 0.66 |
| | | 10 | **0.80** | 0.58 | 0.72 |
| | | 20 | **0.83** | 0.70 | 0.76 |
| | | 30 | **0.84** | 0.74 | 0.77 |
| | | 50 | **0.86** | 0.77 | 0.79 |

Bold indicates the best accuracy among the three classifiers.

cancer, as well as spam. This result suggests that the beta kernel is a better algorithm than LASVM for a handful of data points (i.e., less than 50). When the dataset contains more than 50 data points, the beta kernel and LASVM perform comparably, which echoes the findings

**Table 7.** Results of large-scale simulation for Gaussian mixture model

| Model | Acc | TP | TN | Run-time (seconds) |
|---|---|---|---|---|
| Beta kernel | **0.9996** | 1.0000 | 0.9992 | 206 |
| LASVM | 0.9905 | 0.9817 | 0.9999 | 797 |
| Naive Bayes | 0.9708 | 1.0000 | 0.9850 | 745 |
| Logistic regression | 0.9965 | 0.9993 | 0.9978 | 8 |

Bold indicates the best accuracy among the four classifiers.

from Table 2 in which the beta kernel and SVM algorithms perform similarly. In the online-learning environment, the beta kernel model outperforms LASVM for the Haberman's survival, satellite, adult, transfusion, breast cancer, and tornado datasets after 30 or 50 iterations (which correspond to 150 to 500 data points). Although not shown in the table, the beta kernel runs more quickly than the RVM and LASVM.

### 3.6. Large-Scale Simulation

Another test explores how well the beta kernel model performs for large datasets. As in Section 3.3, a Gaussian mixture model is used, in which the positive and negative class are drawn from a multivariate normal distribution with unit covariance matrices. The simulation uses 10 attributes. The mean for each attribute corresponding to the positive class is 1.0, and the mean for each attribute corresponding to the negative class is $-1.0$. We create a highly imbalanced set so that the chance of having a positively labeled point is 0.1%. The training set has 1 000 000 data points, and the testing set has 10 000 data points. The beta kernel model and the LASVM process this training set by using 1000 batches where each batch has 1000 data points. The logistic regression and naive Bayes algorithms offered by Matlab can process the million data points in one batch. The large-scale simulation is repeated 50 times.

As depicted in Table 7, the beta kernel model outperforms the other three algorithms for this large dataset, using the acccuracy measure described earlier, $\text{Acc} = \sqrt{\text{TP} * \text{TN}}$. The accuracies for the beta kernel, LASVM, and logistic regression are very similar to each other, however. The beta kernel model runs more slowly than logistic regression because the beta kernel requires generating two kernel matrices $\mathbf{K}$, one for the positively labeled training data points and the other for the negatively labeled training data points. The beta kernel model runs more quickly than the other two kernel-based algorithms: the LASVM and naive Bayes. Matlab's logistic regression runs so quickly because it does not need to divide the training set into smaller batches.

The computational complexity of the beta kernel model depends on the algorithm to generate the kernel matrix. If

we use the beta kernel model to predict only one unknown data point, the complexity of generating the RBF kernel matrix is at worst $\mathcal{O}(dm)$, where $d$ is the number of attributes and $m$ is the size of the training set. After the kernel matrix is generated, the complexity of updating the $\alpha$ and $\beta$ parameters is $\mathcal{O}(m)$. Although the beta kernel model can handle large datasets if the memory to store the kernel matrices is sufficient, it may not be cost-effective to use the beta kernel model if non-kernel algorithms such as logistic regression perform reasonably well.

### 4. CONCLUSIONS

This article has explored the usefulness of the beta kernel model and compared the model's accuracy with the RVM (a binary classification algorithm based on Gaussian distributions), the SVM, naive Bayes, and logistic regression. The beta kernel model relies on the well-known beta-binomial Bayesian formula, and deploying a kernel function as a measure of similarity between two different data points enables us to apply these updating techniques to classification problems. Incorporating weighting parameters or beginning with a nonuniform prior can help the model correctly classify imbalanced datasets. The model can be generalized to a Dirichlet kernel model for multiclass classification problems, and future research can compare the accuracy of the Dirichlet kernel model with other multiclass classifiers.

The extensive numerical testing of the beta kernel model with the RVM, SVM, naive Bayes, and logistic regression indicates that the beta kernel model may have some advantages that can be exploited for classification problems. The beta kernel model performs similarly to the SVM, a weighted SVM, and logistic regression for the 10 datasets in which the minority class composes between 7 and 44% of the data. The beta kernel model consistently performs better than the RVM and naive Bayes. If the user desires a probabilistic data mining tool, the beta kernel model may be a superior choice to the RVM. When the minority class comprises only 5% of the data, the beta kernel model is usually more accurate than the RVM and SVM, even though the latter two classifiers were improved using under-sampling, over-sampling, and SMOTE. This suggests that the beta kernel model should be an important tool for classifying heavily imbalanced datasets. The online learning experiment reveals that the beta kernel model consistently outperforms the RVM and LASVM (an incremental learning version of the SVM) if 50 or fewer data points are available, and the model frequently performs better than the RVM and LASVM even if more data are available. Finally, the beta kernel model calculates posterior probabilities very quickly and runs faster than the

RVM and SVM, both of which rely on solving optimization problems.

As this article represents, to our knowledge, the first extensive analysis and testing of the beta kernel model, we believe the model can potentially become a useful tool in machine learning. The beta kernel model provides similar accuracies for classifying datasets where the number in each class is relatively the same, and the model carries other advantages, such as fast computation times. If the dataset is heavily imbalanced, the beta kernel model may be the most accurate. If the data arrive incrementally, the model easily and quickly updates to incorporate the new data and can be relatively accurate with just a few data points. Future work could explore if the beta kernel model can be combined with the SVM, RVM, and logistic regression to improve overall accuracy. We intend to apply the beta kernel model to existing problems and demonstrate how some of these benefits can aid a decision maker.

## ACKNOWLEDGMENTS

## APPENDIX

We seek to prove that the uniform beta prior with a weighted kernel likelihood will label an unknown data point as a beta prior that derives from the training set and is updated with an unweighted kernel likelihood.

The uniform prior ($\alpha = \beta = 1$) and weighted kernel likelihood labels an unknown point $\mathbf{x}_i$ in the positive class if and only if Eq. A1 holds.

$$\frac{1 + \frac{m_-}{m} \sum k_{ij}^+}{2 + \frac{m_-}{m} \sum k_{ij}^+ + \frac{m_+}{m} \sum k_{ij}^-} > 0.5 \qquad (A1)$$

We assign $\sum k_{ij}^+ = \sum_{\{j|y_j=1\}} k(\mathbf{x}_i, \mathbf{x}_j)$ and $\sum k_{ij}^- = \sum_{\{j|y_j=-1\}} k(\mathbf{x}_i, \mathbf{x}_j)$ for notational convenience.

If the prior is formulated so that the prior's expected value is the proportion of positively labeled points in the training set, then an unknown point will be positively classified if and only if Eq. A2 holds.

$$\frac{\sum k_{ij}^+}{\sum k_{ij}^+ + \sum k_{ij}^-} > \frac{m_+}{m} \qquad (A2)$$

We want to show that Eqs. A1 and A2 are equivalent. We begin with A1.

$$1 + \frac{m_-}{m} \sum k_{ij}^+ > 1 + 0.5 \left( \frac{m_-}{m} \sum k_{ij}^+ + \frac{m_+}{m} \sum k_{ij}^- \right)$$

$$m_- \sum k_{ij}^+ > 0.5 \left( m_- \sum k_{ij}^+ + m_+ \sum k_{ij}^- \right)$$

$$0.5 m_- \sum k_{ij}^+ > 0.5 m_+ \sum k_{ij}^- \qquad (A3)$$

$$m_- \sum k_{ij}^+ + m_+ \sum k_{ij}^+ > m_+ \sum k_{ij}^- + m_+ \sum k_{ij}^+$$

$$m \sum k_{ij}^+ > m_+ \left( \sum k_{ij}^- + \sum k_{ij}^+ \right)$$

The last line is equivalent to A2.

## REFERENCES

[1] J. Shawe-Taylor and N. Cristianini, Kernel Methods for Pattern Analysis, Cambridge, Cambridge University Press, 2004.

[2] N. Cristianini and J. Shawe-Taylor, An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods, Cambridge, Cambridge University Press, 2000.

[3] T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, New York, Springer Science, 2001.

[4] B. Schölkopf and A. J. Smola, Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond, Cambridge, MIT Press, 2002.

[5] M. Seeger, Bayesian model selection for support vector machines, Gaussian processes and other kernel classifiers, In S. A. Solla, T. K. Leen, and K.-R. Müller, eds. Advances in Neural Information Processing Systems, Cambridge, MIT Press, 2000, 603–609.

[6] M. E. Tipping, Sparse Bayesian learning and the relevance vector machine, J Mach Learn Res 1 (2001), 211–244.

[7] C. M. Bishop and M. E. Tipping, Bayesian regression and classification, In Advances in Learning Theory: Methods, Models and Applications, J. A. K. Suykens, G. Horváth, S. Basu, C. Micchelli, and J. Vandewalle, eds. Amsterdam, the Netherlands, IOS Press, 2003, 267–288.

[8] M. A. T. Figueiredo, Adaptive sparseness using Jeffreys prior, In Neural Information Processing Systems, Vol. 14, T. Dietterich, S. Becker, and Z. Ghahramani, eds. Cambridge, MIT Press, 2002, 697–704.

[9] B. K. Mallick, D. Ghosh, and M. Ghosh, Bayesian classification of tumours by using gene expression data, J Roy Statl Soc Part B 67(2) (2005), 219–234.

[10] Z. Zhang, G. Dai, and M. I. Jordan, Bayesian generalized kernel mixed models, J Mach Learn Res 12 (2011), 111–139.

[11] L. Montesano and M. Lopes, Learning grasping affordances from local visual descriptors, In Proceedings of the 8th IEEE international conference on development and learning, Shanghai, China, 2009.

[12] E. Parzen, On estimation of a probability density function and mode, Ann Math Stat 33(3) (1962), 1065–1076.

[13] M. Maalouf and T. B. Trafalis, Kernel logistic regression using truncated Newton method, In C. H. Dagli, D. L. Enke, K. M. Bryden, H. Ceylen, and M. Gen, eds. Intelligent Engineering Systems Through Artificial Neural Networks, Vol. 18. New York, ASME Press, 2008, 455–462.

[14] M. E. Tipping and A. C. Faul, Fast marginal likelihood maximisation for sparse Bayesian models, In Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics, C. M. Bishop and B. J. Frey, eds. Key West, FL, 2003. http://www.miketipping.com/papers.htm (accessed April 4, 2011).

[15] G. M. Foody, RVM-based multi-class classification of remotely sensed data, Int J Remote Sens 29(6) (2008), 1817–1823.

[16] S.-C. Huang and T.-K. Wu, Combining wavelet-based feature extractions with relevance vector machines for stock index forecasting, Expert Syst 25(2) (2008), 133−149.

[17] B. Saha, K. Goebel, and J. Christophersen, Comparison of prognostic algorithms for estimating remaining useful life of batteries, Trans Inst Meas Control 31(3/4) (2009), 293−308.

[18] A. K. Gupta and S. Nadarajah, eds. Handbook of Beta Distribution and Its Applications, New York, Marcel Dekker, 2004.

[19] M. Mason and M. Lopes, Robot self-initiative and personalization by learning through repeated interactions, In Proceedings of the 6th ACM/IEEE International Conference on Human-Robot Interaction (HRI'11), Lausanne, Switzerland, 2011. http://flowers.inria.fr/mlopes/myrefs/11-hri.pdf (accessed April 5, 2011).

[20] G. King and L. Zeng, Logistic regression in rare events data, Polit Anal 9 (2001), 137−163.

[21] S. Wang, W. Jiang, and K.-L. Tsui, Adjusted support vector machines based on a new loss function, Ann Oper Res 174 (2010), 83−101.

[22] M. Maalouf and T. B. Trafalis, Robust weighted kernel logistic regression in imbalanced and rare events data, Comput Stat Data Anal 55 (2011), 168−183.

[23] K. Morik, P. Brockhausen, and T. Joachims, Combining statistical learning with a knowledge-based approach-a case study in intensive care monitoring, In Proceedings of the 16th international conference on machine learning, Bled, Slovenia, 1999. http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.71.5781&rep=rep1&type=pdf (accessed April 4, 2011).

[24] B. X. Wang and N. Japkowicz, Boosting support vector machines for imbalanced data sets, In Foundations of Intelligent Systems: Proceedings of the 17th International Symposium on Methodologies for Intelligent Systems, Toronto, Canada, A. Aijun, S. Matwin, Z. W. Raś, and D. Ślęzak, eds. Berlin, Springer-Verlag, 2008.

[25] B. P. Carlin and T. A. Louis, Bayesian Methods for Data Analysis, (3rd ed.), Boca Raton, FL, CRC Press, 2008.

[26] O. Chapelle, Active learning for Parzen window classifier. In Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics, Barbados, 2005. http://olivier.chapelle.cc/pub/aistats05.pdf (accessed April 4, 2011).

[27] H.-G. Chew, D. J. Crisp, R. E. Bogner, and C.-C. Lim, Target detection in radar imagery using support vector machines with training size biasing, In Proceedings of the 6th International Conference on Control, Automation, Robotics, and Vision, Singapore, 2000. http://kernel-machines.org/papers/upload_11483_ICARCV2000-4.ps (accessed April 5, 2011).

[28] C.-C. Chang and C.-J. Lin, LIBSVM: A library for support vector machines. Software available at http://www.csie.ntu.edu.tw/~cjlin/libsvm (accessed October 7, 2010).

[29] Matlab, NaiveBayes.fit. Version R2012b. Mathworks, 2012.

[30] M. A. Little, P. E. McSharry, E. J. Hunter, J. Spielman, and L. O. Ramig, Suitability of dysphonia measurements for telemonitoring of Parkinson's disease, IEEE Trans Biomed Eng 56(4) (2009), 1015−1022.

[31] S. J. Haberman, Generalized residuals for log-linear models. In Proceedings of the 9th International Biometrics Conference, Boston, 1976, 104−122.

[32] I. Guyon, S. R. Gunn, A. Ben-Hur, and G. Dror, Result analysis of the NIPS 2003 feature selection challenge. In Neural Information Processing System, 2004. http://books.nips.cc/papers/files/nips17/NIPS2004_0194.pdf (accessed May 4, 2012).

[33] R. Kohavi, Scaling up the accuracy of Naive-Bayes classifiers: a decision-tree hybrid. In Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, E. Simoudis, J. Han, and U. Fayyad eds. 1996. 202−207. http://robotics.stanford.edu/~ronnyk/nbtree.pdf (accessed January 11, 2014).

[34] I.-C. Yeh, K.-J. Yang, and T.-M. Ting, Knowledge discovery on RFM model using Bernoulli sequence, Expert Syst Appl 36(3) (2009), 5866−5871.

[35] W. N. Street, O. L. Mangasarian, and W. H. Wolberg, An inductive learning approach to prognostic prediction, In Proceedings of the Twelfth International Conference on Machine Learning, A. Prieditis and S. Russell, eds. San Francisco, CA: Morgan Kaufmann, 1995, 522−530.

[36] K. Bache and M. Lichman, UCI Machine Learning Repository School of Information and Computer Science, Irvine, CA, University of California, 2013. http://archive.ics.uci.edu/ml (accessed January 11, 2014).

[37] U. Alon, N. Barkai, D. A. Notterman, K. Gish, S. Ybarra, D. Mack, and A. J. Levine, Broad patterns of gene expression revealed by clustering of tumor and normal colon tissues probed by oligonucleotide arrays, Proc Natl Acad Sci USA 96 (1999), 6745−6750. http://genomics-pubs.princeton.edu/oncology/affydata/index.html (accessed April 4, 2011).

[38] T. B. Trafalis, I. Adrianto, and M. B. Richman, Active learning with support vector machines for tornado prediction, In Computational Science-CCS 2007: Proceedings of the 7th International Conference on Computational Science, Beijing, China, Y. Shi, G. D. van Albada, J. Dongarra, & P. M. A. Sloot, eds. Berlin: Springer-Verlag, 2007, 1130−1137.

[39] T. Kudo and Y. Matsumoto, Fast methods for kernel-based text analysis, In Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics, E. W. Hinrichs and D. Roth, eds. Stroudsburg, PA, Association for Computational Linguistics, 2003, 24−31. http://acl.ldc.upenn.edu/acl2003/main/pdfs/Kudo.pdf (accessed January 6, 2014).

[40] J. Moon, T. Shon, J. Seo, J. Kim, and J. Seo, An approach for spam e-mail detection with support vector machine and n-gram indexing, In Computer and Information Sciences-ISCIS 2004: 19th International Symposium, C. Aykanat, T. Dayar, and İ. Korpeoğu, eds. Kemer-Antalya, Turkey and Berlin: Springer, 2004, 351−362.

[41] M. Kubat, R. Holte, and S. Matwin, Learning when negative examples abound, In Machine Learning, ECML-97: Proceedings of the 9th European conference on machine learning, M. van Someren and G. Widmer, eds. Heidelberg, Springer, 1997, 146−153.

[42] R. Akbani, S. Kwek, and N. Japkowicz, Applying support vector machines to imbalanced datasets, In Machine Learning: ECML 2004, J.-F. Boulicaut, F. Esposito, F. Giannotti, and D. Pedreschi, eds. Berlin, Springer, 2004, 39−50.

[43] G. E. A. P. A. Batista, R. C. Prati, and M. C. Monard, A study of the behavior of several methods for balancing machine learning training data, SIGKDD Explor 6(1) (2004), 20−29.

[44] S. Visa and A. Ralescu, Issues in mining imbalanced data sets-a review paper, In Paper Presented at the Proceedings of Midwest Artificial Intelligence and Cognitive Science Conference (MAICS '05), Dayton, 2005.

[45] N. V. Chawla, Data mining for imbalanced datasets: an overview, In Data Mining and Knowledge Discovery Handbook, (2nd ed.), O. Maimon and L. Rokach, eds. New York, Springer Science, 2010, 875−886.

[46] M. A. Maloof, Learning when data sets are imbalanced and when costs are unequal and unknown. In Paper presented at the International Conference on Machine Learning Workshop on Learning from Imbalanced Data Sets, Washington, DC, 2003.

[47] Y. Tang, Y.-Q. Zhang, N. V. Chawla, and S. Krasser, SVMs modeling for highly imbalanced classification, IEEE Trans Syst Man Cybern Part B: Cybern 39(1) (2009), 281–288.

[48] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, SMOTE: synthetic minority over-sampling technique, J Artif Intell Res 16 (2002), 321–357.

[49] N. V. Chawla, A. Lazarevic, L. O. Hall, and K. W. Bowyer, SMOTEBoost: Improving prediction of the minority class in boosting, In Paper presented at the 7th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD), Dubrovnik, Croatia, 2003.

[50] J. Van Hulse, T. M. Khoshgoftaar, and A. Napolitano, Experimental Perspectives on Learning from Imbalanced Data. In Paper Presented at the 24th International Conference on Machine Learning, Corvallis, OR, 2007, 935–942.

[51] M. Opper, A Bayesian approach to online learning, In On-Line Learning in Neural Networks, D. Saad, ed. New York, Cambridge University Press, 1998, 363–378.

[52] S. A. Solla and O. Winther, Optimal perceptron learning: An online Bayesian approach, In On-Line Learning in Networks, D. Saad, ed. New York, Cambridge University Press, 379–399, 1998.

[53] L. Csató and M. Opper, Sparse representation for Gaussian process models, In NIPS 2000, Vol. 13, T. K. Leen, T. G. Dietterich, and V. Tresp, eds. MIT Press, Cambridge, MA, 2001.

[54] K. M. A. Chai, H. T. Ng, and H. L. Chieu, Bayesian online classifiers for text classification and filtering, In Paper presented at the SIGIR, Tampere, Finland, August 11–15, 2002.

[55] N. Nikolaev and P. Tino, Sequential relevance vector machine learning from time series, In Paper presented at the International Joint Conference on Neural Networks, Montreal, Canada, July 31–August 4, 2005.

[56] A. Bordes, S. Ertekin, J. Weston, and L. Bottou, Fast kernel classifiers with online and active learning, J Mach Learn Res 6 (2005), 1579–1619.

[57] C. Gentile, A new approximate maximal margin classification algorithm, J Mach Learn Res 2 (2001), 213–242.

[58] J. Kivinen, A. J. Smola and R. C. Williamson, Online learning with kernels, IEEE Trans Signal Process 52(8) (2004), 2165–2176.

[59] R. Jin, S. C. H. Hoi, and T. Yang, Online multiple kernel learning: algorithms and mistake bounds, Algorithmic Learn Theory 6331 (2010), 390–404.

[60] M. Revow, Twonorm data set. Delve datasets, 1996. http://www.cs.toronto.edu/~delve/data/twonorm/desc.html (accessed May 8, 2012).

[61] M. E. Tipping, SparseBayes for Matlab, version 2, 2009. http://www.vectoranomaly.com/downloads/downloads.htm (accessed April 4, 2011).

[62] P. H. Gosselin, F. Precioso, and S. Philipp-Foliguet, Incremental kernel learning for active image retrieval without global dictionaries, Pattern Recognit 44(10–11) (2011), 2244–2254.