

Toward Instance-dependent Label Noise Tolerant Classification: A Probabilistic Approach

Jakramate Bootkrajang · Jeerayut
Chaijaruwanich

Received: date / Accepted: date

Abstract Learning from labelled data is becoming more and more challenging due to inherent imperfection of training labels. Existing label-noise tolerant learning machines were primarily designed to tackle class-conditional noise which occurs at random, independently from input instances. However, relatively less attention was given to a more general type of label noise which is influenced by input features. In this paper, we try to address the problem of learning a classifier in the presence of instance-dependent label noise by developing a novel label noise model which is expected to capture the variation of label noise rate within a class. This is accomplished by adopting a probability density function of a mixture of Gaussians to approximate the label flipping probabilities. Experimental results demonstrate the effectiveness of the proposed method over existing approaches.

Keywords instance-dependent label noise · classification · logistic regression

1 Introduction

Classical machinery of learning a classifier relies on a perfect set of labelled data. Unfortunately, it is becoming more difficult to guarantee the correctness of the training labels. In many classification tasks, labelling errors originate from the lack of knowledge or resources to annotate the input instances [1–3]. In addition, the quality of labels is being compromised in order to acquire large amount of training data quickly and cheaply [4–6].

In general, there are three forms of label noises: class-conditional label noise, instance-dependent label noise and adversarial label noise [7,8]. Class-conditional noise is the noise which occurs randomly and independently from

Department of Computer Science, Chiang Mai University, Chiang Mai, Thailand, 50200
Tel.: +66-053-943-414
fax: +66-053-892-281
E-mail: jakramate.b@cmu.ac.th (corresponding author), jeerayut.c@cmu.ac.th

the input feature. The noise is often referred to as random noise because labels of random instances within a class could be flipped subjected to same label flipping probability. Instance-dependent noise is the noise which is in parts influenced by input features. The dependency can be of many forms. One of the most natural ones is the case where instances that live closer to the decision boundary have higher chance of being mislabelled. Consequently, the flipping probabilities for labels in the same class are then not necessarily equal. Adversarial noise is an extreme noise that is intentionally injected into the dataset to perturb the classification algorithm the most [9]. A good example of this type of noise can be found in spam email classification problem where spams were specially crafted to look like normal email and to mislead the algorithm as much as possible.

Among the three types of noises, random label noise is probably the one that received most attention during the past. Accordingly, methods for dealing with random label noise are rather well developed. Notable model-based approaches for dealing with random label noise include robust linear discriminant analysis [10] robust Kernel Fisher Discriminant [11,12] robust Logistic Regression [13,14], and robust kernel logistic regression [15] to name just a few. Furthermore, theoretical properties and its effects on learning algorithm are relatively well understood [16–20].

Compared to random labelling error, the effects of instance-dependent label noise on learning machine is somewhat less understood and methods for tackling such noise are still lacking. Part of the reason might be that there can be many types of dependencies between input instance and label noise. As such, modelling the occurrence of instance-dependent label is more difficult than the random noise case since any modelling assumption might be bound to specific dependencies only. One of the earlier attempts to address instance-dependent noise is a work in [21] where a model of label noise conditioned on the distance of input from the decision boundary was studied. The work focused on theoretical aspect of the model in the presence of such noise but provided no algorithmic solution for learning the model. Recently, [22] proposed a generalisation of the latent variable label noise model [11] which was originally formulated for random noise to cover the cases where noise does not occur at random. The work assumed label flipping is more likely when input instance lives closer to decision boundary and proposed to use a probability density function of the exponential distribution to model the noise. An approach under similar assumption have also been studied in [23] with one subtle difference that an unnormalised Gaussian distribution and the Laplace distribution were instead chosen to model the distribution of label noise. The authors of [23] coined the term ‘Boundary Condition Noise’ (BCN) to refer to this particular type of instance-dependent noise.

One of the difficulties in modelling instance-dependent noise is the fact that distribution of label noise varies from case to case. For example, the models found in [23,22] might comfortably accommodate BCN but might not be suitable for noise which occurs further from the decision boundary. This

prompts a new development of more flexible model which is able to capture wider range of instance-dependent label noises.

In this work, we will investigate a new approach to modelling instance-dependent label noise. We will employ a probability density function of a mixture of Gaussians in order to approximate label noise distribution within a class. The new label noise model will later be incorporated into a standard logistic regression classification model. The increased flexibility of the new label noise model should capture the noise rate better and hence improved classification performance can be expected.

We consider the following points to be our contributions.

- We propose a novel instance-dependent label noise model based on a mixture of Gaussians.
- We develop a new logistic regression classifier employing the new noise model and devise an algorithm to learn the classifier jointly with estimating the label noise model.
- We extensively evaluate the proposed method and demonstrate that the approach is superior to existing approaches in various situations.

The rest of the paper is organised as follows. Section 2 briefly reviews the latent variable model and its generalisation. A novel label noise function based on Gaussian mixture model is proposed in Section 3. Section 4 demonstrates the effectiveness of the proposed method using a wide range of artificial and real-world testbeds. Section 5 concludes the study.

2 Problem setting and backgrounds

A classification task is a task of inferring discrete-valued function $h : \mathbf{x} \rightarrow y$ which maps input vector $\mathbf{x} \in R^m$ to discrete class label $y \in K$ using a finite set of input-label pairs $(\mathbf{x}_n, y_n)_{n=1}^N$ drawn i.i.d from some joint distribution \mathcal{D} . In this study we will focus on binary classification problems where $y \in \{0, 1\}$. In an idealised setting \mathbf{x} is paired with its *true* class assignment. However in many real-world classification problems the correctness of the label cannot be guaranteed and we instead have access to $S = (\mathbf{x}_n, \tilde{y}_n)_{n=1}^N \sim \tilde{\mathcal{D}}$. Here, \tilde{y} denotes an *observed* (and possibly noisy) class label. The question is can we still learn h using this noisy set of labelled data?

One of the early probabilistic approaches to address the problem of learning in the presence of label noise is through the use of a so-called *latent variable model* which initially developed for a generative classifier [11]. A work in [13] later considers the model in discriminative classifier perspective. Essentially, the latent variable model seeks to explain the posterior probability of the observed label using a weighted average of the posterior probabilities of the true class label.

$$p(\tilde{y} = j | \mathbf{x}_n, \phi) = \sum_{k=0}^1 p(\tilde{y} = j | y = k) p(y = k | \mathbf{x}_n, \phi) \quad (1)$$

Here, ϕ denotes parameters of a classification model, e.g., the weight vector of a logistic regression model. Hereinafter, we shall use \tilde{P}_n^j for the posterior probability of the observed class j , and P_n^k for the posterior probability of the true class k to simplify the notation. The term, $p(\tilde{y} = j|y = k)$ in Eq.(1) represents a probability that a label was flipped from class k into class j . Clearly, the probability is independent from the input vector. Therefore, the model in Eq.(1) is well-suited for class-conditional label noise, i.e., random label noise.

Arguably, the constraint imposed by Eq.(1) is rather too rigid. In reality, noise is expected to occur more in the region of maximum confusion, i.e., near the decision boundary or region with too few examples. Previous study has empirically shown that the argument is indeed observable in real-life [23]. They found that the majority of disagreements between experts and non-expert labellers in crowdsourced labelling task occurred near the decision boundary. The work in [22] also acknowledged the limitation of the latent variable model and pioneered the generalisation of the model by taking possible dependencies between flipping probability and input instance into consideration using:

$$\tilde{P}_n^j = \sum_{k=0}^1 p(\tilde{y} = j|y = k, \mathbf{x}_n, \phi) P_n^k \quad (2)$$

Observe that in the generalised label noise model the flipping probability term now depends on both the input vector and the parameter of the classification model. A work in [22] considers the use of constrained exponential distribution function to model the flipping probability. Using logistic regression as a classification model the flipping probability of each instance can then be expressed by $p(\tilde{y}|y = k, \mathbf{x}_n, \mathbf{w}) = \frac{1}{\alpha_k} \exp(-\mathbf{w}^T \mathbf{x}_n / \alpha_k)$. In a similar manner, [23] used the Laplace distribution and an unnormalised zero-mean Gaussian distribution to capture the noise's dynamics. They observed that the Laplace model generally yielded superior classification performances.

Although, all of the aforementioned noise models have been shown to work well on BCN problems, it is still questionable if those noise models will still work if nature of noise deviate from the BCN assumption. For example, in the case where noise rate peaks somewhere within the class and not precisely at the decision boundary, or when distribution of noise is not strictly Laplace or Gaussian. To address the possible shortcoming, in this paper we will propose a more flexible alternative. The new noise model will be based on a probability density function of Gaussian Mixture Model. We shall then incorporate the new noise model into a logistic regression to yield a label noise robust classifier.

3 Gaussian Mixture Noise Model

The generalised label noise model allows label flipping probability $p(\tilde{y} = j|y = k, \mathbf{x}, \phi)$ to be defined in different ways. Previous studies had proposed various types of distribution functions to model the noise rate. However, the predictive

Table 1: Summary of the notations used in the manuscript.

Symbol	Description
\mathbf{x}	The input variables
y	The true label
\tilde{y}	The observed label
P_n^k	The true posterior of class k of \mathbf{x}_n
\tilde{P}_n^k	The observed posterior of class k of \mathbf{x}_n
\mathcal{G}	Gaussian mixture model
C	The number of mixing components
μ_i	The mean of the i -th component
σ_i	The standard deviation of the i -th component
ω_i	The mixing weight of the i -th component
τ_n	The normalised distance of \mathbf{x}_n from decision boundary
\mathbf{w}	The weight vector of logistic regression
λ_i	The regularisation parameter of σ_i

performance of the classifier could be compromised, when the true distribution of label noise cannot be adequately approximated by the chosen function. In this section we shall explore a new label noise model. The approach makes use of Gaussian Mixture Model (GMM) to approximate the label flipping probability. In general, a univariate GMM is expressed by:

$$\mathcal{G}(\tau) = \sum_{i=1}^C \omega_i \mathcal{N}(\tau; \mu_i, \sigma_i^2) \quad (3)$$

with a constraint that $\sum_{i=1}^C \omega_i = 1$. Here C represents the number of mixing components. We will use $\mathcal{G}(\tau)$ to model a chance that label will flip from one class to the other. The chance is expressed as a function of a random variable τ which represents a distance of an example \mathbf{x} from decision boundary induced by a classifier. We note that for more sophisticated treatment of label noise, one can model the noise level as a function of instance's position. In that case, multivariate GMM could be employed in place of the univariate GMM. In the current study, however, we chose the distance from decision boundary because the measure adequately reflects our noise assumption while also being easier to work with. In the case of logistic regression parametrised by the weight vector \mathbf{w} , the distance is defined as $\tau_n = \mathbf{x}_n^T \mathbf{w} / \|\mathbf{w}\|$. To further account for asymmetric label flipping, i.e., the case where one class flips to the other but not vice versa, we will model the label noise associated with instances classified by the model as positive and negative using separate mixtures. In the case of logistic regression this can be easily differentiated by observing the polarity of τ_n . We denote the distance of instances which fall on the *positive* side of the decision boundary by τ_n^1 and those that fall on the *negative* side by τ_n^0 . To this

end, we define the label flipping probabilities to be:

$$p(\tilde{y}_n = 1 | y_n = 0, \mathbf{x}_n, \mathbf{w}) = \mathcal{G}^{01}(\tau_n^0) \quad (4)$$

$$p(\tilde{y}_n = 0 | y_n = 0, \mathbf{x}_n, \mathbf{w}) = 1 - \mathcal{G}^{01}(\tau_n^0) \quad (5)$$

$$p(\tilde{y}_n = 0 | y_n = 1, \mathbf{x}_n, \mathbf{w}) = \mathcal{G}^{10}(\tau_n^1) \quad (6)$$

$$p(\tilde{y}_n = 1 | y_n = 1, \mathbf{x}_n, \mathbf{w}) = 1 - \mathcal{G}^{10}(\tau_n^1) \quad (7)$$

where

$$\mathcal{G}^{01}(\tau_n^0) = \sum_{i=1}^C \frac{\omega_i^{01}}{\sigma_i^{01} \sqrt{2\pi}} \exp \left(- \frac{(\tau_n^0 - \mu_i^{01})^2}{2(\sigma_i^{01})^2} \right) \quad (8)$$

$$\mathcal{G}^{10}(\tau_n^1) = \sum_{i=1}^C \frac{\omega_i^{10}}{\sigma_i^{10} \sqrt{2\pi}} \exp \left(- \frac{(\tau_n^1 - \mu_i^{10})^2}{2(\sigma_i^{10})^2} \right) \quad (9)$$

Here, $\Theta = \{\omega_i^{10}, \omega_i^{01}, \mu_i^{10}, \mu_i^{01}, \sigma_i^{10}, \sigma_i^{01}\}_{i=1}^C$ is a set of noise model's parameters that need to be estimated from the data.

3.1 Regularisations on the standard deviations

By definition, the value of label flipping probability must be between zero and one. However, the density function in Eq.(3) which we adopted to model the flipping probability can output any positive real. One way to guarantee that Eqs.(4)-(7) never exceed one is to control the value of σ_i 's in the two mixtures. To do that we first derive a bound on the maximum of $\mathcal{G}(\tau)$.

$$\mathcal{G}(\tau) = \sum_{i=1}^C \frac{\omega_i}{\sigma_i \sqrt{2\pi}} \exp \left(- \frac{(\tau - \mu_i)^2}{2\sigma_i^2} \right) \quad (10)$$

$$\leq \sum_{i=1}^C \frac{\omega_i}{\sigma_i \sqrt{2\pi}} \quad (11)$$

$$\leq \frac{1}{\sqrt{2\pi}} \sum_{i=1}^C \frac{\omega_i}{\sigma_i} \quad (12)$$

$$\leq \frac{1}{\sqrt{2\pi}} \sum_{i=1}^C \frac{1}{\sigma_i} \quad (13)$$

$$\leq \frac{1}{\sqrt{2\pi}} \frac{C}{\min_i \sigma_i} \quad (14)$$

Since the exponential function in Eq.(10) attains its maximum of one when $\tau = \mu$, the value of $\mathcal{G}(\tau)$ is bounded by Eq.(11). By observing that $0 \leq \omega_i \leq 1$, we further bound the ratio ω_i/σ_i in the summation by $1/\sigma_i$. It is then not difficult to verify that Eq.(14) follows from Eq.(13). We note that the bound is quite loose especially for large C . This might affect the working of the

model that employs a GMM with too many mixing components in the sense that its σ_i can be unreasonably large. Large σ_i translates to flat PDF and hence indicating that the model is ignoring the label noise.

So, to ensure that label noise function gives out number between zero and one we require that $\mathcal{G}(\tau) \leq \frac{1}{\sqrt{2\pi}} \frac{C}{\min_i \sigma_i} \leq 1$, or

$$\min_i \sigma_i \geq \frac{C}{\sqrt{2\pi}} \quad (15)$$

We will introduce a regularisation to restrict the value of the smallest σ_i to be greater than $C/\sqrt{2\pi}$. Since during the parameters learning process we do not know which σ_i will be the smallest we thus enforce the constraint on all σ_i . For this purpose, we will employ a log-barrier function of the form, $\log(\sigma_i - C/\sqrt{2\pi})$, to enforce such constraint on each of the σ_i 's in the two mixtures. Essentially, the log-barrier discourages σ_i from getting too close to the bound.

3.2 Objective function

We will now incorporating the proposed label noise function into the standard logistic regression classifier. To this end, let us first write down the objective function of the Gaussian mixture model based robust logistic regression (GMMLR) as well as all of the regularisations in the form of a penalised log-likelihood:

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \Theta) = & \sum_{n=1}^N \left\{ \tilde{y}_n \log \left([1 - \mathcal{G}_n^{10}] P_n^1 + \mathcal{G}_n^{01} P_n^0 \right) \right. \\ & \left. + (1 - \tilde{y}_n) \log \left(\mathcal{G}_n^{10} P_n^1 + [1 - \mathcal{G}_n^{01}] P_n^0 \right) \right\} \\ & + \sum_{j \in \{01, 10\}} \sum_{i=1}^C \lambda_i^j \log \left(\sigma_i^j - \frac{C}{\sqrt{2\pi}} \right) \end{aligned} \quad (16)$$

where as before $\Theta = \{\omega_i^{10}, \omega_i^{01}, \mu_i^{10}, \mu_i^{01}, \sigma_i^{10}, \sigma_i^{01}\}_{i=1}^C$ denotes a set of parameters to the noise model. Here we simplify some notation by writing $\mathcal{G}_n^{jk} = \mathcal{G}^{jk}(\tau_n^j)$. The first term in Eq.(16) is the objective function for the logistic regression. The second term is a regularisation term imposing constraints on the standard deviations.

3.3 Model selection

It is apparent that our model formulation involves a number of regularisation parameters. A traditional way of picking good value of those is by performing cross validation. However, the technique requires a lot of computational efforts.

To facilitate the selection of λ_i , we propose a simple heuristic to select a good value of the regularisation parameters by:

$$\lambda_i = \frac{1}{(\sigma_i - C/\sqrt{2\pi})} \quad (17)$$

The motivation behind this heuristic is to discourage σ_i from taking extreme values. First, notice that Eq.(16) can be maximised by increasing the value of σ_i . But that is not always favourable because large σ_i translates to a flat Gaussian indicating that the model believes there is no label noise in the dataset. In such case, the robust model will behave like a traditional non-robust classifier. By setting λ_i using Eq.(17), we are balancing the gain in the likelihood value as σ_i increases and hence preventing it from being too large. The heuristic also helps enforcing the constraint that σ_i must be greater than $C/\sqrt{2\pi}$. Observe when σ_i is getting close to $C/\sqrt{2\pi}$, e.g., $(\sigma_i - C/\sqrt{2\pi} < 1)$, the log-barrier function will be negative. At the same time λ_i will be large resulting in the undesirable decrement of the objective function value. In effect, this prevents σ_i from getting too small.

3.4 Learning the model

Once the formulation is finalised, we now move on to model optimisation. There are basically two sets of parameters in concern. The first is \mathbf{w} , the weight vector of the logistic regression model and the second is Θ , the parameters to the label noise model. Since there is no close form solution, we resort to gradient descend methodology for estimating the parameters. In particular, we adopted a variant of Newton's method called quasi-Newton method by which the calculation of the Hessian can be facilitated. We implemented the machinery of learning using publicly available optimisation package [24]. In the following we will derive the gradients of the objective function w.r.t to each of the parameters.

First, the gradient of the objective function w.r.t \mathbf{w} is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_{n=1}^N \left(\frac{\tilde{y}_n}{\tilde{P}_n^1} - \frac{1 - \tilde{y}_n}{\tilde{P}_n^0} \right) (1 - \mathcal{G}_n^{10} - \mathcal{G}_n^{01}) P_n^1 P_n^0 \mathbf{x}_n \quad (18)$$

Next, the derivatives of Eq.(16) w.r.t the means of \mathcal{G}^{10} and \mathcal{G}^{01} are calculated,

$$\frac{\partial \mathcal{L}}{\partial \mu_i^{10}} = \sum_{n=1}^N \left(\frac{1 - \tilde{y}_n}{\tilde{P}_n^0} - \frac{\tilde{y}_n}{\tilde{P}_n^1} \right) \left(\frac{(\tau_n^1 - \mu_i^{10})}{(\sigma_i^{10})^2} \right) \omega_i^{10} P_n^1 \mathcal{N}_i^{10} \quad (19)$$

$$\frac{\partial \mathcal{L}}{\partial \mu_i^{01}} = \sum_{n=1}^N \left(\frac{\tilde{y}_n}{\tilde{P}_n^1} - \frac{1 - \tilde{y}_n}{\tilde{P}_n^0} \right) \left(\frac{(\tau_n^1 - \mu_i^{01})}{(\sigma_i^{01})^2} \right) \omega_i^{01} P_n^0 \mathcal{N}_i^{01} \quad (20)$$

where \mathcal{N}_i^{10} and \mathcal{N}_i^{01} denote the i -th Gaussian of \mathcal{G}^{10} and \mathcal{G}^{01} , respectively. Likewise, the gradients of the objective w.r.t σ_i are found to be:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \sigma_i^{10}} &= \sum_{n=1}^N \left(\frac{1 - \tilde{y}_n}{\tilde{P}_n^0} - \frac{\tilde{y}_n}{\tilde{P}_n^1} \right) \left(\frac{(\tau_n^1 - \mu_i^{10})^2}{(\sigma_i^{10})^3} - \frac{1}{\sigma_i^{10}} \right) \omega_i^{10} P_n^1 \mathcal{N}_i^{10} \\ &\quad + \frac{\lambda_i}{\sigma_i^{10} - \frac{C}{\sqrt{2\pi}}} \end{aligned} \quad (21)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \sigma_i^{01}} &= \sum_{n=1}^N \left(\frac{\tilde{y}_n}{\tilde{P}_n^1} - \frac{1 - \tilde{y}_n}{\tilde{P}_n^0} \right) \left(\frac{(\tau_n^0 - \mu_i^{01})^2}{(\sigma_i^{01})^3} - \frac{1}{\sigma_i^{01}} \right) \omega_i^{01} P_n^0 \mathcal{N}_i^{01} \\ &\quad + \frac{\lambda_i}{\sigma_i^{01} - \frac{C}{\sqrt{2\pi}}} \end{aligned} \quad (22)$$

For the mixing weights, we proceed the same way. Using $\omega_i = u_i^2$ parameterisation, we first optimise for u_i by:

$$\frac{\partial \mathcal{L}}{\partial u_i^{10}} = \sum_{n=1}^N \left(\frac{1 - \tilde{y}_n}{\tilde{P}_n^0} - \frac{\tilde{y}_n}{\tilde{P}_n^1} \right) 2u_i^{10} P_n^1 \mathcal{N}_i^{10} - 2\alpha_i u_i^{10} \quad (23)$$

$$\frac{\partial \mathcal{L}}{\partial u_i^{01}} = \sum_{n=1}^N \left(\frac{\tilde{y}_n}{\tilde{P}_n^1} - \frac{1 - \tilde{y}_n}{\tilde{P}_n^0} \right) 2u_i^{01} P_n^0 \mathcal{N}_i^{01} - 2\alpha_i u_i^{01} \quad (24)$$

We later recover ω_i by squaring the optimised u_i . Additionally, we adopted the projection on simplex method by [25] to ensure that the mixing weights sum to unity. Algorithm 1 summarises the steps to learn the proposed ‘‘Gaussian Mixture Model-based robust Logistic Regression’’ (GMMLR) model

Algorithm 1 The steps to learn the proposed GMMLR model.

Require: Set of training data $(\mathbf{x}_n, y_n)_{n=1}^N$

Initialise $\mathbf{w} \sim \text{uniform}[-0.1, 0.1]$,

Initialise $\mu_i^{10} = i/C$ and $\mu_i^{01} = -i/C$ for $i = 1 : C$,

Initialise σ_i^{10} and σ_i^{01} with value greater than $C/\sqrt{2\pi}$

while Iteration < MaxIteration **do**

 Calculate regularisation parameter λ_i , using Eq.(17)

 Calculate distance from decision boundary $\tau_n^k = \mathbf{w}^T \mathbf{x}_n / \|\mathbf{w}\|$

 Calculate flipping probabilities using Eqs.(4)-(7)

 Update \mathbf{w} using Eq.(18)

 Update μ_i^{jk} using Eq.(19) and Eq.(20)

 Update σ_i^{jk} using Eq.(21) and Eq.(22)

 Update ω_i^{jk} using Eq.(23) and Eq.(24)

end while

Ensure: Optimised weight vector, \mathbf{w} . Optimised $\mu_i^{jk}, \sigma_i^{jk}, \omega_i^{jk}$.

Table 2: Characteristics of the datasets used in this study. Top 11 datasets are from the UCI repository. The bottom 5 datasets are datasets which genuinely contains label noise. The ground truths for all of the datasets except LULC are available. For LULC, the dataset is divided into training and testing sets apriori. The training labels are noisy although the exact mislabelling rate is unknown. The testing labels are clean.

UCI dataset	# Instances	# Features	Remarks
australian	690	14	-
biodegradation	1055	41	-
boston	506	13	-
diabetes	768	8	-
fertility	100	9	-
german	1000	24	-
ionosphere	351	34	-
magic04	19020	10	-
musk	476	166	-
spambase	4601	57	-
usps46	11000	256	class 4 vs class 6
Realworld dataset	# Instances	# Features	Remarks
breast [26]	49	7129	9 mislabellings
colon [27]	62	2000	9 mislabellings
leukaemia [28]	72	7129	1 mislabelling
websearch [15]	1030	1318	183 mislabellings
LULC [3]	10845	28	forest vs grass

4 Experiments

4.1 Experimental protocol

The datasets used in the experiments fall into two categories: data with no report of labelling errors and real-world data that genuinely contains label noise according to the literature. The first group of data contain 11 datasets from UCI repository [29]. For the second group, we used colon cancer dataset [27], breast cancer dataset [26], leukaemia dataset [28], websearch dataset [15] and Land Use/Land Cover (LULC) dataset [3], all of which have been reported to contain label noise. Table 2 summarises the characteristics of the datasets used in this study.

Since only part of the datasets in our test-suite originally contains label noise, we had to simulate the existence of label noise in the remaining datasets. To generate instance-dependent noise we first train a logistic regression on clean training data. The resulting weight vector, together with parameters of respective noise distributions are then used to calculate label flipping probabilities of instances in the training set. We adopt PDFs of Gaussian distribution $\mathcal{N}(\mu = 1, \sigma^2 = 2)$, Gamma distribution $\Gamma(shape = 1, scale = 2)$ and mixture of Gaussians $0.5 * \mathcal{N}(\mu = 0, \sigma^2 = 0.5) + 0.5 * \mathcal{N}(\mu = 2, \sigma^2 = 0.5)$ as label noise generating functions. The gamma noise distribution is used to simulate boundary condition noise which can often be observed in practice [23]. Noise generated by Gaussian and mixture of Gaussians are less common but the two

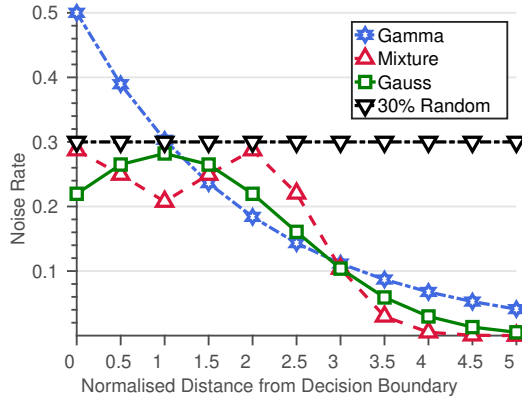


Fig. 1: Illustration of label noise rates studied in this work.

distributions will serve as testbeds for examining the flexibility of the noise model. Still, one could imagine that such complex label noise may exist in crowdsourced datasets, where labellers of different backgrounds are asked to perform the labelling task. Presumably each of the labeller who is responsible for a ‘subset’ of the dataset might introduce labelling errors to the dataset at different mislabelling rates. Along the same line, such complex label noise has also been found in automate product categories tagging task where labels of instances in specific data space are subjected to different label flipping probabilities [5]. In addition to instance-dependent label noises, we also investigated the model’s generality in dealing with random label noise at 30% level. The noise rates corresponding the four noise generating functions are illustrated in Fig 1. These artificial label noises will be referred to as *gaussian noise*, *gamma noise*, *mixture noise* and *random noise* respectively.

In the case where label noise is already presented in the data, we did not inject artificial noise any further. Since the ground truths are available for these set of data we can always evaluate the learned model using noise-free test data. We divided the data using 80/20 percentage of training and testing set split except for the case of breast dataset where we used 90/10 splitting due to high-dimensional and low sample size nature of the data.

4.2 Effect of the number of mixing components

We will begin with a set of experiments which seeks to better understand the working of the GMMLR model. Firstly, we want to see how C , the number of mixing components, affects the model’s predictive performance. We will compare GMMLR utilising 1, 3, 5, 7 and 9 mixing components named GMMLR-1, GMMLR-3, GMMLR-5, GMMLR-7 and GMMLR-9 on 11 UCI datasets corrupted by the four types of label noises described above. We performed 20 experiment repetitions using 80/20 training/testing data splitting for each

Table 3: Mean classification errors (%) and standard errors of GMMLR model with varying number of mixing components under *mixture noise* averaged over 20 runs. Boldface entry highlights the best mean rank.

Dataset	GMMLR-1	GMMLR-3	GMMLR-5	GMMLR-7	GMMLR-9
australian	14.49 \pm 3.51	14.93 \pm 3.38	14.42 \pm 3.14	13.95 \pm 2.76	14.06 \pm 2.83
biodeg	17.18 \pm 2.80	16.97 \pm 2.66	17.35 \pm 3.18	17.46 \pm 3.47	17.87 \pm 4.90
boston	15.39 \pm 4.34	15.59 \pm 4.52	15.78 \pm 4.49	16.23 \pm 4.50	16.37 \pm 4.41
diabetes	25.71 \pm 2.93	25.00 \pm 2.68	25.26 \pm 2.66	25.26 \pm 2.70	25.26 \pm 2.71
fertility	19.75 \pm 7.69	19.75 \pm 10.06	18.75 \pm 9.44	18.75 \pm 9.44	18.75 \pm 9.44

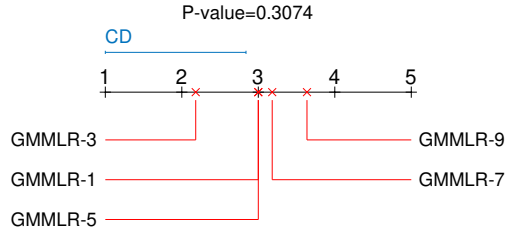


Fig. 2: Nemenyi post-hoc test for comparing five GMMLR models with varying number of mixing components under *mixture noise*. In this case, the predictive performances were not statistically different as tested by Friedman test at 0.05 level.

Table 4: Mean classification errors (%) and standard errors GMMLR model with varying number of mixing components under *gaussian noise* averaged over 20 runs. Boldface entry highlights the best mean rank.

Dataset	GMMLR-1	GMMLR-3	GMMLR-5	GMMLR-7	GMMLR-9
australian	15.62 \pm 3.26	15.04 \pm 4.10	15.11 \pm 4.17	15.25 \pm 4.18	15.25 \pm 4.21
biodeg	17.27 \pm 3.16	17.13 \pm 3.15	17.13 \pm 3.28	17.54 \pm 3.17	17.84 \pm 3.29
boston	14.17 \pm 3.12	14.66 \pm 3.43	14.80 \pm 4.08	16.32 \pm 4.46	16.96 \pm 4.52
diabetes	25.45 \pm 4.81	25.32 \pm 4.85	25.88 \pm 4.60	25.91 \pm 4.56	25.88 \pm 4.58
fertility	18.75 \pm 10.11	18.50 \pm 10.01	18.50 \pm 9.75	18.75 \pm 9.72	18.50 \pm 9.75
german	28.00 \pm 3.57	28.93 \pm 4.59	28.70 \pm 4.27	28.70 \pm 4.29	28.70 \pm 4.29
ionosphere	15.77 \pm 3.89	16.13 \pm 4.27	15.99 \pm 4.07	16.13 \pm 4.07	16.06 \pm 4.17
magic04	21.27 \pm 0.70	21.31 \pm 0.73	21.63 \pm 1.07	21.84 \pm 1.12	21.99 \pm 1.13
musk	25.78 \pm 5.31	25.94 \pm 5.62	25.99 \pm 4.88	25.89 \pm 5.05	26.15 \pm 5.39
spambase	8.12 \pm 0.90	8.12 \pm 0.90	8.31 \pm 0.82	8.52 \pm 0.81	8.97 \pm 1.01
usps46	1.92 \pm 1.20	1.91 \pm 1.21	1.80 \pm 1.08	1.81 \pm 1.14	1.82 \pm 1.16
Meanrank	2.3182	2.5000	2.5455	3.7273	3.9091

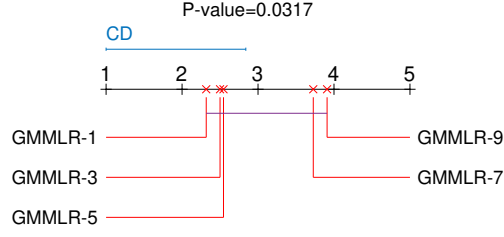


Fig. 3: Nemenyi post-hoc test for comparing GMMLR with varying number of mixing components under *gaussian noise*. Although, p-value from Friedman test indicates significant differences in performance, the post-hoc test shows that the differences are still within critical value.

noise types and computed the average classification errors and their standard errors.

The average classification errors incurred by GMMLR-1 to GMMLR-9 in the presence of *mixture noise*, *gaussian noise*, *random noise* and *gamma noise* are presented in Table 3 - 6, respectively. In each table, we ranked the performance of the five variants of GMMLR on each dataset, and summarised their mean ranks in the bottom row of the respective table. We performed Friedman test at 0.05 level as described in [30] followed by Nemenyi post-hoc test in order to compare the performances of the five classifiers. The graphical representation of the Nemenyi test are placed right after its corresponding table. We decided to draw the graphs for all of the experiments regardless of whether the p-value from the Friedman test is smaller than 0.05 or not.

From the results, there seem to be some significant differences in the performances of the five variants of GMMLR in all noise types except for the case of mixture noise, as first tested by the Friedman test. That is, under *gaussian noise* (Table 4 and Fig.3), *gamma noise* (Table 6 and Fig.5) and *random noise* (Table 5 and Fig.4) the p-values are smaller than 0.05 and so a post-hoc test was necessary. However, Nemenyi post-hoc test showed that groups of classifiers with significantly different performances can be formed only in *random noise* case. From Fig.4, we see that GMMLR-3 and GMMLR-1 is significantly better than GMMLR-9 while there is not enough evidence to tell if GMMLR-3 is better than GMMLR-1, GMMLR-5 or GMMLR-7.

Despite the statistical tests showing that all variants of GMMLR are not statistically different in most of the cases tested, we can still observe some general trends from the graphical representations of Nemenyi test. From Fig.2-4, we notice that GMMLR-1, GMMLR-3, and GMMLR-5, albeit not statistically different, were ranked higher than GMMLR-7 and GMMLR-9. The slightly lower predictive performances exhibited by GMMLR with large C could be due to the bound on the standard deviation in Eq.(14), where the s.d. is excessively constrained for large C . It can be understood that when the s.d. is large, as enforced by the bound, GMMLR becomes less aware of label noise. As such, the performance of GMMLR with larger C might be compromised.

Table 5: Mean classification errors (%) and standard errors of GMMLR models with varying number of mixing components under *random noise* averaged over 20 runs. Boldface entry highlights the best mean rank.

Dataset	GMMLR-1	GMMLR-3	GMMLR-5	GMMLR-7	GMMLR-9
australian	16.45 \pm 3.70	16.41 \pm 3.89	17.36 \pm 4.10	17.68 \pm 3.93	17.68 \pm 3.93
biodeg	20.47 \pm 4.19	20.52 \pm 3.43	20.95 \pm 3.54	21.14 \pm 4.03	21.42 \pm 4.10
boston	19.36 \pm 6.89	19.17 \pm 7.19	19.85 \pm 6.92	20.34 \pm 6.51	20.78 \pm 6.02
diabetes	25.97 \pm 4.40	25.94 \pm 4.20	25.94 \pm 4.18	26.14 \pm 3.87	26.10 \pm 3.91
ecoli	24.95 \pm 11.50	22.75 \pm 10.45	22.75 \pm 10.45	22.75 \pm 10.45	22.75 \pm 10.45

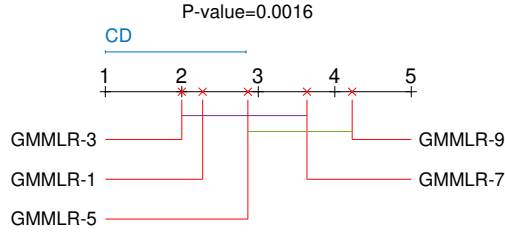


Fig. 4: Nemenyi post-hoc test for comparing GMMLR with varying number of mixing components under *random noise*. The post-hoc test groups the classifiers into two groups where we see that GMMLR-3 is better than GMMLR-9. However, there is not enough evidence to say if any of GMMLR-1,3,5,7 is better than the others.

In terms of the running time, it is apparent as shown in Table 7 that a model with smaller number of mixing components runs faster. We see that the time required to train the model increases, rather linearly, as the number of mixing components increases. So, in light of the predictive performance and the running time, it seems that GMMLR with $C = 3$ is preferable to other variants of GMMLR.

4.3 Effects of standard deviation regularisation

Next, we will investigate the influence and the importance of regularising the value of the standard deviations of the Gaussians in the mixture. The standard deviation indicates the spread of each Gaussian and hence directly controls the value of the flipping probability. By formulation, the flipping probability must be between zero and one. Figure 6 illustrates one of the estimated label flipping probabilities obtained from GMMLR with s.d. regulariser (top) and without s.d. regulariser (bottom). We see that when s.d. is too small the output from the PDF of the GMM is no longer below one. Although, this is perfectly fine in

Table 6: Mean classification errors (%) and standard errors of GMMLR with varying number of mixing components under *gamma noise* averaged over 20 runs. Boldface entry highlights the classifier with the best mean rank.

Dataset	GMMLR-1	GMMLR-3	GMMLR-5	GMMLR-7	GMMLR-9
australian	16.92 \pm 4.97	18.22 \pm 4.77	19.38 \pm 5.50	19.93 \pm 5.65	19.93 \pm 5.65
biodeg	16.37 \pm 4.62	16.23 \pm 4.65	17.16 \pm 4.17	17.84 \pm 3.73	18.27 \pm 3.81
boston	17.25 \pm 5.18	17.25 \pm 5.33	18.73 \pm 5.21	19.17 \pm 5.09	19.36 \pm 4.95
diabetes	23.31 \pm 4.10	24.03 \pm 4.62	25.36 \pm 4.01	25.45 \pm 4.06	25.45 \pm 4.06
fertility	31.75 \pm 13.89	29.75 \pm 14.82	29.75 \pm 14.82	30.00 \pm 14.51	29.75 \pm 14.82
...					19
...					14
...					4
...					8
...					7
...					5
...					

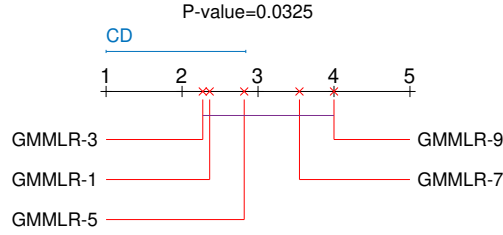


Fig. 5: Nemenyi post-hoc test for comparing GMMLR with varying number of mixing components under *gamma noise*. Although, p-value from Friedman test indicates differences in performance, the post-hoc test shows that the differences are still within critical value.

Table 7: Running times (seconds) of variants of GMMLR learning from Gamma noise on Intel Core-i5 3.2GHz machine averaged over 20 runs.

Dataset	GMMLR-1	GMMLR-3	GMMLR-5	GMMLR-7	GMMLR-9
australian	0.69 \pm 0.12	1.54 \pm 0.27	2.75 \pm 0.20	4.22 \pm 0.35	5.83 \pm 0.38
biodeg	0.90 \pm 0.07	2.06 \pm 0.18	3.62 \pm 0.26	5.30 \pm 0.49	7.29 \pm 0.61
boston	0.63 \pm 0.08	1.50 \pm 0.24	2.71 \pm 0.25	3.99 \pm 0.27	5.44 \pm 0.32
diabetes	0.61 \pm 0.11	1.39 \pm 0.20	2.68 \pm 0.27	4.02 \pm 0.25	5.57 \pm 0.21
fertility	0.31 \pm 0.07	0.79 \pm 0.08	1.48 \pm 0.03	2.13 \pm 0.03	2.86 \pm 0.04
german	0.73 \pm 0.13	1.64 \pm 0.24	3.07 \pm 0.26	4.55 \pm 0.27	6.31 \pm 0.28
ionosphere	0.63 \pm 0.04	1.44 \pm 0.16	2.55 \pm 0.18	3.68 \pm 0.29	4.93 \pm 0.32
magic04	5.10 \pm 0.20	12.35 \pm 0.77	22.57 \pm 1.65	34.67 \pm 2.90	48.18 \pm 3.61
musk	0.74 \pm 0.08	1.70 \pm 0.17	2.94 \pm 0.24	4.27 \pm 0.32	5.69 \pm 0.40
spambase	2.16 \pm 0.10	4.54 \pm 0.22	7.84 \pm 0.58	11.70 \pm 0.82	16.11 \pm 1.00
usps46	1.91 \pm 0.12	4.26 \pm 0.20	6.85 \pm 0.32	9.64 \pm 0.46	12.96 \pm 0.52

standard statistical usage where the probability is taken to be the integral of the PDF, it might not always work for our case where the output from the PDF is directly used as label flipping probability. When the flipping probability is not between zero and one, the term on the left hand side of Eq.(2) is no longer a valid posterior probability in which case the classification performance might suffer.

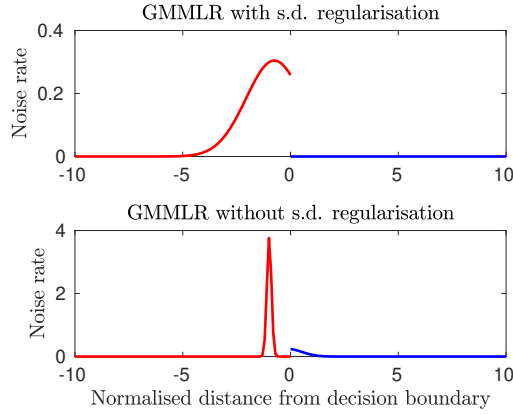


Fig. 6: Illustration of estimated noise rate. Observe that, without regularisation, too small s.d. would produce the flipping probability that violates the probabilistic constraint.

To validate our speculation, we investigated the effect of standard deviation regularisers λ_i , by comparing GMMLR *with* regulariser (GMMLR-reg) and GMMLR *without* regulariser (GMMLR-noreg). We opted for GMMLR with 3 mixtures as it achieved relatively good performance in the previous experiment. In this experiment, GMMLR-reg selects λ_i using the proposed heuristic in Eq.(17). We tested the two variants on the 11 UCI datasets corrupted with the four types of noises mentioned above. Again, we repeated the experiment for 20 repetitions for each label noise configurations and computed the average classification errors and their standard errors. The results are summarised in Table 8.

From the results, we observe that the regularisation on the s.d. is quite important and turning it off results in slightly higher classification errors. To test if such differences are significant, we then employed Wilcoxon ranksum test at 0.05 level to compare the performances of the two GMMLRs. Surprisingly, we only found that the differences were significant, e.g., p-value < 0.05 in two cases namely in the spambase dataset corrupted by *mixture noise* and biodeg dataset corrupted by *gaussian noise*. The reason why the majority of the performance differences of the two GMMLRs are not statistically significant could be that the estimated s.d. of GMMLR-noreg were still sometimes legitimate (the resulting PDF outputs values between zero and one). However, there were also times when the estimated s.d. of GMMLR-noreg was unusable which can be observed through the relatively larger error variances compared to those of GMMLR-reg in most of the datasets.

Despite the fact that we do not have enough evidence to say that GMMLR-reg is statistically better than GMMLR-noreg, we still see that GMMLR-reg achieved lower error rates in general and in real world classification problems there is a good chance that GMMLR with s.d. regulariser will perform better.

Table 8: Mean classification errors (%) and standard errors of GMMLR without standard deviation regulariser (GMMLR-noreg) and GMMLR which employs the regulariser (GMMLR-reg). Boldface entries are the best ones (not necessarily statistically significant). Entry marked with ‘-’ is worse than GMMLR-reg, with ‘+’ is better than GMMLR-reg and with ‘ \approx ’ is equivalent to GMMLR-reg as tested by Wilcoxon rank sum test at 0.05 level. The last row counts the cases where GMMLR-reg is better than, equivalent to, and worse than GMMLR-noreg.

Dataset	Mixture noise		Gaussian noise	
	GMMLR-noreg	GMMLR-reg	GMMLR-noreg	GMMLR-reg
australian	14.75 \pm 2.73 \approx	14.93 \pm 3.38	15.58 \pm 4.48 \approx	15.04 \pm 4.10
biodeg	19.53 \pm 6.58 \approx	16.97 \pm 2.66	20.14 \pm 4.79 -	17.13 \pm 3.15
boston	17.11 \pm 4.83 \approx	15.59 \pm 4.52	17.55 \pm 5.61 \approx	14.66 \pm 3.43
diabetes	25.39 \pm 2.82 \approx	25.00 \pm 2.68	26.27 \pm 4.30 \approx	25.32 \pm 4.85
fertility	19.50 \pm 10.37 \approx	19.75 \pm 10.06	21.00 \pm 9.12 \approx	18.50 \pm 10.01
german	27.73 \pm 3.39 \approx	26.82 \pm 2.72	28.80 \pm 4.51 \approx	28.93 \pm 4.59
ionosphere	17.25 \pm 6.93 \approx	15.42 \pm 4.56	16.76 \pm 4.81 \approx	16.13 \pm 4.27
magic04	22.87 \pm 2.64 \approx	21.83 \pm 1.04	22.32 \pm 2.16 \approx	21.31 \pm 0.73
musk	28.54 \pm 7.16 \approx	26.56 \pm 5.09	28.39 \pm 9.85 \approx	25.94 \pm 5.62
spambase	11.78 \pm 4.47 -	8.50 \pm 1.67	9.19 \pm 2.24 \approx	8.12 \pm 0.90
usps46	4.12 \pm 7.08 \approx	2.19 \pm 1.28	4.17 \pm 9.81 \approx	1.91 \pm 1.21
Win-Tie-Loss	1-10-0	n/a	1-10-0	n/a

Dataset	Random noise		Gamma noise	
	GMMLR-noreg	GMMLR-reg	GMMLR-noreg	GMMLR-reg
australian	17.25 \pm 3.77 \approx	16.41 \pm 3.89	20.00 \pm 9.74 \approx	18.22 \pm 4.77
biodeg	23.08 \pm 5.12 \approx	20.52 \pm 3.43	19.10 \pm 4.88 \approx	16.23 \pm 4.65
boston	17.94 \pm 4.74 \approx	19.17 \pm 7.19	18.92 \pm 4.71 \approx	17.25 \pm 5.33
diabetes	26.14 \pm 4.18 \approx	25.94 \pm 4.20	27.11 \pm 8.31 \approx	24.03 \pm 4.62
fertility	21.25 \pm 10.99 \approx	22.75 \pm 10.45	31.00 \pm 14.56 \approx	29.75 \pm 14.82
german	29.85 \pm 6.18 \approx	28.27 \pm 4.47	31.18 \pm 6.15 \approx	29.48 \pm 4.62
ionosphere	22.39 \pm 5.97 \approx	20.85 \pm 5.04	17.25 \pm 4.88 \approx	18.24 \pm 4.63
magic04	22.62 \pm 1.60 \approx	21.87 \pm 1.20	21.03 \pm 1.32 \approx	20.97 \pm 0.71
musk	29.48 \pm 9.07 \approx	27.60 \pm 4.82	29.48 \pm 9.64 \approx	23.49 \pm 7.73
spambase	11.28 \pm 1.65 \approx	10.36 \pm 2.05	9.32 \pm 1.60 \approx	9.07 \pm 1.53
usps46	11.67 \pm 10.83 \approx	10.65 \pm 4.62	3.98 \pm 4.86 \approx	3.88 \pm 0.83
Win-Tie-Loss	0-11-0	n/a	0-11-0	n/a

4.4 Cross validation versus automatic selection of λ_i

In this work we have proposed a simple heuristic for determining λ_i , the regularisation parameters of the standard deviations, instead of performing the traditional cross validation (CV). We want to study if this heuristic finds a good set of regularisation parameters compared to those found by performing a grid search using the traditional CV. To be of any usefulness, the heuristic should be able to select a “good” GMMLR model as measured by the model’s predictive performance, and also within shorter period of learning time compared to CV. To test the ability of the proposed heuristic in selecting good GMMLR model, we will again take GMMLR-3, which seems to work best from the previous experiments and compare it with a GMMLR-3 for which λ_i was selected by the traditional CV. We will refer to the latter as GMMLR-3-CV. In particular, GMMLR-3-CV selected its regularisation parameters (λ_i) from the set $\{2^i\}$, $i \in \{\pm 1, \pm 3, \pm 5\}$ via 3 folds cross validation.

Table 9: Mean classification errors (%) and standard errors of GMMLR where its regularisers were selected using cross validation (GMMLR-3-CV) and GMMLR-3 which selects regularisers via the proposed heuristic. Bold-face entries are the best ones (not necessarily statistically significant). The performance obtained from using the heuristic is similar to what we would get from performing time-consuming cross validation. The differences are not statistically significant as tested by Wilcoxon rank sum test at 0.05 level.

Dataset	Mixture noise		Gaussian noise	
	GMMLR-3-CV	GMMLR-3	GMMLR-3-CV	GMMLR-3
australian	14.53 \pm 3.29 \approx	14.93 \pm 3.38	15.36 \pm 4.08 \approx	15.04 \pm 4.10
biodeg	17.06 \pm 2.86 \approx	16.97 \pm 2.66	17.16 \pm 3.18 \approx	17.13 \pm 3.15
boston	15.93 \pm 4.46 \approx	15.59 \pm 4.52	14.26 \pm 4.19 \approx	14.66 \pm 3.43
diabetes	25.23 \pm 2.62 \approx	25.00 \pm 2.68	25.91 \pm 4.43 \approx	25.32 \pm 4.85
fertility	19.00 \pm 8.83 \approx	19.75 \pm 10.06	19.50 \pm 10.50 \approx	18.50 \pm 10.01
german	27.32 \pm 2.40 \approx	26.82 \pm 2.72	28.57 \pm 4.14 \approx	28.93 \pm 4.59
ionosphere	15.35 \pm 4.50 \approx	15.42 \pm 4.56	16.13 \pm 3.97 \approx	16.13 \pm 4.27
magic04	22.18 \pm 1.52 \approx	21.83 \pm 1.04	21.70 \pm 1.49 \approx	21.31 \pm 0.73
musk	26.35 \pm 5.16 \approx	26.56 \pm 5.09	25.62 \pm 4.91 \approx	25.94 \pm 5.62
spambase	8.68 \pm 1.65 \approx	8.50 \pm 1.67	8.13 \pm 0.94 \approx	8.12 \pm 0.90
usps46	2.28 \pm 1.40 \approx	2.19 \pm 1.28	1.67 \pm 1.00 \approx	1.91 \pm 1.21

Dataset	Random noise		Gamma noise	
	GMMLR-3-CV	GMMLR-3	GMMLR-3-CV	GMMLR-3
australian	17.07 \pm 3.70 \approx	16.41 \pm 3.89	16.63 \pm 5.06 \approx	18.22 \pm 4.77
biodeg	20.81 \pm 3.47 \approx	20.52 \pm 3.43	16.97 \pm 4.04 \approx	16.23 \pm 4.65
boston	19.31 \pm 7.46 \approx	19.17 \pm 7.19	17.55 \pm 5.43 \approx	17.25 \pm 5.33
diabetes	26.46 \pm 3.94 \approx	25.94 \pm 4.20	25.19 \pm 4.07 \approx	24.03 \pm 4.62
fertility	22.75 \pm 12.40 \approx	22.75 \pm 10.45	29.00 \pm 14.01 \approx	29.75 \pm 14.82
german	28.48 \pm 4.95 \approx	28.27 \pm 4.47	31.15 \pm 5.65 \approx	29.48 \pm 4.62
ionosphere	21.20 \pm 4.37 \approx	20.85 \pm 5.04	17.96 \pm 4.26 \approx	18.24 \pm 4.63
magic04	22.43 \pm 1.96 \approx	21.87 \pm 1.20	21.59 \pm 2.32 \approx	20.97 \pm 0.71
musk	27.19 \pm 5.23 \approx	27.60 \pm 4.82	23.33 \pm 7.38 \approx	23.49 \pm 7.73
spambase	10.34 \pm 2.02 \approx	10.36 \pm 2.05	9.15 \pm 1.62 \approx	9.07 \pm 1.53
usps46	10.90 \pm 4.24 \approx	10.65 \pm 4.62	3.47 \pm 1.09 \approx	3.88 \pm 0.83

The results from the experiment are summarised in Table 9. Statistically, the differences in the error rates incurred by the two variants of GMMLR are negligible. This is because the p-values from Wilcoxon rank sum test are all greater than 0.05 level. The results can be used as an evidence to validate the usefulness of the proposed heuristic for automatic selection of λ_i . Combining this observation with the fact that CV requires much longer time to perform the model selection, it is natural to conclude that the proposed heuristic is useful for selecting λ_i as it can speed up the learning process without sacrificing the classification performance.

4.5 Comparisons with other classifiers on artificial label noises

We have seen from the previous experiments about the effect of the number of mixing components, the benefits of regularising the standard deviation and the usefulness of the proposed automatic regularisation parameters selection and come to the conclusion that GMMLR with $C = 3$, with s.d. regularisation and with the automated model selection is likely to perform best compared to other variations of GMMLR model. In this section, we will investigate the

comparative performance of the GMMLR with other existing label noise robust classifiers, with traditional classifier which is unaware of label noise as well as with the current state-of-the-art algorithm for classification problems.

Specifically, we compared the proposed method to the generalised Logistic Regression (gLR) [22], Logistic regression with Laplace noise model (Laplace-BCN) [23] which were developed for boundary condition noise and robust Logistic Regression (rLR) [13] which targets random label noise. In the case of rLR, instead of using the EM algorithm as in [13], we employed a gradient-based algorithm [14] which was shown to converge faster for optimising the model. We additionally included a traditional logistic regression (LR) and an SVM with linear kernel (SVM-Linear) ¹ as baseline models. The C parameter of the SVM were chosen using 5-fold cross validation from values in $\{2^i\}$ where $i \in \{-10, -8, -6, \dots, 6, 8, 10\}$. We decided not to include a more common SVM with radial basis function (rbf) kernel in this study since the nonlinearity of the model can yield extreme positive/negative effects on the performance e.g., serious overfitting on noisy linearly separable data and incomparably better accuracy on nonlinear datasets. As such, it is difficult to understand the intrinsic label noise robustness of the model. Still, we believe that the SVMs employing nonlinear kernels could be well compared with the kernelised version of the above robust classifiers in the future study.

The comparison with the traditional classifiers will validate the practical usefulness of label noise modelling. The comparison with gLR and LaplaceBCN will demonstrate whether the added flexibility of the GMMLR is beneficial while the comparison with rLR should highlight the inadequacies of existing latent variable model in the presence of instance-dependent label noises.

We will now present the predictive performances as measured by classification errors on 11 UCI datasets subjected to the four artificial label noises mentioned above. We note that, unlike the previous set of experiments, here we performed 100 experiment repetitions for each label noise configurations in order to obtain more reliable statistics. Table 10 summarises the results for *gamma noise* (top) and *gaussian noise* (bottom) while the results for *mixture noise* and *random noise* are presented in Table 11.

Let us first compare the proposed method to the traditional LR. From the results we see that the proposed GMMLR model exhibited more robustness towards label noise than the traditional classifier, which is unaware of mislabelling, in the majority of the datasets and under all noise types. For each of the noise types, GMMLR was statistically significantly better than LR on at least 4 datasets. Further GMMLR was never worse than LR in all of the cases. These observations well demonstrate the benefit of having the noise model built into the classifier while also highlighting the inadequacy of the traditional LR classifier in the presence of label noise.

Compared to rLR which assumes uniform noise rate, GMMLR tends to perform better in the cases where noise distributions are not uniform, i.e., *mixture noise*, *gaussian noise* and *gamma noise*. This signifies the necessity

¹ We used LIBLINEAR [31] in this study.

of having more flexible noise model. We observe further that in random noise cases, rLR sometime outperformed GMMLR, which is somehow expected. Still, GMMLR did quite well under random noise and this demonstrates the generality of the proposed model. Further, within the group of classifiers adopting instance-dependent noise model, LaplaceBCN and gLR showed impressive performances compared to GMMLR in the case where the noise matches their assumptions, e.g., *gamma noise*. However, when noise distribution deviates from their assumptions, for example *gaussian noise*, the number of times that the two classifiers were better than GMMLR was reduced.

We observed that SVM-Linear is quite robust to *gamma noise* having more winning counts than GMMLR. We think that this might due to the inherent robustness of SVM towards noise near decision boundary. As by formulation, SVM can neglect some hard-to-classify points in the confusion region through the use of slack variables and can instead selects points on the boundary of the confusion region to be the support vectors. So, in the presence of label noises near decision boundary, SVM seems to be able to select meaningful support vectors and thus the orientation of its decision boundary is not much altered. However, SVM-Linear seems to be troubled most by random noise. This is also understandable since if all of the examples are subjected to label flipping, it is then harder to ensure the correct orientation of the decision boundary since the selected support vectors might not be the optimal ones.

The last row of each tables provides a summary of the number of times that GMMLR won, tied or lose to its peers as tested by Wilcoxon ranksum test at 0.05 level. We see that GMMLR is the top performer in the non-random noise scenarios based on the winning counts. In the presence of 30% random label noise, we observe that the performances of GMMLR are comparable to rLR, and are generally better than those of LaplaceBCN and SVM. We also notice that under random noise, gLR seems to do better than LaplaceBCN. We speculate that the modified exponential distribution used by gLR has heavier tails than the noise distribution used by LaplaceBCN. It was then able to approximate uniform label noise better. Still, by formulation, the GMMLR is more flexible and so better classification performances were observed in the random noise case.

Table 10: Mean errors (%) and standard errors of GMMLR compared to existing traditional and robust classifiers. Boldface entries (also marked by ‘-’ sign) are statistically worse than GMMLR, Underline entries (‘+’ sign) are statistically better than GMMLR while normal entries (‘ \approx ’ sign) are statistically equivalent to GMMLR as tested by Wilcoxon rank sum test at 0.05 level. The last row counts the cases where GMMLR is better than, equivalent to, and worse than the method in that column.

Dataset	Classification error (%) under <i>gamma noise</i>					
	LR	gLR	rLR	LaplaceBCN	SVM-Linear	GMMLR
australian	19.04 \pm 4.77 -	19.14 \pm 4.80 -	18.21 \pm 5.61 \approx	18.94 \pm 6.49 \approx	<u>15.23</u> \pm <u>3.33</u> +	17.51 \pm 4.44
biodeg	18.69 \pm 4.12 -	18.53 \pm 4.39 \approx	19.56 \pm 8.21 \approx	16.01 \pm 3.53 +	17.73 \pm 3.83 \approx	17.19 \pm 3.72
boston	20.09 \pm 4.60 -	19.67 \pm 4.75 -	20.33 \pm 8.40 \approx	19.76 \pm 8.68 \approx	20.66 \pm 5.11 -	17.85 \pm 5.35
diabetes	27.32 \pm 4.69 -	27.32 \pm 4.68 -	27.42 \pm 5.12 -	25.53 \pm 4.48 \approx	28.57 \pm 5.46 -	25.66 \pm 4.73
fertility	29.65 \pm 13.66 \approx	29.95 \pm 13.40 \approx	40.40 \pm 20.01 -	40.15 \pm 17.12 -	<u>20.30</u> \pm <u>15.07</u> +	30.05 \pm 13.49
german	29.64 \pm 4.48 \approx	29.66 \pm 4.48 \approx	32.23 \pm 8.29 -	<u>27.39</u> \pm <u>4.06</u> +	30.75 \pm 3.89 -	29.02 \pm 4.60
ionosphere	17.23 \pm 4.81 \approx	17.99 \pm 5.35 \approx	18.89 \pm 5.55 -	17.23 \pm 10.60 \approx	16.55 \pm 5.02 \approx	17.21 \pm 4.69
magic04	23.72 \pm 3.32 -	21.41 \pm 0.97 \approx	21.22 \pm 0.76 \approx	<u>20.86</u> \pm <u>0.70</u> +	24.53 \pm 6.22 -	21.22 \pm 0.88
musk	30.87 \pm 6.24 -	26.01 \pm 6.48 -	25.04 \pm 5.97 \approx	32.36 \pm 9.56 -	23.97 \pm 5.29 \approx	23.60 \pm 5.12
spambase	10.60 \pm 1.40 -	8.19 \pm 1.58 +	8.36 \pm 1.70 \approx	9.51 \pm 1.66 -	11.84 \pm 4.54 -	8.60 \pm 1.33
usps46	8.12 \pm 1.58 -	6.09 \pm 2.39 -	3.27 \pm 1.67 -	3.96 \pm 1.40 -	<u>0.32</u> \pm <u>0.33</u> +	2.67 \pm 1.24
Win-Tie-Loss	8-3-0	5-5-1	5-6-0	4-4-3	5-3-3	n/a
Dataset	Classification error (%) under <i>gaussian noise</i>					
	LR	gLR	rLR	LaplaceBCN	SVM-Linear	GMMLR
australian	15.04 \pm 3.00 \approx	14.99 \pm 3.09 \approx	15.36 \pm 3.15 \approx	14.65 \pm 3.09 \approx	14.91 \pm 3.07 \approx	14.94 \pm 2.88
biodeg	17.99 \pm 3.79 -	17.51 \pm 3.75 \approx	17.18 \pm 3.89 \approx	16.13 \pm 3.30 \approx	17.24 \pm 3.68 \approx	16.38 \pm 3.04
boston	17.12 \pm 4.39 -	16.40 \pm 3.99 -	15.41 \pm 3.65 \approx	14.42 \pm 3.43 \approx	17.07 \pm 4.16 -	15.02 \pm 3.26
diabetes	25.42 \pm 4.10 \approx	25.39 \pm 4.11 \approx	25.24 \pm 3.62 \approx	24.69 \pm 3.70 \approx	24.78 \pm 3.91 \approx	24.96 \pm 3.84
fertility	18.65 \pm 7.97 \approx	18.75 \pm 8.11 \approx	30.45 \pm 15.73 -	27.05 \pm 12.91 -	<u>13.60</u> \pm <u>7.22</u> +	19.30 \pm 8.20
german	26.70 \pm 3.39 \approx	26.68 \pm 3.40 \approx	28.14 \pm 5.41 -	25.79 \pm 3.34 \approx	<u>27.06</u> \pm <u>3.49</u> \approx	26.25 \pm 3.31
ionosphere	16.70 \pm 5.72 \approx	16.28 \pm 5.70 \approx	16.69 \pm 5.51 \approx	23.94 \pm 23.86 \approx	14.75 \pm 4.52 \approx	15.96 \pm 5.23
magic04	22.40 \pm 1.21 -	22.74 \pm 1.23 -	21.65 \pm 0.73 -	21.48 \pm 0.57 \approx	23.15 \pm 2.60 -	21.32 \pm 0.55
musk	30.25 \pm 6.75 -	25.45 \pm 6.62 -	25.02 \pm 6.12 -	31.65 \pm 10.78 -	22.94 \pm 4.57 \approx	22.27 \pm 5.09
spambase	10.72 \pm 1.62 -	8.41 \pm 1.58 \approx	8.23 \pm 1.35 \approx	8.21 \pm 1.00 +	11.24 \pm 2.63 -	8.47 \pm 1.00
usps46	5.55 \pm 2.96 -	2.94 \pm 1.83 -	1.49 \pm 1.00 \approx	2.49 \pm 1.25 -	<u>0.21</u> \pm <u>0.21</u> +	1.25 \pm 0.73
Win-Tie-Loss	6-5-0	4-7-0	4-7-0	3-7-1	3-6-2	n/a

Table 11: Mean errors (%) and standard errors of GMMLR compared to existing traditional and robust classifiers. Boldface entries (also marked by ‘-’ sign) are statistically worse than GMMLR, Underline entries (‘+’ sign) are statistically better than GMMLR while normal entries (‘ \approx ’ sign) are statistically equivalent to GMMLR as tested by Wilcoxon rank sum test at 0.05 level. The last row counts the cases where GMMLR is better than, equivalent to, and worse than the method in that column.

Dataset	Classification error (%) under <i>mixture noise</i>					
	LR	gLR	rLR	LaplaceBCN	SVM-Linear	GMMLR
australian	15.12 \pm 3.30 \approx	15.22 \pm 3.37 \approx	15.62 \pm 3.32 \approx	15.49 \pm 6.27 \approx	14.77 \pm 3.06 \approx	15.20 \pm 3.20
biodeg	19.28 \pm 4.09 -	18.55 \pm 4.12 \approx	18.14 \pm 4.51 \approx	16.58 \pm 2.78 +	18.57 \pm 3.81 \approx	17.58 \pm 3.18
boston	16.81 \pm 4.15 -	15.83 \pm 3.96 \approx	15.86 \pm 3.88 \approx	15.44 \pm 7.52 \approx	16.90 \pm 3.87 -	15.45 \pm 3.46
diabetes	24.32 \pm 3.48 \approx	24.33 \pm 3.52 \approx	24.43 \pm 4.20 \approx	24.05 \pm 3.43 \approx	24.27 \pm 3.56 \approx	24.05 \pm 3.51
fertility	19.25 \pm 9.93 \approx	19.30 \pm 9.77 \approx	26.10 \pm 12.18 -	24.65 \pm 12.76 -	<u>15.15 \pm 9.36 +</u>	19.50 \pm 9.76
german	27.11 \pm 3.15 \approx	27.11 \pm 3.13 \approx	27.43 \pm 4.21 \approx	25.82 \pm 2.56 \approx	28.09 \pm 3.39 -	26.32 \pm 2.82
ionosphere	17.73 \pm 5.82 \approx	17.32 \pm 5.33 \approx	17.59 \pm 5.25 \approx	31.04 \pm 27.96 \approx	<u>15.21 \pm 4.90 +</u>	16.85 \pm 5.48
magic04	22.23 \pm 1.29 -	23.76 \pm 2.30 -	22.25 \pm 1.03 -	21.95 \pm 0.91 -	23.48 \pm 4.45 -	21.51 \pm 0.76
musk	32.08 \pm 6.45 -	28.68 \pm 6.24 -	28.21 \pm 6.64 -	32.57 \pm 8.99 -	26.54 \pm 5.79 \approx	25.51 \pm 5.60
spambase	11.79 \pm 2.31 -	12.65 \pm 5.91 \approx	11.57 \pm 4.60 \approx	8.89 \pm 1.36 +	12.79 \pm 3.10 -	9.37 \pm 1.15
usps46	3.92 \pm 2.78 -	3.09 \pm 2.21 -	1.70 \pm 1.31 \approx	2.53 \pm 1.48 -	<u>0.23 \pm 0.25 +</u>	1.33 \pm 1.02
Win-Tie-Loss	6-5-0	3-8-0	3-8-0	4-5-2	4-4-3	n/a
Dataset	Classification error (%) under <i>random noise</i>					
	LR	gLR	rLR	LaplaceBCN	SVM-Linear	GMMLR
australian	16.53 \pm 3.40 \approx	16.66 \pm 3.51 \approx	17.32 \pm 5.59 \approx	16.32 \pm 4.06 \approx	15.17 \pm 3.21 \approx	15.69 \pm 3.26
biodeg	21.70 \pm 5.04 \approx	21.24 \pm 6.49 \approx	21.11 \pm 7.56 \approx	<u>18.59 \pm 4.19 +</u>	20.79 \pm 4.57 \approx	20.39 \pm 4.50
boston	20.73 \pm 4.74 -	19.33 \pm 4.75 \approx	19.96 \pm 7.10 \approx	17.80 \pm 5.22 \approx	20.30 \pm 4.79 -	18.24 \pm 4.84
diabetes	27.86 \pm 4.48 \approx	27.81 \pm 4.47 \approx	26.14 \pm 4.49 \approx	26.31 \pm 3.78 \approx	29.29 \pm 5.09 -	26.89 \pm 4.23
fertility	25.60 \pm 12.40 \approx	25.60 \pm 12.48 \approx	35.95 \pm 19.47 -	36.50 \pm 16.94 -	<u>19.50 \pm 15.15 +</u>	25.80 \pm 12.14
german	29.40 \pm 4.13 \approx	29.38 \pm 4.14 \approx	29.35 \pm 5.50 \approx	27.32 \pm 3.18 \approx	30.31 \pm 4.03 -	28.31 \pm 3.72
ionosphere	21.00 \pm 5.71 \approx	21.61 \pm 7.43 \approx	22.46 \pm 8.63 \approx	25.58 \pm 19.87 \approx	18.97 \pm 5.38 \approx	19.85 \pm 5.97
magic04	23.83 \pm 2.30 -	<u>21.01 \pm 0.87 +</u>	24.06 \pm 4.92 \approx	22.38 \pm 6.39 -	24.45 \pm 3.41 -	22.13 \pm 1.24
musk	33.90 \pm 7.66 -	30.41 \pm 7.06 -	29.60 \pm 7.53 -	34.96 \pm 9.29 -	28.86 \pm 6.49 \approx	27.16 \pm 8.02
spambase	13.71 \pm 1.93 -	11.22 \pm 5.21 \approx	10.88 \pm 4.42 +	11.14 \pm 2.30 \approx	14.17 \pm 4.84 -	11.15 \pm 1.42
usps46	16.08 \pm 3.06 -	13.20 \pm 6.91 -	7.61 \pm 5.38 \approx	14.67 \pm 9.51 -	<u>0.73 \pm 0.49 +</u>	7.57 \pm 3.95
Win-Tie-Loss	5-6-0	2-8-1	2-8-1	4-6-1	5-4-2	n/a

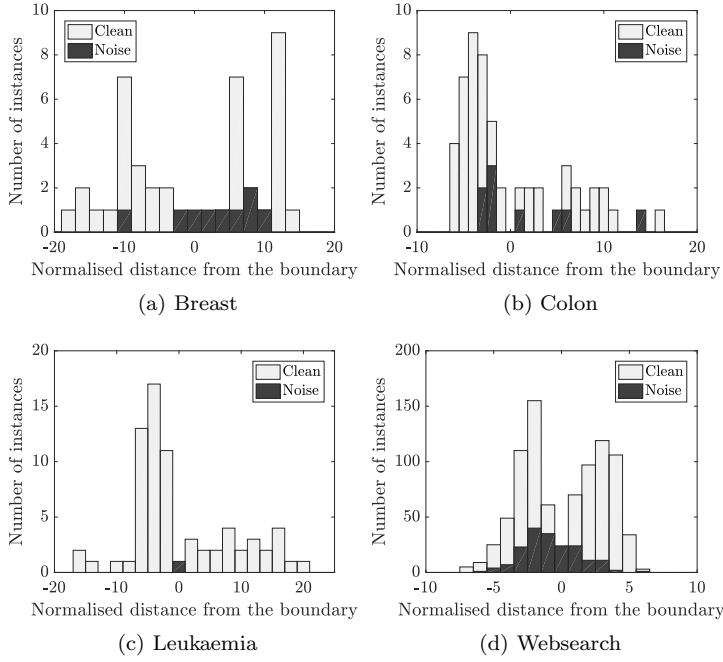


Fig. 7: Characteristics of the label noises in 4 real-world datasets.

4.6 Comparisons with other classifiers on noisy real-world datasets

Real-world dataset can be challenging since we have no control over the characteristics of the noise. In this section we will investigate the algorithms' robustness against real label noise. We use five datasets from various application domains and probably with different sources of labelling errors. For colon cancer dataset the task is to classify the microarray samples as normal tissue (N) or tumour (T). According to the literature [27], the samples T2, T30, T33, T36, T37, N8, N12, N34, N36 may be mislabelled. For breast cancer data, the aim is to distinguish between estrogen positive and estrogen negative observations. Biological evidence [26] suggests that the samples 11, 14, 16, 31, 33, 40, 43, 45, 46 are mislabelled. The objective of leukaemia data analysis is to discriminate between Acute Lymphoblastic Leukaemia (ALL) and Acute Myeloid Leukaemia (AML). Unlike the previous two datasets, the only one suspect mislabelling in this dataset was identified by consensus of previous label noise detection algorithms. The websearch dataset is an image classification datasets where the images were gathered quickly via image search engine. The annotations were also sloppily done by taken the query keyword as the class label directly. Obviously, some of the images returned by the search engine might not relate to the keyword (class label) at all. The websearch dataset

is composed of two classes constructed using keywords ‘bike’ and ‘not bike’. The ground truths for this dataset were manually identified [15]. Lastly, LULC dataset is a land usage classification dataset. The images were obtained from satellite imagery while labels of the land patches were extracted from crowd-sourced OpenStreetMap (OSM) database. We selected two out of six classes, namely ‘forest’ and ‘grass’, which are most often mistaken for one another according to the analysis in [3] for our experiment. The LULC dataset was divided into training and testing sets apriori by the creator [3]. The approximate noise rate was found to be around 40% but the exact noisy labels were not identified. The testing labels, however, are all clean. For each of the datasets except breast and LULC, we trained the classifiers on 80% of the data and held out 20% for testing. We used 90/10 train/test splitting for breast dataset due to the small sample size and high dimensional nature of the data. In the case of breast, colon, leukaemia and websearch datasets, mislabelled instances that were split into test set were corrected before testing using the available ground truths. For LULC data, we sampled 80% of the instances in the training set for training, and evaluated the model on the whole test set. Since the test labels of LULC dataset are clean, we did not perform any label correction.

Since the ground truth labels are available (except for LULC dataset), we first summarised the number of correct labels versus noisy labels as a function of their distances from the decision boundary in order to get a better idea of how noise distributions might look like. The plots generated using the weight vector of LR trained in noise-free data are depicted in Fig.7. We observe that the noises spread rather randomly in the case of breast and colon cancer data. Although, according to the literature [26] label noise in breast dataset is more symmetric. That is, the ratio of mislabelling in class ALL and AML are quite similar, and previous theoretical studies seem to suggest that symmetric label noise is relatively harmless [16]. On the other hand, the noise in colon cancer is somewhat asymmetric and we expect that the benefit of label noise modelling, especially the one assuming uniform random noise, will be more obvious in this case. The leukaemia dataset contains only one mislabelling and it should pose little problem to all of the algorithms. Still, it is interesting to see how much label noise aware algorithms would be fooled into believing that some hard-to-classify samples are mislabelled. Lastly, we see that majority of label noise in websearch data occurs near the decision boundary, in which case LaplaceBCN and SVM-linear might be able to do well in this case.

We now present the classification results obtained from 100 experiment repetitions in Table 12. Boldface entries highlight the cases where GMMLR was significantly better, while underlined entries marks the case where GMMLR was worse than its peers as tested by Wilcoxon ranksum test at 0.05 level.

Let us first consider the classifiers’ performances in the cases of breast cancer, colon cancer and leukaemia datasets. Statistically, the differences in predictive performances of all models against GMMLR under these three datasets were not significant except for the case of colon cancer data where we spot impressive performance from rLR. We were quite surprised how well the traditional classifiers (LR and SVM) performed on these noisy datasets, and at

Table 12: Mean classification errors (%) and standard errors of the six algorithms on real-world datasets that genuinely contains label noise. Boldface entries highlight the cases where GMMLR was significantly better. Underlined entries mark the cases where GMMLR was worse than its peers as tested by Wilcoxon ranksum test at 0.05 level.

Dataset	Classification model		
	LR	gLR	rLR
breast	15.80 \pm 13.42 \approx	15.40 \pm 12.98 \approx	16.00 \pm 12.71 \approx
colon	10.83 \pm 9.29 \approx	10.67 \pm 9.41 \approx	<u>7.25 \pm 7.55 +</u>
leukaemia	4.07 \pm 4.91 \approx	4.20 \pm 5.25 \approx	<u>4.13 \pm 5.17 \approx</u>
websearch	20.53 \pm 3.12 -	20.22 \pm 3.16 -	18.50 \pm 2.94 \approx
LULC	33.33 \pm 1.08 -	34.27 \pm 1.33 -	34.24 \pm 1.36 -
Win-Tie-Loss	2-3-0	2-3-0	1-3-1

Dataset	Classification model		
	LaplaceBCN	SVM-Linear	GMMLR
breast	15.80 \pm 13.42 \approx	14.00 \pm 13.18 \approx	15.40 \pm 13.29
colon	11.08 \pm 8.95 \approx	10.08 \pm 8.40 \approx	10.25 \pm 8.93
leukaemia	4.20 \pm 5.16 \approx	4.13 \pm 5.08 \approx	4.20 \pm 5.16
websearch	18.67 \pm 2.95 \approx	<u>16.20 \pm 2.72 +</u>	18.91 \pm 2.84
LULC	32.58 \pm 1.27 -	46.12 \pm 1.00 -	31.94 \pm 0.91
Win-Tie-Loss	1-4-0	1-3-1	n/a

the same time wondered why label noise modelling did not help. We speculate firstly that the sample size and data dimensionality might affect the learning of label noise model such that the estimated noise rates do not reflect the real flipping probabilities when the number of training instances is limited. Indeed, by observing the label flipping probabilities learned by GMMLR in these datasets, we found that they were quite close to zero, making the classifier behaves much like the traditional LR. Meanwhile, rLR which is the simplest label noise model seems to required less training instances compared to other robust models, and hence was able to perform well under colon cancer dataset. Among the three datasets, colon cancer is probably less demanding due to its larger sample size and relatively low dimensionality. This is also where we started to see some slight improvement of GMMLR over LR.

For websearch dataset, we see that SVM-Linear was the top performer. GMMLR was worse than SVM-Linear but not by a large margin. Its performances were comparable to rLR and LaplaceBCN, the models which assume boundary conditioned noise. The improvements of GMMLR over traditional LR was well apparent, and this once again demonstrates the benefit of having the label noise model built into the traditional classifier. Interestingly, SVM-Linear was quite robust to noise around the decision boundary. We think that the slack variables might have helped masking out the negative effect of instances that were hard to classify (which could be the result of flipped labels) to some extent.

In the case where training examples are abundance as in LULC dataset, it turned out that GMMLR can better estimate the noise rates and hence an improved classification accuracy can be expected. We see that GMMLR

outperformed its peers in this particular case. We were surprised by the twisted performance of the SVM which was rather impressive in the previous four datasets. The imbalance of class ‘forest’ (7431 instances) and class ‘grass’ (446 instances) might added up to the already severe noise level (estimated to be around 40%). The last row of Table 12 summarises the win-tie-loss counts of GMMLR against other classifiers. We observe that GMMLR was never worse than the traditional LR, gLR and LaplaceBCN, while its winning counts are comparable to SVM-Linear and rLR. Overall, the empirical evidences seem to suggest that it is still a very challenging task to come up with a general label noise model which can do well on all types of label noise. A particular model tends to work best only when the nature of label noise matches its noise assumption. Still, through the series of experiments, we have witnessed that the proposed GMMLR model was quite flexible compared to the existing ones. It was able to counteract the negative effects of various instance-dependent label noises. It also performed reasonable well in the presence of random label noise.

5 Conclusion

In this paper, we have proposed a novel label noise model which is based on a mixture of Gaussians. We incorporated the new noise model into a logistic regression and developed a Gaussian Mixture Model based robust Logistic Regression classifier. Experimental results demonstrated that the added flexibility of the noise function helps counteracting negative effects of wide range of label noises including class-conditional label noise and instance-dependent label noises. Since the current work is still limited to binary classification problems, it is then interesting to investigate the multi-class version of the proposed model in the future.

Acknowledgements

The authors would like to thank anonymous reviewers for constructive comments. This research is financially supported by the Thailand Research Fund (Grant number: MRG59080235). Department of Computer Science, Faculty of Science at Chiang Mai University provides research and computing facilities.

References

1. Eyal Beigman and Beata Beigman Klebanov. Learning with annotation noise. In *ACL 2009, Proceedings of the 47th Annual Meeting of the Association for Computational Linguistics, 2-7 August 2009, Singapore*, pages 280–287, 2009.
2. Aleksander Kolcz and Gordon V. Cormack. Genre-based decomposition of email class noise. In *SIGKDD’09*, pages 427–436, 2009.
3. Brian A. Johnson and Kotaro Iizuka. Integrating openstreetmap crowdsourced data and landsat time-series imagery for rapid land use/land cover (lulc) mapping: Case study of the laguna de bay area of the philippines. *Applied Geography*, 67:140 – 149, 2016.

4. Rion Snow, Brendan O'Connor, Daniel Jurafsky, and Andrew Y. Ng. Cheap and fast - but is it good? evaluating non-expert annotations for natural language tasks. In *EMNLP*, pages 254–263, 2008.
5. Dan Shen, Jean-David Ruvini, and Badrul Sarwar. Large-scale item categorization for e-commerce. In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management*, CIKM '12, pages 595–604, New York, NY, USA, 2012. ACM.
6. Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive noisy labeled data for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2691–2699, 2015.
7. Benoît Frénay and Michel Verleysen. Classification in the Presence of Label Noise: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*, 25(5):845–869, 2014.
8. Aditya Krishna Menon, Brendan van Rooyen, and Nagarajan Natarajan. Learning from binary labels with instance-dependent corruption. *arXiv preprint arXiv:1605.00751*, 2016.
9. Battista Biggio, Blaine Nelson, and Pavel Laskov. Support vector machines under adversarial label noise. In *ACML*, volume 20 of *JMLR Proceedings*, pages 97–112. JMLR.org, 2011.
10. Raj S. Chhikara and Jim McKeon. Linear discriminant analysis with misallocation in training samples. *Journal of the American Statistical Association*, 79(388):899–906, 1984.
11. Neil D. Lawrence and Bernhard Schölkopf. Estimating a kernel fisher discriminant in the presence of label noise. In *ICML'01*, pages 306–313. Morgan Kaufmann, 2001.
12. Yunlei Li, Lodewyk F.A. Wessels, Dick de Ridder, and Marcel J.T. Reinders. Classification in the presence of class noise using a probabilistic kernel Fisher method. *Pattern Recognition*, 40(12):3349–3357, 2007.
13. Vikas C. Raykar, Shipeng Yu, Linda H. Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, and Linda Moy. Learning from crowds. *Journal of Machine Learning Research*, 11:1297–1322, 2010.
14. Jakramate Bootkrajang and Ata Kabán. Label-noise robust logistic regression and its applications. In *ECML-PKDD'12*, pages 143–158, 2012.
15. Jakramate Bootkrajang and Ata Kabán. Learning kernel logistic regression in the presence of class label noise. *Pattern Recognition*, 47(11):3641–3655, 2014.
16. Gábor Lugosi. Learning with an unreliable teacher. *Pattern Recognition*, 25:79–87, 1992.
17. Philip M. Long and Rocco A. Servedio. Random classification noise defeats all convex potential boosters. *Machine Learning*, 78(3):287–304, March 2010.
18. Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *NIPS'13*, pages 1196–1204, 2013.
19. Naresh Manwani and P. S. Sastry. Noise tolerance under risk minimization. *IEEE T. Cybernetics*, 43(3):1146–1151, 2013.
20. Aritra Ghosh, Naresh Manwani, and P. S. Sastry. Making risk minimization tolerant to label noise. *Neurocomputing*, 160:93–107, 2015.
21. Peter A. Lachenbruch. Discriminant analysis when the initial samples are misclassified ii: Non-random misclassification models. *Technometrics*, 16(3):pp. 419–424, 1974.
22. Jakramate Bootkrajang. A generalised label noise model for classification in the presence of annotation errors. *Neurocomputing*, 192:61–71, 2016.
23. Jun Du and Zhihua Cai. Modelling class noise with symmetric and asymmetric distributions. In *AAAI*, pages 2589–2595, 2015.
24. Mark Schmidt. minfunc: unconstrained differentiable multivariate optimization in matlab, 2005.
25. Yunmei Chen and Xiaojing Ye. Projection onto a simplex. *arXiv preprint arXiv:1101.6081*, 2011.
26. Mike West, Carrie Blanchette, Holly Dressman, Erich Huang, Seiichi Ishida, Rainer Spang, Harry Zuzan, John A. Olson Jr., Jeffrey R. Marks, and Joseph R. Nevins. Predicting the clinical status of human breast cancer by using gene expression profiles. *Proc. Natl. Acad. Sci. USA*, 98(20):pp. 11462–11467, 2001.

27. U. Alon, N. Barkai, D. A. Notterman, K. Gishdagger, S. Ybarradagger, D. Mackdagger, and A. J. Levine. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc. Natl. Acad. Sci. USA*, 96(12):6745–6750, 1999.
28. T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, and C. D. Bloomfield. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 286:531–537, 1999.
29. Andrew Frank, Arthur Asuncion, et al. Uci machine learning repository, 2010.
30. Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, December 2006.
31. Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008.