# Accepted Manuscript

Learning from Multiple Annotators: Distinguishing Good from Random Labelers
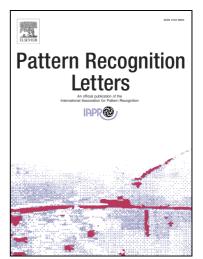
Filipe Rodrigues, Francisco Pereira, Bernardete Ribeiro

# Learning from Multiple Annotators: Distinguishing Good from Random Labelers

Filipe Rodrigues[1], Francisco Pereira[2] and Bernardete Ribeiro[1]

[1]*Centre for Informatics and Systems of the University of Coimbra (CISUC)*
*Department of Informatics Engineering, University of Coimbra*
*3030-290 Coimbra, Portugal*
*Tel.: +351 239790056*
*{fmpr,bribeiro}@dei.uc.pt*

[2]*Singapore-MIT Alliance for Research and Technology (SMART)*
*1 CREATE Way, Singapore 138602*
*Tel.: +65 93233653*
*camara@smart.mit.edu*

## Abstract

With the increasing popularity of online crowdsourcing platforms such as Amazon Mechanical Turk (AMT), building supervised learning models for datasets with multiple annotators is receiving an increasing attention from researchers. These platforms provide an inexpensive and accessible resource that can be used to obtain labeled data, and in many situations the quality of the labels competes directly with those of experts. For such reasons, much attention has recently been given to annotator-aware models. In this paper, we propose a new probabilistic model for supervised learning with multiple annotators where the reliability of the different annotators is treated as a latent variable. We empirically show that this model is able to achieve state of the art performance, while reducing the number of model parameters, thus avoiding a potential overfitting. Furthermore, the proposed model is easier to implement and extend to other classes of learning problems such as sequence

labeling tasks.

## 1. Introduction

Crowdsourcing (Howe, 2008) is rapidly changing the way datasets are built. With the development of crowdsourcing platforms such as Amazon Mechanical Turk (AMT)[1], it is becoming increasingly easier to obtain labeled data for a wide range of tasks covering different areas such as Computer Vision, Natural Language Processing, Speech Recognition, etc. The attractiveness of these platforms comes not only from their low cost and accessibility, but also from the surprisingly good quality of the labels obtained, which in many cases competes directly with those of "experts" (Snow et al., 2008). Furthermore, by distributing the workload among multiple annotators, labeling tasks can be completed much faster.

The current trend of social web, where citizens' participation is growing in many forms, has come to stay, and information is being produced at a massive rate. This information can take many forms: document tags, opinions, product ratings, user clicks, contents, etc. These new sources of data also motivate the development of new machine learning approaches for learning from multiple sources.

On another perspective, there are tasks for which ground truth labels simply cannot be obtained due to their highly subjective nature. Consider

---

[1]http://www.mturk.com

2

for instance the tasks of sentiment analysis, movie rating or keyphrase extraction. These tasks are subjective in nature and hence no absolute gold standard can be defined. In such cases the only attainable goal is to build a model that captures the *wisdom of the crowds* (Surowiecki, 2004) as well as possible. For such tasks crowdsourcing platforms like AMT become a natural solution. However, the large amount of labeled data needed to compensate for the heterogeneity of annotators' expertise can rapidly rise its actual cost beyond acceptable values. Since different annotators have different levels of expertise, it is important to consider how *reliable* the annotators are when learning from their answers, and a parsimonious solution needs to be designed that is able to deal with such real world constraints (e.g. annotation cost) and heterogeneity.

Even in situations where a ground truth can be obtained, it may be too costly. For example, in Medical Diagnosis, determining whether a patient has cancer may require a biopsy, which is an invasive procedure, and thus should only be used as a last resource. On the other hand, it is rather easy for a diagnostician to consult its colleagues for their opinions before making a decision. Therefore, although there is no crowdsourcing involved here, there are still multiple experts, with different levels of expertise, providing their own (possibly wrong) opinions, from which we have to be able to learn from.

Many approaches have recently been proposed that deal with this increasingly important problem of supervised learning from multiple annotators in different paradigms: classification (Raykar et al., 2009; Yan et al., 2011), regression (Groot et al., 2011), ranking (Wu et al., 2011), etc. However, most of the work developed so far is centered on the *unknown* true labels of

3

the data, for which noisy versions are provided by the various annotators. Therefore, there has been a tendency to include these *unobserved* true labels as latent variables in a probabilistic framework, which, as we demonstrate, is not necessarily the best option. Furthermore, this choice of latent variables hinders a natural extension of these approaches to structured prediction problems such as sequence labeling tasks due to combinatorial explosion of possible outcomes of the latent variables. Contrarily to these approaches, we argue that the focus should be on the annotators, and that including the also *unknown* reliabilities of the annotators as latent variables can be preferable, since it not only leads to simpler models that are less prone to overfitting, but also bypasses the problem of the high number of possible labelings to marginalize over.

In this paper, we propose a new probabilistic model that explores these ideas, and explicitly handles the annotators' reliabilities as latent variables. We empirically show, using both simulated annotators and human annotators from AMT, that for many tasks the new model can be competitive with the state of the art methods, and can even significantly outperform previous approaches under certain conditions. Although we focus on multi-class Logistic Regression as the base classifier, the proposed model is simple and generic enough to be implemented with other classifiers. Furthermore the extension to structured prediction problems such as sequence labeling tasks can be much easier than with latent ground truth models (e.g. Raykar et al. (2010); Yan et al. (2011)).

The remainder of this paper is organized as follows: Section 2 provides the reader with an overview of state of the art; Section 3 clarifies the problem

4

⁸⁶ with latent ground truth models; Section 4 presents the proposed model, and ⁸⁷ Section 5 compares the results obtained by this model with two majority ⁸⁸ voting baselines and a state of the art approach; the article will end with a ⁸⁹ short discussion and conclusions (Section 6).

## 2. State of the art

⁹¹ There is considerable work on estimating ground truth labels from the ⁹² responses of multiple annotators. Most of the early important works were in ⁹³ the fields of Biostatistics and Epidemiology. In 1979, Dawid and Skene (1979) ⁹⁴ proposed an approach for estimating the error rates of multiple patients ⁹⁵ (annotators) given their responses (labels) to multiple medical questions. ⁹⁶ However, like most of the early works with multiple annotators, this work ⁹⁷ only focused on estimating the unobserved ground truth labels. Only later, ⁹⁸ researchers started paying more attention to the specific problem of learning ⁹⁹ a classifier from the multiple annotator's data. In 1995, Smyth et al. (1995) ¹⁰⁰ proposed a similar approach to the one from Dawid and Skene (1979) to ¹⁰¹ estimate the ground truth from the labels of multiple experts, which was ¹⁰² then used to train a classifier. As with previous works, the authors employed ¹⁰³ a model where the unknown true labels were treated as latent variables.

¹⁰⁴ More recently, with the increasing popularity of AMT and other crowd- ¹⁰⁵ sourcing and work-recruiting platforms, researchers started recognizing the ¹⁰⁶ importance of the problem of learning from the labels of multiple non-expert ¹⁰⁷ annotators. The researchers' interest grew even further with works such as ¹⁰⁸ (Snow et al., 2008) and (Novotney and Callison-Burch, 2010), which show ¹⁰⁹ that, for many tasks, learning from multiple non-experts can be as good as

5

110 learning from an expert.

111    With the rising interest in crowdsourcing as a source of labeled data,
112 more challenging approaches for learning from multiple annotators started
113 to appear. In 2009, Raykar et al. (2009) proposed an innovative probabilis-
114 tic approach where the unknown ground truth labels and the classifier are
115 learnt jointly. By handling the unobserved ground truth labels as latent vari-
116 ables, the authors are able to find the maximum likelihood parameters for
117 their model by iteratively estimating the posterior distribution of the ground
118 truth labels and then using this estimate to determine the qualities of the
119 annotators and the parameters of a Logistic Regression model. Unlike most
120 of the previous works, this approach also has the advantage of relaxing the
121 requirement of repeated labeling, i.e. the same instance being annotated by
122 multiple annotators. Later works then relaxed other assumptions made by
123 the authors. For example, Yan et al. (2010) relaxed the assumption that
124 the quality of the labels provided by the annotators does not depend on the
125 instance they are labeling.

126    This main line of work also inspired many variations and extensions in the
127 past couple of years. Groot et al. (2011) proposed an extension of Gaussian
128 processes to do regression in a multiple annotator setting. In the field of
129 ranking, Wu et al. (2011) presented an approach to learn how to rank from
130 the opinions of multiple annotators. In an active learning setting, Yan et al.
131 (2011) proposed an approach for multiple annotators by providing answers to
132 the following questions: what instance should be selected to be labeled next
133 and which annotators should label it? On a different perspective, in (Donmez
134 et al., 2010) the authors propose the use of a particle filter to model the

6

135  time-varying accuracies of the different annotators. Despite the plausibility

136  of their assumptions, i.e. it is legitimate to assume that the quality of the

137  labels provided by an annotator will vary with time, the results obtained

138  showed only a small improvement on the performance of their model through

139  the inclusion of this time dependance.

140      The approaches above mentioned typically treat the *unknown* ground

141  truth labels as latent variables and build a model on that basis. We argue

142  that explicitly handling the reliabilities of the annotators as latent variables,

143  as opposed to the true labels, in a fashion that slightly resembles a mixture of

144  experts (Jacobs et al., 1991; Bishop, 2006), brings many attractive advantages

145  and can, under certain conditions, outperform latent ground truth models.

146  **3. The problem with latent ground truth models**

147      In order to help motivate the proposed model, we now introduce a typ-

148  ical class of approaches for learning from multiple annotators, in which the

149  *unknown* true labels are treated as latent variables (e.g. Raykar et al. (2009,

150  2010); Yan et al. (2010)).

151      Let $y_i^r$ be the label assigned to instance $\mathbf{x}_i$ by the $r^{th}$ annotator, and let

152  $y_i$ be the true (unobserved) label for that instance. Contrarily to a typical

153  classification problem with a single annotator, in a setting with $R$ annotators,

154  a dataset $\mathcal{D}$ with size $N$ consists of a set of labels $\{y_i^1, y_i^2, ..., y_i^R\}$ for each of

155  the $N$ instances $\mathbf{x}_i$.

156      In general, the class of models we refer to as "latent ground truth mod-

157  els" tend to assume the following generative process: for each instance $\mathbf{x}_i$

158  there is an *unobserved* true label $y_i$, and each of the different annotators in-
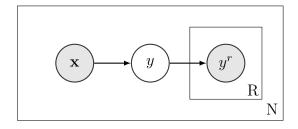
7

Figure 1: Plate representation of general latent ground truth model.

159 dependently provides its own version $(y_i^r)$ of this true label, which in practice
160 corresponds to an approximation to the real label $y_i$. Figure 1 depicts such
161 a model in plate notation. Shaded nodes represent observed variables, and
162 non-shaded nodes represent unobserved (latent) variables.

If besides the dataset $\mathcal{D} = \{y_i^1, ..., y_i^R, \mathbf{x}_i\}_{i=1}^N$ we were given the true labels
$\mathcal{Y} = \{y_i\}_{i=1}^N$ as well, the likelihood for this model would take the form

$$p(\mathcal{D}, \mathcal{Y}) = \prod_{i=1}^N \Big( p(y_i|\mathbf{x}_i) \prod_{r=1}^R p(y_i^r|y_i) \Big). \tag{1}$$

163 Since we do not actually observe the true labels $y_i$ we must treat them as
164 latent variables and marginalize them out of the likelihood, and this leads
165 us to the first problem with this approach: although this marginalization is
166 not difficult for classification problems where the number of classes $(K)$ is
167 small, for other types of problems like sequence labeling tasks (or any task
168 with structured outputs), marginalizing over the output space is intractable
169 in general (Sutton, 2012). If we consider, for example, the tasks of part-
170 of-speech (POS) tagging or Named Entity Recognition (NER), which are
171 usually handled as a sequence labelling problems, it is easy to see that the
172 number of possible label sequences grows exponentially with the length of
173 the sentence, deeming the marginalization over the output space intractable.

8

<sub>174</sub> The second problem with this class of models is related with the prob-

<sub>175</sub> ability $p(y_i^r|y_i)$, which for a classification problem with $K$ classes requires a

<sub>176</sub> $K \times K$ table of parameters for each annotator. Even though this approach

<sub>177</sub> allows to capture certain biases in the annotators answers, like for example

<sub>178</sub> the tendency to confuse two classes, in practice, on a crowdsourcing platform

<sub>179</sub> like AMT, each annotator only labels a rather small set of instances. There-

<sub>180</sub> fore, under such conditions, having a model with so many parameters for

<sub>181</sub> the reliability of the annotators can easily lead to overfitting. Consider, for

<sub>182</sub> example, a classification problem with 10 classes. Such a problem requires a

<sub>183</sub> total of 100 parameters (a $10 \times 10$ probability table) to model the expertise

<sub>184</sub> of a single annotator. To effectively learn such a number of parameters, each

<sub>185</sub> annotator would be required to label a large number of instances, at least in

<sub>186</sub> the order of the thousands, something that is both unrealistic and hard to

<sub>187</sub> control in a crowdsourcing platform.

<sub>188</sub> Taking these issues into consideration, we developed a new probabilis-

<sub>189</sub> tic model for learning from multiple annotators, which we present in the

<sub>190</sub> following section.

<sub>191</sub> **4. Proposed model**

<sub>192</sub> *4.1. Maximum likelihood estimator*

Given a dataset $\mathcal{D} = \{y_i^1, ..., y_i^R, \mathbf{x}_i\}_{i=1}^N$ with $N$ instances and $R$ different

annotators, and assuming that the instances are independent and identically

distributed (i.i.d.), the likelihood is given by

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{N} p(y_i^1, ..., y_i^R|\mathbf{x}_i, \theta) \tag{2}$$
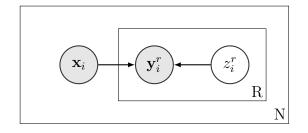
9

Figure 2: Plate representation of the proposed model.

193  where $\theta$ denotes the model parameters.

Let us now assume the following generative process of the annotators' labels: when the annotators are asked to provide a label to a given instance $\mathbf{x}_i$, they flip a biased coin, and based on the outcome of those coin flips, they decide whether or not to provide the correct label. This intuition amounts to introducing a binary random variable $z_i^r$, whose value indicates whether the $r^{th}$ annotator labeled the $i^{th}$ instance correctly or not. Hence, $z_i^r \sim Bernoulli(\pi_r)$, where $\pi_r$ is the accuracy of the $r^{th}$ annotator, and

$$p(z_i^r | \pi_r) = (\pi_r)^{z_i^r} (1 - \pi_r)^{1 - z_i^r}. \tag{3}$$

194  The expectation of this Bernoulli random variable $\mathbb{E}[z_i^r] = p(z_i^r = 1)$ can be
195  interpreted as the probability of an annotator providing a correct label or, in
196  other words, as an indicator of how reliable an annotator is. For the sake of
197  simplicity, we assume that an unreliable annotator provides labels according
198  to some random model $p_{\text{Rand}}(y_i^r = k | \mathbf{x}_i)$.

199  Figure 2 shows a plate representation of this generative model. Notice
200  that the variables $z_i^r$ are not observed in this model, hence their nodes are
201  not shaded in the figure.

If we were told the true values for $\mathcal{Z} = \{z_i^1, ..., z_i^R\}_{i=1}^N$, and assuming

10

the annotators make their decisions independently of the each other, the complete-data likelihood could then be factored as

$$p(\mathcal{D}, \mathcal{Z}|\theta) = \prod_{i=1}^{N} \prod_{r=1}^{R} p(z_i^r|\pi_r) \; p(y_i^r|\mathbf{x}_i, z_i^r, \mathbf{w}) \qquad (4)$$

where $\theta = \{\boldsymbol{\pi}, \mathbf{w}\}$ are the model parameters. The values of $\boldsymbol{\pi} = \{\pi_r\}_{r=1}^{R}$ correspond to the parameters of the $R$ Bernoulli distributions (one for each annotator). In turn, $\mathbf{w}$ are the weights of a Logistic Regression model.

Following the generative process described above, we can now define $p(y_i^r|\mathbf{x}_i, z_i^r, \mathbf{w})$ as

$$p(y_i^r|\mathbf{x}_i, z_i^r, \mathbf{w}) = \left(p_{\mathrm{LogReg}}(y_i^r|\mathbf{x}_i, \mathbf{w})\right)^{z_i^r} \left(p_{\mathrm{Rand}}(y_i^r|\mathbf{x}_i)\right)^{1-z_i^r} \qquad (5)$$

where $p_{\mathrm{LogReg}}(y_i^r|\mathbf{x}_i, \mathbf{w})$ denotes the likelihood of the label provided by the $r^{th}$ annotator for the instance $\mathbf{x}_i$ according to a multi-class Logistic Regression model with parameters $\mathbf{w}$, which for a classification task with $K$ classes is given by

$$p_{\mathrm{LogReg}}(y_i^r = k|\mathbf{x}_i, \mathbf{w}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x}_i)}{\sum_{k'=1}^{K} \exp(\mathbf{w}_{k'}^T \mathbf{x}_i)}. \qquad (6)$$

Similarly, $p_{\mathrm{Rand}}(y_i^r|\mathbf{x}_i)$ denotes the likelihood of the label $y_i^r$ according to a random model, which we assume to be uniformly distributed. Hence,

$$p_{\mathrm{Rand}}(y_i^r = k|\mathbf{x}_i) = \frac{1}{K}. \qquad (7)$$

To summarize, this is akin to saying that if $z_i^r = 1$ then the label provided by the $r^{th}$ annotator ($y_i^r$) fits a Logistic Regression model, which is assumed to capture the correct (true) labeling process. Conversely, if $z_i^r = 0$ then $y_i^r$ is assumed to be drawn from a random model where all the classes are equiprobable.

11

Since we do not actually observe the set $\mathcal{Z}$ we must treat the variables $z_i^r$ as latent and marginalize them out of the likelihood by summing over all its possible outcomes. The (observed) data likelihood then becomes

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{N} \prod_{r=1}^{R} \sum_{z_i^r \in \{0,1\}} p(z_i^r|\pi_r) \, p(y_i^r|\mathbf{x}_i, z_i^r, \mathbf{w}). \tag{8}$$

Making use of equations 3 and 5, this expression can be further simplified, giving

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{N} \prod_{r=1}^{R} \Big( \pi_r \, p_{\mathrm{LogReg}}(y_i^r|\mathbf{x}_i, \mathbf{w}) + (1 - \pi_r) \, p_{\mathrm{Rand}}(y_i^r|\mathbf{x}_i) \Big). \tag{9}$$

210 Our goal is then to estimate the maximum likelihood parameters $\theta_{\mathrm{ML}}$,
211 which are found by determining $\theta_{\mathrm{ML}} = \arg\max_{\theta} \ln p(\mathcal{D}|\theta)$.

212 At this point, it is important to note that extending this approach to
213 sequence labeling problems, or any kind of structured prediction problems
214 in general, could be as simple as replacing in equation 5 the probabilities
215 $p_{\mathrm{LogReg}}(y_i^r|\mathbf{x}_i, \mathbf{w})$ and $p_{\mathrm{Rand}}(y_i^r|\mathbf{x}_i)$ with their sequence labeling counterparts,
216 which for $p_{\mathrm{LogReg}}(\cdot)$ could be an Hidden Markov Model (HMM) or a Condi-
217 tional Random Field (CRF), and updating the remaining equations accord-
218 ingly.

219 *4.2. Expectation-Maximization*

220 As with other latent variable models, we rely on the Expectation-
221 Maximization (EM) algorithm (Dempster et al., 1977) to optimize this oth-
222 erwise intractable maximization problem. The EM algorithm is an iterative
223 method for finding maximum likelihood solutions for probabilistic models
224 with latent variables, and consist of two steps: the E-step and M-step. In

12

<sup>225</sup> the E-step the posterior distribution of the latent variables is computed based
<sup>226</sup> on the current model parameters. This posterior distribution is then used
<sup>227</sup> to estimate the new model parameters (M-step). These two steps are then
<sup>228</sup> iterated until convergence.

<sup>229</sup>    If we observed the complete dataset $\{\mathcal{D}, \mathcal{Z}\}$ then the loglikelihood func-
<sup>230</sup> tion would simply take the form $\ln p(\mathcal{D}, \mathcal{Z}|\theta)$. Since we only have access to
<sup>231</sup> the "incomplete" dataset $\mathcal{D}$, our state of the knowledge about the values
<sup>232</sup> of $\mathcal{Z}$ (the reliabilities of the annotators) can be given by the posterior dis-
<sup>233</sup> tribution $p(\mathcal{Z}|\mathcal{D}, \theta)$. Therefore, instead of the complete data loglikelihood,
<sup>234</sup> we consider its expected value under the posterior distribution of the latent
<sup>235</sup> variable $p(\mathcal{Z}|\mathcal{D}, \theta)$, which corresponds to the E-step of the EM algorithm.
<sup>236</sup> Hence, in the E-step we use the current parameter values $\theta^{old}$ to find the
<sup>237</sup> posterior distribution of the latent variables in $\mathcal{Z}$. We then use this poste-
<sup>238</sup> rior distribution to find the expectation of the complete-data loglikelihood
<sup>239</sup> evaluated for some general parameter values $\theta$. This expectation is given by

$$
\begin{aligned}
&\mathbb{E}_{p(\mathcal{Z}|\mathcal{D},\theta_{old})}\Big[\ln p(\mathcal{D}, \mathcal{Z}|\theta)\Big] \\
&= \sum_{\mathcal{Z}} p(\mathcal{Z}|\mathcal{D}, \theta_{old}) \ln p(\mathcal{D}, \mathcal{Z}|\theta) \\
&= \sum_{i=1}^{N} \sum_{r=1}^{R} \sum_{z_i^r \in \{0,1\}} p(z_i^r|y_i^r, \mathbf{x}_i, \theta_{old}) \ln \Big( p(z_i^r|\pi_r) \ p(y_i^r|\mathbf{x}_i, z_i^r, \mathbf{w}) \Big).
\end{aligned} \tag{10}
$$

The posterior distribution of the latent variables $z_i^r$ (denoted by $\gamma(z_i^r)$)

13

can be estimated using the Bayes theorem giving

$$
\begin{aligned}
\gamma(z_i^r) &= p(z_i^r = 1 | y_i^r, \mathbf{x}_i, \theta^{old}) \\
&= \frac{p(z_i^r = 1 | \pi_r^{old}) \; p(y_i^r | \mathbf{x}_i, z_i^r = 1, \mathbf{w}^{old})}{p(z_i^r = 1 | \pi_r^{old}) \; p(y_i^r | \mathbf{x}_i, z_i^r = 1, \mathbf{w}^{old}) + p(z_i^r = 0 | \pi_r^{old}) \; p(y_i^r | \mathbf{x}_i, z_i^r = 0, \mathbf{w})} \\
&= \frac{\pi_r^{old} \; p_{\mathrm{LogReg}}(y_i^r | \mathbf{x}_i, \mathbf{w}^{old})}{\pi_r^{old} \; p_{\mathrm{LogReg}}(y_i^r | \mathbf{x}_i, \mathbf{w}^{old}) + (1 - \pi_r^{old}) \; p_{\mathrm{Rand}}(y_i^r | \mathbf{x}_i)}
\end{aligned}
\tag{11}
$$

240 where we also made use of equations 3 and 5.

The expected value of the complete data loglikelihood then becomes

$$
\begin{aligned}
\mathbb{E}_{p(\mathcal{Z}|\mathcal{D}, \theta_{old})}\Big[ \ln p(\mathcal{D}, \mathcal{Z}|\theta) \Big] = \sum_{i=1}^{N} \sum_{r=1}^{R} \gamma(z_i^r) \ln \Big( \pi_r \; p_{\mathrm{LogReg}}(y_i^r | \mathbf{x}_i, \mathbf{w}) \Big) \\
+ (1 - \gamma(z_i^r)) \ln \Big( (1 - \pi_r) \; p_{\mathrm{Rand}}(y_i^r | \mathbf{x}_i) \Big).
\end{aligned}
\tag{12}
$$

In the M-step of the EM algorithm we maximize this expectation with respect to the model parameters $\theta$, obtaining new parameter values $\theta^{new}$ given by

$$
\theta^{new} = \arg \max_{\theta} \mathbb{E}_{p(\mathcal{Z}|\mathcal{D}, \theta_{old})}\Big[ \ln p(\mathcal{D}, \mathcal{Z}|\theta) \Big].
\tag{13}
$$

241 The EM algorithm can then be summarized as follows:

242 **E-step** Compute the posterior distribution of the latent variables $z_i^r$ by mak-
243 ing use of equation 11.

**M-step** Estimate the new model parameters $\theta^{new} = \{\boldsymbol{\pi}^{new}, \mathbf{w}^{new}\}$ given by

$$
\mathbf{w}^{new} = \arg \max_{\mathbf{w}} \sum_{i=1}^{N} \sum_{r=1}^{R} \gamma(z_i^r) \ln p_{\mathrm{LogReg}}(y_i^r | \mathbf{x}_i, \mathbf{w})
\tag{14}
$$

$$
\widehat{\mathcal{Y}}^{new} = \arg \max_{\widehat{\mathcal{Y}}} p_{\mathrm{LogReg}}(\widehat{\mathcal{Y}} | \mathcal{X}, \mathbf{w}^{new})
\tag{15}
$$

$$
\pi_r^{new} = accuracy_r = \frac{\#\{i : y_i^r = \widehat{y}_i\}}{N_r}
\tag{16}
$$

14

where $N_r$ denotes the number of instances labeled by annotator $r$. In order to optimize equation 14 we use limited-memory BFGS (Liu and Nocedal, 1989). The first order derivate is given by

$$\nabla_{\mathbf{w}} = \sum_{i=1}^{N}\sum_{r=1}^{R}\left(\gamma(z_i^r)\sum_{k=1}^{K}\left(t_{ik}^r - p_{\text{LogReg}}(y_i = k|\mathbf{x}_i, \mathbf{w})\right)\mathbf{x}_i\mathbf{x}_i^T\right) \quad (17)$$

where $\mathbf{t}_i^r$ is a vector representation of $y_i^r$ in a 1-of-$K$ coding scheme, thus $t_{ik}^r$ would be 1 when $k$ corresponds to the label provided by the $r^{th}$ annotator and 0 otherwise.

Notice that this is very similar to the typical training of a multi-class Logistic Regression model. However, in this case, the contributions of the labels provided by each annotator to the loglikelihood are being weighted by her reliability, or in other words, by how likely it is for her to be correct. This makes our proposed approach quite easy to implement in practice.

## 5. Experiments

The proposed Multiple-Annotator Logistic Regression (MA-LR)[2] model was evaluated using both multiple-annotator data with simulated annotators and data manually labelled using AMT. The model was compared with the multi-class extension of the model proposed by Raykar et al. (2009, 2010), which is a latent ground truth model, and with two majority voting baselines:

- Soft Majority Voting (MVsoft): this corresponds to a multi-class Logistic Regression model trained with the *soft* probabilistic labels resultant from the voting process.

---

[2]Source code is available at: http://amilab.dei.uc.pt/fmpr/malr.tar.gz

15

261     • Hard Majority Voting (MVhard): this corresponds to a multi-class
262         Logistic Regression model trained with the most voted labels resultant
263         from the voting process (i.e. the most voted class for a given instance
264         gets "1" and the others get "0").

265 In all experiments the EM algorithm was initialized with majority voting.

266 *5.1. Simulated annotators*

267      With the purpose of comparing the presented approaches in different
268 classification tasks we used six popular benchmark datasets from the UCI
269 repository[3] - a collection of databases, domain theories, and data generators
270 that are used by the machine learning community for the empirical analysis
271 of machine learning algorithms. Since these datasets do not have labels from
272 multiple annotators, the latter were simulated from the ground truth using
273 two different methods. The first method, denoted "label flips", consists in
274 randomly flipping the label of an instance with a given uniform probability
275 $p(flip)$ in order to simulate an annotator with an average reliability of $(1 -$
276 $p(flip))$. The second method, referred to as "model noise", seeks simulating
277 annotators that are more consistent in their opinions, and can be summarized
278 as follows. First, a multi-class Logistic Regression model is trained on the
279 original training set. Then, the resulting weights $\mathbf{w}$ are perturbed, such that
280 the classifier consistently "fails" in a coherent fashion throughout the test set.
281 To do so, the values of $\mathbf{w}$ are standardized, and then random "noise" is drawn
282 from a Gaussian distribution with zero mean and $\sigma^2$ variance and added
283 to the weights $\mathbf{w}$. These weights are then "unstandardized" (by reversing

---

[3]http://archive.ics.uci.edu/ml/index.html

16

Table 1: Details of the UCI datasets

| Dataset | Num. Instances | Num. Features | Num. Classes |
| --- | --- | --- | --- |
| Annealing | 798 | 38 | 6 |
| Image Segmentation | 2310 | 19 | 7 |
| Ionosphere | 351 | 34 | 2 |
| Iris | 150 | 4 | 3 |
| Parkinson's | 197 | 23 | 2 |
| Wine | 178 | 13 | 3 |

<sub>284</sub> the standardization process previously used), and the modified multi-class
<sub>285</sub> Logistic Regression model is re-applied to the training set in order to simulate
<sub>286</sub> an annotator. The quality of this annotator will vary depending on the value
<sub>287</sub> of $\sigma^2$ used.

<sub>288</sub> Since in practice each annotator only labels a small subset of all the in-
<sub>289</sub> stances in the dataset, we introduce another parameter in this annotator
<sub>290</sub> simulation process: the probability $p(label)$ of an annotator labeling an in-
<sub>291</sub> stance.

<sub>292</sub> Table 1 describes the UCI datasets used in these experiments. Special care
<sub>293</sub> was taken in choosing datasets that correspond to real data and that were
<sub>294</sub> among the most popular ones in the repository and, consequently, among
<sub>295</sub> the Machine Learning community. Datasets that were overly unbalanced,
<sub>296</sub> i.e. with too many instances of some classes and very few instances of oth-
<sub>297</sub> ers, were avoided. Despite that, the selection process was random, which
<sub>298</sub> resulted in a rather heterogeneous collection of datasets: with different sizes,

17

<sup>299</sup> dimensionalities and number of classes.

<sup>300</sup>    Figures 3 and 4 show the results obtained using 5 simulated annotators
<sup>301</sup> with different reliabilities using distinct simulation methods: "label flips"
<sup>302</sup> and "model noise" respectively. Although not all the results (i.e. using both
<sup>303</sup> simulation methods on all the six datasets) are presented here, we note that
<sup>304</sup> the omitted results are similar to those shown. Hence, to avoid redundancy
<sup>305</sup> and preserve brevity, only a random subset of these are presented. All the
<sup>306</sup> experiments use 10-fold cross-validation. Due to the stochastic nature of the
<sup>307</sup> simulation process of the annotators, each experiment was repeated 30 times
<sup>308</sup> and the average results were collected. The plots on the left show the root
<sup>309</sup> mean squared error (RMSE) between the estimated annotators accuracies
<sup>310</sup> and their actual accuracies evaluated against the ground truth. The plots
<sup>311</sup> on the center and on the right show, respectively, the trainset and testset
<sup>312</sup> accuracies. Note that here, unlike in "typical" supervised learning tasks,
<sup>313</sup> trainset accuracy is quite important since it indicates how well the models
<sup>314</sup> are estimating the *unobserved* ground truth labels from the opinions of the
<sup>315</sup> multiple annotators.

<sup>316</sup>    From a general perspective on the results of figures 3 and 4 we can con-
<sup>317</sup> clude that both methods for learning from multiple annotators (MA-LR and
<sup>318</sup> Raykar) tend to outperform the majority voting baselines under most condi-
<sup>319</sup> tions. Not surprisingly, as the value of $p(label)$, and consequently the average
<sup>320</sup> number of instances labeled by each annotator, decreases, both the trainset
<sup>321</sup> and testset accuracies of all the approaches decrease or stay roughly the same.
<sup>322</sup> As expected, a higher trainset accuracy usually translates in a higher testset
<sup>323</sup> accuracy and a better approximation of the annotators accuracies (i.e. lower
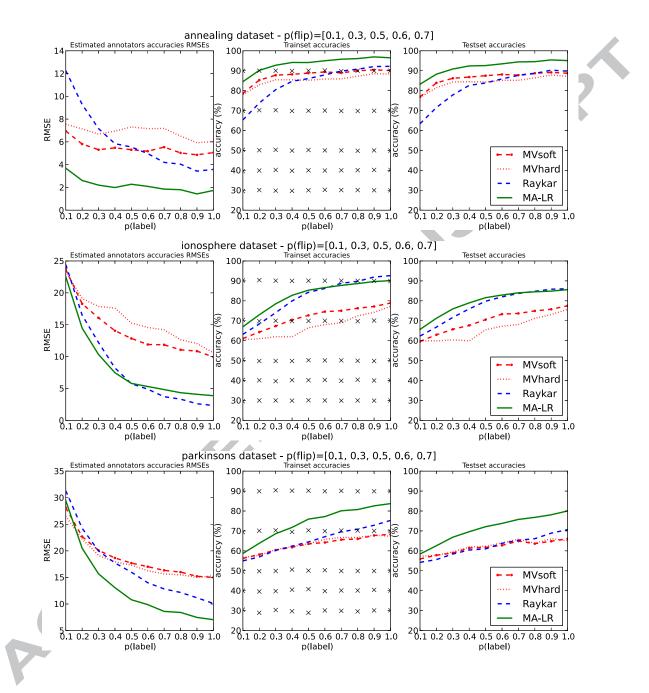
18

Figure 3: Results for the Annealing, Ionosphere and Parkinsons datasets using the "label flips" method for simulating annotators. The "x" marks indicate the average true accuracies of the simulated annotators.
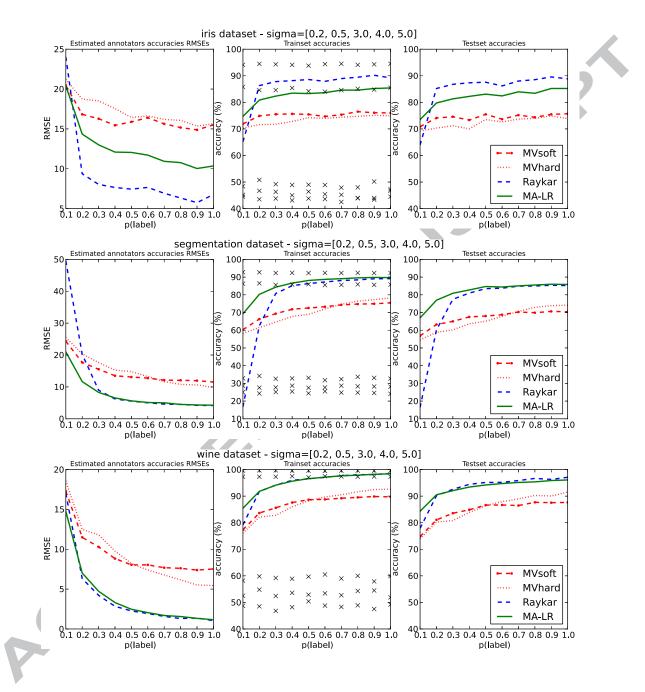
19

Figure 4: Results for the Iris, Segmentation and Wine datasets using the "model noise" method for simulating annotators. The "x" marks indicate the average true accuracies of the simulated annotators.

324 RMSE), since the approximation of the ground truth is also better.

325 A more careful analysis of the results reveals that, contrarily to the model
326 by Raykar et al. (2009, 2010), the proposed model (MA-LR) is less prone
327 to overfitting when the number of instances labeled by each annotator de-
328 creases. This is a direct consequence of the number of parameters used to
329 model the annotators expertise. While the model by Raykar et al. (2009,
330 2010) uses a $K \times K$ confusion matrix for each annotator, making a total of
331 $RK^2$ parameters, the proposed model only has $R$ parameters. However, it is
332 important to note that there is a tradeoff here, since the model by Raykar et
333 al. can capture certain biases in the annotators answers by keeping a $K \times K$
334 confusion matrix for each annotator, which is not possible with the MA-LR
335 model. Notwithstanding, in practice, on crowdsourcing platforms like AMT,
336 the number of instances labeled by each annotator is usually low. Hence, we
337 believe that the proposed model is preferable in most situations. Further-
338 more, our experimental results show that even when the number of instances
339 labeled by each annotator is high, the MA-LR model can achieve similar or
340 even better results than the model by Raykar et al. (2009, 2010).

341 *5.2. Amazon Mechanical Turk*

342 In order to assess the performance of the proposed model in learning from
343 the labels of multiple non-expert human annotators and compare it with the
344 other approaches, two experiments were conducted using AMT: sentiment
345 polarity and music genre classification[4].

346 The sentiment polarity experiment was based on the sentiment analysis

---

[4]Datasets are available at: http://amilab.dei.uc.pt/fmpr/mturk-datasets.tar.gz

21

347   dataset introduced by Pang and Lee (2005), which corresponds to a collection
348   of more than ten thousand sentences extracted from the movie review website
349   RottenTomatoes[5]. These are labeled as positive or negative depending on
350   whether they were marked as "fresh" or "rotten" respectively. From this
351   collection, a random subset of 5000 sentences were selected and published on
352   Amazon Mechanical Turk for annotation. Given the sentences, the workers
353   were asked to provide the sentiment polarity (positive or negative). The
354   remaining 5428 sentences were kept for evaluation.

355   For the music genre classification experiment, the audio dataset intro-
356   duced by Tzanetakis and Cook (2002) was used. This dataset consists of
357   a thousand samples of songs with 30 seconds of length and divided among
358   10 different music genres: classical, country, disco, hiphop, jazz, rock, blues,
359   reggae, pop and metal. Each of the genres has 100 representative samples.
360   A random 70/30 train/test split was performed on the dataset, and the 700
361   training samples were published on AMT for classification. In this case, the
362   workers were required to listen to a 30-second audio excerpt and classify it
363   as one of the 10 genres enumerated above.

364   On both experiments, the AMT workers were required to have an *HIT*
365   *approval rate* - an AMT quality indicator that reflects the percentage of
366   accepted answers of a worker - of 95%, which ensures some reliability on the
367   quality of the answers.

368   Table 2 shows some statistics about the answers of the AMT workers for
369   both datasets. Figure 5 further explore the distributions of the number of

---

[5]http://www.rottentomatoes.com/

22

Table 2: Statistics of the answers of the AMT workers for the two experiments performed. Note that the worker accuracies correspond to trainset accuracies.

|  | Sentiment polarity | Music genre |
|---|---|---|
| Number of answers collected | 27747 | 2946 |
| Number of workers | 203 | 44 |
| Avg. answers per worker ($\pm$ std) | $136.68 \pm 345.37$ | $66.93 \pm 104.41$ |
| Min. answers per worker | 5 | 2 |
| Max. answers per worker | 3993 | 368 |
| Avg. worker accuracy ($\pm$ std) | $77.12 \pm 17.10\%$ | $73.28 \pm 24.16\%$ |
| Min. worker accuracy | 20% | 6.8% |
| Max. worker accuracy | 100% | 100% |

370 answers provided by each annotator and their accuracies for the sentiment
371 polarity and music genre datasets. The figure reveals a highly skewed dis-
372 tribution of number of answers per worker, which support our intuition that
373 on this kind of crowdsourcing platforms each worker tends to only provide
374 a small number of answers, with only a couple of workers performing high
375 quantities of labelings.

376    Standard preprocessing and features extraction techniques were performed
377 on both experiments. In the case of the sentiment polarity dataset, the stop-
378 words were removed and the remaining words were reduced to their root by
379 applying a stemmer. This resulted in a vocabulary with size 8919, which still
380 makes a bag-of-words representation computationally expensive. Hence, La-
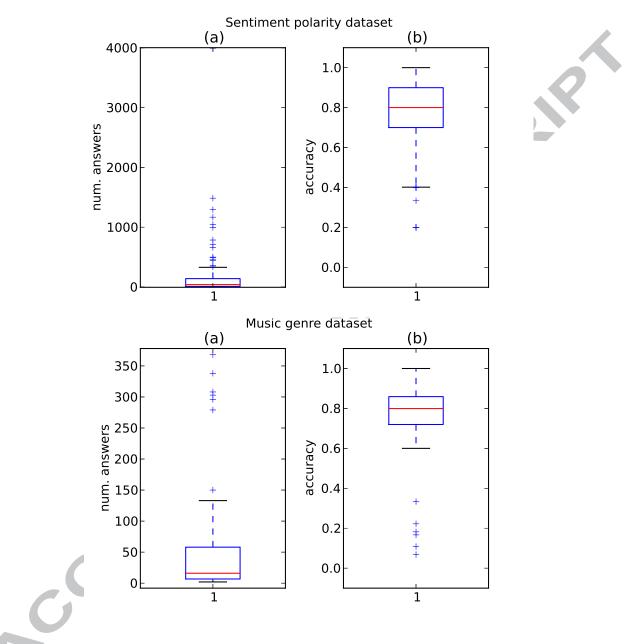381 tent Semantic Analysis (LSA) was used to further reduce the dimensionally

23

Figure 5: Boxplots for the number of answers (a) and for the accuracies (b) of the AMT workers for the sentiment polarity (top) and music genre (bottom) datasets.

Table 3: Trainset and testset accuracies for the different approaches on the datasets obtained from AMT.

|  | Sentiment polarity | | Music genre | |
|---|---|---|---|---|
| Method | Train acc. | Test acc. | Train acc. | Test acc. |
| MVsoft | 80.70% | 71.65% | 67.43% | 60.33% |
| MVhard | 79.68% | 70.27% | 67.71% | 59.00% |
| Raykar | 49.91% | 48.67% | 9.14% | 12.00% |
| Raykar (w/prior) | 84.92% | 70.78% | 71.86% | 63.00% |
| MA-LR | 85.40% | 72.40% | 72.00% | 64.00% |

382 of the dataset to 1200 features.

383     Regarding the music genre dataset, we used Marsyas[6], a standard music
384 information retrieval tool, to extract a collection of commonly used features
385 in this kind of tasks (Tzanetakis and Cook, 2002). These include means and
386 variances of timbral features, time-domain Zero-Crossings, Spectral Centroid,
387 Rolloff, Flux and Mel-Frequency Cepstral Coefficients (MFCC) over a texture
388 window of 1 second. A total of 124 features were extracted. The details on
389 these features fall out of the scope of this article. The interested reader is
390 redirected to the appropriate literature (e.g. Aucouturier and Pachet (2003);
391 Tzanetakis and Cook (2002)).

392     Table 3 presents the results obtained by the different methods on the sen-
393 timent polarity and music genre datasets. As expected, the results indicate
394 that both annotator-aware methods are clearly superior when compared to

---

[6]http://marsyasweb.appspot.com

25

³⁹⁵ the majority voting baselines. Also, notice that due to the fact some anno-
³⁹⁶ tators only label a very small portion of instances, the "standard" model by
³⁹⁷ Raykar et al. (2009, 2010) performs very poorly (as bad as a random classi-
³⁹⁸ fier) due to overfitting. In order to overcome this, a prior had to be imposed
³⁹⁹ on the probability distribution that controls the quality of the annotators.
⁴⁰⁰ In the case of the sentiment polarity task, a $Beta(1, 1)$ prior was used, and
⁴⁰¹ for the music genre task we applied a symmetric Dirichlet with parameter
⁴⁰² $\alpha = 1$. Despite the use of a prior, the model by Raykar et al. (2009, 2010)
⁴⁰³ still performs worse than the proposed MA-LR model, which takes advan-
⁴⁰⁴ tage of its single quality parameter per annotator to produce better estimates
⁴⁰⁵ of the annotators' reliabilities. These results are coherent with our findings
⁴⁰⁶ with the simulated annotators, which highlights the quality of the proposed
⁴⁰⁷ model.

## 6. Conclusions and Future Work

⁴⁰⁹ In this paper we presented a new probabilistic model for supervised multi-
⁴¹⁰ class classification from multiple annotator data. Unlike previous approaches,
⁴¹¹ in this model the reliabilities of the annotators are treated as latent variables.
⁴¹² This design choice results in a model with various attractive characteristics,
⁴¹³ such as: its easy implementation and extension to other classifiers, the nat-
⁴¹⁴ ural extension to structured prediction problems (as opposed to the com-
⁴¹⁵ monly used latent ground truth models), and the ability to overcome the
⁴¹⁶ overfitting to which more complex models of the annotators expertise are
⁴¹⁷ susceptible as the number of instances labeled by each annotator decreases.
⁴¹⁸ We empirically showed, using both simulated annotators and human-

26

labeled data from Amazon Mechanical Turk, that under most conditions, the proposed approach achieves comparable or even better results when compared to a state of the art model (Raykar et al., 2009, 2010) despite its much smaller set of parameters to model the annotators expertise. In fact, it turned out that this reduced number of parameters plays a key role in making the model less prone to overfitting.

Future work will explore the behavior of the proposed model when we relax the assumption that the reliability of the annotators does not depend on the instances that they are labeling, similarly to what is done in Yan et al. (2010). Furthermore, the generalization to sequence labeling tasks will also be investigated.

## References

Aucouturier, J., Pachet, F., 2003. Representing musical genre: A state of the art. Journal of New Music Research 32 (1), 83–93.

Bishop, C. M., 2006. Pattern Recognition and Machine Learning. Springer.

Dawid, A. P., Skene, A. M., 1979. Maximum likelihood estimation of observer error-rates using the EM algorithm. Journal of the Royal Statistical Society. Series C 28 (1), 20–28.

Dempster, A. P., Laird, N. M., Rubin, D. B., 1977. Maximum likelihood from incomplete data via the em algorithm. Journal of the Royal Statistical Society, Series B 39 (1), 1–38.

Donmez, P., Schneider, J., Carbonell, J., 2010. A probabilistic framework to

441 learn from multiple annotators with time-varying accuracy. In: Proc. of

442 the SIAM Int. Conf. on Data Mining. pp. 826–837.

443 Groot, P., Birlutiu, A., Heskes, T., 2011. Learning from multiple annotators

444 with gaussian processes. In: Proc. of the 21st Int. Conf. on Artificial Neural

445 Networks. Vol. 6792. pp. 159–164.

446 Howe, J., 2008. Crowdsourcing: Why the Power of the Crowd is Driving the

447 Future of Business, 1st Edition. Crown Publishing Group, New York, NY,

448 USA.

449 Jacobs, R., Jordan, M. I., Nowlan, S., Hinton, G., 1991. Adaptive mixtures

450 of local experts. Neural Computation, 79–87.

451 Liu, D. C., Nocedal, J., 1989. On the limited memory BFGS method for large

452 scale optimization. Mathematical Programming 45, 503–528.

453 Novotney, S., Callison-Burch, C., 2010. Cheap, fast and good enough: Au-

454 tomatic speech recognition with non-expert transcription. In: Proc. of the

455 Human Language Technologies. Association for Computational Linguis-

456 tics, pp. 207–215.

457 Pang, B., Lee, L., 2005. Seeing stars: Exploiting class relationships for sen-

458 timent categorization with respect to rating scales. In: Proc. of the ACL.

459 pp. 115–124.

460 Raykar, V., Yu, S., Zhao, L., Jerebko, A., Florin, C., Valadez, G., Bogoni,

461 L., Moy, L., 2009. Supervised learning from multiple experts: whom to

28

462    trust when everyone lies a bit. In: Proc. of the 26th Int. Conf. on Machine
463    Learning. pp. 889–896.

464    Raykar, V., Yu, S., Zhao, L., Valadez, G., Florin, C., Bogoni, L., Moy, L.,
465    2010. Learning from crowds. Journal of Machine Learning Research, 1297–
466    1322.

467    Smyth, P., Fayyad, U., Burl, M., Perona, P., Baldi, P., 1995. Inferring ground
468    truth from subjective labelling of venus images. In: Advances in Neural
469    Information Processing Systems. pp. 1085–1092.

470    Snow, R., O'Connor, B., Jurafsky, D., Ng, A., 2008. Cheap and fast - but is it
471    good?: evaluating non-expert annotations for natural language tasks. In:
472    Proc. of the Conf. on Empirical Methods in Natural Language Processing.
473    pp. 254–263.

474    Surowiecki, J., 2004. The Wisdom of Crowds : Why the many are smarter
475    than the few and how collective wisdom shapes business, economies, soci-
476    eties, and nations. Doubleday.

477    Sutton, C., 2012. An introduction to conditional random fields. Foundations
478    and Trends® in Machine Learning 4 (4), 267–373.

479    Tzanetakis, G., Cook, P., 2002. Musical genre classification of audio signals.
480    Speech and Audio Processing, IEEE Transactions on 10 (5), 293–302.

481    Wu, O., Hu, W., Gao, J., 2011. Learning to rank under multiple annotators.
482    In: Proc. of the 22nd Int. Joint Conf. on Artificial Intelligence. pp. 1571–
483    1576.

29

Yan, Y., Rosales, R., Fung, G., Dy, J., 2011. Active learning from crowds. In: Proc. of the 28th Int. Conf. on Machine Learning. pp. 1161–1168.

Yan, Y., Rosales, R., Fung, G., Schmidt, M., Valadez, G., Bogoni, L., Moy, L., Dy, J., 2010. Modeling annotator expertise: Learning when everybody knows a bit of something. Journal of Machine Learning Research 9, 932–939.

# Highlights:

We propose a new probabilistic model for learning with multiple annotators.

The reliability of the different annotators is treated as a latent variable.

Model is able to achieve state of the art performance (or superior).

Reduced number of model parameters is able to avoid overfitting.

Model is easier to implement and extend to other classes of learning problems.