

# Semi-Supervised Learning Applied to Text Classification

António Rodrigues

Faculdade de Engenharia da Universidade do Porto  
Rua Dr. Roberto Frias S/N, 4200-465 Porto, Portugal

up200400437@fe.up.pt

## Abstract

*The work presented in this paper explores the field of semi-supervised learning, applied to a particular problem: text classification. We explore a generative model for text classification — Multinomial Naive Bayes (MNB) — and from there we advance to a semi-supervised setting by combining it with the Expected Maximization (EM) algorithm, studying some techniques to augment it. We finally present results from our own and third-party implementations of such models, by using the well-known 20 News-groups dataset.*

## 1. Introduction

Semi-supervised learning (SSL) is a branch of machine learning that makes use of unlabeled data in an attempt to improve the performance of purely supervised learning methods in cases where labeled data is ‘hard to get’ and scarce, when compared to unlabeled data [1, 9, 10]. Such scenarios are rather frequent, e.g. in application areas such as text classification, image recognition, web content analysis. Here we focus on the study of specific semi-supervised classification tasks, even though other sub-fields such as semi-supervised regression [1] exist.

The work presented in this paper explores the field of semi-supervised learning, applied to a particular problem: text classification. We start by exploring a generative model for text classification — Multinomial Naive Bayes (MNB) [5] — still in its fully-supervised form. Based on that model, we then advance to a semi-supervised setting by combining it with the Expected Maximization (EM) [6] algorithm, studying some techniques to augment it. We finally present results from our own and third-party implementations of such models, over the well-known 20 News-groups dataset [3].

## 2. Semi-Supervised Learning

In purely Supervised Learning problems, one is given a labeled set of observations  $\mathbf{X}_\ell = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{|\mathbf{X}_\ell|}\}$ <sup>1</sup> together with their class correspondences  $\mathbf{Y} = \{y_1, y_2, \dots, y_{|\mathbf{Y}|}\}$ , taken from the joint distribution  $P(\mathbf{X}, \mathbf{Y})$ . The goal is to find a function that maps any given observation to a class,  $\mathbf{x} \rightarrow y_i \in \mathbf{Y}$ , such that the classification error is minimized.

Considering an SSL setting, the classification problem is similar — a labeled set  $\mathbf{X}_\ell$  is also given, the goal is to estimate a relation  $\mathbf{x} \rightarrow y_i \in \mathbf{Y}$  — but now with the addition of an unlabeled set  $\mathbf{X}_u = \{\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_{|\mathbf{X}_u|}\}$  (i.e. without class label correspondences), where usually  $|\mathbf{X}_u| \gg |\mathbf{X}_\ell|$ , with samples now taken from the marginal distribution  $P(\mathbf{X})$ . Starting with assumptions about the distribution of the data  $P(\mathbf{X})$  (e.g. treating it as a mixture distribution), SSL methods tackle the problem by supplying  $\mathbf{X}_u$  to unsupervised learning techniques, in order to reach good estimates of ‘hidden’ (also referred to as ‘latent’) class-conditional probabilities  $P(\mathbf{X}|\mathbf{Z})$ . The difference is that now the number and type of cases for  $\mathbf{Z}$  are pre-determined by the class correspondences of the labeled dataset,  $\mathbf{Y}$ . The assignment of each ‘latent’ component of the distribution to elements in  $\mathbf{Y}$  can then be guided with the help of the labeled dataset,  $\mathbf{X}_\ell$ .

The success of SSL techniques is heavily dependent on the model’s assumptions about the distribution of  $P(\mathbf{X})$ , with severe negative impact on performance in the case of inappropriate model-to-problem matchings [10, 9]. Chapelle et al. [1] list the key types of assumptions often made by SSL methods:

- **Smoothness:** If two points  $x_1$  and  $x_2$  are close to each other in a high density region, then so should be the respective labels  $y_1$  and  $y_2$ , i.e. the labeling function is assumed to be *smoother* in high density regions.
- **Clustering or low density separation:** Points belonging to the same cluster are likely to belong to the same

<sup>1</sup>We will henceforth use the notation  $|S|$  to represent the size of a given set  $S$ .

class. Following a somewhat complementing idea to ‘smoothness’, decision boundaries should lie in low-density regions.

- **Manifold structures** [1]
- **Transduction** [1]

Several classes of SSL methods exist, each considering specialized cases of one (or more) of the aforementioned assumptions: Self-Training, Co-Training [9], Generative Models, Low-Density Separation methods, Graph-Based Methods [1], among others [9].

### 2.1. Semi-Supervised Generative Models

In their way to estimate  $P(Y|X)$ , Generative Models (GMs) first find a  $y_i \rightarrow \mathbf{x}$  mapping by modeling the class-conditional distributions  $P(X|Y)$ , instead of directly going for an estimation of the posteriors. In other words, GMs go through the ‘trouble’ of estimating how data is generated by each element in  $Y$  in order to classify an observation.

GMs treat  $P(X|Y)$  as a mixture model, where components belong to a family of parametric distributions (e.g. Gaussian, Multinomial, etc.), governed by a set of parameters  $\theta = \{\theta_1, \theta_2, \dots, \theta_n\}$ . Unsupervised learning methods are applied over the unlabeled data  $\mathbf{X}_u$  in order to find estimates for those parameters,  $\hat{\theta}$ <sup>2</sup>, and identify the components of the mixture model, usually with a 1:1 correspondence between the ‘latent’ clusters and each class in  $Y$  (Nigam et al. [6] propose both 1:1 and N:1 correspondence alternatives, Zhu [9] also points out that such assumptions should be carefully considered in order not to incur in model incorrectness). The posterior probabilities  $P(Y = y_i | X = \mathbf{x}, \hat{\theta})$  can then be determined via Bayes’ Rule:

$$P(y_i | \mathbf{x}, \hat{\theta}) = \frac{P(y_i | \hat{\theta}) P(\mathbf{x} | y_i, \hat{\theta})}{\sum_{i=1}^{|Y|} P(y_i | \hat{\theta}) P(\mathbf{x} | y_i, \hat{\theta})}$$

The training of a classifier consists in determining the parameter estimates that maximize the likelihood (Maximum Likelihood estimation) of both labeled and unlabeled data:

$$\prod_{\mathbf{x}_i \in \mathbf{X}_\ell} P(y_i) P(\mathbf{x}_i | y_i, \theta) \times \prod_{\mathbf{x}_i \in \mathbf{X}_u} \sum_{k=1}^{|Y|} P(y_k) P(\mathbf{x}_i | y_k, \theta) \quad (1)$$

or via Maximum a Posteriori (MAP) estimation of  $P(\theta | \mathbf{X}_\ell, \mathbf{X}_u)$ . These estimates are usually tackled via gradient descent methods [1] or using Expectation Maximization (EM) algorithms [6], as with the case described in Section 3.4.

<sup>2</sup>We represent estimations of a variable  $x$  by topping it with an ‘^’ sign,  $\hat{x}$ .

## 3. Methodology

The work presented here follows the consulted background literature about supervised and semi-supervised text classification [5, 6, 7], reflecting our understanding of the material.

### 3.1. Notation

Text classification mainly deals with categorizing a set of text articles into some topic, according to a set of features such as the words contained in them. Here we present the notation to formalize these concepts.

An article  $a$ , belonging to a topic (i.e. a label/class)  $t$ , is represented by an array  $a = \{w_1, w_2, \dots, w_{|\mathcal{D}|}, t\}$ , i.e. a list of  $|\mathcal{D}|$  features and its label. In this case, each feature typically corresponds to the number of occurrences of the word  $w_i$ , the  $i$ -th word in a dictionary  $\mathcal{D}$ , within an article  $a_i$ . Other options for representing word frequency exist, such as Term Frequency and Inverse Document Frequency (TF-IDF) [7].

As an example, consider an article set composed by two instances,  $\mathcal{A} = \{a_1, a_2\}$ , each one of them with the following contents:

$a_1$ : You like potato, I like potato.

$a_2$ : I say tomato, you say tomato.

In this case, the dictionary  $\mathcal{D}$  would contain the following words<sup>3</sup>:

$\mathcal{D} = \{\text{You, like, potato, I, say, tomato}\}$

The articles  $a_1$  and  $a_2$  could then be codified (using the sparse ‘bag-of-words’ approach [8]) as:

$a_1 = \{1, 2, 2, 1, 0, 0\}$

$a_2 = \{1, 0, 0, 1, 2, 2\}$

### 3.2. Generative Models for Text Classification

McCallum et al. [5] present two generative models for text classification, (1) a Multivariate Bernoulli event model and (2) a Multinomial event model: the first case considers multivariate Bernoulli as the parametric distributions describing each mixture component, only capturing the (non-)occurrence of word events in articles, while the second case considers Multinomial distributions, now capturing the quantity of word events. The authors state that the multinomial event model generally outperforms the multivariate

<sup>3</sup>As we will see in Section 5.1, in practice, the words in  $\mathcal{D}$  might be converted to a common format, e.g. all words converted to lowercase, stripped of punctuation, etc.

Bernoulli model, specially when considering large dictionary sizes [5]. For the remainder of this paper, we focus our attention on the Multinomial event model.

Despite our focus on the Multinomial model, being GMs, both models assume (1) that an article  $a$  is generated according to a mixture model, encompassing several mixture components  $c_j \in \mathcal{C} = \{c_1, c_2, \dots, c_{|\mathcal{C}|}\}$ . The shape of each component distribution is governed by a set of parameters  $\theta$ . It is also assumed (2) there is a 1:1 correspondence between the components in  $\mathcal{C}$  and topics of articles  $t_i$ .

One can look at the process of ‘generating’ an article  $a_i$  in the following manner: (1) Selecting a component  $c_j$  from the mixture model, with probability  $P(c_j|\theta)$ ; (2) letting  $c_j$  generate  $a_i$  according to its own distribution  $P(a_i|c_j, \theta)$ . This results in the probability of an article  $a_i$  being generated by a component  $c_j$ :

$$P(c_j|\theta)P(a_i|c_j, \theta) \quad (2)$$

As different components in  $\mathcal{C}$  can contribute to origin a similar article  $a_i$ , the probability of finding an article  $a_i$  is obtained by marginalizing expression 2 over all the components in  $\mathcal{C}$ :

$$P(a_i|\theta) = \sum_{j=1}^{|\mathcal{C}|} P(c_j|\theta)P(a_i|c_j, \theta) \quad (3)$$

### 3.3. Multinomial Naive Bayes

The Multinomial Naive Bayes (MNB) model described by McCallum et al. in [5] captures the frequency of words in articles and applies a ‘bag-of-words’ approach for article representation.

The second step in the process of generating an article  $a_i$  — related to the  $P(a_i|c_j, \theta)$  term in expression 2 — can be further expressed as a sequence of  $|a_i|$  draws (with replacement) of words  $w_k$  from the dictionary  $\mathcal{D}$ . Besides the assumptions mentioned in Section 3.2, the MNB model considers that (1) the length of an article  $|a_i|$  (word count) is independent of the topic/mixture model component  $c_j$ ; and (2) the appearances of words in an article are conditionally independent from each other, given an article topic: the so-called ‘Naive Bayes’ assumption. The class conditional probability of an article  $a_i$  can then be thought of as a multinomial distribution over words, with  $a_i$  independent trials, in the form:

$$P(a_i|c_j, \theta) = |a_i|! \prod_{k=1}^{|\mathcal{D}|} \frac{P(w_k|c_j, \theta)^{N_{i,k}}}{N_{i,k}!} \propto \prod_{k=1}^{|\mathcal{D}|} P(w_k|c_j, \theta)^{N_{i,k}} \quad (4)$$

where  $N_{i,k}$  is the number of times word  $w_k$  of a dictionary  $\mathcal{D}$  appears in an article  $a_i$ . In practice, as the multinomial coefficient  $\frac{|a_i|!}{N_{i,1}! \dots N_{i,|\mathcal{D}|}!}$  does not depend on the mixture components  $c_j$ , it is often ignored when the purpose is to maximize the likelihood  $P(a_i|c_j, \theta)$  [6, 2, 8].

The set of mixture model parameters  $\theta$  to be estimated during the training phase of the classifier consists of (1) each of the class conditional probabilities of words  $\hat{\theta}_{w_k|c_j} \equiv P(w_k|c_j, \hat{\theta})$ ; and (2) the prior probabilities for each topic/mixture model component  $\hat{\theta}_{c_j} \equiv P(c_j|\hat{\theta})$ .

These estimates are obtained by Maximum a Posteriori (MAP) estimation of the parameters, i.e. we find the values of  $\theta$  that maximize the posterior probability  $P(\theta|\mathcal{A})$ :

$$P(\theta|\mathcal{A}) \propto P(\theta) \times P(\mathcal{A}|\theta) \quad (5)$$

Here we follow the same representation of  $P(\theta)$  used by Nigam et al. [6], a Dirichlet distribution, with  $\sigma = 2$ . The probability  $P(\mathcal{A}|\theta)$  is equal to  $\prod_{i=1}^{|\mathcal{A}|} P(a_i|\theta)$ , as one assumes the article generation events are independent between each other, given the mixture model parameters  $\theta$ .

The results of the MAP estimation reduce to ‘counting problems’. Specifically,  $\hat{\theta}_{w_k|c_j}$  is given by the ratio of appearances of a word  $w_k$  within all articles  $a_i$  belonging to a component  $c_j$  vs. the total number of word events for  $c_j$ :

$$\hat{\theta}_{w_k|c_j} \equiv P(w_k|c_j, \hat{\theta}) = \frac{\alpha + \sum_{i=1}^{|\mathcal{A}|} N_{i,k} P(t_i = c_j|a_i)}{\alpha|\mathcal{D}| + \sum_{s=1}^{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{A}|} N_{i,s} P(t_i = c_j|a_i)} \quad (6)$$

where  $P(c_j|a_i) \in \{0, 1\}$  depending on the label of an article  $a_i$  and  $N_{i,k}$  is the number of occurrences of the word  $w_k$  in an article  $a_i$ . Notice the inclusion of ‘smoothing priors’  $\alpha$ , used to avoid probabilities equal to zero in the lack of particular word events for a component  $c_j$ . In [5, 6] the authors use  $\alpha = 1$ , which is designated by Laplace smoothing.

The parameters  $\hat{\theta}_{c_j}$  are given by the ratio of the articles belonging to a component  $c_j$  vs. the total number of articles  $\mathcal{A}$ :

$$\hat{\theta}_{c_j} \equiv P(c_j|\hat{\theta}) = \frac{\alpha + \sum_{i=1}^{|\mathcal{A}|} P(t_i = c_j|a_i)}{\alpha|\mathcal{C}| + \sum_{i=1}^{|\mathcal{A}|} P(t_i = c_j|a_i)} \quad (7)$$

The final expression for the posterior probabilities  $P(c_j|a_i, \theta)$ , i.e. the probability of the class given an article, is obtained via Bayes’ Rule:

$$P(c_j|a_i, \hat{\theta}) = \frac{P(c_j|\hat{\theta})P(a_i|c_j, \hat{\theta})}{P(a_i|\hat{\theta})} \propto P(c_j|\hat{\theta})P(a_i|c_j, \hat{\theta}) \quad (8)$$

Expression 8 can be progressively expanded by replacing each one of the terms by the corresponding expressions, given in equations 3, 4, 6 and 7.

Despite its unrealistic simplifications and assumptions (data generated by mixture model; 1:1 correspondence between mixture components and topics; conditional independence of word events; article length independence), Naive Bayes classifiers have proven to provide fair classification performance [5, 6]. Due to their simplicity, these are well suited for text classification tasks, where the number of features is usually large (i.e. dictionary sizes often reaching orders of thousands of words [5]).

### 3.4. Semi-Supervised Learning via Expectation Maximization (EM)

The MNB text classifier described in Section 3.3 falls into the scope of supervised learning, only taking labeled data into account. Despite its fair performance when trained with large amounts of labeled data, Nigam et al. [6] notice how it suffers when faced with small-sized datasets and show some advantages (regarding classification accuracy) of expanding such models to the semi-supervised learning scope, i.e. considering both labeled and unlabeled data.

As with the case of the MNB classifier, the parameters  $\theta$  are obtained by MAP estimation, i.e. maximizing expression 5. However, we should note that now the training dataset is composed by the labeled and unlabeled subsets  $\mathcal{A}^\ell$  and  $\mathcal{A}^u$ . Following the Semi-Supervised GM ideas introduced in Section 3.2, the expression for  $P(\mathcal{A}|\theta)$  becomes:

$$P(\mathcal{A}|\theta) = \prod_{a_i \in \mathcal{A}^u} \sum_{j=1}^{|\mathcal{C}|} P(c_j|\theta) P(a_i|c_j, \theta) \times \prod_{a_i \in \mathcal{A}^\ell} P(t_i = c_j|\theta) P(a_i|t_i = c_j, \theta) \quad (9)$$

Note that for the set of labeled articles,  $\mathcal{A}^\ell$ , one already knows the true topic/mixture model component  $c_j$  which generated each article  $a_i$ , hence there is not the need of referring to all components in  $\mathcal{C}$ . However, for the unlabeled set  $\mathcal{A}^u$ , each component  $c_j$  has a contribution to the generation of each article  $a_i$  which must be taken into account.

As described in [6] (and as with the case of the MNB classifier), expression 9 can be passed to logarithmic form, with the maximization of  $P(\theta)P(\mathcal{A}|\theta)$  being accomplished by solving the system of partial derivatives of  $\log(P(\theta)P(\mathcal{A}|\theta))$ . Here we use the same nomenclature as

that used in [6],  $\ell(\theta|\mathcal{A}) \equiv \log(P(\theta)P(\mathcal{A}|\theta))$ :

$$\begin{aligned} \ell(\theta|\mathcal{A}) &= \log(P(\theta)) \\ &+ \sum_{a_i \in \mathcal{A}^u} \log \left( \sum_{j=1}^{|\mathcal{C}|} P(c_j|\theta) P(a_i|c_j, \theta) \right) \\ &+ \sum_{a_i \in \mathcal{A}^\ell} \log(P(t_i = c_j|\theta) P(a_i|t_i = c_j, \theta)) \end{aligned} \quad (10)$$

The log of sums over all mixture components  $\mathcal{C}$  for all  $a_i \in \mathcal{A}^u$  makes the problem computationally intractable [6]. This is where the EM algorithm comes into play, providing an iterative process to obtain a MAP estimation of the parameters  $\theta$ , including the unlabeled data [6]. We now proceed with a brief description of the overall MNB + EM process, without detailing the inner works of EM (refer to [6] for further details):

1. Train a MNB classifier with the labeled dataset  $\mathcal{A}^\ell$  only. Find the estimated parameters  $\hat{\theta}$  using expressions 6 and 7.
2. **EM's E Step:** Use the classifier governed by the current  $\hat{\theta}$  parameters to estimate the contribution of each mixture model component  $c_j$  to the generation of each article in the unlabeled dataset  $\mathcal{A}^u$ , i.e. determine  $P(c_j|a_i, \hat{\theta}) \forall a_i \in \mathcal{A}^u$  using expression 8.
3. **EM's M Step:** Re-estimate the  $\hat{\theta}$  parameters at the light of the new  $P(c_j|a_i, \hat{\theta})$  for both  $\mathcal{A}^\ell$  and  $\mathcal{A}^u$ , using expressions 6 and 7. Note that now  $P(c_j|a_i, \hat{\theta})$  varies between 0 and 1 for  $\mathcal{A}^u$  (as opposed to  $P(c_j|a_i, \hat{\theta}) \in \{0, 1\} \forall a_i \in \mathcal{A}^\ell$ ).
4. Evaluate  $\ell(\theta|\mathcal{A})$  using expression 10. If  $\Delta\ell(\theta|\mathcal{A}) > T$ ,  $T$  being a convergence threshold, return to step 2. Else, accept the classifier governed by the current  $\hat{\theta}$  as the final solution.

Nigam et al. [6] note that while this simple MNB + EM combination performs well over datasets containing small amounts of labeled data vs. large amounts of unlabeled data (e.g. with differences within the range of 1000 vs. 10000 [6]), in the case of the 20 Newsgroups dataset [3]), it may decrease the classification accuracy of an MNB classifier in the presence of large labeled datasets. As with other semi-supervised learning problems [10], these reductions in performance are due to violations of the model assumptions, previously stated in Sections 3.2 and 3.3.

#### 3.4.1 Extensions to the EM Algorithm: EM- $\lambda$

Nigam et al. [6] propose two extensions to the EM approach described in Section 3.4 which attempt to cope with violations of some MNB model assumptions: the EM-Multiple

and EM- $\lambda$  techniques. In the first case, the idea is to tackle violations of the 1:1 mixture component-to-topic correspondence assumption, allowing N:1 correspondences. The idea is that some topics may be separated into sub-topics (e.g. ‘football’ and ‘cricket’ in ‘sports’), with co-occurrences of words that may be better captured by multiple multinomial distributions. The unsupervised learning component of the problem is now detached from a particular number of ‘soft clusters’, which may be now determined via cross-validation [6].

We describe the second case – EM- $\lambda$  – with more detail. Notice that when  $|\mathcal{A}^u| \gg |\mathcal{A}^\ell|$ , the influence of  $\mathcal{A}^\ell$  in the maximization of expression 10 is negligible, i.e. EM will be essentially be performing unsupervised clustering [6]. The role of  $\mathcal{A}^\ell$  would then be limited to provide the initial parameter estimates  $\hat{\theta}$  and provide the number and topic correspondences for the ‘latent’ variables of the mixture model. This may result in poor classification accuracy if the distribution of the data does not follow the GM’s assumptions.

One solution is to reduce the influence of the unlabeled data in expression 10 by a factor  $\lambda$  with  $0 \leq \lambda \leq 1$ . The difference between EM- $\lambda$  and the method shown in Section 3.4 is in the M-Step, with equations 6 and 7 altered to include the  $\lambda$  factors:

$$\begin{aligned} \hat{\theta}_{w_k|c_j} &\equiv P(w_k|c_j, \hat{\theta}) \\ &= \frac{\alpha + \sum_{i=1}^{|\mathcal{A}|} \Lambda(i) N_{i,k} P(t_i = c_j|a_i)}{\alpha|\mathcal{D}| + \sum_{s=1}^{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{A}|} N_{i,s} P(t_i = c_j|a_i)} \end{aligned} \quad (11)$$

$$\hat{\theta}_{c_j} \equiv P(c_j|\hat{\theta}) = \frac{\alpha + \sum_{i=1}^{|\mathcal{A}|} \Lambda(i) P(t_i = c_j|a_i)}{\alpha|\mathcal{C}| + |\mathcal{A}^\ell| + \lambda|\mathcal{A}^u|} \quad (12)$$

where the function  $\Lambda(i)$  defines the  $\lambda$  weighing factor to apply, whether the an article  $a_i$  belongs to  $\mathcal{A}^\ell$  or  $\mathcal{A}^u$ :

$$\Lambda(i) = \begin{cases} \lambda & \text{if } a_i \in \mathcal{A}^u \\ 1 & \text{if } a_i \in \mathcal{A}^\ell \end{cases} \quad (13)$$

If we consider the extreme values allowed for  $\lambda$ , when  $\lambda = 0$  the influence of the unlabeled data is mitigated, resulting in a MNB classifier; with  $\lambda = 1$  the contribution of the unlabeled dataset is maximum, corresponding to the basic MNB + EM approach described in Section 3.4. In [6], the value of  $\lambda$  is chosen as that which maximizes the leave-one-out cross-validation classification accuracy.

## 4. Implementation

This section provides noteworthy details on our implementation of the methods introduced in Section 3.

### 4.1. Multinomial Naive Bayes

The computation of the class conditional probabilities  $P(a_i|c_j, \theta)$  via expression 4 involves the multiplication of  $|\mathcal{D}|$  probabilities  $P(w_k|c_j, \theta)$ . Considering that text classification tasks work with dictionaries composed by thousands of words (e.g.  $\sim 60000$  for the 20 Newsgroups dataset [3]), the multiplication result may end up at zero due to floating point precision underflow. To overcome this issue, the same computations can be performed in the logarithmic domain, with expression 4 becoming:

$$\log(P(a_i|c_j, \theta)) \propto \sum_{k=1}^{|\mathcal{D}|} N_{i,k} \times \log(P(w_k|c_j, \theta)) \quad (14)$$

### 4.2. Expectation Maximization (EM)

As we get the values of  $P(c_j|a_i, \theta)$  in logarithmic form, the floating point underflow problem may also arise when computing the ‘log-of-sums’ component in expression 10. In order to circumvent it, we apply the ‘Log-Sum-Exp’ (LSE) trick, i.e. considering:

$$\begin{aligned} \log \sum_{j=1}^{|\mathcal{C}|} P(c_j|a_i, \theta) &= \log \sum_{j=1}^{|\mathcal{C}|} e^{\log(P(c_j|a_i, \theta))} \\ &= m + \log \sum_{j=1}^{|\mathcal{C}|} e^{\log(P(c_j|a_i, \theta)) - m} \end{aligned} \quad (15)$$

where  $m$  is the maximum value of  $\log(P(c_j|a_i, \theta))$ , for each  $a_i$ . In our case we also set  $P(c_j|a_i, \theta) = 0$  when  $\log(P(c_j|a_i, \theta)) - m < p$ , i.e. we discard the posteriors which, even after LSE, are still smaller than a threshold  $e^p$ , e.g. and therefore considered too small to impact the final result. We also use LSE after EM’s E step, before applying expression 6 over  $\mathcal{A}^u$ .

We only run each EM’s M step (see Section 3), over the unlabeled data  $\mathcal{A}^u$ , since the  $\hat{\theta}$  values for  $\mathcal{A}^\ell$  are previously calculated in step 1 and do not change over the iterations (the respective  $P(c_j|a_i, \hat{\theta})$  values do not change during the E step, as these are given).

## 5. Experiments

This section describes the experiments we conducted to evaluate our and third party implementation of the Multinomial Naive Bayes models and respective semi-supervised extensions described in Section 3.

### 5.1. Experimental Setup

We start with a brief description of the dataset used in our experiments, followed by the plan of experiments whose results are shown in Section 5.2.



Partition	# Articles	# Unique Words
Training	11260	53485
Test	7502	60698

Table 1. Characteristics of version of 20news-bydate dataset used in experiments.

### 5.1.1 Dataset

We work with the 20 Newsgroups dataset [3], which consists in a collection of approx. 19000 web forum posts, (almost) evenly separated across 20 different topics<sup>4</sup>. Although separated from each other, some of the 20 classes are related to each other as sub-topics of a larger category (e.g. `rec.sport.baseball` and `rec.sport.baseball` as sub-topics of `rec.sport.*`). The posts were collected over a period of several months in 1993 [6].

For our experiments, we consider a modified version of the 20news-bydate dataset<sup>5</sup>. The 20news-bydate provides a dataset sorted by date, with separate training and test subsets (test subset composed by later posts), comprising a total of 18846 newsgroups articles. Using the Rainbow tools<sup>6</sup>, developed by the authors of [5, 6], we have modified the dataset to (1) remove newsgroups headers (a common practice since the article headers include the name of the newsgroup they belong to, which would make classification trivial), (2) remove 524 common ‘stop words’ (e.g. ‘the’, ‘of’) and (3) remove words which occur only 1 time. These feature selection approaches are made in previous works on text classification [5, 6, 2, 8] and are mainly used as attempts to reduce the data dimensionality (i.e. size of the dictionary,  $|\mathcal{D}|$ ). We have noticed that (2) was necessary to achieve similar MNB accuracies to those reported in [5, 6]. The characteristics of the modified 20news-bydate are summarized in Table 1.

### 5.1.2 Experimental Protocol

The protocol is structured to test the different methods presented in Section 3, in the same order.

#### A) Multinomial Naive Bayes

We test our implementation of the MNB classifier, together with Rainbow’s, evaluating accuracy and confusion between actual and predicted classifications under different conditions. We evaluate total and per-class accuracy  $Acc$  in %, according to expression 16:

$$Acc = \frac{\text{\# correctly classified articles}}{\text{total \# classified articles}} \quad (16)$$

The MNB classifiers are tested over the test set (see Table 1), for different sample sizes  $|\mathcal{A}^\ell| = \{20, 50, 100, 500, 1000, 5000, 7520\}$  from the training set. For each case, we randomly sample  $|\mathcal{A}^\ell|/20$  articles from each of the 20 topics and show the average results over 20 runs of this procedure.

#### B) Multinomial Naive Bayes + EM

We test our implementation of the MNB + EM combination, together with Rainbow’s, again evaluating accuracy and confusion between actual and predicted classifications under different conditions. The selection of  $\mathcal{A}^\ell$  and  $\mathcal{A}^u$  sets is done as follows:

1. We first select a  $\mathcal{A}^\ell$  sample of size  $|\mathcal{A}^\ell| = \{20, 40, 100, 300, 500, 800, 1000, 1260\}$  from the training subset, again  $|\mathcal{A}^\ell|/20$  articles from each of the 20 topics.
2. We use the remaining  $11260 - |\mathcal{A}^\ell|$  as unlabeled data  $\mathcal{A}^u$ .

The values of  $|\mathcal{A}^\ell|$  are chosen so that  $|\mathcal{A}^u| \geq 10000$ , for comparison with the results from [6]. Again, we show the average results over 20 runs of this procedure. As a baseline for direct comparison we use the results for the fully supervised MNB classifier, trained with each different  $\mathcal{A}^\ell$  set.

#### C) Multinomial Naive Bayes + EM- $\lambda$

We test our implementation of the MNB + EM- $\lambda$  combination, in a similar way to that of the basic EM case. Here we additionally test several values of  $\lambda = \{0.01, 0.02, 0.1, 0.2, 0.5, 1.0\}$ , evaluating accuracy and confusion between actual and predicted classifications.

## 5.2. Experimental Results

We present the results for the experiments described in the previous section. Due to problems with our implementation of EM and EM- $\lambda$ <sup>7</sup>, we only compare our implementation of MNB in MATLAB to that of the Rainbow framework. Nevertheless, we provide the results obtained with Rainbow, using the modified version of the 20news-bydate dataset.

In Figure 1 we present the results of the supervised MNB classifier. The accuracy values for the Rainbow’s version closely follow those reported in [6]. For low values of  $|\mathcal{A}^\ell|$ , our implementation of MNB (which follows the methodology provided in Section 3.3) provides worse results than Rainbow’s, with closer values for  $|\mathcal{A}^\ell| > 5000$ .

<sup>4</sup>A list of topics is available in <http://qwone.com/~jason/20Newsgroups/>.

<sup>5</sup>Available in <http://qwone.com/~jason/20Newsgroups/>.

<sup>6</sup>Available in <http://www.cs.cmu.edu/~mccallum/bow/rainbow/>. Due to compiling issues with the original versions of Rainbow, we have used the following patched version which compiles with gcc version 4.6.3: <https://github.com/brendano/bow/>.

<sup>7</sup>Errors in the M-step of the EM algorithm. The code is available in <http://paginas.fe.up.pt/~up200400437/mlproject.zip> for reference

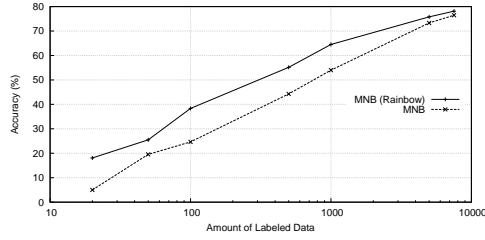


Figure 1. Results for test A: MNB, using Rainbow’s framework and our implementation in MATLAB.

The results for test B, shown in Figure 2, partly reproduce the results found in [6], since for the lower  $|A^\ell|$  the average results of MNB + EM surpass those for MNB. For  $|A^\ell| > 40$ , MNB provides better accuracies, nevertheless one should note we use a larger testing set (7502 vs. 4000) and a different method for choosing  $A^\ell$  and  $A^u$ . MNB + EM accuracy strongly varies during the 20 runs, consistently providing higher maximums and lower minimums than MNB, for all  $|A^\ell|$  (see Figure 3). The maximum values for MNB + EM approach those reported in [6]. Figures 4 and 5 show four cases of confusion matrices for test B,  $A^\ell = 1260$ , for both MNB and MNB + EM methods. One can clearly verify EM’s clustering nature on the right side of Figure 3, with a significant number of false predictions for articles belonging to `comp.*` subgroups, mostly classified as belonging to the class `comp.graphics`. As expected, the main focus of confusion occur for clusters of subgroups, namely `comp.*`, `sci.*` (with a considerable number of articles belonging to several `comp.*` subgroups being misclassified as `sci.crypt`) and `talk.*`, for both MNB and MNB + EM cases.

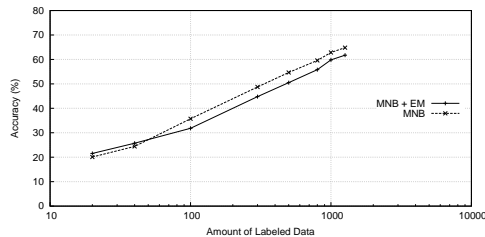


Figure 2. Results for test B: MNB + EM, using Rainbow’s framework. The results of MNB for the same  $A^\ell$  are also given.

Figure 6 shows the results for the EM- $\lambda$  method, for sev-

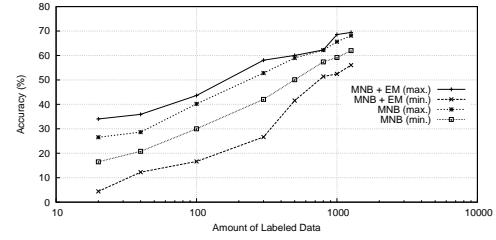


Figure 3. Results for test B: MNB + EM, using Rainbow’s framework. The results of MNB for the same  $A^\ell$  are also given.

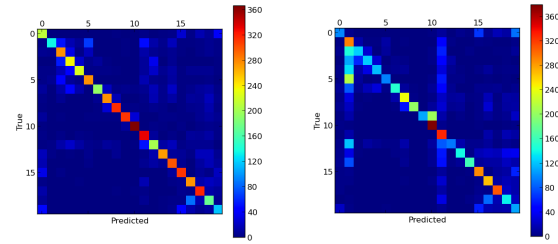


Figure 4. Results for test B: Confusion matrices for best (68.12 %) and worst (61.99 %) results of MNB,  $|A^\ell| = 1260$ .

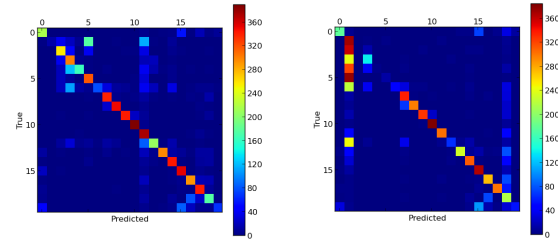


Figure 5. Results for test B: Confusion matrices for best (left side, 69.44 %) and worst (right side, 56.04 %) results of MNB + EM,  $|A^\ell| = 1260$ .

eral values of  $L = 1/\lambda$ <sup>8</sup>. We verified that for low values of  $|A^\ell|$ , lower values of  $L$  (and therefore higher  $\lambda$ ) provide better results, with the top accuracy verified for  $\lambda = 1$ . As  $|A^\ell|$  increases, higher accuracy values tend to be favored by lower values of  $\lambda$ . Nevertheless, the best results for EM- $\lambda$  keep occurring for  $\lambda = 1$ , which seems somewhat counter-intuitive.

<sup>8</sup>We use Rainbow’s option `--em-unlabeled-normalizer`, which determines the number of unlabeled articles it takes to equal a labeled article.

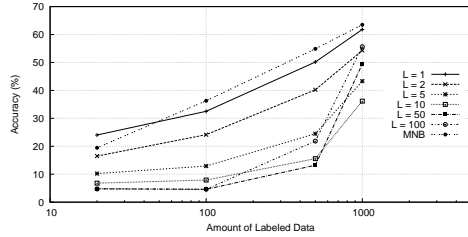


Figure 6. Results for test C: MNB + EM- $\lambda$ , using Rainbow’s framework. The results of MNB for the same  $A^\ell$  are also given.

## 6. Related Work

The study of techniques for text classification presented here is mostly based on the work by McCallum et al. [5] and Nigam et al. [6], extensively covered in Section 3. Within the field of supervised learning, starting from the MNB model presented in [5], Rennie et al. [7] propose a series of transforms, based on Information Retrieval techniques, which improve the performance of MNB, resulting in an enhanced method designated by Transformed Weight-normalized Complement Naive Bayes (TWCNB). In the semi-supervised scope, Mann et al. [4] proposed a general approach to semi-supervised learning, but tested on text classification, designated by Expectation Regularization (XR) [4]. XR is based on exponential-family parametric models, normally trained by MAP estimation. XR adds it with a second term, which attempts to minimize the difference between the predicted posteriors  $P(Y|X)$  for unlabeled data and estimated pre-known values for  $P(Y|X)$ , obtained from empirical data or the labeled data. The method is tested against MNB and MNB + EM, being outperformed by the two in the Simulated/Real/Aviation/Auto (SRAA) text classification task, similar to 20 Newsgroups task.

## 7. Conclusions

In this work we study the field of semi-supervised learning, applied to the field of text classification. We closely follow the work in [5, 6], reporting out understanding of the presented models.

We follow the methods presented by Nigam et al. [6] to implement our versions of the MNB and MNB + EM (including EM- $\lambda$ ) approaches. We take the well-known 20 Newsgroups dataset [3], particularly a version designated by 20news-bydate, to evaluate the studied semi-supervised models in the context of a real-world dataset. The 20news-bydate dataset is manipulated using the

Rainbow tools — developed by the authors of the base literature reviewed in this paper — in order to pre-process it according to commonly used procedures and divide it into training and test partitions, also creating unlabeled samples for the semi-supervised setting.

Our version of MNB is compared with Rainbow’s, producing lower accuracy values for small amounts of labeled data, approaching Rainbow’s accuracy rates for higher amounts of labeled data ( $|A^\ell| > 5000$ ). Due to problems with our implementation of MNB + EM, making its results unsuitable for comparison, we use Rainbow’s implementation of MNB + EM and EM- $\lambda$  over the 20news-bydate dataset to experimentally validate the previously studied semi-supervised methods. We verify that, on average MNB + EM performs better than MNB for  $|A^\ell| < 100$ , with  $11220 \leq |A^u| \leq 11240$ . We also verify that, despite its large variations, for at least one of the test runs performed for every  $|A^\ell|$ , MNB + EM surpasses MNB. Regarding EM- $\lambda$ , while smaller weights for unlabeled data seem to favor accuracy for larger  $|A^\ell|$ , the best results are always obtained for  $\lambda = 1$ . For  $|A^\ell| \geq 1000$ , the results seem intuitive, as the accuracy approaches that of MNB for smaller values of  $\lambda$  (i.e. smaller weights given to the unlabeled data).

## References

- [1] O. Chapelle, B. Schölkopf, and A. Zien. *Semi-Supervised Learning*. Adaptive Computation and Machine Learning. Mit Press, 2010. 1, 2
- [2] A. M. Kibriya, E. Frank, B. Pfahringer, and G. Holmes. Multinomial Naive Bayes for Text Categorization Revisited. In *Proceedings of the 17th Australian Joint Conference on Advances in Artificial Intelligence*, AI’04, pages 488–499, Berlin, Heidelberg, 2004. Springer-Verlag. 3, 6
- [3] K. Lang. Newsweeder: Learning to Filter Netnews. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 331–339, 1995. 1, 4, 5, 6, 8
- [4] G. S. Mann and A. McCallum. Simple, Robust, Scalable Semi-Supervised Learning via Expectation Regularization. In *Proceedings of the 24th international conference on Machine learning - ICML ’07*, pages 593–600, New York, New York, USA, 2007. ACM Press. 8
- [5] A. McCallum and K. Nigam. A Comparison of Event Models for Naive Bayes Text Classification. In *AAAI-98 Workshop on Learning for Text Categorization*, pages 41–48. AAAI Press, 1998. 1, 2, 3, 4, 6, 8
- [6] K. Nigam, A. K. McCallum, S. Thrun, and T. Mitchell. Text Classification from Labeled and Unlabeled Documents using EM. *Machine Learning - Special Issue on Information Retrieval*, 39(2-3):103–134, 2000. 1, 2, 3, 4, 5, 6, 7, 8
- [7] J. D. M. Rennie, L. Shih, J. Teevan, and D. R. Karger. Tackling the Poor Assumptions of Naive Bayes Text Classifiers. In *Proceedings of the Twentieth International Conference on Machine Learning*, pages 616–623, 2003. 2, 8
- [8] J. Su, J. Sayyad-Shirabad, and S. Matwin. Large Scale Text Classification using Semi-Supervised Multinomial Naive



Bayes. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011*, pages 97–104, 2011. [2](#), [3](#), [6](#)

[9] X. Zhu. Semi-Supervised Learning Literature Survey. Technical Report 1530, Computer Sciences, University of Wisconsin -Madison, 2008. [1](#), [2](#)

[10] X. Zhu and A. Goldberg. *Introduction to Semi-Supervised Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2009. [1](#), [4](#)







