

Quick review of mixture models and EM algorithm

- Let \mathbf{X} be a set of multivariate data points (vectors):

$$\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}.$$

- General expression for finite Gaussian mixture model:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (1)$$

- **Assumptions:**

- Each data point is generated by one of the K Gaussian components, with probability π_k .
- There is a one-to-one mapping between Gaussian components and classes.

- **Goal:** Find $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, k=1, \dots, K$ that maximizes likelihood of training data $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. This allows us to calculate, for each data point \mathbf{x}_i , degree of responsibility of each Gaussian component for \mathbf{x}_i .

EM algorithm

1. Initialize (randomly or some other way) the means $\boldsymbol{\mu}_k$, covariances $\boldsymbol{\Sigma}_k$, and mixing coefficients π_k , and evaluate initial value of log likelihood.

$$\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{n,k}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}, n = 1, \dots, N \text{ and } k = 1, \dots, K$$

3. **M step.** Re-estimate the parameters using the current responsibilities.

$$\boldsymbol{\mu}_k^{new} = \frac{1}{\underbrace{\sum_{n=1}^N \gamma(z_{n,k})}_{\text{called } N_k \text{ in Bishop book}}} \sum_{n=1}^N \gamma(z_{n,k}) \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{\sum_{n=1}^N \gamma(z_{n,k})} \sum_{n=1}^N \gamma(z_{n,k}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\boldsymbol{\pi}_k^{new} = \frac{\sum_{n=1}^N \gamma(z_{n,k})}{N}$$

4. Evaluate the log likelihood with the new parameters

$$\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

and check for convergence of either the parameters or the log likelihood. If not converged, return to step 2.

Text classification from labeled and unlabeled documents using EM

K. Nigam et al., *Machine Learning*, 2000

- Big problem with text classification: need labeled data.
- What we have: lots of unlabeled data.
- Question of this paper: Can unlabeled data be used to increase classification accuracy?
- I.e.: Any information implicit in unlabeled data? Any way to take advantage of this implicit information?

General idea: A version of EM algorithm

- Train a classifier with small set of available labeled documents.
- Use this classifier to assign probabilistically-weighted class labels to unlabeled documents by calculating expectation of missing class labels.
- Then train a new classifier using all the documents, both originally labeled and formerly unlabeled.
- Iterate.

Probabilistic framework

- Assumes data are generated with Gaussian mixture model
- Assumes one-to-one correspondence between mixture components and classes.
- “These assumptions rarely hold in real-world text data”

Probabilistic framework

Let $C = \{c_1, \dots, c_K\}$ be the classes / mixture components

Let $\theta = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\} \cup \{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\} \cup \{\pi_1, \dots, \pi_K\}$ be the mixture parameters.

Assumptions: A document d_i is created by first selecting a mixture component according to the mixture weights π_j , then having this selected mixture component generate a document according to its own parameters, with distribution

$$p(d_i | c_j; \theta).$$

- Likelihood of document d_i :

$$p(d_i | \theta) = \sum_{j=1}^k \pi_k p(d_i | c_j; \theta)$$

- Now, we will apply EM to a naive Bayes Classifier

Recall naive Bayes classifier: Assume each attribute is conditionally independent, given c_j .

Let $\mathbf{x} = (a_1, a_2, \dots, a_D)$

We have:

$$p(a_1, a_2, \dots, a_D | c_j) = p(a_1 | c_j) p(a_2 | c_j) \cdots p(a_D | c_j)$$

$$p(c_j | \mathbf{x}) = p(c_j) \prod_i p(a_i | c_j), i = 1, \dots, D; j = 1, \dots, K$$

To “train” naive Bayes from labeled data, estimate

$$p(c_j) \text{ and } p(a_i | c_j), j = 1, \dots, K; i = 1, \dots, D$$

These values are estimates of the parameters in θ . Call these values $\hat{\theta}$.

Note that Naive Bayes can be thought of as a generative mixture model.

Document d_i is represented as a vector of word frequencies $(w_1, \dots, w_{|V|})$, where V is the vocabulary (all known words).

There is an assumed probability distribution over words associated with each class, parameterized by θ .

We need to estimate $\hat{\theta}$ to determine what probability distribution document $d_i = (w_1, \dots, w_{|V|})$ is most likely to come from.

Applying EM to Naive Bayes

- We have a small number of labeled documents S_{labeled} and a large number of unlabeled documents, $S_{\text{unlabeled}}$.
- The initial parameters $\hat{\theta}$ are estimated from the labeled documents S_{labeled} .
- **Expectation step:** The resulting classifier is used to assign probabilistically-weighted class labels $p(c_j | \mathbf{x})$ to each unlabeled document $\mathbf{x} \in S_{\text{unlabeled}}$.
- **Maximization step:** Re-estimate $\hat{\theta}$ using $p(c_j | \mathbf{x})$ values for $\mathbf{x} \in S_{\text{unlabeled}} \cup S_{\text{labeled}}$.
- Repeat until $p(c_j | \mathbf{x})$ or $\hat{\theta}$ has converged.

Augmenting EM

What if basic assumptions (each document generated by one component; one-to-one mapping between components and classes) do not hold?

They tried two things to deal with this:

(1) Weighting unlabeled data less than labeled data

(2) Allow multiple mixture components per class:

A document may be comprised of several different sub-topics, each best captured with a different word distribution.

Data

- 20 UseNet newsgroups
- Web pages (WebKB)
- Newswire articles (Reuters)

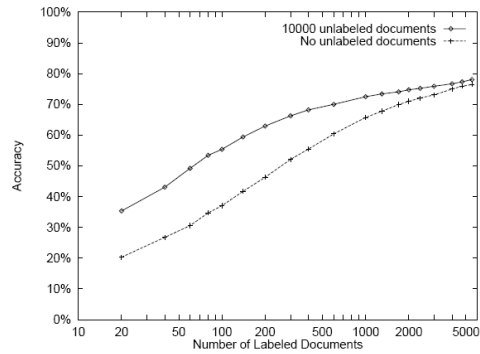


Figure 2. Classification accuracy on the 20 Newsgroups data set, both with and without 10,000 unlabeled documents. With small amounts of training data, using EM yields more accurate classifiers. With large amounts of labeled training data, accurate parameter estimates can be obtained without the use of unlabeled data, and the two methods begin to converge.

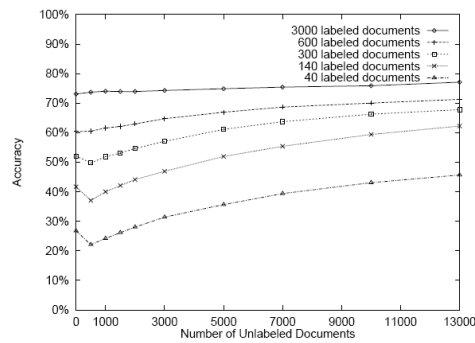


Figure 3. Classification accuracy while varying the number of unlabeled documents. The effect is shown on the 20 Newsgroups data set, with 5 different amounts of labeled documents, by varying the amount of unlabeled data on the horizontal axis. Having more unlabeled data helps. Note the dip in accuracy when a small amount of unlabeled data is added to a small amount of labeled data. We hypothesize that this is caused by extreme, almost 0 or 1, estimates of component membership, $P(c_j|d_i, \hat{\theta})$, for the unlabeled documents (as caused by naive Bayes' word independence assumption).

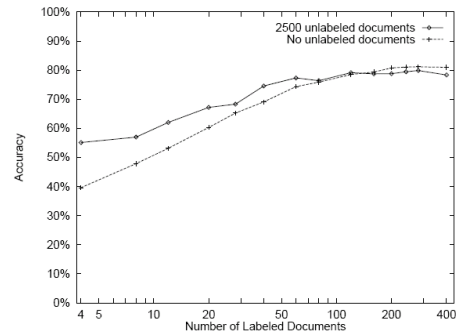


Figure 4. Classification accuracy on the WebKB data set, both with and without 2500 unlabeled documents. When there are small numbers of labeled documents, EM improves accuracy. When there are many labeled documents, however, EM degrades performance slightly—indicating a misfit between the data and the assumed generative model.

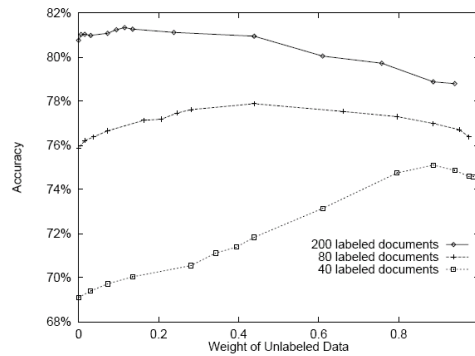


Figure 5. The effects of varying λ , the weighting factor on the unlabeled data in EM- λ . These three curves from the WebKB data set correspond to three different amounts of labeled data. When there is less labeled data, accuracy is highest when more weight is given to the unlabeled data. When the amount of labeled data is large, accurate parameter estimates are attainable from the labeled data alone, and the unlabeled data should receive less weight. With moderate amounts of labeled data, accuracy is better in the middle than at either extreme. Note the magnified vertical scale.

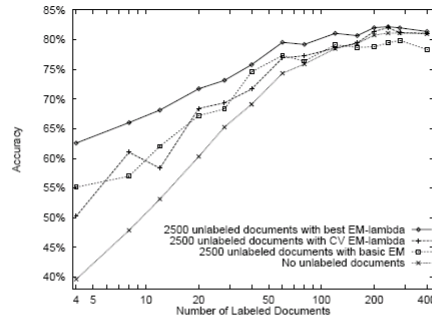


Figure 6. Classification accuracy on the WebKB data set, with modulation of the unlabeled data by the weighting factor λ . The top curve shows accuracy when using the best value of λ . In the second curve, λ is chosen by cross-validation. With small amounts of labeled data, the results are similar to basic EM; with large amounts of labeled data, the results are more accurate than basic EM. Thanks to the weighting factor, large amounts of unlabeled data no longer degrades accuracy, as it did in Figure 4, and yet the algorithm retains the large improvements with small amounts of labeled data. Note the magnified vertical axis to facilitate the comparisons.

Table 4. Precision-recall breakeven points showing performance of binary classifiers on Reuters with traditional naive Bayes (NB1), multiple mixture components using just labeled data (NB*), basic EM (EM1) with labeled and unlabeled data, and multiple mixture components EM with labeled and unlabeled data (EM*). For NB* and EM*, the number of components is selected optimally for each trial, and the median number of components across the trials used for the negative class is shown in parentheses. Note that the multi-component model is more natural for Reuters, where the negative class consists of many topics. Using both unlabeled data and multiple mixture components per class increases performance over either alone, and over naive Bayes.

Category	NB1	NB*	EM1	EM*	EM* vs NB1	EM* vs NB*
acq	69.4	74.3 (4)	70.7	83.9 (10)	+14.5	+9.6
corn	44.3	47.8 (3)	44.6	52.8 (5)	+8.5	+5.0
crude	65.2	68.3 (2)	68.2	75.4 (8)	+10.2	+7.1
earn	91.1	91.6 (1)	89.2	89.2 (1)	-1.9	-2.4
grain	65.7	66.6 (2)	67.0	72.3 (8)	+6.3	+5.7
interest	44.4	54.9 (5)	36.8	52.3 (5)	+7.9	-2.6
money-fx	49.4	55.3 (15)	40.3	56.9 (10)	+7.5	+1.6
ship	44.3	51.2 (4)	34.1	52.5 (7)	+8.2	+1.3
trade	57.7	61.3 (3)	56.1	61.8 (3)	+4.1	+0.5
wheat	56.0	67.4 (10)	52.9	67.8 (10)	+11.8	+0.4