

Bayesian Machine Learning: Assignment 3

Seung-Hoon Na

April 18, 2018

1 Generative models for document classification

1.1 Document representation

In this assignment, you will implement generative models for document classification. Each document is represented by two schemes – 1) **binary vector** based on word occurrences, 2) **word frequency vector** based on word frequencies. Suppose that D is a given document and its corresponding document vector is $\mathbf{x} = [x_1, \dots, x_V]$ where V is the number of all selected words. Each scheme for document representation is defined as follows:

1. **Binary vector:** $x_i = I(w_i \in D)$ where w_i indicates i -th word. So, when w_i occurs in D , $x_i = 1$. Otherwise, $x_i = 0$.
2. **Word frequency vector:** $x_i = c(w_i, D)$ where $c(w_i, D)$ means the number of times w_i occurs in D .

1.2 Generative classifiers using Bayesian naive Bayes

In generative approaches for classification task, the **class posterior** $P(y = c|\mathbf{x}, \theta)$ is computed from 1) the **class prior** $P(y = c) = \pi_c$ and 2) the **class-conditional density**, $P(\mathbf{x}|y = c, \theta)$.

Our goal is to implement different prediction methods and compare their classification accuracies, for each document representation scheme.

1.3 Dataset

Please use the following standard datasets.

1. **Reuters-21578**

<http://www.daviddlewis.com/resources/testcollections/reuters21578/>

2. **20 newsgroups**

<http://qwone.com/~jason/20Newsgroups/>

If necessary, please split data randomly to training/dev/test sets, with the ratios of 85%, 5%, 10%, respectively. (The full cross validation is not required in this assignment).

2 Task 1: Generative classifiers using Bayesian naive Bayes - Binary vector case

Assuming that a document is represented by a **binary vector**, **implement the following three classifiers and compare them in terms of classification accuracy** (using python codes).

1. **MLE**: MLE approximation for Bayesian naive Bayes(section 3.5). For prediction, implement Eq. (3.69) where $\hat{\pi}_c$ and θ are MLE.
2. **MAP**: MAP approximation for Bayesian naive Bayes (section 3.5). For prediction, implement Eq. (3.69) where $\hat{\pi}_c$ and θ are MAP estimate.
3. **Full posterior predictive** for Bayesian naive Bayes (section 3.5). For prediction, implement Eq. (3.66).

For binary case, refer to Algorithm 3.1 and 3.2. for parameter fitting and prediction, respectively.

You should submit the followings:

1. Source codes for data preprocessing and preprocessed data.
2. Source codes for model fitting and prediction.
3. Trained model files
4. Comparison results
5. README (i.e. describe how to train and test your models)

3 Task 2: Generative classifiers using Bayesian naive Bayes - Word frequency vector case

Now, assume that a document is represented by a **word frequency vector**.

For this word frequency scheme, **implement the following three classifiers and compare them in terms of classification accuracy** (using python codes).

Hint. Use Eq. (3.78) for the class-conditional density and derive its Bayesian naive Bayes prediction formula, similar to Eq. (3.66) for the binary case.

1. **MLE**: MLE approximation for Bayesian naive Bayes(section 3.5). For prediction, implement Eq. (3.69) where $\hat{\pi}_c$ and θ are MLE.
2. **MAP**: MAP approximation for Bayesian naive Bayes (section 3.5). For prediction, implement Eq. (3.69) where $\hat{\pi}_c$ and θ are MAP estimate.
3. **Full posterior predictive** for Bayesian naive Bayes (section 3.5). For prediction, implement Eq. (3.66).

For binary case, refer to Algorithm 3.1 and 3.2. for parameter fitting and prediction, respectively.

You should submit the followings:

1. Required formulas for MLE, MAP, Full posterior predictive methods for word frequency case.
2. Source codes for data preprocessing and preprocessed data.
3. Source codes for model fitting and prediction.
4. Trained model files
5. Comparison results
6. README (i.e. describe how to train and test your models)

4 Task 3: Applying feature selection using mutual information

Read Section 3.5.4, apply the MI-based feature selection to the “Full posterior predictive models” implemented in Task 1 and Task2.

You should submit the followings:

1. Required formulas for MI-based feature selection.
2. Source codes for feature selection
3. Trained model files
4. Comparison results (b/w original model and the feature selected model)
5. README (i.e. describe how to train and test your models)