

Computational Lab IV

Issued: Friday, November 3

Due: Friday, November 17, 10pm

Spam Filtering

We are going to build a spam filter. Let $C \in \{\text{spam}, \text{ham}\}$ be a random variable that represents the type of an email message; non-spam emails are commonly referred to as “ham.” We assume that each email is independent of all other emails. Let $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ be the set of all words used in all emails. Let $Y_i \in \{0, 1\}$ be a random variable that takes on value 1 if the email contains word \mathbf{w}_i and 0 if it does not. Our goal is to build a classifier for identifying spam from a training data set of emails that have been labeled as spam or ham.

Code and Data

The file `spam_filtering.zip` contains the functions:

<code>naivebayes.py</code>	Your code goes here.
<code>logistic.py</code>	Your code goes here.
<code>util.py</code>	Python functions/classes that will be useful as you write your classifier. Please do not modify.

and the data¹

<code>data/spam/</code>	Spam emails.
<code>data/ham/</code>	Non-spam (i.e., ham) emails.
<code>data/testing/</code>	Test data for you to classify.

Part I: Naïve Bayes Classifier

In building this classifier, we model a word’s presence in the email as independent of all other words, given the type of email. Specifically, our (generative) model is

$$p_{C,Y_1,\dots,Y_n}(c, y_1, \dots, y_n; \boldsymbol{\theta}) = p_C(c; s) \prod_{i=1}^n p_{Y_i|C}(y_i|c; p_i, q_i),$$

¹The data used in this assignment comes from a preprocessed version of the Enron email database. See V. Metsis, I. Androutsopoulos and G. Paliouras, “Spam Filtering with Naïve Bayes—Which Naïve Bayes?” in *Proc. Conf. Email and Anti-Spam (CEAS-2006)*, (Mountain View, CA), 2006.

with

$$p_C(c; s) = \begin{cases} s & c = \text{spam} \\ 1 - s & c = \text{ham}, \end{cases}$$

and, for $i = 1, \dots, n$,

$$p_{Y_i}(y_i|c; p_i, q_i) = \begin{cases} q_i^{y_i}(1 - q_i)^{1-y_i} & y_i \in \{0, 1\}, c = \text{spam} \\ p_i^{y_i}(1 - p_i)^{1-y_i} & y_i \in \{0, 1\}, c = \text{ham}, \end{cases}$$

where $\theta = (s, p_1, \dots, p_n, q_1, \dots, q_n)$ are the model parameters (all with values between 0 and 1) to be learned from the training set. In this part of the lab, we only work with functions in `naivebayes.py`.

- (a) *Training*. First, we estimate the parameters of the model.
 - i) Determine the maximum likelihood (ML) estimates of the parameters $(s, p_1, \dots, p_n, q_1, \dots, q_n)$, based on a training set consisting of k emails. This part has no coding.
 - ii) Implement the function `get_counts` that counts the number of files each word occurs in.
 - iii) Implement the function `get_log_probabilities` that computes the log of a smoothed frequency for each word. See Recitation 13 notes for more details on Laplace smoothing.
 - iv) Implement the function `learn_distributions`. This function takes two lists of filenames, one for spam and one for ham, and computes ML estimates of parameters $(s, p_1, \dots, p_n, q_1, \dots, q_n)$.
- (b) *Testing*. Implement the function `classify_message`, which computes the MAP estimate of the type of an email from the test dataset. The function should use Bayes' rule with the ML parameter estimates you obtained in the training. You can test your code by running

```
python naivebayes.py data/testing/ data/spam/ data/ham/
```

How well does your classifier perform?

- (c) Estimating the spam frequency s from data depends heavily on the number of examples in your training set. In the real world, it is often difficult to find good training examples for ham, since nobody wants to give out their private email for the world to read. As a result, spam datasets often have many more spam examples than ham examples. By setting s to reflect your belief of how often you get spam emails, we can adjust how much the algorithm favors catching spam at the expense of falsely flagging a ham message. Try setting s to a few values, for example, from 0.0001 to 0.9999, and briefly explain what happens to your performance as s increases and decreases.

Part II: Logistic Regression Classifier

We next use a logistic model for posterior distribution of email type C , with the simple features $g_i(y_i) = y_i$ for $i = 1, \dots, n$; specifically, our (discriminative) model is

$$p_{C|Y_1, \dots, Y_n}(c|y_1, \dots, y_n; \boldsymbol{\theta}) = \begin{cases} \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i y_i\right) & c = \text{spam} \\ 1 - \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i y_i\right) & c = \text{ham}, \end{cases}$$

with

$$\sigma(u) = \frac{1}{1 + e^{-u}}$$

denoting the sigmoid function defined in class, and where $\boldsymbol{\theta} = (\theta_0, \dots, \theta_n)$ are the model parameters to be learned from the training set.

Given training data in the form of k i.i.d. labeled data examples

$$(c, y_{1, \dots, n})^{(1, \dots, k)} = (c, y_{1, \dots, n})^{(1)}, \dots, (c, y_{1, \dots, n})^{(k)},$$

we learn the parameters of the model above by maximizing the component of the (log) likelihood used in discriminative modeling, or alternatively, by minimizing the negative log-likelihood, viz.,

$$\begin{aligned} \varphi\left(\boldsymbol{\theta}; (c, y_{1, \dots, n})^{(1, \dots, k)}\right) &= - \sum_{j=1}^k \ln p_{C|Y_1, \dots, Y_n}(c^{(j)}|y_1^{(j)}, \dots, y_n^{(j)}; \boldsymbol{\theta}) \\ &= - \sum_{j=1}^k \left[\mathbb{1}(c^{(j)} = \text{spam}) \ln \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i y_i^{(j)}\right) \right. \\ &\quad \left. + \mathbb{1}(c^{(j)} = \text{ham}) \ln \left(1 - \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i y_i^{(j)}\right)\right) \right], \end{aligned} \tag{1}$$

where to obtain the last equality we have used that a $B(p)$ distribution over $\mathcal{Z} = \{\mathbf{a}_1, \mathbf{a}_2\}$ with $p_Z(\mathbf{a}_1) = p$ can be expressed in the convenient form

$$p_Z(z) = p^{\mathbb{1}(z=\mathbf{a}_1)}(1-p)^{\mathbb{1}(z=\mathbf{a}_2)}.$$

We carry out this minimization numerically through a procedure referred to as *gradient descent*. Briefly, the procedure evaluates the gradient of the function $\varphi(\cdot)$ with respect to parameters $\boldsymbol{\theta}$ and takes a step towards lower value of the function along its gradient. When the value of the function does not increase with subsequent steps, we have arrived at a (local) maximum. We have provided an implementation of the gradient descent for you to use. In this part of the lab, we only work with functions in `logistic.py`.

(d) *Training.* We first estimate the model parameters.

- i) Implement the function `extract_features` that extracts the binary features y_1, \dots, y_n from the input email message as defined at the start of the problem.
- ii) Implement the function `logistic_eval` that computes

$$\sum_{j=1}^k \ln p_{C|Y_1, \dots, Y_n}(c^{(j)} | y_1^{(j)}, \dots, y_n^{(j)}; \boldsymbol{\theta}).$$

Note that the cost function $\varphi(\cdot)$ in (1) is just the negative of this quantity.

- iii) Obtain an expression for the derivative of the cost function

$$\varphi\left(\boldsymbol{\theta}; (c, y_{1, \dots, n})^{(1, \dots, k)}\right) \quad (2)$$

with respect to each of the constituent parameters $\theta_0, \dots, \theta_n$, including its derivation in your writeup. Then implement the function `logistic_derivative` that computes this derivative.

Hint: You can take the derivative for one data point at a time, and then sum the derivatives for all data points. Also, you may find it useful to use the following properties of the sigmoid function from class:

$$\sigma(-u) = 1 - \sigma(u), \quad \text{and} \quad \frac{d}{du} \sigma(u) = \sigma(u) \sigma(-u).$$

- iv) Implement the function `train_logistic`. This function takes two lists of filenames, one for spam and one for ham, and computes the ML estimates of parameters $\theta_0, \dots, \theta_n$. Set the `SHOW_LOSS_PLOT` global variable at the top of `logistic.py` to `True` to observe the behavior of the cost function (2) over iterations of the gradient descent procedure. Include the plot in your writeup.
- (e) *Testing.* Implement the function `classify_message` in `logistic.py`, which computes the MAP estimate of the email's type with the ML parameter estimates you obtained in training. You can test your code by running

```
python logistic.py data/testing/ data/spam/ data/ham/
```

How well does your classifier perform? How well does it perform compared to your previous classifier from `naivebayes.py`?

(Note that we've set `logistic.py` to use only 100 examples from each of spam and ham for training, governed by the `NUM_EXAMPLES` global variable at the top

of `logistic.py`. It's very slow to train on the full data, and the results aren't much different, so it's fine to train your model on just this smaller set.)

For the writeup, report the performance for `NUM_EXAMPLES=100`.

(Optional) Note that the performance of your classifier also depends on the `learning_rate` parameter for the training (defined in `optimize_theta` method). Feel free to experiment with different values. How are the training time and convergence of parameters be affected? For the writeup, please use the given `learning_rate` of 0.05.

- (f) Implement an alternative function `extract_features` with a different choice of features. For example, you could choose to use word *frequencies* instead of just binary labels for whether a word appears in the file, as we did previously. Describe your implementation of `extract_features`, and compare the results with those from the previous implementation. What do you think accounts for the difference in performance, if any?

What to hand in: Upload to Gradescope your writeup as a single PDF file and completed files `naivebayes.py` and `logistic.py`. The writeup should include your answers for all parts. Scans/images of handwritten work are fine but please write neatly — we can't grade what we can't decipher!

Gradescope will run a number of tests which are visible to you. These are to make sure your code runs correctly in the Gradescope system. Please make sure you pass these tests as we will use the autograder to check your submission.