

Homework 3: Support Vector Machines*

Pattern Recognition and Machine Learning

September 21, 2016

Instructions

- **Submission:** 本次作业需提交的内容包括实验报告文档和代码, 实验报告文档主要描述原理和结果, 请每位同学压缩打包, 压缩包命名格式为“学号-班级-姓名-PRMLHW3.zip”. 请各班负责人(班长、学委等)统一将本班的所有作业收集并按班级和作业名的格式统一打包(如计算机1301班的命名为“计算机1301-PRMLHW3.zip”), 在deadline前发到yxliang@csu.edu.cn. 未按时提交的则在deadline后自行按指定的命名格式打包后发到yxliang@csu.edu.cn.
- **Late homework policy:** 本次作业提交的deadline为2016.12.13 24:00, 约占平时成绩的1/3, 迟交则按每天0.8的比例进行计算, 如本来作业能打10分, 迟交1天则算8分, 迟交2天则算6.4分, 以此类推.
- **Collaboration policy:** Homeworks must be done individually, except where otherwise noted in the assignments. “Individually” means each student must hand in their own answers, and each student must write and use their own code in the programming parts of the assignment. It is acceptable for students to collaborate in figuring out answers and to help each other solve the problems, though you must in the end write up your own solutions individually, and you must list the names of students you discussed this with. We will be assuming that you will be taking the responsibility to make sure you personally understand the solution to any work arising from such collaboration.

Part I

Support Vector Machines

In the first half of this exercise, you will be using support vector machines (SVMs) with various example 2D datasets. Experimenting with these datasets will help you gain an intuition of how SVMs work and how to use a Gaussian kernel with SVMs. In the next half of the exercise, you will be using support vector machines to build a spam classifier. The provided script, `ex6.m`, will help you step through the first half of the exercise.

1 Linear SVM

We will begin by with a 2D example dataset which can be separated by a linear boundary. The script `ex6.m` will plot the training data (Figure 1). In this dataset, the positions of the positive examples (indicated with +) and the negative examples (indicated with o) suggest a natural separation indicated by the gap. However, notice that there is an outlier positive example + on the far left at about (0;1; 4:1). As part of this exercise, you will also see how this outlier affects the SVM decision boundary.

*本次作业内容来自Andrew Ng's mlClass <http://www.ml-class.org/course/class/index/>.

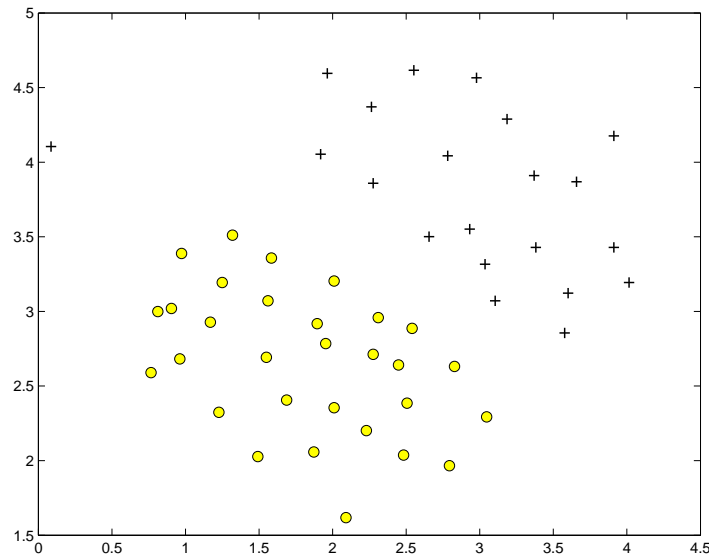


Figure 1: Example Dataset 1

In this part of the exercise, you will try using different values of the C parameter with SVMs. Informally, the C parameter is a positive value that controls the penalty for misclassified training examples. A large C parameter tells the SVM to try to classify all the examples correctly. C plays a role similar to $1/\lambda$, where λ is the regularization parameter that we were using previously for logistic regression.

The next part in `ex6.m` will run the SVM training (with $C = 1$) using SVM software that we have included with the starter code, `svmTrain.m`¹. When $C = 1$, you should find that the SVM puts the decision boundary in the gap between the two datasets and misclassifies the data point on the far left (Figure 2).

Implementation Note: Most SVM software packages (including `svmTrain.m`) automatically add the extra feature $x_0 = 1$ for you and automatically take care of learning the intercept term θ_0 . So when passing your training data to the SVM software, there is no need to add this extra feature $x_0 = 1$ yourself. In particular, in Octave your code should be working with training examples $x \in \mathbb{R}^n$ (rather than $x \in \mathbb{R}^{n+1}$); for example, in the first example dataset $x \in \mathbb{R}^2$.

Your task is to try different values of C on this dataset. Specifically, you should change the value of C in the script to $C = 100$ and run the SVM training again. When $C = 100$, you should find that the SVM now classifies every single example correctly, but has a decision boundary that does not appear to be a natural fit for the data (Figure 3).

2 SVM with Gaussian Kernels

In this part of the exercise, you will be using SVMs to do non-linear classification. In particular, you will be using SVMs with Gaussian kernels on datasets that are not linearly separable.

2.1 Gaussian Kernel

To find non-linear decision boundaries with the SVM, we need to first implement a Gaussian kernel. You can think of the Gaussian kernel as a similarity function that measures the “distance” between a pair of examples,

¹In order to ensure compatibility with Octave, we have included this implementation of an SVM learning algorithm. However, this particular implementation was chosen to maximize compatibility, and is not very efficient. If you are training an SVM on a real problem, especially if you need to scale to a larger dataset, we strongly recommend instead using a highly optimized SVM toolbox such as LIBLINEAR <https://www.csie.ntu.edu.tw/~cjlin/liblinear/> or LIBSVM <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.

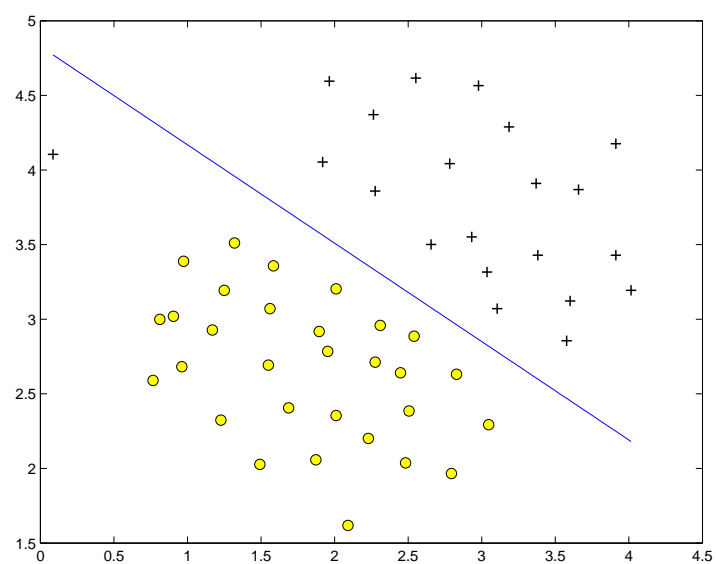


Figure 2: SVM Decision Boundary with $C = 1$ (Example Dataset 1)

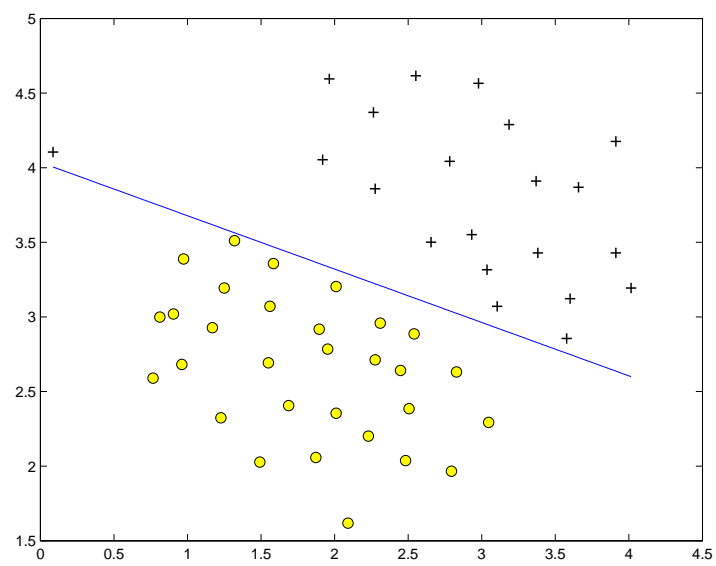


Figure 3: SVM Decision Boundary with $C = 100$ (Example Dataset 1)

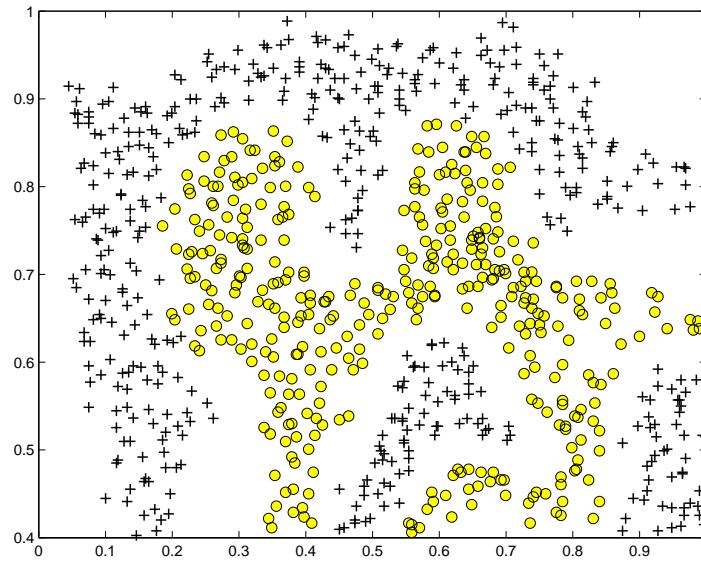


Figure 4: Example Dataset 2

$(x^{(i)}, x^{(j)})$. The Gaussian kernel is also parameterized by a bandwidth parameter, σ , which determines how fast the similarity metric decreases (to 0) as the examples are further apart.

You should now complete the code in `gaussianKernel.m` to compute the Gaussian kernel between two examples, $(x^{(i)}, x^{(j)})$. The Gaussian kernel function is defined as:

$$K(x^{(i)}, x^{(j)}) = \exp\left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{2\sigma^2}\right).$$

Once you've completed the function `gaussianKernel.m`, the script `ex6.m` will test your kernel function on two provided examples and you should expect to see a value of 0.324652.

2.2 Example Dataset 2

The next part in `ex6.m` will load and plot dataset 2 (Figure 4). From the Figure, you can observe that there is no linear decision boundary that separates the positive and negative examples for this dataset. However, by using the Gaussian kernel with the SVM, you will be able to learn a non-linear decision boundary that can perform reasonably well for the dataset.

If you have correctly implemented the Gaussian kernel function, `ex6.m` will proceed to train the SVM with the Gaussian kernel on this dataset. (Figure 5) shows the decision boundary found by the SVM with a Gaussian kernel. The decision boundary is able to separate most of the positive and negative examples correctly and follows the contours of the dataset well.

2.3 Example Dataset 3

In this part of the exercise, you will gain more practical skills on how to use a SVM with a Gaussian kernel. The next part of `ex6.m` will load and display a third dataset (Figure 6). You will be using the SVM with the Gaussian kernel with this dataset.

In the provided dataset, `ex6data3.mat`, you are given the variables `X`, `y`, `Xval`, `yval`. The provided code in `ex6.m` trains the SVM classifier using the training set `(X, y)` using parameters loaded from `dataset3Params.m`.

Your task is to use the cross validation set `Xval`, `yval` to determine the best C and σ parameter to use. You should write any additional code necessary to help you search over the parameters C and σ . For both C and σ , we suggest trying values in multiplicative steps (e.g., 0:01; 0:03; 0:1; 0:3; 1; 3; 10; 30). Note that you should try all possible pairs of values for C and σ (e.g., $C = 0:3$ and $\sigma = 0:1$). For example, if you try each of

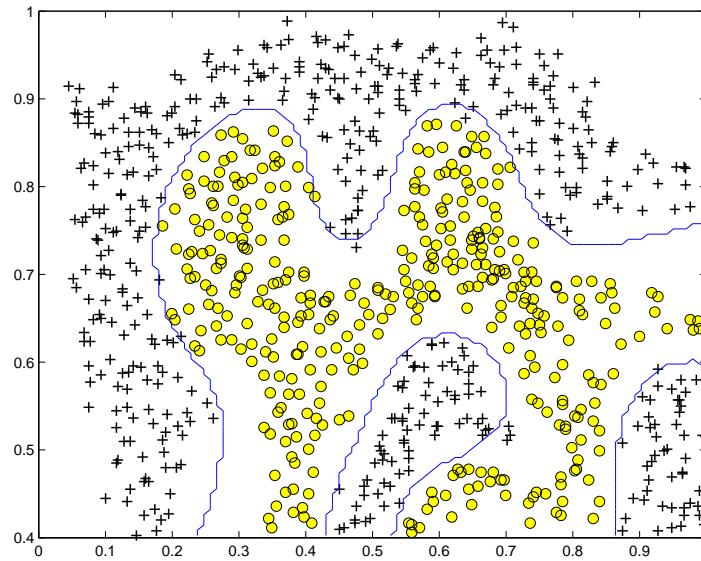


Figure 5: SVM (Gaussian Kernel) Decision Boundary (Example Dataset 2)

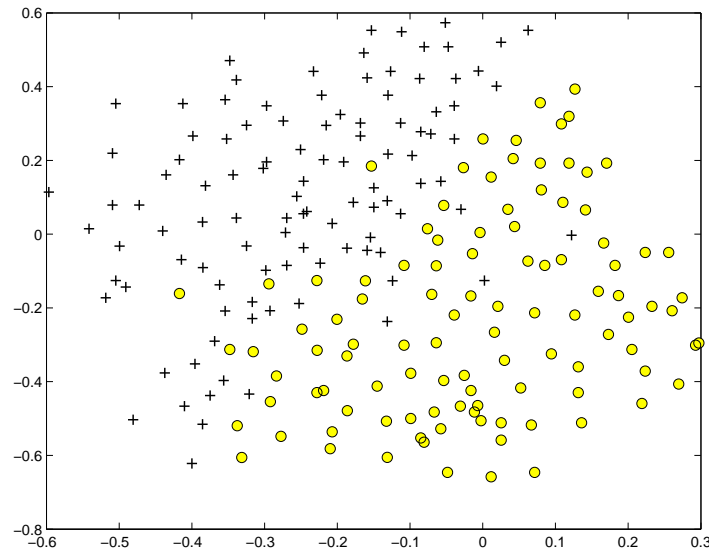


Figure 6: Example Dataset 3

the 8 values listed above for C and for σ^2 , you would end up training and evaluating (on the cross validation set) a total of $8^2 = 64$ different models.

After you have determined the best C and σ parameters to use, you should modify the code in `dataset3Params.m` filling in the best parameters you found. For our best parameters, the SVM returned a decision boundary shown in Figure 7.

Implementation Note: When implementing cross validation to select the best C and σ parameter to use, you need to evaluate the error on the cross validation set. Recall that for classification, the error is defined as the fraction of the cross validation examples that were classified incorrectly. In Octave, you can compute this error using `mean(double(predictions ~= yval))`, where `predictions` is a vector containing all the predictions from the SVM, and `yval` are the true labels from the cross validation set. You can use the `svmPredict` function to generate the predictions for the cross validation set.

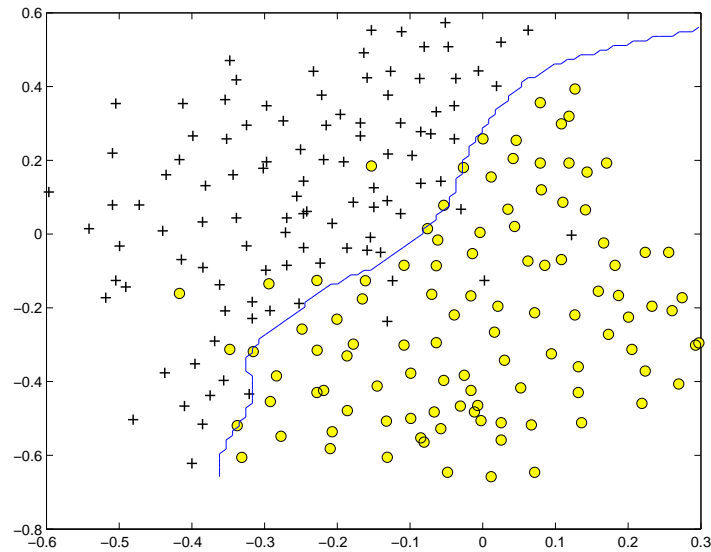


Figure 7: SVM (Gaussian Kernel) Decision Boundary (Example Dataset 3)

Part II

Spam Classification

Many email services today provide spam filters that are able to classify emails into spam and non-spam email with high accuracy. In this part of the exercise, you will use SVMs to build your own spam filters.

You will be training a classifier to classify whether a given email, x , is spam ($y = 1$) or non-spam ($y = 0$). In particular, you need to convert each email into a feature vector $x \in \mathbb{R}^n$. The following parts of the exercise will walk you through how such a feature vector can be constructed from an email.

Throughout the rest of this exercise, you will be using the script `ex6_spam.m`. The dataset included for this exercise is based on a subset of the **SpamAssassin** Public Corpus². For the purpose of this exercise, you will only be using the body of the email (excluding the email headers).

3 Preprocessing Emails

```
> Anyone knows how much it costs to host a web portal ?
>
Well, it depends on how many visitors youre expecting. This can be anywhere from less than 10 bucks a
month to a couple of $100. You should checkout http://www.rackspace.com/ or perhaps Amazon EC2 if
youre running something big..
To unsubscribe yourself from this mailing list, send an email to: groupname-unsubscribe@egroups.com
```

Figure 8: Sample Email

Before starting on a machine learning task, it is usually insightful to take a look at examples from the dataset. Figure 8 shows a sample email that contains a URL, an email address (at the end), numbers, and dollar amounts. While many emails would contain similar types of entities (e.g., numbers, other URLs, or other email addresses), the specific entities (e.g., the specific URL or specific dollar amount) will be different in almost every email. Therefore, one method often employed in processing emails is to “normalize” these values, so that all URLs are treated the same, all numbers are treated the same, etc. For example, we could replace each URL in the email with the unique string “httpaddr” to indicate that a URL was present.

²<http://spamassassin.apache.org/publiccorpus/>

```
anyon know how much it cost to host a web portal well it depend on how mani visitor your
expect thi can be anywher from less than number buck a month to a coupl of dollarnumb
you should checkout httpaddr or perhap amazon ecnumb if your run someth big to unsubscrib
yourself from thi mail list send an email to emailaddr
```

Figure 9: Preprocessed Sample Email

This has the effect of letting the spam classifier make a classification decision based on whether *any* URL was present, rather than whether a specific URL was present. This typically improves the performance of a spam classifier, since spammers often randomize the URLs, and thus the odds of seeing any particular URL again in a new piece of spam is very small.

In `processEmail.m`, we have implemented the following email preprocessing and normalization steps:

- **Lower-casing:** The entire email is converted into lower case, so that capitalization is ignored (e.g., `IndIcaTE` is treated the same as `Indicate`).
- **Stripping HTML:** All HTML tags are removed from the emails. Many emails often come with HTML formatting; we remove all the HTML tags, so that only the content remains.
- **Normalizing URLs:** All URLs are replaced with the text “`httpaddr`”.
- **Normalizing Email Addresses:** All email addresses are replaced with the text “`emailaddr`”.
- **Normalizing Numbers:** All numbers are replaced with the text “`number`”.
- **Normalizing Dollars:** All dollar signs (\$) are replaced with the text “`dollar`”.
- **Word Stemming:** Words are reduced to their stemmed form. For example, “discount”, “discounts”, “discounted” and “discounting” are all replaced with “discount”. Sometimes, the Stemmer actually strips off additional characters from the end, so “include”, “includes”, “included”, and “including” are all replaced with “includ”.
- **Removal of non-words:** Non-words and punctuation have been removed. All white spaces (tabs, newlines, spaces) have all been trimmed to a single space character.

The result of these preprocessing steps is shown in Figure 9. While preprocessing has left word fragments and non-words, this form turns out to be much easier to work with for performing feature extraction.

```
1 aa
2 ab
3 abil
...
86 anyon
...
916 know
...
1898 zero
1899 zip
```

Figure 10: Vocabulary List

```
86 916 794 1077 883 370
1699 790 1822 1831 883 431
1171 794 1002 1893 1364
592 1676 238 162 89 688
945 1663 1120 1062 1699 375
1162 479 1893 1510 799 1182
1237 810 1895 1440 1547 181
1699 1758 1896 688 1676 992
961 1477 71 530 1699 531
```

Figure 11: Word Indices for Sample Email

3.1 Vocabulary List

After preprocessing the emails, we have a list of words (e.g., Figure 9) for each email. The next step is to choose which words we would like to use in our classifier and which we would want to leave out.

For this exercise, we have chosen only the most frequently occurring words as our set of words considered (the vocabulary list). Since words that occur rarely in the training set are only in a few emails, they might cause the model to overfit our training set. The complete vocabulary list is in the file `vocab.txt` and also shown in Figure 10. Our vocabulary list was selected by choosing all words which occur at least a 100 times in the spam corpus, resulting in a list of 1899 words. In practice, a vocabulary list with about 10,000 to 50,000 words is often used.

Given the vocabulary list, we can now map each word in the preprocessed emails (e.g., Figure 9) into a list of word indices that contains the index of the word in the vocabulary list. Figure 11 shows the mapping for the sample email. Specifically, in the sample email, the word “anyone” was first normalized to “anyon” and then mapped onto the index 86 in the vocabulary list.

Your task now is to complete the code in `processEmail.m` to perform this mapping. In the code, you are given a string `str` which is a single word from the processed email. You should look up the word in the vocabulary list `vocabList` and find if the word exists in the vocabulary list. If the word exists, you should add the index of the word into the word indices variable. If the word does not exist, and is therefore not in the vocabulary, you can skip the word.

Once you have implemented `processEmail.m`, the script `ex6_spam.m` will run your code on the email sample and you should see an output similar to Figures 9 & 11.

Octave Tip: In Octave, you can compare two strings with the `strcmp` function. For example, `strcmp(str1, str2)` will return 1 only when both strings are equal. In the provided starter code, `vocabList` is a “cell-array” containing the words in the vocabulary. In Octave, a cell-array is just like a normal array (i.e., a vector), except that its elements can also be strings (which they can’t in a normal Octave matrix/vector), and you index into them using curly braces instead of square brackets. Specifically, to get the word at index `i`, you can use `vocabList{i}`. You can also use `length(vocabList)` to get the number of words in the vocabulary.

4 Extracting Features from Emails

You will now implement the feature extraction that converts each email into a vector in \mathbb{R}^n . For this exercise, you will be using $n = \#$ words in vocabulary list. Specifically, the feature $x_i \in \{0, 1\}$ for an email corresponds to whether the i -th word in the dictionary occurs in the email. That is, $x_i = 1$ if the i -th word is in the email and $x_i = 0$ if the i -th word is not present in the email.

Thus, for a typical email, this feature would look like:

$$x = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^n.$$

You should now complete the code in `emailFeatures.m` to generate a feature vector for an email, given the word indices.

Once you have implemented `emailFeatures.m`, the next part of `ex6_spam.m` will run your code on the email sample. You should see that the feature vector had length 1899 and 45 non-zero entries.

5 Training SVM for Spam Classification

After you have completed the feature extraction functions, the next step of `ex6_spam.m` will load a preprocessed training dataset that will be used to train a SVM classifier. `spamTrain.mat` contains 4000 training examples of spam and non-spam email, while `spamTest.mat` contains 1000 test examples. Each original email was processed using the `processEmail` and `emailFeatures` functions and converted into a vector $x^{(i)} \in \mathbb{R}^{1899}$.

After loading the dataset, `ex6_spam.m` will proceed to train a SVM to classify between spam ($y = 1$) and non-spam ($y = 0$) emails. Once the training completes, you should see that the classifier gets a training accuracy of about 99.8% and a test accuracy of about 98.5%.

6 Top Predictors for Spam

To better understand how the spam classifier works, we can inspect the parameters to see which words the classifier thinks are the most predictive of spam. The next step of `ex6_spam.m` finds the parameters with the largest positive values in the classifier and displays the corresponding words (Figure 12). Thus, if an email contains words such as “guarantee”, “remove”, “dollar”, and “price” (the top predictors shown in Figure 12), it is likely to be classified as spam.

our click remov guarante visit basenumb dollar will price pleas nbsp most lo ga dollarnumb

Figure 12: Top predictors for spam email

7 Try your own emails

Now that you have trained a spam classifier, you can start trying it out on your own emails. In the starter code, we have included two email examples (`emailSample1.txt` and `emailSample2.txt`) and two spam examples (`spamSample1.txt` and `spamSample2.txt`). The last part of `ex6_spam.m` runs the spam classifier over the first spam example and classifies it using the learned SVM. You should now try the other examples we have provided and see if the classifier gets them right. You can also try your own emails by replacing the examples (plain text files) with your own emails.

8 Build your own dataset

In previous exercises, we provided a preprocessed training set and test set. These datasets were created using the same functions (`processEmail.m` and `emailFeatures.m`) that you now have completed. For this exercise, you will build your own dataset using the original emails from the [SpamAssassin Public Corpus](#).

Your task in this exercise is to download the original files from the public corpus and extract them. After extracting them, you should run the `processEmail`³ and `emailFeatures` functions on each email to extract a feature vector from each email. This will allow you to build a dataset X , y of examples. You should then randomly divide up the dataset into a training set, a cross validation set and a test set.

While you are building your own dataset, we also encourage you to try building your own vocabulary list (by selecting the high frequency words that occur in the dataset) and adding any additional features that you think might be useful.

Finally, we also suggest trying to use highly optimized SVM toolboxes such as [LIBSVM](#).

³The original emails will have email headers that you might wish to leave out. We have included code in `processEmail` that will help you remove these headers.