

Project 1: Structural Risk Minimization



"One person's spam is another person's dinner."
-- ancient German wisdom

Introduction

In this project you will be building an email spam filter. First perform an "svn update" in your svn root directory.

The code for this project (`project1`) consists of several files, some of which you will need to read and understand in order to complete the assignment, and some of which you can ignore.

Files you'll edit:

<code>partners.txt</code>	If you work in a group, this file should contain the two wustlkeys of you and your group partner. These should be in two separate lines (the first two lines). There should be nothing else in this file. Please make sure that your partner also puts your wustlkey in his/her <code>partners.txt</code> file, as project partnerships must be reciprocal . If you don't have a partner then just leave the file alone.
<code>ridge.m</code>	Computes the ridge regression loss and gradient
<code>hinge.m</code>	Computes the hinge loss and gradient
<code>grdescent.m</code>	Performs gradient descent
<code>linearmodel.m</code>	Returns the predictions for a weight vector and a data set.
<code>logistic.m</code>	Computes the logistic regression loss and gradient.
<code>spamupdate.m</code>	(optional) Allows you to update the spam filter when you make a mistake.
<code>trainspamfilter.m</code>	Trains your spam filter and saves the final weight vector in a file w0.mat .

Files you want to look at and maybe change:

tokenizedata.m	Calls tokenize.py to tokenize the raw data into vectorial format.
tokenize.py	A simple python script that turns raw emails into bag of word vectors.
example_tests.m	Describes several unit tests to find obvious bugs in your implementation.

Files you might want to look at:

valsplit.m	This function takes the data and splits it into 80% training (xTr,yTr) and 20% validation (xTv,yTv). The splitting is not random but by time (i.e. the training data consists of emails that were received before the validation data.)
spamfilter.m	Loads in the file w0.mat and applies the corresponding spam filter on whatever test set you pass on as argument.
spamdemo.m	Runs your classifier on some sample emails, and shows you the ones it misclassifies.

How to submit: You can commit your code with subversion, with the command line

```
svn commit -m "some insightful comment"
```

where you should substitute "some meaningful comment" with something that describes what you did. You can submit as often as you want until the deadline. Please be aware that the last submission determines your grade.

Evaluation: Your code will be autograded for technical correctness. Please *do not* change the names of any provided functions or classes within the code, or you will wreak havoc on the autograder. However, the correctness of your implementation -- not the autograder's output -- will be the final judge of your score. If necessary, we will review and grade assignments individually to ensure that you receive due credit for your work.

Regrade Requets: We will set up a forum/goolge doc for regrade requests in due course.

Academic Dishonesty: We will be checking your code against other submissions in the class for logical redundancy. If you copy someone else's code and submit it with minor changes, we will know. These cheat detectors are quite hard to fool, so please don't try. We trust you all to submit your own work only; *please* don't let us down. If you do, we will pursue the strongest consequences available to us.

Getting Help: You are not alone! If you find yourself stuck on something, contact the course TAs for help. Office hours and [Piazza](#) are there for your support; please use the appropriate tags (**project1** and/or **autograder**). If you can't make any of our office hours, let us know and we can schedule an alternative time. We want these projects to be rewarding and instructional, not frustrating and demoralizing. But, we don't know when or how to help unless you ask.

Computing derivatives

Before you dive into the programming part of this assignment you will need to derive the gradients for several loss functions. **You do not have to hand this part in, but save your derivations as these are part of written homework 2.**

Derive the gradient function for each of the following loss functions with respect to the weight vector w . Write down the gradient update (with stepsize c).

(Note that: $\|w\|_2^2 = w^\top w$ and λ is a non-negative constant.)

1. Ridge Regression: $\mathcal{L}(w) = \sum_{i=1}^n (w^\top x_i - y_i)^2 + \lambda \|w\|_2^2$
2. Logistic Regression: ($y_i \in \{+1, -1\}$): $\mathcal{L}(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^\top x_i))$
3. Hinge loss: ($y_i \in \{+1, -1\}$): $\mathcal{L}(w) = \sum_{i=1}^n \max(1 - y_i(w^\top x_i + b), 0) + \lambda \|w\|_2^2$

Building an email spam filter

You will now implement these functions and their gradient updates. The file **data_train.mat** in **data** contains the pre-processed email data, where emails are represented as bag-of-words vectors. You will only need this file for now, but to compete in the spam filtering competition you might want to use **data_train** folder in **data**, which contains the raw email data, so that you can invent your own features.

Type in

```
>> addpath('data')
>> load data_train
>> valsplit
>> whos
```

Variables in the current scope:					
Name	Size	Bytes	Class	At	
X	1024x5000	20022664	double	s	
Y	1x5000	40000	double	s	
xTr	1024x4000	16425048	double	s	
xTv	1024x1000	3597624	double	s	
yTr	1x4000	32000	double	s	
yTv	1x1000	8000	double	s	

This should generate a training data set **xTr**, **yTr** and a validation set **xTv**, **yTv** for you. (You can check which variables are in your environment with the **whos** command.)

It is now time to implement your classifiers. We will always use gradient descent, but with various loss functions.

1. Implement the function **ridge.m** which computes the loss and gradient for a particular data set **xTr**, **yTr** and a weight vector **w**. Make sure you don't forget to incorporate your regularization constant λ . You can check your gradient with the following expression (it checks the difference between the interpolated and the actual function value):

```
err= checkgrad('ridge', zeros(size(xTr,1),1), 1e-5
```

The last four input variables are passed on to the function specified as the first argument. $1e-5$ represents 10^{-5} in MATLAB. The return value should be very small (less than 10^{-8}). (You can also use the **checkgrad** function for the later functions in this assignment and throughout your entire life whenever you have to implement functions and their gradients.)

2. Implement the function **grdescent.m** which performs gradient descent. Make sure to include the tolerance variable to stop early if the norm of the gradient is less than the tolerance value (you can use the function **norm(x)**). When the norm of the gradient is tiny it means that you have arrived at a minimum.
The first parameter of **grdescent** is a function which takes a weight vector and returns loss and gradient. In Octave you can make inline functions e.g. with the following code (first line):

```
f=@(w) ridge(w,xTr,yTr,0.1);
w=grdescent(f,zeros(size(xTr,1),1),1e-6,1000);
```

You can choose what kind of step-size you implement (e.g. constant, decreasing, line search,...). [HINT: Personally, I increase the stepsize by a factor of 1.01 each iteration where the loss goes down, and decrease it by a factor 0.5 if the loss went up. ... if you are smart you also undo the last update in that case to make sure the loss decreases every iteration.]

3. Write the (almost trivial) function `linearmodel` which returns the predictions for a vector `w` and a data set `xTv`.
4. Now call:

```
>> trainspamfilter(xTr,yTr);
>> spamfilter(xTv,yTv,0.1);
False positive rate: 0.78%
True positive rate: 77.39%
AUC: 98.50%
```

The first command trains a spam filter with ridge regression and saves the resulting weight vector in `w0.mat`.

The second command will run your spam filter with the weights in `w0.mat` over the validation data set.

The outputs of `spamfilter.m` are:

- *false positive rate (fpr)* (how many emails you accidentally classify as spam).
 - *true positive rate (tpr)* (how many spam emails you catch).
 - *area under the curve (AUC)*, which different from tpr and fpr is independent of the cut-off threshold (the last argument into `spamfilter.m`). As the name suggests, it computes the area of the [ROC curve](#) and is a good measure to compare spam filters.
5. Now implement the function `hinge.m`, which is the equivalent to ridge but with the hinge loss. Take a look at `trainspamfilter.m`. You can change it to use the logistic loss instead of ridge regression to train the classifier.
 6. Now implement the function `logistic.m`, which is the equivalent to ridge but with the log-loss (logistic regression). [By default the logistic loss does not take a regularization constant, but feel free to incorporate regularization if you want to.]
 7. Now, run `vis_rocs` to see if your algorithms all work. You might have to fiddle with the `STEPSIZE` parameter at the very top (maybe set it to something very small initially (e.g. `1e-08`) and work yourself up). If you want to, change the `trainspamfilter.m` to a different loss function with different parameters.

Hints

Tests. To test your code you can implement and run `example_tests.m`, which describes and partially implements several example unit tests. Those tests are a subset of what we will use in the autograder to grade your submission.

Feature Extraction (Quality Evaluation)

30% of the grade for your project 1 submission will be assigned by how well your Spam classifier performs on a secret test set of emails. If you want to improve your classifier modify `tokenize.py` or `trainspamfilter`. You must train your weight vector locally and commit the changes to it (the file `w0` is written by `trainspamfilter` and contains the weight vector learned by your algorithm. The autograder will evaluate this classifier on emails from the same authors as the ones in your dataset, but emails that arrived later on. You can also modify `spamupdate.m` see task 1 below for tips to get you started. Also consider changing the default threshold in

spamfilter.m which is currently set to 0.3. We will use your modified tokenizers to tokenize the test emails as well.

1. **(Optional)** you can implement the function spamupdate to make small gradient steps during test time (basically you still correct the classifier after you made a mistake).
2. **(Optional)** If you take a look at the script tokenizedata.m, you can get an idea of how the tokenization is done. You can modify this if you want to change how the tokenization is done. For example, by default the data uses $2^{10} = 1024$ dimensional features. You could change this by increasing 10 to 11. You could also change the tokenize.py function (e.g. to include bigrams). A common trick is e.g. to remove stopwords. A full list is [here](#). You can also include bi-grams or feature re-weighting with [TFIDF](#)

After making modifications, you can call the function tokenizedata to generate data_train.mat.

```
>> tokenizedata
Running python tokenization script ...
Loading data ...
>> ls *.mat
data_train.mat
```

You must commit data_train.mat (in addition to tokenizedata.m and tokenize.py) to your repository if you want the autograder to use your tokenization in place of ours.

Note to Windows users: You must have Python 2.7 installed and accessible from one of the directories in your PATH environment variable in order for the tokenize script to work. If you have already installed Python, open command prompt and type "python". If the Python 2.7 interpreter comes up, you should be all set. If it does not, and you have installed Python 2.7, you can add your Python directory to your path environment variable by following [these instructions](#). The tokenizer script is not compatible with Python 3.

Credits: Project adapted from Kilian Weinberger (Thanks for sharing!).