

Problem Set 3

*Handed Out: March 22, 2017**Official Due Date: April 4, 2017**Extended Due Date: April 10, 2017*

- Please submit your solutions via your CCIS `github` account.
- Materials associated with this problem set are available at <https://github.ccs.neu.edu/cs6140-03-spring2017/materials>.
- I encourage you to discuss the homework with other members of the class. The goal of the homework is for you to learn the course material. However, you should write your own solution.
- Please keep your solution brief, clear, and legible. If you are feeling generous, I would *really* appreciate typed solutions (and if you plan on publishing CS/Math/Engineering research, this is actually a good exercise) – see the source material if you would like to use \LaTeX to do this.
- I encourage you to ask questions before class, after class, via email, or the Piazza QA section. However, please do not start e-mailing me questions the night before the homework is due. ☺

1. [Classification via Linear Programming – 15 points]

Let \mathcal{S} represent a training dataset of size m where each $\mathbf{x}_i \in \mathbb{R}^d$ is a d -dimensional vector and $y_i \in \{-1, 1\}$ is the corresponding label. Let $\mathbf{w} \in \mathbb{R}^d$ be the learned weight vector and θ be the threshold value such that we predict $h(\mathbf{x}_i) = 1$ if and only if $\mathbf{w}^T \mathbf{x}_i + \theta \geq 0$. Else, $h(\mathbf{x}_i) = -1$.¹

Instead of the SVM formulation, consider the following linear program formulation:²

$$\min \quad z = \xi \tag{1}$$

$$\text{s.t.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + \theta) \geq 1 - \xi \quad \forall (\mathbf{x}_i, y_i) \in \mathcal{S} \tag{2}$$

$$\xi \geq 0 \tag{3}$$

A (restricted) general form of a linear program specification can be stated as

$$\min \quad z = \mathbf{c}^T \mathbf{x} \tag{4}$$

$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \geq \mathbf{b} \tag{5}$$

where \mathbf{c} is often referred to as the *cost vector* and z as the *objective function*.³

¹Note that this is precisely the setting you are used to by this point.

²This is essentially a simplified variant of K.P. Bennett and O.L. Mangasarian, Robust Linear Programming Discrimination of Two Linearly Inseparable Sets. Optimization Methods and Software 1, 1992, 23-34. (<http://research.cs.wisc.edu/techreports/1991/TR1054.pdf>)

³This formulation also closely parallels the linear programming representation found in *Pattern Classification*, by Duda, Hart, and Stork.

Please rewrite the linear program given by Equations (1-3) in terms of \mathbf{c} , \mathbf{x} , \mathbf{A} , and \mathbf{b} . Note that given this formulation, you can use Matlab, Xpress, CVXOPT, GLPK, or many other LP solvers (including open-source) to solve this linear discriminant formulation.⁴

2. [Online Learning – 85 points]

The purpose of this problem set is to compare the respective performance of a few well-known online learning algorithms by comparing their performance on a synthetic dataset. While the assignments thus far have placed an emphasis on using *real* data, simulations are often better suited for understanding specific properties of learning algorithms.

Specifically, we are going to explore learning a linear function resulting from training on data generated by an *l-of-m-of-n* function. An *l-of-m-of-n* is defined over a n -dimensional Boolean vector where there is a defined subset of m attributes such that $f(\mathbf{x}) = 1$ if and only if at least l of these m attributes are active in \mathbf{x} . Note that this is a linear function,⁵ meaning that WINNOW and PERCEPTRON are able to represent the target concept (i.e., the target concept is contained in the hypothesis space).

Algorithms

Specifically, you must implement two online learning algorithms: PERCEPTRON and WINNOW (with and without margin in both cases) as defined below for Boolean vectors $\mathbf{x} \in \{0,1\}^n$. Note that these are very closely related algorithms with the most significant different being the update rules.

Algorithm 1 PERCEPTRON

Input: Training data \mathcal{S} , Number of rounds T , Learning rate η , Estimated margin γ

```

 $\mathbf{w} \leftarrow \mathbf{0}$ ;  $b \leftarrow 0$ ;  $t \leftarrow 0$ 
while  $t < T$  do
  for all  $(\mathbf{x}, y) \in \mathcal{S}$  do
    if  $y(\mathbf{w} \cdot \mathbf{x} + b) \leq \gamma$  then
       $\mathbf{w} \leftarrow \mathbf{w} + \eta y \mathbf{x}$ 
       $b \leftarrow b + \eta y$ 
    end if
  end for
   $t \leftarrow t + 1$ 
end while

```

Output: Learned hypothesis \mathbf{w}

⁴Not that I am asking you to do this – although it is useful to do at some point in your life...

⁵If you don't see this immediately, spend some time convincing yourself.

Algorithm 2 WINNOW

Input: Training data \mathcal{S} , Number of rounds T , Learning rate η , Threshold θ , Estimated margin γ

```
w  $\leftarrow$  1;  $t \leftarrow 0$ 
while  $t < T$  do
  for all  $(\mathbf{x}, y) \in \mathcal{S}$  do
    if  $y(\mathbf{w} \cdot \mathbf{x} - \theta) \leq \gamma$  then
       $\mathbf{w} \leftarrow \mathbf{w} \circ \eta^{y\mathbf{x}}$ 
    end if
  end for
   $t \leftarrow t + 1$ 
end while
```

Output: Learned hypothesis \mathbf{w}

Some Notes

While it is difficult for me to preemptively guess which problems you will encounter, here are some notes that I have observed to be confusing for some.

- For PERCEPTRON without margin, if we set $[\mathbf{w} \ b] = \mathbf{0}$ (i.e., the concatenation of \mathbf{w} and b), the learning rate has no effect. Therefore, we will just set $\eta = 1$.
- For PERCEPTRON with margin, the relationship between η and γ are closely related – updates with a large η will more likely make $y(\mathbf{w} \cdot \mathbf{x} + b) > \gamma$. We will set $\gamma = 1$ and experiment with $\eta = \{1.5, 0.25, 0.03, 0.005, 0.001\}$.
- Since l -of- m -of- n functions are monotone, we can use the simplest version of WINNOW as shown, noting that we **do not** update θ . The parameters we will consider includes $\eta = \{1.1, 1.01, 1.005, 1.0005, 1.0001\}$.
- WINNOW with margin is somewhat sensitive. Therefore, while we will experiment with varying $\eta = \{1.1, 1.01, 1.005, 1.0005, 1.0001\}$ and $\gamma = \{2.0, 0.3, 0.04, 0.006, 0.001\}$.

Generating Data

Data is generated as follows. First, the label is generated from a Bernoulli distribution with $\mu = 0.5$.⁶

- For each positive example, select a number of active features from the interval $l \leq a \leq m$. Secondly, select randomly from a uniform distribution a attributes amongst x_1, \dots, x_m and set to 1. Set the other $m - a$ attributes to 0. Set the remainder of the $n - m$ attributes (e.g., x_{m+1}, \dots, x_n) to 1 uniformly with a probability of 0.5.
- For each negative example, select a number of active features from the interval $0 \leq a \leq l - 1$. Secondly, select randomly from a uniform distribution a attributes amongst x_1, \dots, x_m and set to 1. Set the other $m - a$ attributes to 0. Set the remainder of the $n - m$ attributes (e.g., x_{m+1}, \dots, x_n) to 1 uniformly with a probability of 0.5.

You have been provided with the Java file `GenerateLMN.java` which can be used to generate this data in LIBSVM format (sparse representation). Note that we are manipulating x_1, \dots, x_m to represent the relevant portion of the function out of convenience (as there is no loss of generality). However, you should *not* use this information to influence your learning algorithms.

Tuning Parameters

While learning algorithms aren't generally as sensitive as it may seem in this assignment, one purpose of this assignment is to learn how to set hyper-parameters (e.g., η, γ , etc.).⁷

To set the hyper-parameters, we will set aside two distinct subsamples of the training data, \mathcal{D}_1 and \mathcal{D}_2 , each consisting of 10% of the training data. For each set of hyper-parameters settings, training the desired algorithm on \mathcal{D}_1 by running the algorithm twenty times over the data. Based on the resulting learned parameters, evaluate this model on \mathcal{D}_2 and record the resulting accuracy.

Choosing the set of hyper-parameters that achieve the best performance on \mathcal{D}_2 , conduct the experiments described in each section. Note that you will be testing the outer product of hyper-parameters settings, meaning that two hyper-parameters with five values each will result in testing 25 sets of hyper-parameters. For your own sanity, I suggest writing code such that this can be automated.

⁶Yes, you too can now make a fair coin flip seem sophisticated.

⁷Note that the process I describe also works for fixed-size data (although there are other reasonable protocols); for this specific problem, we could actually do this differently since we can generate an arbitrary amount of data.

Experiments

(a) [20 points] Counting Mistakes

The first set of experiments is to evaluate the specified learning algorithms with two target function configurations: (a) $l = 10, m = 100, n = 500$ and (b) $l = 10, m = 100, n = 1000$.

In each case, your experiments should consist of the following steps:

- i. Generate a *clean*⁸ 50,000 instance dataset for each specified $\{l, m, n\}$ configuration.
- ii. Fill in the Table 1 with the best performing hyper-parameters:

Algorithm	$n = 500$			$n = 1000$		
	Params	Acc(\mathcal{D}_2)	M	Params	Acc(\mathcal{D}_2)	M
PERCEPTRON($\gamma = 0$)						
PERCEPTRON($\gamma > 0$)						
WINNOW($\gamma = 0$)						
WINNOW($\gamma > 0$)						

Table 1: Tuning Hyper-parameters for Experiment 1

- iii. For each of the four algorithms, run the algorithm over the generated dataset (of 50,000 instances) once and keep track of the number of mistakes, M , that each algorithm makes. Note that a mistake is different than an update (as the margin PERCEPTRON/WINNOW updates for many cases where no mistake is made).
- iv. Plot the cumulative number of mistakes M versus the number of instances ($0 \leq N \leq 50,000$) observed. For each configuration, plot all four curves on the same graph such that they can be easily compared. In each case, the x -axis should be the number of instances observed, N , and the y -axis should be the number of mistakes made M . You should have two graphs (one for each dataset) with four curves on each. Please clearly label your graphs.⁹

(b) [30 points] Learning Curves

The second set of experiments is to construct learning curves. We will begin by setting $l = 10, m = 20$ and vary n such that $40 \leq n \leq 200$ in increments of 40. For each of the 5 different functions, begin by generating 50,000 *clean* instances with the specified $\{l, m, n\}$ configuration.

As in the previous section, fill in Table 2. To run the experiment,

⁸meaning *not* noisy

⁹If you are having memory issues, you can consider plotting the cumulative error every 10 or 100 instances.

Algorithm		PERCEPTRON	PERCEPTRON(γ)	WINNOWER	WINNOWER(γ)
$n = 40$	Params				
	Acc(\mathcal{D}_2)				
	M				
$n = 80$	Params				
	Acc(\mathcal{D}_2)				
	M				
$n = 120$	Params				
	Acc(\mathcal{D}_2)				
	M				
$n = 160$	Params				
	Acc(\mathcal{D}_2)				
	M				
$n = 200$	Params				
	Acc(\mathcal{D}_2)				
	M				

Table 2: Tuning Hyper-parameters for Experiment 2

- i. Present an example to the learning algorithm.
- ii. Again, keep track of the number of mistakes, M , the algorithm makes.

The way we will measure convergence is that we will let the algorithm run until S consecutive examples are presented such that no mistakes are made (i.e, the current hypothesis *survives* for 1000 instances). Note that we are able to do this since we know that the algorithms can learn the target function. Once S instances are encountered, record M at this point and halt. For each algorithm, plot a curve of M (on the y -axis) as a function of n (on the x -axis) such that there are four curves, each determined by the five points associated with each value of n , on a single plot.

Note that you may have to play a bit with the value of S . I would use $S = 1000$ as a good starting point. However, if this is too large such that the algorithm does not halt before $N = 50,000$, you may have to lower this value. Note that you must use the same value of S for all experiments for valid comparisons.

(c) [35 points] **Batch Performance**

In this case, we will run the online algorithms in batch mode to gain understanding of the relative performance in more common scenarios. These experiments should be conducted as follows:

- i. For a given $\{l, m, n\}$ configuration, generate a *noisy* 50,000 instance training set and *clean* 10,000 instance testing set. We will flip each label with probability 0.05 and each attribute with probability 0.001 using `GenerateLMN.java`.

- ii. Optimize the hyper-parameters as in previous sections. Note that technically, since we are running online algorithms in batch mode, T is another hyper-parameter – however, we will just set $T = 20$ as you have probably had enough by now (although if you find something interesting, I would be most receptive).
- iii. Using these hyper-parameters, train the model with the 50,000 training examples and evaluate your model on the testing data. Report the accuracy of each respective learning algorithm.

Repeat this experiment with the following three $\{l, m, n\}$ configurations.

- $l = 10, m = 100, n = 1000$
- $l = 10, m = 500, n = 1000$
- $l = 10, m = 1000, n = 1000$

For each $\{l, m, n\}$ configuration, the same training and testing data should be used for all appropriate comparisons. This can either be accomplished by using files or setting the random seed appropriately. You should generate results similar to Table 3.

Algorithm	$m = 100$			$m = 500$			$m = 1000$		
	Params	Acc(\mathcal{D}_2)	Acc(Test)	Params	Acc(\mathcal{D}_2)	Acc(Test)	Params	Acc(\mathcal{D}_2)	Acc(Test)
PERC.									
PERC.(γ)									
WINNOW									
WINNOW(γ)									

Table 3: Tuning Hyper-parameters for Experiment 3

What to submit

- A detailed, yet concise report. In addition to the requested information, summarize the findings. Discuss differences in performance of the algorithms and attempt to explain why. Were the results consistent across all experiments? If possible, try to make the report interesting. Note that these experiments were (heavily) influenced by [Kivinen, Warmuth, and Auer; The Perceptron Algorithm versus Winnow: Linear versus Logarithmic Mistake Bounds When Few Input Variables are Relevant, Artificial Intelligence, pg 325-343, 1997] if you would like to do some reading.
- Two plots from the first experiment and one plot from the second experiment. Please clearly label these.
- One table for the third experiment.
- Tables associated with each experiment.

- Source code. Submit all relevant files via **github**. You are free to use the programming language of your choice. However, please include a **README** file that provides instructions on compiling and running your program. Of course, a shell script would be greatly appreciated.
- Use your CCIS github repository to submit all relevant files. You are free to use the programming language of your choice, but please attempt to conform to the instructions above. To be safe, try submitting something **before** the assignment deadline.
- The code you submit must be your own. If you find/use information about specific algorithms from the Web, etc., be sure to cite the source(s) clearly in your source code.

Appendix: Linear Programming

In this appendix, we will walk through a simple linear programming example.¹⁰ If you want to read more on the topic, a good reference is *Linear Programming: Foundations and Extensions* by Vanderbei. Some classic texts include *Linear Programming* by Chvatal; and *Combinatorial Optimization: Algorithms and Complexity* by Papadimitrou and Steiglitz (in particular, the beginning of chapter 2 may be helpful). A widely available (albeit incomplete) reference is *Introduction to Algorithms* by Cormen, Leiserson, Rivest, and Stein.

Example: Consider the following problem.¹¹ You are given a choice of three foods, namely eggs at a cost of \$0.10 an egg, pasta at a cost of \$0.05 a bowl, and yogurt at a cost of \$0.25 a cup. An egg (t_1) provides 3 portions of protein, 1 portion of carbohydrates, and 2 portions of fat. A bowl of pasta (t_2) provides 1 portion of protein, 3 portions of carbohydrates, and no portions of fat. A cup of yogurt (t_3) provides 2 portions of protein, 2 portions of carbohydrates, and 1 portion of fat. You are required to consume at least 7 portions of protein and 9 portions of carbohydrates per day and are not allowed to consume more than 4 portions of fat. In addition, you obviously may not consume a negative amount of any food. The objective now is to find the cheapest combination of foods that still meet your daily nutritional requirements.

¹⁰Note that SVM uses *quadratic* programming (QP) which means that the objective function includes a quadratic term.

¹¹The problem closely resembles an instantiation of the *diet problem* by G. J. Stigler, *The Cost of Subsistence*, 1945.

This can be written as the following linear program:

$$z = 0.1t_1 + 0.05t_2 + 0.25t_3 \rightarrow \min \quad (6)$$

$$3t_1 + t_2 + 2t_3 \geq 7 \quad (7)$$

$$t_1 + 3t_2 + 2t_3 \geq 9 \quad (8)$$

$$-2t_1 - t_3 \geq -4 \quad (9)$$

$$t_i \geq 0 \quad \forall i \quad (10)$$

Note that inequality (9) of the LP is equivalent to $2t_1 + t_3 < 4$. This corresponds to:

$$\mathbf{A} = \begin{pmatrix} 3 & 1 & 2 \\ 1 & 3 & 2 \\ -2 & 0 & -1 \end{pmatrix} \quad \vec{b} = \begin{pmatrix} 7 \\ 9 \\ -4 \end{pmatrix} \quad \vec{c} = \begin{pmatrix} 0.1 \\ 0.05 \\ 0.25 \end{pmatrix} \quad \vec{t} = \begin{pmatrix} t_1 \\ t_2 \\ t_3 \end{pmatrix}$$

To solve this program using Matlab:

```
c = [0.1; 0.05; 0.25];
A = [3 1 2; 1 3 2; -2 0 -1];
b = [7; 9; -4];
lowerBound = zeros(3, 1); % this constrains t >= 0
[t, z] = linprog(c, -A, -b, [], [], lowerBound)
```

The results of this linear program show that to meet your nutritional requirements at a minimum cost, you should eat 1.5 eggs, 2.5 bowls of pasta, and no cups of yogurt, and the cost for such a diet is \$0.275.¹²

¹²Note that this is not intended to be actual nutritional advice. Your mileage may vary.