

# CS489/698: Introduction to Machine Learning

## Homework 1

Due: 11:59 pm, September 26, 2017, submit on LEARN.

Include your name, student number and session!

Submit your writeup in pdf and all source code in a zip file (with proper documentation). Write a script for each programming exercise so that the TAs can easily run and verify your results. Make sure your code runs!

[Text in square brackets are hints that can be ignored.]

### Exercise 1: Perceptron and Winnow (50 pts)

**Convention:** All algebraic operations, when applied to a vector or matrix, are understood to be element-wise (unless otherwise stated).

**Algorithm 1:** The perceptron algorithm.

---

**Input:**  $X \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \{-1, 1\}^n$ ,  $\mathbf{w} = \mathbf{0}_d$ ,  $b = 0$ ,  $\text{max\_pass} \in \mathbb{N}$   
**Output:**  $\mathbf{w}, b, \text{mistake}$

```

1 for  $t = 1, 2, \dots, \text{max\_pass}$  do
2    $\text{mistake}(t) \leftarrow 0$ 
3   for  $i = 1, 2, \dots, n$  do
4     if  $y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \leq 0$  then
5        $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$  //  $\mathbf{x}_i$  is the  $i$ -th row of  $X$ 
6        $b \leftarrow b + y_i$ 
7        $\text{mistake}(t) \leftarrow \text{mistake}(t) + 1$ 

```

---

**Algorithm 2:** The winnow algorithm.

---

**Input:**  $X \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \{-1, 1\}^n$ ,  $\mathbf{w} = \frac{1}{d+1} \mathbf{1}_d$ ,  $b = \frac{1}{d+1}$ , step size  $\eta > 0$ ,  $\text{max\_pass} \in \mathbb{N}$   
**Output:**  $\mathbf{w}, b, \text{mistake}$

```

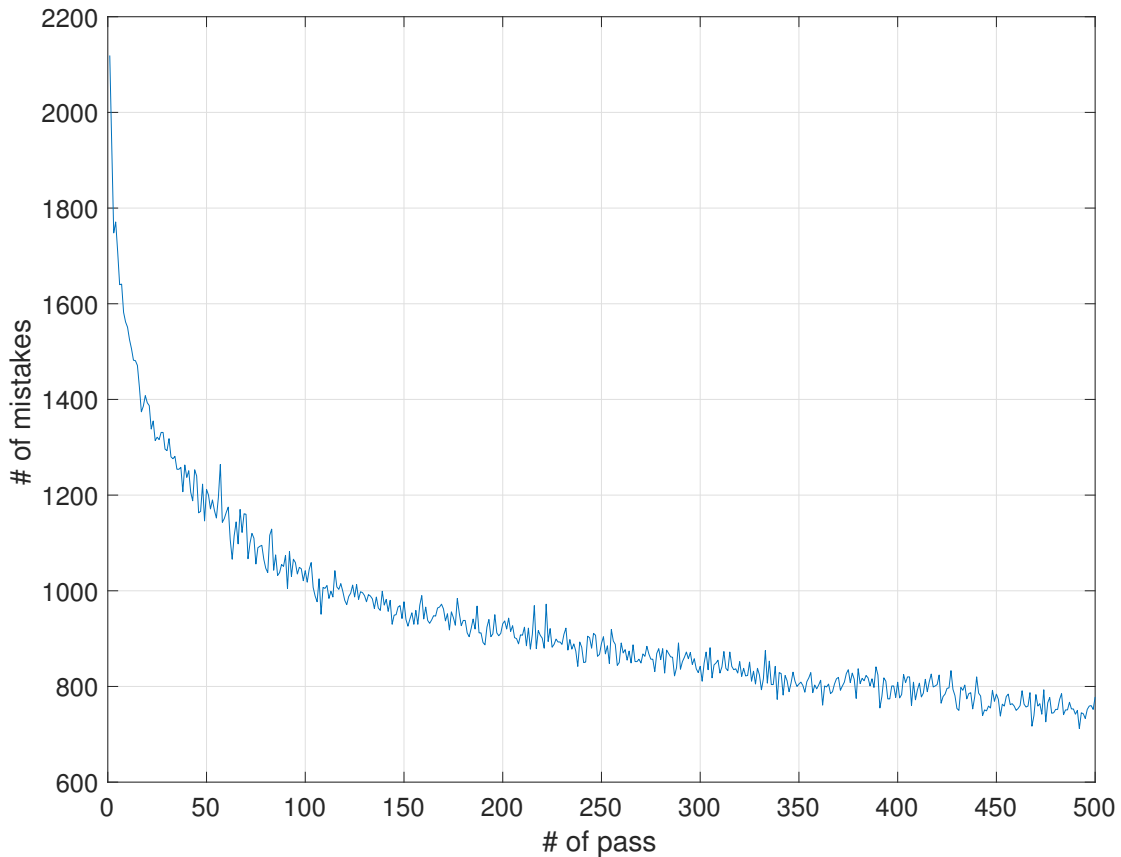
1 for  $t = 1, 2, \dots, \text{max\_pass}$  do
2    $\text{mistake}(t) \leftarrow 0$ 
3   for  $i = 1, 2, \dots, n$  do
4     if  $y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \leq 0$  then
5        $\mathbf{w} \leftarrow \mathbf{w} \odot \exp(\eta y_i \mathbf{x}_i)$  //  $\odot$  is the element-wise product
6        $b \leftarrow b \exp(\eta y_i)$  // element-wise exponential
7        $s \leftarrow b + \sum_{i=1}^d w_i$  // normalize
8        $\mathbf{w} \leftarrow \mathbf{w} / s$ 
9        $b \leftarrow b / s$ 
10     $\text{mistake}(t) \leftarrow \text{mistake}(t) + 1$ 

```

---

- (10 pts) Implement the perceptron in Algorithm 1. Your implementation should take input as  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \{-1, 1\}^n$ , an initialization of the hyperplane parameters  $\mathbf{w} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$ , and the maximum number of passes of the training set [suggested  $\text{max\_pass} = 500$ ]. Run your perceptron algorithm on the [spambase](#) dataset (available on [course website](#)), and plot the number of mistakes ( $y$ -axis) w.r.t. the number of passes ( $x$ -axis).

**Ans:** The number of mistakes should drop as the number of passes of the training data increases.



2. (5 pts) Run your implementation on **spambase** again, but this time moving the updates (line 5 and line 6 in Algorithm 1) outside of the IF-clause, i.e., we update the weight vectors even when perceptron predicts correctly. We count the number of mistakes as before, i.e., only when perceptron makes a mistake. Plot again the number of mistakes w.r.t. the number of passes.

**Ans:** The number of mistakes stays as a constant 2788 (the number of negative instances).

3. (5 pts) Prove the following claim: If there exist  $\mathbf{w}^*$  and  $b^*$  such that

$$\begin{cases} \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* \geq 0, & \text{if } y_i = 1 \\ \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* < 0, & \text{if } y_i = -1 \end{cases}, \quad (1)$$

then there exist  $\mathbf{w}^*$  and  $b^*$  such that

$$\begin{cases} \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* > 0, & \text{if } y_i = 1 \\ \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* < 0, & \text{if } y_i = -1 \end{cases}. \quad (2)$$

This confirms again that we can be indifferent about the value of  $\text{sign}(0)$ .

**Ans:** Since we only have a finite training set, let

$$\epsilon = \min_{i: y_i = -1} -\langle \mathbf{x}_i, \mathbf{w}^* \rangle.$$

Then, from our assumption (1) we know  $b^* < \epsilon$ . Let  $\alpha = (\epsilon - b^*)/2 > 0$ . We claim that  $\mathbf{w}^* = \mathbf{w}^*$ ,  $b^* = b^* + \alpha$  satisfy (2). Indeed,

$$\begin{cases} \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* + \alpha \geq 0 + \alpha = \alpha > 0, & \text{if } y_i = 1 \\ \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* = \langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^* + \alpha < b^* - \epsilon + \alpha = -\alpha < 0, & \text{if } y_i = -1 \end{cases}$$

4. (10 pts) Implement winnow (Algorithm 2), the multiplicative version of perceptron. Tune and report a few step sizes  $\eta$  and plot the number of mistakes winnow makes on the **spambase** dataset (w.r.t. the number of passes for each step size  $\eta$  you tried). Explain the effect of the step size. [A good step size should be inversely proportional to the maximum absolute value in  $X$ . You may want to normalize  $X$  so that its maximum absolute value is, say, 1.]

Ans: Winnow will mistake all (2788) negative instances, since both  $X$  and  $\mathbf{w}$  and  $b$  are nonnegative. The step size has no effect here.

5. (10 pts) It is known that if there exist **nonnegative**  $\mathbf{w}^*$  and  $b^*$  such that  $y_i(\langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^*) > 0$  for all  $i$ , then winnow converges after at most  $O(\frac{2\log(d+1)}{\gamma^2})$  mistakes, where

$$\gamma = \max_{\substack{\mathbf{w} \geq \mathbf{0}, b + \sum_i w_i = 1}} \min_i \frac{y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b)}{\max_i \max_j \max\{|X_{ij}|, 1\}}.$$

Recall that in perceptron, there is no nonnegativity assumption. Find a simple transformation of the data  $(X, \mathbf{y})$  such that if there exist  $\mathbf{w}^*$  and  $b^*$  so that  $y_i(\langle \mathbf{x}_i, \mathbf{w}^* \rangle + b^*) > 0$  for all  $i$ , then there exist nonnegative  $\mathbf{w}^* \geq \mathbf{0}$  and  $b^* \geq 0$  on the **transformed** data that separate the two classes. Thus, the additional nonnegativity assumption in winnow is really not a big deal. Plot the number of mistakes that winnow makes on the transformed **spambase** w.r.t. the number of passes (with a similar step size as in the previous exercise). [Each real number  $w$  can be written as the difference of two nonnegative numbers  $w_+ := \max\{w, 0\}$  and  $w_- := \max\{-w, 0\}$ . Consider duplicating each point  $\mathbf{x}$  with some version of itself.]

Ans: As hinted, any real vector  $\mathbf{w}$  can be written as the difference of two nonnegative real vectors:

$$\mathbf{w} = \mathbf{w}_+ - \mathbf{w}_-, \text{ where } \mathbf{w}_+ = \max\{\mathbf{w}, \mathbf{0}\}, \mathbf{w}_- = \max\{-\mathbf{w}, \mathbf{0}\}.$$

As a result,

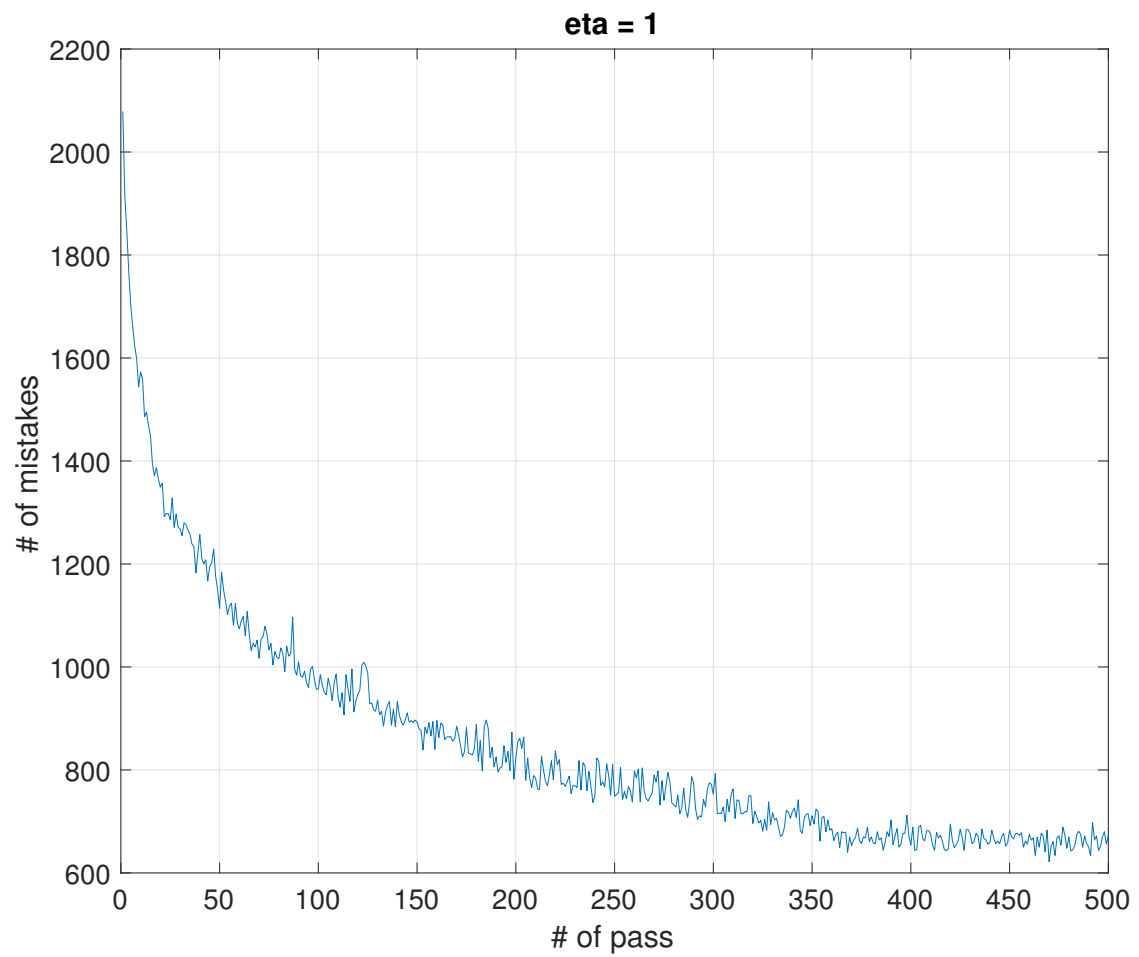
$$\langle \mathbf{x}, \mathbf{w}^* \rangle + b^* = \langle \mathbf{x}, \mathbf{w}_+^* - \mathbf{w}_-^* \rangle + b_+^* - b_-^* = \left\langle \begin{bmatrix} \mathbf{x} \\ 1 \\ -\mathbf{x} \\ -1 \end{bmatrix}, \begin{bmatrix} \mathbf{w}_+^* \\ b_+^* \\ \mathbf{w}_-^* \\ b_-^* \end{bmatrix} \right\rangle.$$

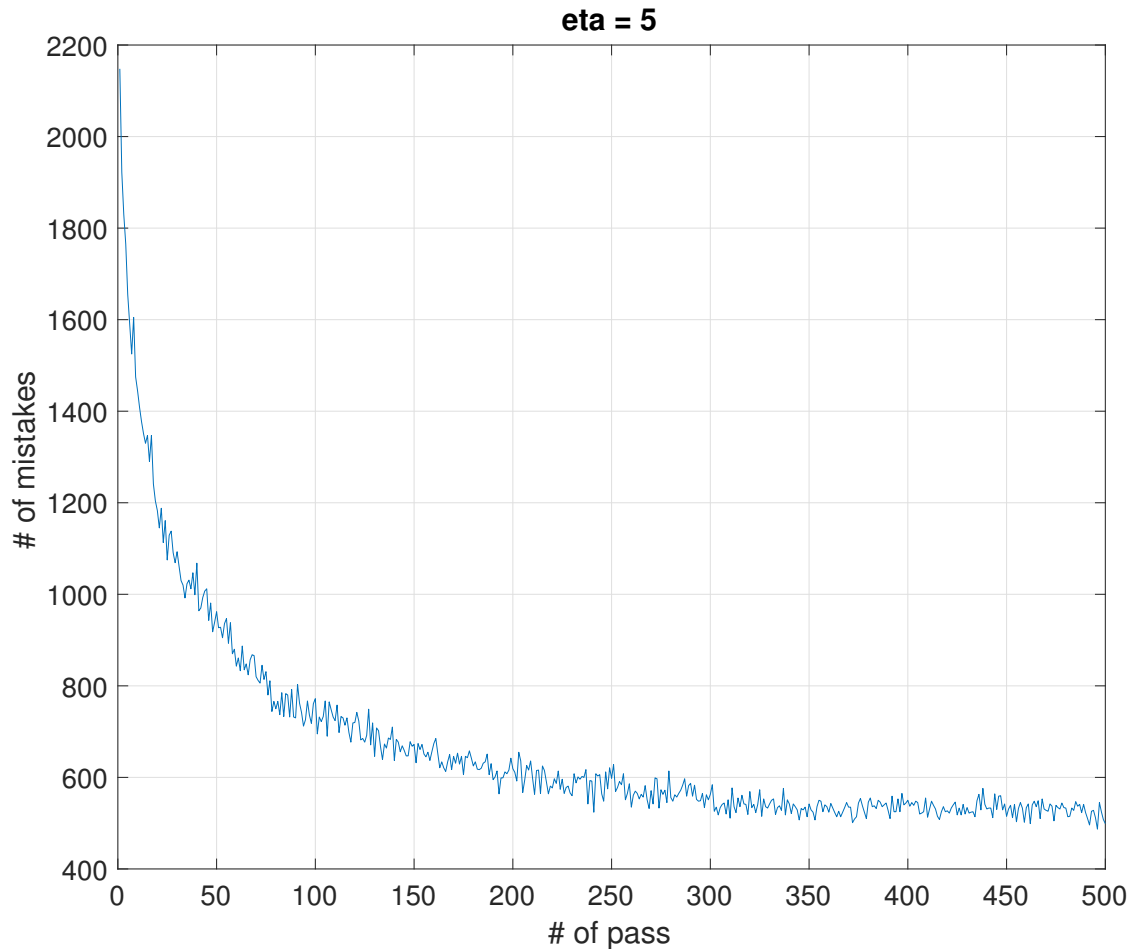
Therefore, if we define

$$\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \\ -\mathbf{x} \\ -1 \end{bmatrix}, \mathbf{w}^* = \begin{bmatrix} \mathbf{w}_+^* \\ b_+^* \\ \mathbf{w}_-^* \\ b_-^* \end{bmatrix}, b^* = 0,$$

then  $\mathbf{w}^*$  and  $b^*$  are nonnegative, and they separate the two classes  $(\tilde{\mathbf{x}}, y)$ .

In our implementation of winnow, we normalize the data so that  $\|X\|_\infty = 1$ .





6. (10 pts) Compare what you learned about perceptron with the above result of winnow. What are the respective pros and cons of the two algorithms? [Recall that the margin  $\gamma$  in perceptron can be written as:  $\gamma = \max_{b^2 + \sum_i w_i^2 = 1} \min_i \frac{y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b)}{\max_j \sqrt{1 + \sum_j |X_{ij}|^2}}$ , with which the perceptron converges after at most  $O(1/\gamma^2)$  mistakes.] Modify the `spambase` dataset by adding 100 irrelevant features (say, random numbers from the uniform distribution on  $[-1, 1]$ ) [`rand` in `Matlab` or `Julia`, and `numpy.random.uniform` in `python`] and run both perceptron and winnow (you may want to use the variant in Ex 1.5). Plot the number of mistakes that each algorithm makes this time (again w.r.t. the number of passes).

**Ans:** The number of mistakes made by perceptron and winnow are in general not comparable, because they depend on different definitions of the margin  $\gamma$ . First, both algorithms are online, and they make fewer mistakes if the data is “more” separated. However, we note that the perceptron bound does not have any explicit dependence on the dimension while the winnow bound scales with  $\log(d)$ . To apply winnow, we usually need to apply the “balancing” trick in Ex 1.5, hence doubling the dimension of the problem, which is another slight disadvantage. On the other hand, if the true separating hyperplane depends only on a few number of features, then the winnow bound is much more appealing (because we are taking the maximum instead of the  $\ell_2$  norm of all features.)

### Exercise 2: Linear Regression and Regularization (50 pts)

**Convention:** The Matlab notation  $X_{:,j}$  means the  $j$ -th column of  $X$  while  $X_{i,:}$  means the  $i$ -th row of  $X$ . We

use `argmin` to denote the set of minimizers of a minimization problem.

---

**Algorithm 3:** Alternating minimization for Lasso.

---

**Input:**  $X \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{w} = \mathbf{0}_d$ ,  $\lambda \geq 0$   
**Output:**  $\mathbf{w}$

```

1 repeat
2   for  $j = 1, \dots, d$  do
3      $w_j \leftarrow \underset{z \in \mathbb{R}}{\operatorname{argmin}} \frac{1}{2} \|X_{:j}z + \sum_{k \neq j} X_{:k}w_k - \mathbf{y}\|_2^2 + \lambda|z|$  // fix all  $\mathbf{w}$  but optimize  $w_j$  only
4 until convergence
```

---

- (10 pts) Implement ridge regression that we discussed in class:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2. \quad (3)$$

Your implementation should take input as  $X \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \mathbb{R}^n$ ,  $\lambda \geq 0$ , and maybe an initializer for  $\mathbf{w} \in \mathbb{R}^d$ . [To solve a linear system  $A\mathbf{x} = \mathbf{b}$ , use `A\b` in Matlab or Julia, and `numpy.linalg.solve` in python.] Test your algorithm on the Boston [housing](#) dataset (to predict the median house price, i.e.,  $y$ ). Train and test splits are provided on [course website](#). Find the best  $\lambda^*$  in the range  $[0, 100]$  with increment 10 using 10-fold cross validation. Report the mean square error on the training set ( $\frac{1}{n} \|X\mathbf{w} - \mathbf{y}\|_2^2$ ), validation set (averaged over 10-fold) and test set [different normalization constant  $n$ ] for each candidate  $\lambda$ . Report the percentage of nonzeros in  $\mathbf{w}$  (for each  $\lambda$ ).

**Ans:** Simply solve the linear system  $(X^\top X + \lambda I)\mathbf{w} = X^\top \mathbf{y}$ .

- (10 pts) Randomly choose a sample pair  $(\mathbf{x}, y)$  from your training set. Multiply  $\mathbf{x}$  by  $10^6$  and/or  $y$  by  $10^3$  and run ridge regression again (with the same cross-validation procedure to choose  $\lambda$  as above). Report the mean square error on the training set ( $\frac{1}{n} \|X\mathbf{w} - \mathbf{y}\|_2^2$ ), validation set (averaged over 10-fold) and test set [different normalization constant  $n$ ] for each candidate  $\lambda$ .

**Ans:** Here the errors are supposedly much larger, due to the one and only extreme outlier.

- (10 pts) Add 1000 irrelevant columns to  $X$  (both training and test splits), with each entry an *iid* sample from the standard normal distribution (`randn` in Matlab or Julia, and `numpy.random.standard_normal` in python). Run ridge regression again (with the same cross-validation procedure to choose  $\lambda$  as above). Report the mean square error on the training set ( $\frac{1}{n} \|X\mathbf{w} - \mathbf{y}\|_2^2$ ), validation set (averaged over 10-fold) and test set [different normalization constant  $n$ ] for each candidate  $\lambda$ . Report the percentage of nonzeros in the *last 1000 entries* in  $\mathbf{w}$  (i.e., weights corresponding to the added irrelevant features).

**Ans:** Most (and perhaps all) entries of the weight vector should be nonzero.

- (10 pts) The Lasso replaces the (squared) 2-norm penalty in ridge regression with the 1-norm:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1. \quad (4)$$

Implement the alternating minimization algorithm for Lasso (Algorithm 3). You might find the following fact useful:

$$\operatorname{sign}(w) \cdot \max\{0, |w| - \lambda\} = \underset{z \in \mathbb{R}}{\operatorname{argmin}} \frac{1}{2} (z - w)^2 + \lambda|z|, \quad (5)$$

which is known as the soft-thresholding operator. You should try to perform step 3 of Algorithm 3 in  $O(n)$  time and space. Stop the algorithm when the change of  $\mathbf{w}$  drops below some tolerance `tol`, say  $10^{-3}$ . [Any idea to make your implementation even more efficient?] Report the mean square error of Lasso on the training set ( $\frac{1}{n} \|X\mathbf{w} - \mathbf{y}\|_2^2$ ), validation set (averaged over 10-fold) and test set [different normalization constant  $n$ ] for each candidate  $\lambda$ . Report the percentage of nonzeros in  $\mathbf{w}$  (for each  $\lambda$ ). [A rule of thumb here is that  $\lambda = O(\sqrt{\log(d)/n})$ .]

**Ans:** Here the key is to maintain the residual vector explicitly. We start with  $\mathbf{r} = X\mathbf{w} - \mathbf{y} = -\mathbf{y}$  (because we initialize with  $\mathbf{w} = \mathbf{0}$ ). In each iteration we do the following:

$$\mathbf{r} \leftarrow \mathbf{r} - X_{:j}w_j \quad (6)$$

$$w_j \leftarrow \underset{z}{\operatorname{argmin}} \frac{1}{2}\|X_{:j}\|_2^2 z^2 + \langle X_{:j}, \mathbf{r} \rangle z + \lambda|z| \quad (7)$$

$$= -\operatorname{sign}(\langle X_{:j}, \mathbf{r} \rangle) \max\{0, |\langle X_{:j}, \mathbf{r} \rangle| - \lambda\} / \|X_{:j}\|_2^2 \quad (8)$$

$$\mathbf{r} \leftarrow \mathbf{r} + X_{:j}w_j \quad (9)$$

[A good practice is to randomly perturb all indices and then go through the indices one by one sequentially.]

Clearly, the above iteration costs  $O(n)$ . We can further improve the efficiency by pre-compute  $\|X_{:j}\|_2$  for all  $j$ .

5. (10 pts) Run Lasso on the housing dataset that you modified in Ex2.3, with  $\lambda$  cross-validated as before. Report the mean square error on the training set ( $\frac{1}{n}\|X\mathbf{w} - \mathbf{y}\|_2^2$ ), validation set (averaged over 10-fold) and test set [different normalization constant  $n$ ] for each candidate  $\lambda$ . Report the percentage of nonzeros in the *last 1000 entries* in  $\mathbf{w}$  (i.e., weights corresponding to the added irrelevant features). Comparing with your results in Ex2.3, what can you conclude? [Mean square error, time complexity, sparsity, etc. ]

**Ans:** Lasso should yield a much sparser weight, hopefully with many (or even most) of the last 1000 entries being 0.