

# Machine Learning and Computational Statistics, Spring 2016

## Homework 2: Lasso

**Due: Tuesday, February 16, 2016, at 6pm (Submit via NYU Classes)**

**Instructions:** Your answers to the questions below, including plots and mathematical work, should be submitted as a single PDF file. You may include your code inline or submit it as a separate file. You may either scan hand-written work or, preferably, write your answers using software that typesets mathematics (e.g. L<sup>A</sup>T<sub>E</sub>X, L<sup>A</sup>T<sub>E</sub>X, or MathJax via iPython).

### 1 Preliminaries

#### 1.1 Dataset construction

Start by creating a design matrix for regression with  $m = 150$  examples, each of dimension  $d = 75$ . We will choose a true weight vector  $\theta$  that has only a few non-zero components:

1. Let  $X \in \mathbf{R}^{m \times d}$  be the “design matrix,” where the  $i$ ’th row of  $X$  is  $x_i \in \mathbf{R}^d$ . Construct a random design matrix  $X$  using `numpy.random.rand()` function.
2. Construct a true weight vector  $\theta \in \mathbf{R}^{d \times 1}$  as follows: Set the first 10 components of  $\theta$  to 10 or -10 arbitrarily, and all the other components to zero.
3. Let  $y = (y_1, \dots, y_m)^T \in \mathbf{R}^{m \times 1}$  be the response. Construct the vector  $y = X\theta + \epsilon$ , where  $\epsilon$  is an  $m \times 1$  random noise vector generated using `numpy.random.randn()` with mean 0 and standard deviation 0.1.
4. Split the dataset by taking the first 80 points for training, the next 20 points for validation, and the last 50 points for testing.

Note that we are not adding an extra feature for the bias in this problem. By construction, the true model does not have a bias term.

#### 1.2 Experiments with Ridge Regression

By construction, we know that our dataset admits a sparse solution. Here, we want to evaluate the performance of ridge regression (i.e.  $\ell_2$ -regularized linear regression) on this dataset.

1. Run ridge regression on this dataset. Choose the  $\lambda$  that minimizes the square loss on the validation set. For the chosen  $\lambda$ , examine the model coefficients. Report on how many components with true value 0 have been estimated to be non-zero, and vice-versa (don’t worry if they are

all nonzero). Now choose a small threshold (say  $10^{-3}$  or smaller), count anything with magnitude smaller than the threshold as zero, and repeat the report. (For running ridge regression, you may either use your code from HW1, or you may use `scipy.optimize.minimize` (see the demo code provided for guidance). For debugging purposes, you are welcome, even encouraged, to compare your results to what you get from `sklearn.linear_model.Ridge`.)

## 2 Coordinate Descent for Lasso (a.k.a. The Shooting algorithm)

The Lasso optimization problem can be formulated as

$$\hat{w} = \arg \min_{w \in \mathbf{R}^d} \sum_{i=1}^m (h_w(x_i) - y_i)^2 + \lambda \|w\|_1,$$

where  $h_w(x) = w^T x$ , and  $\|w\|_1 = \sum_{i=1}^d |w_i|$ . Since the  $\ell_1$ -regularization term in the objective function is non-differentiable, it's not clear how gradient descent or SGD could be used to solve this optimization problem.

Another approach to solving optimization problems is coordinate descent, in which at each step we optimize over one component of the unknown parameter vector, fixing all other components. The descent path so obtained is a sequence of steps each of which is parallel to a coordinate axis in  $\mathbf{R}^d$ , hence the name. It turns out that for the Lasso optimization problem, we can find a closed form solution for optimization over a single component fixing all other components. This gives us the following algorithm:

---

**Algorithm 13.1:** Coordinate descent for lasso (aka shooting algorithm)

---

```

1 Initialize  $\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$ ;
2 repeat
3   for  $j = 1, \dots, D$  do
4      $a_j = 2 \sum_{i=1}^n x_{ij}^2$ ;
5      $c_j = 2 \sum_{i=1}^n x_{ij} (y_i - \mathbf{w}^T \mathbf{x}_i + w_j x_{ij})$ ;
6      $w_j = \text{soft}(\frac{c_j}{a_j}, \frac{\lambda}{a_j})$ ;
7 until converged;
```

---

(Source: Murphy, Kevin P. Machine learning: a probabilistic perspective. MIT press, 2012.)

The “soft thresholding” function is defined as

$$\text{soft}(a, \delta) = \text{sign}(a) (|a| - \delta)_+.$$

NOTE: Algorithm 13.1 does not account for the case that  $a_j = c_j = 0$ , which occurs when the  $j$ th column of  $X$  is the identically 0. One can either eliminate the column (as it cannot possibly help the solution), or you can set  $w_j = 0$  in that case, as that is, in fact, the coordinate minimizer in that case. Note that Murphy is suggesting to initialize the optimization with the ridge regression solution, though this is not necessary.

The solution should have a sparsity pattern that is similar to the ground truth. Estimators that preserve the sparsity pattern (with enough training data) are said to be “**sparsistent**”<sup>1</sup> (sparse + consistent). Formally, an estimator  $\hat{\beta}$  of parameter  $\beta$  is said to be consistent if the estimator  $\hat{\beta}$  converges to the true value  $\beta$  in probability as our sample size goes to infinity. Analogously, if we define the support of a vector  $\beta$  as the indices with non-zero components, i.e.  $\text{Supp}(\beta) = \{j \mid \beta_j \neq 0\}$ , then an estimator  $\hat{\beta}$  is said to be sparsistent if as the number of samples becomes large, the support of  $\hat{\beta}$  converges to the support of  $\beta$ , or  $\lim_{m \rightarrow \infty} P[\text{Supp}(\hat{\beta}_m) = \text{Supp}(\beta)] = 1$ .

There are a few tricks that can make selecting the hyperparameter  $\lambda$  easier and faster. First, you can show that for any  $\lambda > 2\|X^T(y - \bar{y})\|_\infty$ , the estimated weight vector  $\hat{w}$  is entirely zero, where  $\bar{y}$  is the mean of values in the vector  $y$ , and  $\|\cdot\|_\infty$  is the infinity norm (or supremum norm), which is the maximum absolute value of any component of the vector. Thus we need to search for an optimal  $\lambda$  in  $[0, \lambda_{\max}]$ , where  $\lambda_{\max} = 2\|X^T(y - \bar{y})\|_\infty$ . (Note: This expression for  $\lambda_{\max}$  assumes we have an unregularized bias term in our model. That is, our decision functions are  $h_{w,b}(x) = w^T x + b$ . For the experiments, you can exclude the bias term, in which case  $\lambda_{\max} = 2\|X^T y\|_\infty$ .)

Second, we can make use of the fact that when  $\lambda$  and  $\lambda'$  are close, so are the corresponding solutions  $\hat{w}(\lambda)$  and  $\hat{w}(\lambda')$ . Start by finding  $\hat{w}(\lambda_{\max})$  and initialize the optimization at  $w = 0$ . Next,  $\lambda$  is reduced (e.g. by a constant factor), and the optimization problem is solved using the previous optimal point as the starting point. This is called **warm starting** the optimization. The entire technique of computing a set of solutions for a chain of nearby  $\lambda$ 's is called a **continuation** or **homotopy method**. In the context of finding a good regularization hyperparameter, it may be referred to as a **regularization path** approach. (Lots of names for this!)

## 2.1 Experiments with the Shooting Algorithm

1. Write a function that computes the Lasso solution for a given  $\lambda$  using the shooting algorithm described above. This function should take a starting point for the optimization as a parameter. Run it on the dataset constructed in (1.1), and select the  $\lambda$  that minimizes the square error on the validation set. Report the optimal value of  $\lambda$  found, and the corresponding test error. Plot the validation error vs  $\lambda$ . [Don't use the homotopy method in this part, as we want to measure the speed improvement of homotopy methods in part 3. Also, no need to vectorize the calculations until part 4, where again we'll compare the speedup. In any case, having two different implementations of the same thing is a good way to check your work.]
2. Analyze the sparsity of your solution, reporting how many components with true value zero have been estimated to be non-zero, and vice-versa.
3. Implement the homotopy method described above. Compare the runtime for computing the full regularization path (for the same set of  $\lambda$ 's you tried in the first question above) using the homotopy method compared to the basic shooting algorithm.
4. The algorithm as described above is not ready for a large dataset (at least if it has been implemented in basic Python) because of the implied loop over the dataset (i.e. where we sum over the training set). By using matrix and vector operations, we can eliminate the loops. This is called “vectorization” and can lead to dramatic speedup in languages such as Python, Matlab, and R. Derive matrix expressions for computing  $a_j$  and  $c_j$ . (Hint: A matlab version of this

---

<sup>1</sup>Li, Yen-Huan, et al. “Sparsistency of  $l_1$ -Regularized  $M$ -Estimators.”

vectorized method can be found here: [http://pmtk3.googlecode.com/svn-history/r1393/trunk/toolbox/Variable\\_selection/lassoExtra/LassoShooting.m](http://pmtk3.googlecode.com/svn-history/r1393/trunk/toolbox/Variable_selection/lassoExtra/LassoShooting.m)). Implement the matrix expressions and measure the speedup to compute the regularization path.

## 2.2 Deriving the Coordinate Minimizer for Lasso

This problem is to derive the expressions for the coordinate minimizers used in the Shooting algorithm. This is often presented using subgradients (e.g. <http://davidrosenberg.github.io/ml2015/docs/2.Lab.subgradient-descent.pdf#page=15>), but here we will walk you through a bare hands approach (which is essentially equivalent).

In each step of the shooting algorithm, we would like to find the  $w_j$  minimizing

$$\begin{aligned} f(w_j) &= \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda |w|_1 \\ &= \sum_{i=1}^n \left[ w_j x_{ij} + \sum_{k \neq j} w_k x_{ik} - y_i \right]^2 + \lambda |w_j| + \lambda \sum_{k \neq j} |w_k|, \end{aligned}$$

where we've written  $x_{ij}$  for the  $j$ th entry of the vector  $x_i$ . This function is strictly convex in  $w_j$ , and thus it has a unique minimum. Furthermore, the only thing keeping  $f$  from being differentiable is the term with  $|w_j|$ . So  $f$  is differentiable everywhere except  $w_j = 0$ . We'll break this problem into 3 cases:  $w_j > 0$ ,  $w_j < 0$ , and  $w_j = 0$ . In the first two cases, we can simply differentiate  $f$  w.r.t.  $w_j$  to get optimality conditions. For the last case, we'll use the following more bare-hands characterization: Since  $f$  is convex, 0 is a minimizer of  $f$  iff

$$\lim_{\varepsilon \downarrow 0} \frac{f(\varepsilon) - f(0)}{\varepsilon} \geq 0 \quad \text{and} \quad \lim_{\varepsilon \downarrow 0} \frac{f(-\varepsilon) - f(0)}{\varepsilon} \geq 0.$$

This is a special case of the optimality conditions described in this slide <http://davidrosenberg.github.io/ml2015/docs/5.Lab.misc.pdf#page=10>, where now the “direction”  $v$  is simply taken to be the scalars 1 and  $-1$ , respectively.

1. First let's get a trivial case out of the way. If  $x_{ij} = 0$  for  $i = 1, \dots, n$ , what is the coordinate minimizer  $w_j$ ? In the remaining questions below, you may assume that  $\sum_{i=1}^n x_{ij}^2 > 0$ .
2. Give an expression for the derivative  $f(w_j)$  for  $w_j \neq 0$ . It will be convenient to write your expression in terms of the following definitions:

$$\begin{aligned} \text{sign}(w_j) &:= \begin{cases} 1 & w_j > 0 \\ 0 & w_j = 0 \\ -1 & w_j < 0 \end{cases} \\ a_j &:= 2 \sum_{i=1}^n x_{ij}^2 \\ c_j &:= 2 \sum_{i=1}^n x_{ij} \left( y_i - \sum_{k \neq j} w_k x_{ik} \right). \end{aligned}$$

3. If  $w_j > 0$  and minimizes  $f$ , then show that  $w_j = -\frac{1}{a_j}(\lambda - c_j)$ . Similarly, if  $w_j < 0$  and minimizes  $f$ , show that  $w_j = \frac{1}{a_j}(\lambda + c_j)$ . Give conditions on  $c_j$  that imply the minimizer  $w_j > 0$  and  $w_j < 0$ , respectively.
4. Derive expressions for the two one-sided derivatives at  $f(0)$ , and show that  $c_j \in [-\lambda, \lambda]$  implies that  $w_j = 0$  is a minimizer.
5. Conclude that the minimizer is given by

$$w_j = \begin{cases} \frac{1}{a_j}(c_j - \lambda) & c_j > \lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ \frac{1}{a_j}(c_j + \lambda) & c_j < -\lambda \end{cases}$$

and show that this is equivalent to the expression given in 2.

### 3 Lasso Properties

#### 3.1 Deriving $\lambda_{\max}$

In this problem we will derive an expression for  $\lambda_{\max}$ . For the first three parts, use the Lasso objective function excluding the bias term i.e,  $L(w) = \|Xw - y\|_2^2 + \lambda \|w\|_1$ . Show that for any  $\lambda \geq 2\|X^T y\|_\infty$ , the estimated weight vector  $\hat{w}$  is entirely zero, where  $\|\cdot\|_\infty$  is the infinity norm (or supremum norm), which is the maximum absolute value of any component of the vector.

1. The one-sided directional derivative of  $f(x)$  at  $x$  in the direction  $v$  is defined as:

$$f'(x; v) = \lim_{h \downarrow 0} \frac{f(x + hv) - f(x)}{h}$$

Compute  $L'(0; v)$ . That is, compute the one-sided directional derivative of  $L(w)$  at  $w = 0$  in the direction  $v$ . [Hint: the result should be in terms of  $X, y, \lambda$ , and  $v$ .]

2. Since the Lasso objective is convex, for  $w^*$  to be a minimizer of  $L(w)$  we must have that the directional derivative  $L'(w^*; v) \geq 0$  for all  $v$ . Starting from the condition  $L'(0; v) \geq 0$ , rearrange terms to get a lower bounds on  $\lambda$ . [Hint: this should be in terms of  $X, y$ , and  $v$ .]
3. Since our lower bounds on  $\lambda$  hold for all  $v$ , we want to compute the maximum lower bound. Compute the maximum lower bound of  $\lambda$  by maximizing the expression over  $v$ . Show that this expression is equivalent to  $\lambda_{\max} = 2\|X^T y\|_\infty$ .
4. [Optional] Show that for  $L(w) = \|Xw + b - y\|_2^2 + \lambda \|w\|_1$ ,  $\lambda_{\max} = 2\|X^T(y - \bar{y})\|_\infty$  where  $\bar{y}$  is the mean of values in the vector  $y$ .

Reference: <http://davidrosenberg.github.io/ml2015/docs/5.Lab.misc.pdf#page=10>.

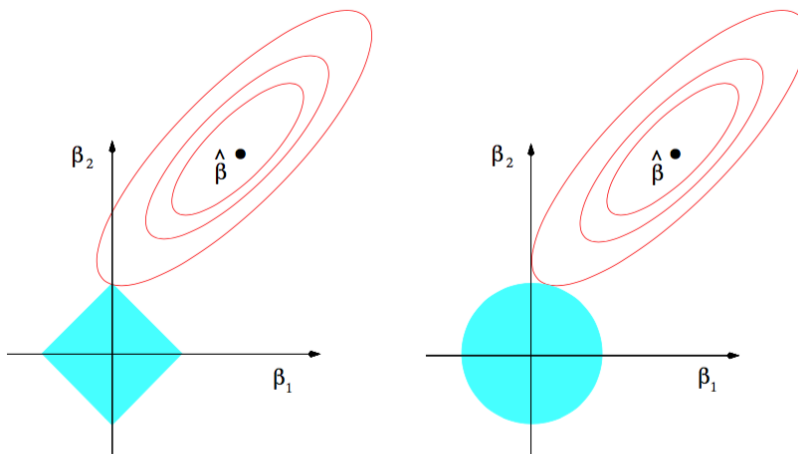
### 3.2 [Optional] Feature Correlation

In this problem, we will examine and compare the behavior of the Lasso and ridge regression in the case of an exactly repeated feature. That is, consider the design matrix  $X \in \mathbf{R}^{m \times d}$ , where  $X_{.i} = X_{.j}$  for some  $i$  and  $j$ , where  $X_{.i}$  is the  $i^{th}$  column of  $X$ . We will see that ridge regression divides the weight equally among identical features, while Lasso divides the weight arbitrarily. In an optional part to this problem, we will consider what changes when  $X_{.i}$  and  $X_{.j}$  are highly correlated (e.g. exactly the same except for some small random noise) rather than exactly the same.

1. [Optional] Derive the relation between  $\hat{\theta}_i$  and  $\hat{\theta}_j$ , the  $i^{th}$  and the  $j^{th}$  components of the optimal weight vector obtained by solving the Lasso optimization problem.  
[Hint: Assume that in the optimal solution,  $\hat{\theta}_i = a$  and  $\hat{\theta}_j = b$ . First show that  $a$  and  $b$  must have the same sign. Then, using this result, rewrite the optimization problem to derive a relation between  $a$  and  $b$ .]
2. [Optional] Derive the relation between  $\hat{\theta}_i$  and  $\hat{\theta}_j$ , the  $i^{th}$  and the  $j^{th}$  components of the optimal weight vector obtained by solving the ridge regression optimization problem.
3. [Optional] What do you think would happen with Lasso and ridge when  $X_{.i}$  and  $X_{.j}$  are highly correlated, but not exactly the same. You may investigate this experimentally or geometrically.

## 4 The Ellipsoids in the $\ell_1/\ell_2$ regularization picture

Recall the famous picture purporting to explain why  $\ell_1$  regularization leads to sparsity, while  $\ell_2$  regularization does not. Here's the instance from Hastie et al's *The Elements of Statistical Learning*:



**FIGURE 3.11.** Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.

(While Hastie et al. use  $\beta$  for the parameters, we'll continue to use  $w$ .)

In this problem we'll show that the level sets of the empirical risk are indeed ellipsoids centered at the empirical risk minimizer  $\hat{w}$ .

Consider linear prediction functions of the form  $x \mapsto w^T x$ . Then the empirical risk for any  $w$ , the empirical risk for  $f(x) = w^T x$  under the square loss is

$$\begin{aligned}\hat{R}_n(w) &= \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 \\ &= \frac{1}{n} (Xw - y)^T (Xw - y).\end{aligned}$$

1. Let  $\hat{w} = (X^T X)^{-1} X^T y$ . Show that  $\hat{w}$  has empirical risk given by

$$\hat{R}_n(\hat{w}) = \frac{1}{n} (-y^T X \hat{w} + y^T y)$$

2. Show that for any  $w$  we have

$$\hat{R}_n(w) = \frac{1}{n} (w - \hat{w})^T X^T X (w - \hat{w}) + \hat{R}_n(\hat{w}).$$

Note that the RHS (i.e. “right hand side”) has one term that's quadratic in  $w$  and one term that's independent of  $w$ . In particular, the RHS does not have any term that's linear in  $w$ . On the LHS (i.e. “left hand side”), we have  $\hat{R}_n(w) = \frac{1}{n} (Xw - y)^T (Xw - y)$ . After expanding the out, you'll have terms that are quadratic, linear, and constant in  $w$ . Completing the square is the tool for rearranging an expression to get rid of the linear terms. The following “completing the square” identity is easy to verify just by multiplying out the expressions on the RHS:

$$x^T M x - 2b^T x = (x - M^{-1}b)^T M (x - M^{-1}b) - b^T M^{-1}b$$

3. Using the expression derived for  $\hat{R}_n(w)$  in 2, give a very short proof that  $\hat{w} = (X^T X)^{-1} X^T y$  is the empirical risk minimizer. That is:

$$\hat{w} = \arg \min_w \hat{R}_n(w).$$

Hint: Note that  $X^T X$  is positive semidefinite and, by definition, a symmetric matrix  $M$  is positive semidefinite iff for all  $x \in \mathbf{R}^d$ ,  $x^T M x \geq 0$ .

4. Give an expression for the set of  $w$  for which the empirical risk exceeds the minimum empirical risk  $\hat{R}_n(\hat{w})$  by an amount  $c > 0$ . This set is an ellipse – what is its center?

## 5 [Optional] Projected SGD via Variable Splitting

In this question, we consider another general technique that can be used on the Lasso problem. We first use the variable splitting method to transform the Lasso problem to a smooth problem with linear inequality constraints, and then we can apply a variant of SGD.

Representing the unknown vector  $\theta$  as a difference of two non-negative vectors  $\theta^+$  and  $\theta^-$ , the  $\ell_1$ -norm of  $\theta$  is given by  $\sum_{i=1}^d \theta_i^+ + \sum_{i=1}^d \theta_i^-$ . Thus, the optimization problem can be written as

$$(\hat{\theta}^+, \hat{\theta}^-) = \arg \min_{\theta^+, \theta^- \in \mathbf{R}^d} \sum_{i=1}^m (h_{\theta^+, \theta^-}(x_i) - y_i)^2 + \lambda \sum_{i=1}^d \theta_i^+ + \lambda \sum_{i=1}^d \theta_i^-$$

such that  $\theta^+ \geq 0$  and  $\theta^- \geq 0$ ,

where  $h_{\theta^+, \theta^-}(x) = (\theta^+ - \theta^-)^T x$ . The original parameter  $\theta$  can then be estimated as  $\hat{\theta} = (\hat{\theta}^+ - \hat{\theta}^-)$ .

This is a convex optimization problem with a differentiable objective and linear inequality constraints. We can approach this problem using projected stochastic gradient descent, as discussed in lecture. Here, after taking our stochastic gradient step, we project the result back into the feasible set by setting any negative components of  $\theta^+$  and  $\theta^-$  to zero.

1. [Optional] Implement projected SGD to solve the above optimization problem for the same  $\lambda$ 's as used with the shooting algorithm. Since the two optimization algorithms should find essentially the same solutions, you can check the algorithms against each other. Report the differences in validation loss for each  $\lambda$  between the two optimization methods. (You can make a table or plot the differences.)
2. [Optional] Choose the  $\lambda$  that gives the best performance on the validation set. Describe the solution  $\hat{w}$  in term of its sparsity. How does the sparsity compare to the solution from the shooting algorithm?