

Machine Learning and Computational Statistics, Spring 2017

Homework 2: Lasso Regression

Saturday, February 11, 2017 Yidi Zhang

1 Introduce

2 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. ℓ_2 -regularized linear regression) on this dataset.

1. Run ridge regression on this dataset. Choose the λ that minimizes the square loss on the validation set. For the chosen λ , 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`.)

Answer 2.1:

```
import numpy as np
from scipy.optimize import minimize

X = np.loadtxt("/Users/twff/Downloads/machine
learning/hw/hw2-lasso/data/X_train.txt")
y = np.loadtxt("/Users/twff/Downloads/machine
learning/hw/hw2-lasso/data/y_train.txt")
X_valid = np.loadtxt("/Users/twff/Downloads/machine
learning/hw/hw2-lasso/data/X_valid.txt")
y_valid = np.loadtxt("/Users/twff/Downloads/machine
learning/hw/hw2-lasso/data/y_valid.txt")

(N,D) = X_train.shape

w = np.random.rand(D,1)

def ridge(Lambda):
```

```

def ridge_obj(theta):
    return ((np.linalg.norm(np.dot(X, theta) - y)**2)/(2*N) +
            Lambda*(np.linalg.norm(theta)**2)
    return ridge_obj

def compute_loss(X, y, Lambda, theta):
    (N,D) = X.shape
    return ((np.linalg.norm(np.dot(X, theta) - y)**2)/(2*N)

for i in range(-8,6):
    Lambda = 10**i;
    w_opt = minimize(ridge(Lambda), w)
    print(Lambda, compute_loss(X_valid, y_valid, Lambda, w_opt.x))

```

When $\lambda = 1e - 09$, it can minimize the square loss, which is 0.0587394466821. $\lambda = 1e - 09$, 0 component with true value 0 has been estimated to be non-zero. And when choosing a small threshold which equals to 10^{-3} , the number is 50.

```

w_opt = minimize(ridge(1e-5), w)
len(w_opt.x[w_opt.x==0])
len(w_opt.x[w_opt.x <= 0.001])

```

When comparing the results with the one using `sklearn.linear_model.Ridge`, the square loss is similar, which is 0.10057515598174918 when $\lambda = 1e - 5$, so the code do not have bug.

```

from sklearn.linear_model import Ridge
clf = Ridge(alpha=1e-5)
clf.fit(X, y)
a = clf.predict(X)#numpy.dot(X,theta)
a=((np.linalg.norm(a - y)**2)/(2*N)
b = compute_loss(X_valid, y_valid, Lambda, w_opt.x)
b

```

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

3.1 Experiments with the Shooting Algorithm

1. Write a function that computes the Lasso solution for a given λ 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 λ that minimizes the square error on the validation set. Report the optimal value of λ found, and the corresponding test error. Plot the validation error vs λ . [Don't use the homotopy method in this part, as we want to measure the speed improvement of homotopy methods in question 3. Also, no need to vectorize the calculations until question 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.]

Answer 3.1.1:

```
def lassoShooting(X, y, lambda_reg, w, ep = 1e-6, iter_max=1000):
    (n,p) = X.shape
    converge = True
    itera = 0

    while converge and (itera + 1) <= iter_max:
        w_old = w.copy()
        for j in range(p):
            a = []
            c = []
            for i in range(n):
                aj = 2 * X[i,j]*X[i,j]
                cj = 2 * X[i,j]*(y[i]-np.dot(w.T,X[i]) + X[i,j]*w[j])
                a.append(aj)
                c.append(cj)
            aj_sum = sum(a)
            cj_sum = sum(c)
            w[j] = soft(cj_sum/aj_sum, lambda_reg/aj_sum)
            converge = (sum(abs(w-w_old)) >= ep)
            #print (sum(abs(w-w_old)))
            itera +=1
        print(itera)
    return w

def soft(a, b):
    if a <= b:
        c = a+b
    elif a > b:
        c = a - b
    else:
        c = 0
    return c
```

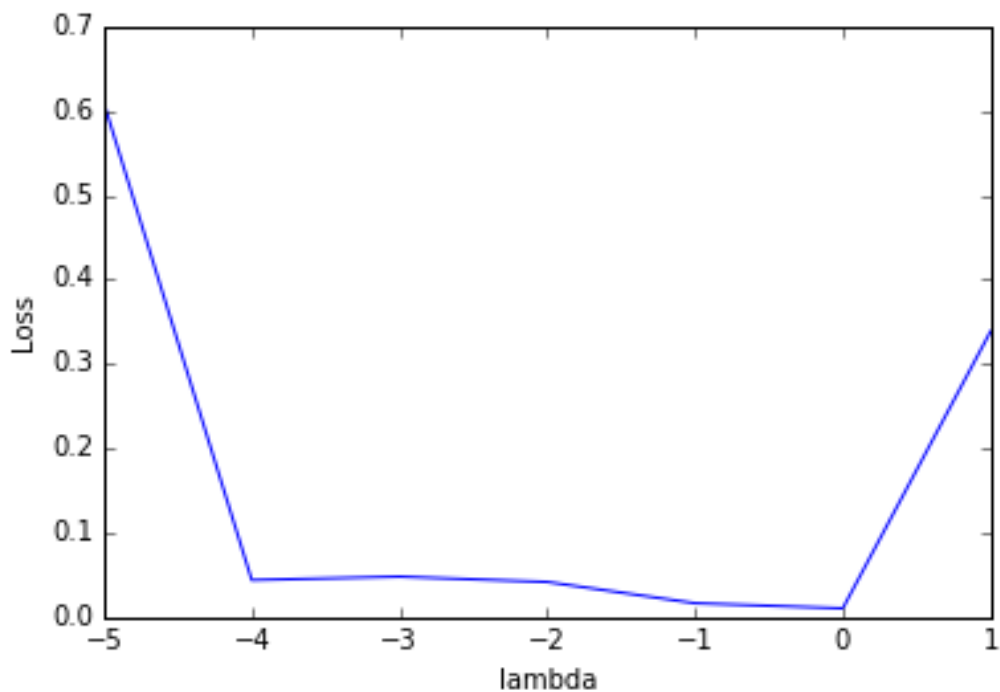


Figure 1: Objective function for each step size.

```
import matplotlib.pyplot as plt
%matplotlib inline

w = np.zeros(X.shape[1])
loss = []
for i in range(-5,2):
    Lambda = 10**i;
    w_opt = lassoShooting(X, y, lambda_reg=Lambda, w=w, ep=1e-3)
    loss.append(compute_loss(X_valid, y_valid, Lambda, w_opt))
    print(Lambda, compute_loss(X_valid, y_valid, Lambda, w_opt))

plt.plot(range(-5,2), loss)
plt.xlabel("lambda")
plt.ylabel("Loss")
plt.savefig("/Users/twff/Downloads/machine learning/hw/hw2-lasso/img311")
```

```
w = np.zeros(X.shape[1])
w_opt = lassoShooting(X, y, lambda_reg=0.4, w=w, ep=1e-3)
```

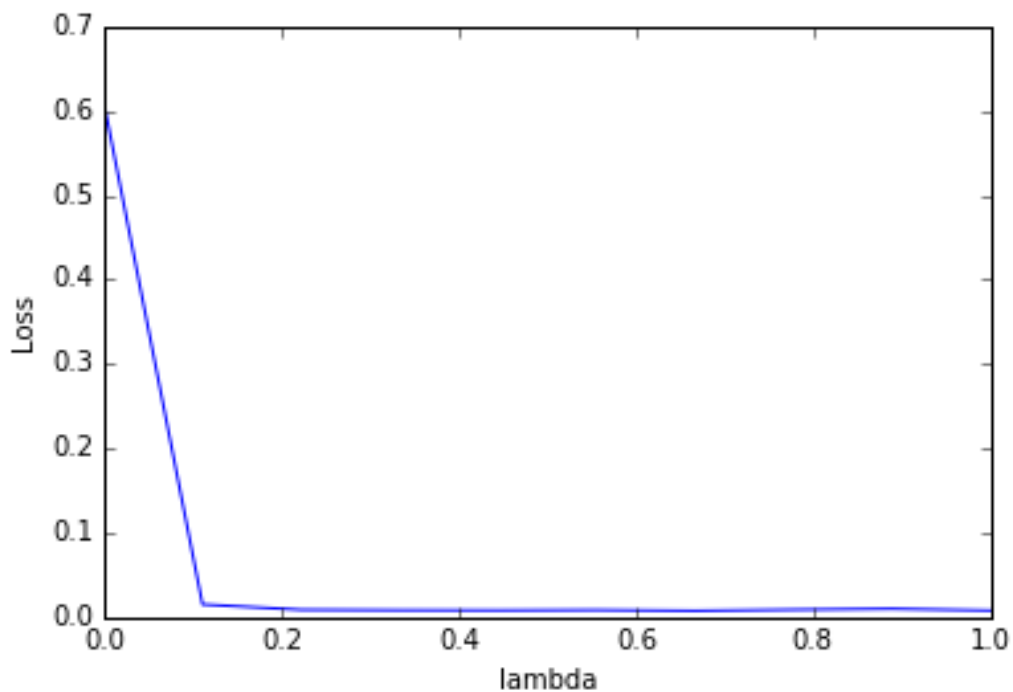


Figure 2: Objective function for each step size.

```
compute_loss(X_test, y_test, 0.4 ,w_opt)
```

And I found the minimum square loss is between 0 and 1. After zooming in, the optimal λ found is 0.4, and the corresponding test error is 0.0093066225296064235.

```
w = np.zeros(X.shape[1])
lambda_reg = np.linspace(0,1,10)
loss = []
for i in range(len(lambda_reg)):
    w_opt = lassoShooting(X, y, lambda_reg=lambda_reg[i], w=w, ep=1e-3)
    loss.append(compute_loss(X_valid, y_valid, lambda_reg[i], w_opt))
    print(lambda_reg[i], compute_loss(X_valid, y_valid, lambda_reg[i],
        w_opt))

plt.plot(lambda_reg, loss)
plt.xlabel("lambda")
plt.ylabel("Loss")
plt.savefig("/Users/twff/Downloads/machine learning/hw/hw2-lasso/img311b")
```

2. Analyze the sparsity¹ of your solution, reporting how many components with true value zero have been estimated to be non-zero, and vice-versa.

Answer 3.1.2:

```
w = np.zeros(X.shape[1])
w_opt = lassoShooting(X, y, lambda_reg=0.4, w=w, ep=1e-3)
compute_loss(X_test, y_test, 0.4, w_opt)
```

0 components with true value zero has been estimated to be non-zero. But if the threshold is 10^{-3} , the number would be 51. In this case, we can see the solution above as a sparsity pattern.

¹One might hope that the solution will 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**” (sparse + consistent). Formally, an estimator $\hat{\beta}$ of parameter β is said to be consistent if the estimator $\hat{\beta}$ converges to the true value β in probability as our sample size goes to infinity. Analogously, if we define the support of a vector β 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 β , or $\lim_{m \rightarrow \infty} P[\text{Supp}(\hat{\beta}_m) = \text{Supp}(\beta)] = 1$.

3. Implement the homotopy method described above. Compare the runtime for computing the full regularization path (for the same set of λ 's you tried in the first question above) using the homotopy method compared to the basic shooting algorithm.

Answer 3.1.3:

```
def lassoShooting_homotopy(X, y, w, ep = 1e-3, pa=0.1, iter_max=2000):
    lambda_max = 2 * np.linalg.norm(np.dot(X.T, y), np.inf)
    (n,p) = X.shape
    converge = True
    itera = 0
    lambda_reg = 1000
    lambda_old = 0
    lambda_list = []
    loss = []
    w = np.zeros(p)
    #print(w)

    while np.abs(lambda_old - lambda_reg)>ep:
        w_prev = w.copy()
        w_new = lassoShooting(X, y, lambda_reg, w, ep = 1e-3,
                               iter_max=iter_max)
        #print(w_new)
        loss.append(compute_loss(X_valid, y_valid, lambda_reg, w_new))
        lambda_list.append(lambda_reg)
        if(sum(abs(w_prev-w_new)) < ep):
            #print(sum(abs(w-w_new)))
            break
        w = w_new.copy()
        lambda_reg = pa*lambda_reg

    return lambda_list, loss
```

The homotopy method time is 124.1882728970013? while the regular method time is 302.8984123150003.

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 vectorized method can be found in the assignment zip file.) Implement the matrix expressions and measure the speedup in computing the regularization path.

Answer 3.1.4:

```
def lassoShooting_matrix(X, y, lambda_reg, w, ep = 1e-3, iter_max=500):
    (n,p) = X.shape
    w_opt = w
    converge = True
    itera = 0

    while converge and (itera+1) <= iter_max:
        w_old = w.copy()
        a = np.dot(X.T, X)
        for j in range(p):
            a = 2*X[:,j].dot(X[:,j])
            c = 2*X[:,j].dot(y) - 2* X.dot(w).dot(X[:,j]) +
                2*X[:,j].dot(X[:,j])*w[j]
            w[j] = soft(c/a, lambda_reg/a)
        converge = (sum(abs(w-w_old)) >= ep)
        #print(sum(abs(w-w_old)))
        itera +=1
    return w
```

For $\lambda = 0.4$, the regular lasso shooting method time is 20.220894887999748, but the vectorization method time is 0.46611224899970694. For computing the regularization path, vectorization time is 9.962034545998904, and regular time is 302.8984123150003.

3.2 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 [derived using subgradients \(slide 15\)](#), but here we will take 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. 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 fact that since $f : \mathbf{R} \rightarrow \mathbf{R}$ 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 [slide 12 here](#), 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 $\sum_{i=1}^n x_{ij}^2 > 0$.

Answer 3.2.1: If $x_{ij} = 0$ for $i = 1, \dots, n$

$$a_j = 2 \sum_{i=1}^n x_{ij}^2 = 0$$

,

$$c_j = 2 \sum_{i=1}^n x_{ij} \left(y_i - \sum_{k \neq j} w_k x_{ik} \right) = 0$$

thus when $a_j = c_j = 0$, we can get $w_j = 0$.

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:

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

Answer 3.2.2:

$$w_j = \text{sign}(w_j) \left(\left| \frac{c_j}{a_j} \right| - \frac{\lambda}{a_j} \right)$$

for $w_j \neq 0$:

$$f(w_j) := \begin{cases} \frac{1}{a_j}(c_j - \lambda) & w_j > 0 \\ \frac{1}{a_j}(c_j + \lambda) & w_j < 0 \end{cases}$$

3. If $w_j > 0$ and minimizes f , show that $w_j = \frac{1}{a_j}(c_j - \lambda)$. Similarly, if $w_j < 0$ and minimizes f , show that $w_j = \frac{1}{a_j}(c_j + \lambda)$. Give conditions on c_j that imply that a minimizer w_j is positive and conditions for which a minimizer w_j is negative.

Answer 3.2.3:

if $w_j > 0$, $w_j = \frac{1}{a_j}(|c_j| - \lambda) = w_j = \frac{1}{a_j}(c_j - \lambda)$, in this case, c_j must be bigger than λ .

if $w_j < 0$, $w_j = -\frac{1}{a_j}(|c_j| - \lambda) = w_j = \frac{1}{a_j}(c_j + \lambda)$, in this case, c_j must be smaller than $-\lambda$.

Give the condition on c_j to imply w_j :

$$f(w_j) := \begin{cases} \frac{1}{a_j}(c_j - \lambda) > 0 & \text{if } c_j > \lambda \\ \frac{1}{a_j}(c_j + \lambda) < 0 & \text{if } c_j < -\lambda \\ w_j = 0 & \text{otherwise} \end{cases}$$

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.

Answer 3.2.4:

When $w_j > 0$, for $c_j > \lambda$, suppose $\varepsilon > 0$

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

When $w_j < 0$, for $c_j < -\lambda$

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

But if $c_j \in [-\lambda, \lambda]$

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

Putting together the preceding expressions, we can conclude that $c_j \in [-\lambda, \lambda]$ implies that $w_j = 0$ is a minimizer.

5. Putting together the preceding results, we conclude the following:

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

Show that this is equivalent to the expression given in 3.

Answer 3.2.5:

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

The derivate of $g(w)$ is:

$$\begin{aligned} \nabla_w g(w) &= 2 \sum_i^n \left(\sum_k^p x_{ik} w - y_i \right) x_{ij} \\ &= -2 \sum_i^n \left(y_i - \sum_k^p x_{ik} w \right) x_{ij} \\ &= -c_j \end{aligned}$$

Then we can get:

$$g(w_{j+1}) = g(w_j) + \lambda \nabla_w g(w) + \frac{a_j}{2}$$

The minimizer of w then is given by:

$$\begin{aligned} \hat{w} &= \text{sign}(w_j) \left(w_j; \frac{\nabla_j g(w) - \lambda}{a_j}, \frac{\nabla_j g(w) + \lambda}{a_j} \right) \\ &= \text{sign}(w_j) \left(w_j; \frac{c_j - \lambda}{a_j}, \frac{c_j + \lambda}{a_j} \right) \end{aligned}$$

Putting together the preceding results, we can conclude that :

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

is equivalent to the expression:

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

4 Lasso Properties

4.1 Deriving λ_{\max}

In this problem we will derive an expression for λ_{\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, λ , and v .]

Answer 4.1.1:

$$\begin{aligned} l'(0; v) &= \lim_{h \downarrow 0} \frac{l(0 + hv) - l(0)}{h}, \\ &= \lim_{h \downarrow 0} \frac{(Xhv - y)(Xhv - y)^T + \lambda hv - y^T y}{h} \\ &= -2X^T v^T y + \lambda \|v\|_1 \end{aligned}$$

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 λ . [Hint: this should be in terms of X, y , and v .]

Answer 4.1.2:

To make $l'(0; v) \geq 0$,

$$\begin{aligned} -2X^T v^T y + \lambda \|v\|_1 &\geq 0 \\ \lambda \|v\|_1 &\geq 2X^T v^T y \\ \lambda &\geq \frac{2X^T v^T y}{\|v\|_1} \end{aligned}$$

Thus, the lower bound of λ is $\frac{2X^T v^T y}{\|v\|_1}$

3. In the previous problem, we get a different lower bound on λ for each choice of v . Compute the maximum lower bound of λ by maximizing the expression over v . Show that this expression is equivalent to $\lambda_{\max} = 2\|X^T y\|_{\infty}$.

Answer 4.1.3:

$$\begin{aligned}\lambda &\geq \frac{2X^T v^T y}{\|v\|_1} \\ \lambda &\geq 2 \frac{v^T}{\|v\|_1} X^T y\end{aligned}$$

Suppose $z^T = \frac{v^T}{\|v\|_1}$, we can get $\|z\|_1 = \|\frac{v^T}{\|v\|_1}\|_1 = \frac{\|v\|_1}{\|v\|_1} = 1$. We let $p = X^T y$, then:

$$\begin{aligned}z^T p &= \sum_{i=1}^n z_i p_i \\ z^T p &\leq \sum_{i=1}^n |z_i| |p_i| \\ &\leq \|z\|_{\infty} \|p\|_{\infty} \\ &\leq \|X^T y\|_{\infty}\end{aligned}$$

Thus we can conclude that $\lambda_{\max} = 2\|X^T y\|_{\infty}$

4. [Optional] Show that for $L(w, b) = \|Xw + b\mathbf{1} - y\|_2^2 + \lambda \|w\|_1$, $\lambda_{\max} = 2\|X^T(y - \bar{y})\|_\infty$ where \bar{y} is the mean of values in the vector y , and $\mathbf{1} \in \mathbf{R}^n$ is a column vector of 1's .

Answer 4.1.4:

for $L(w, b) = \|Xw + b\mathbf{1} - y\|_2^2 + \lambda \|w\|_1$,

$$\begin{aligned} l'(0, 0; v) &= \lim_{h \downarrow 0} \frac{l(w + hv, b + hv) - l(w, b)}{h}, \\ &= \lim_{h \downarrow 0} \frac{((X + 1)hv - y)((X + 1)hv - y)^T + \lambda hv - y^T y}{h} \\ &= -2(X + 1)^T v^T y + \lambda \|v\|_1 \end{aligned}$$

To make $L'(0, 0; v) \geq 0$:

$$\begin{aligned} \lambda &\geq \frac{2(X + 1)^T v^T y}{\|v\|_1} \\ \lambda &\geq 2 \frac{v^T}{\|v\|_1} (X + 1)^T y \end{aligned}$$

Suppose $z^T = \frac{v^T}{\|v\|_1}$, we can get $\|z\|_1 = \left\| \frac{v^T}{\|v\|_1} \right\|_1 = \frac{\|v\|_1}{\|v\|_1} = 1$. We let $p = (X + 1)^T y$, then:

$$\begin{aligned} z^T p &= \sum_{i=1}^n z_i p_i \\ z^T p &\leq \sum_{i=1}^n |z_i| |p_i| \\ &\leq \|z\|_\infty \|p\|_\infty \\ &\leq \|X^T(y - \bar{y})\|_\infty \end{aligned}$$

Thus we can conclude that $\lambda_{\max} = 2\|X^T(y - \bar{y})\|_\infty$

4.2 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. Without a loss of generality, assume the first two columns of X are our repeated features. Partition X and θ as follows:

$$X = \begin{pmatrix} x_1 & x_2 & X_r \end{pmatrix}, \theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_r \end{pmatrix}$$

We can write the Lasso objective function as:

$$\begin{aligned} L(\theta) &= \|X\theta - y\|_2^2 + \lambda \|\theta\|_1 \\ &= \|x_1\theta_1 + x_2\theta_2 + X_r\theta_r - y\|_2^2 + \lambda|\theta_1| + \lambda|\theta_2| + \lambda\|\theta_r\|_1 \end{aligned}$$

With repeated features, there will be multiple minimizers of $L(\theta)$. Suppose that

$$\hat{\theta} = \begin{pmatrix} a \\ b \\ r \end{pmatrix}$$

is a minimizer of $L(\theta)$. Give conditions on c and d such that $(c, d, r^T)^T$ is also a minimizer of $L(\theta)$. [Hint: First show that a and b must have the same sign, or at least one of them is zero. Then, using this result, rewrite the optimization problem to derive a relation between a and b .]

Answer 4.2.1: When having

$$\hat{\theta} = \begin{pmatrix} a \\ b \\ r \end{pmatrix}$$

We can get:

$$L(\hat{\theta}) = \|(a+b)x_{ab} + X_r\theta_r - y\|_2^2 + \lambda(|a| + |b|) + \lambda\|\theta_r\|_1$$

Suppose existing $ab < 0$ and $cd > 0$, so we can get $|a| + |b| \geq |a+b|$, $|c| + |d| = |c+d|$. Besides, since $(c, d, r^T)^T$ and $(a, b, r^T)^T$ are both minimizers, we can get $a+b = c+d = p$, and p is a const. So we can conclude that:

$$\begin{aligned} L(\hat{\theta}) &= \|(a+b)x_{ab} + X_r\theta_r - y\|_2^2 + \lambda(|a| + |b|) + \lambda\|\theta_r\|_1 \\ &\geq \|(c+d)x_{cd} + X_r\theta_r - y\|_2^2 + \lambda(|c| + |d|) + \lambda\|\theta_r\|_1 \end{aligned}$$

Which contraries to hypothesis that $(a, b, r^T)^T$ is a minimizer. So we can conclude that relation between a and b is $ab > 0$ and $a+b = p$, and p is an const.

2. Using the same notation as the previous problem, suppose

$$\hat{\theta} = \begin{pmatrix} a \\ b \\ r \end{pmatrix}$$

minimizes the ridge regression objective function. What is the relationship between a and b , and why?

Answer 4.2.2:

We can write the ridge regression function as:

$$\begin{aligned} L(\theta) &= \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2 \\ &= \|x_1\theta_1 + x_2\theta_2 + X_r\theta_r - y\|_2^2 + \lambda\|\theta\|_2^2 + \lambda\|\theta\|_2^2 + \lambda\|\theta_r\|_2^2 \end{aligned}$$

For minimizer

$$\hat{\theta} = \begin{pmatrix} a \\ b \\ r \end{pmatrix}$$

$$\begin{aligned} L(\hat{\theta}) &= \|X\hat{\theta} - y\|_2^2 + \lambda \|\hat{\theta}\|_2^2 \\ &= \|(a+b)x_{ab} + X_r\theta_r - y\|_2^2 + \lambda(a^2 + b^2) + \lambda\|\theta_r\|_2^2 \end{aligned}$$

Suppose $(a, b, r^T)^T$ and $(c, d, r^T)^T$ are both minimizer, so $a + b = p = c + d$, while p is a const.

Suppose existing $ab < 0$ and $cd > 0$, so we can get $a^2 + b^2 > (a+b)^2$, $c^2 + d^2 < (c+d)^2$. So we can conclude that:

$$\begin{aligned} L(\hat{\theta}) &= \|(a+b)x_{ab} + X_r\theta_r - y\|_2^2 + \lambda(a^2 + b^2) + \lambda\|\theta_r\|_2^2 \\ &\geq \|(c+d)x_{cd} + X_r\theta_r - y\|_2^2 + \lambda(c^2 + d^2) + \lambda\|\theta_r\|_2^2 \end{aligned}$$

Which contraries to hypothesis that $(a, b, r^T)^T$ is a minimizer.

Since $a + b = p$, which is a const, to minimize $a^2 + b^2$, a should equal to b , that is $a = b$. So we can conclude that relation between a and b is $ab > 0$ and $a + b = p$, where p and r is const, and $a = b$

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 theoretically.

Answer 4.2.3:

If the features are highly correlated, the lasso method will only select one of them, that is, for $ab > 0$, one of them would be 0 while the other is const p . For correlated features, ridge method, however, tends to have similar coefficients.

```
for i in range(10):
    X = np.array([0, 1])
    y = np.array([0, 1])
    means = [X.mean(), y.mean()]
    stds = [X.std() / 3, y.std() / 3]
    corr = 0.1*(i+1) # correlation
    covs = [[stds[0]**2, stds[0]*stds[1]*corr],
            [stds[0]*stds[1]*corr, stds[1]**2]]
    m = np.random.multivariate_normal(means, covs, 1000).T
    X = m[0]
    y = m[1]
    ridge = Ridge(alpha=10)
    ridge.fit(m.T,y)
    print ("Ridge model:", pretty_print_linear(ridge.coef_))
```

The results show below:

```
Ridge model: 0.027 * X0 + 0.725 * X1
Ridge model: 0.033 * X0 + 0.726 * X1
Ridge model: 0.068 * X0 + 0.733 * X1
Ridge model: 0.083 * X0 + 0.708 * X1
Ridge model: 0.11 * X0 + 0.693 * X1
Ridge model: 0.15 * X0 + 0.679 * X1
Ridge model: 0.188 * X0 + 0.646 * X1
Ridge model: 0.244 * X0 + 0.583 * X1
Ridge model: 0.313 * X0 + 0.532 * X1
Ridge model: 0.424 * X0 + 0.424 * X1
```

5 [Optional] The Ellipsoids in the ℓ_1/ℓ_2 regularization picture

Recall the famous picture purporting to explain why ℓ_1 regularization leads to sparsity, while ℓ_2 regularization does not. Here's the instance from Hastie et al's *The Elements of Statistical Learning*:

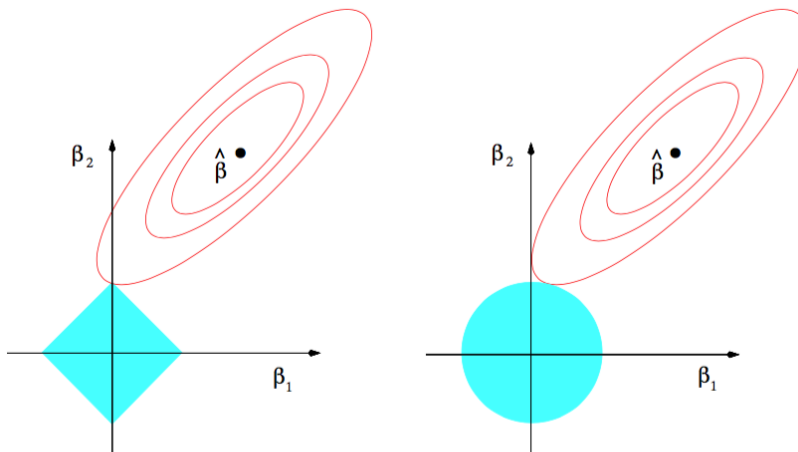


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the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

(While Hastie et al. use β 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 $f(x) = w^T x$ 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. [Optional] 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. [Optional] 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 this 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. [Optional] 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. [Optional] 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 $c > 0$. If X is full rank, then $X^T X$ is positive definite, and this set is an ellipse – what is its center?

6 [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 θ as a difference of two non-negative vectors θ^+ and θ^- , the ℓ_1 -norm of θ 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 θ 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 θ^+ and θ^- to zero.

1. [Optional] Implement projected SGD to solve the above optimization problem for the same λ '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 λ between the two optimization methods. (You can make a table or plot the differences.)

2. [Optional] Choose the λ 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?