

# DS-GA 1003 HW2

Yiyan Chen

February 2018

## 1 Introduction

## 2 Ridge Regression

### 2.1

---

```
1 from ridge_regression import RidgeRegression, plot_prediction_functions,
2     compare_parameter_vectors
3 import setup_problem
4 import numpy as np
5 import pandas as pd
6 from sklearn.model_selection import GridSearchCV, PredefinedSplit
7 from sklearn.metrics import mean_squared_error, make_scorer
8 import matplotlib.pyplot as plt
9 %matplotlib inline
10 lasso_data_fname = "lasso_data.pickle"
11 x_train, y_train, x_val, y_val, target_fn, coeffs_true, featurize =
12     setup_problem.load_problem(lasso_data_fname)
13 X_train = featurize(x_train)
14 X_val = featurize(x_val)
15 X_train_val = np.vstack((X_train, X_val))
16 y_train_val = np.concatenate((y_train, y_val))
17 val_fold = [-1]*len(X_train) + [0]*len(X_val)
18
19 param_grid = [{'l2reg':np.concatenate((10.**np.arange(-3, 1, 0.5), np.arange(1,3,0.3))) }]
20
21 ridge_regression_estimator = RidgeRegression()
22 grid = GridSearchCV(ridge_regression_estimator,
23                     param_grid,
24                     return_train_score=False,
25                     cv = PredefinedSplit(test_fold=val_fold),
26                     refit = True,
27                     scoring = make_scorer(mean_squared_error,
```

```

28                                     greater_is_better = False))
29 grid.fit(X_train_val, y_train_val)
30 df = pd.DataFrame(grid.cv_results_)
31 df['mean_test_score'] = -df['mean_test_score']
32 cols_to_keep = ["param_l2reg", "mean_test_score"]
33 df_toshow = df[cols_to_keep].fillna('-')
34 df_toshow = df_toshow.sort_values(by=["param_l2reg"])
35 df_toshow

```

---

|    | param_l2reg | mean_test_score |
|----|-------------|-----------------|
| 0  | 0.001000    | 0.162705        |
| 1  | 0.003162    | 0.151900        |
| 2  | 0.010000    | 0.141887        |
| 3  | 0.031623    | 0.139648        |
| 4  | 0.100000    | 0.144566        |
| 5  | 0.316228    | 0.152607        |
| 6  | 1.000000    | 0.171068        |
| 8  | 1.000000    | 0.171068        |
| 9  | 1.300000    | 0.179521        |
| 10 | 1.600000    | 0.187993        |
| 11 | 1.900000    | 0.196361        |
| 12 | 2.200000    | 0.204553        |
| 13 | 2.500000    | 0.212530        |
| 14 | 2.800000    | 0.220271        |
| 7  | 3.162278    | 0.229295        |

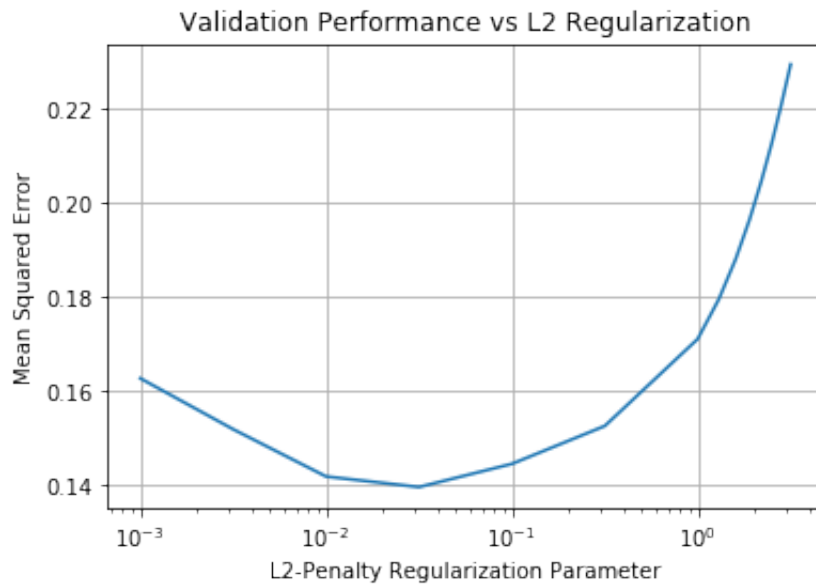
---

```

1 fig, ax = plt.subplots()
2 ax.semilogx(df_toshow["param_l2reg"], df_toshow["mean_test_score"])
3 ax.grid()
4 ax.set_title("Validation Performance vs L2 Regularization")
5 ax.set_xlabel("L2-Penalty Regularization Parameter")
6 ax.set_ylabel("Mean Squared Error")
7 fig.show()

```

---



## 2.2

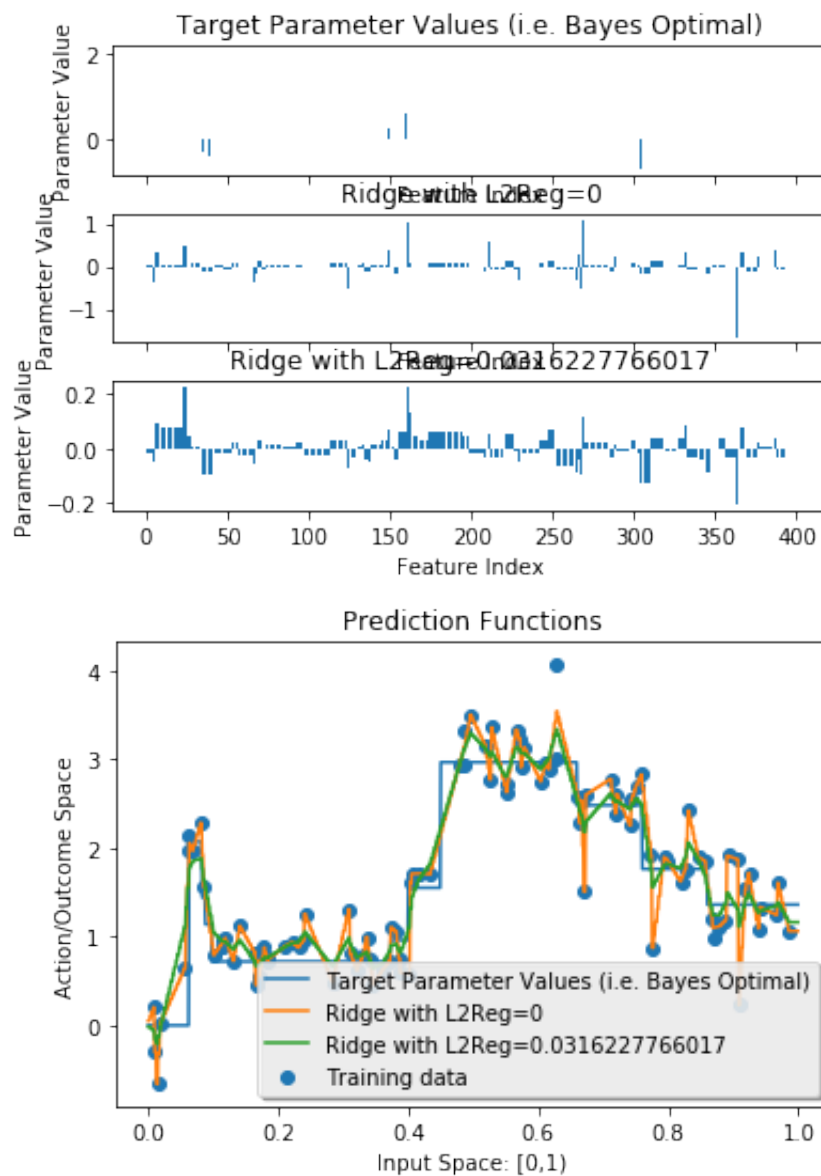
---

```

1 para_final = df_toshow['param_l2reg'][np.argmin(df_toshow['mean_test_score'])]
2 pred_fns = []
3 x = np.sort(np.concatenate([np.arange(0,1,.001), x_train]))
4 name = "Target Parameter Values (i.e. Bayes Optimal)"
5 pred_fns.append({"name":name, "coefs":coefs_true, "preds": target_fn(x) })
6
7 l2regs = [0, para_final]
8 X = featurize(x)
9 for l2reg in l2regs:
10     ridge_regression_estimator = RidgeRegression(l2reg=l2reg)
11     ridge_regression_estimator.fit(X_train, y_train)
12     name = "Ridge with L2Reg="+str(l2reg)
13     pred_fns.append({"name":name,
14                     "coefs":ridge_regression_estimator.w_,
15                     "preds": ridge_regression_estimator.predict(X) })
16
17 plot_prediction_functions(x, pred_fns, x_train, y_train, legend_loc="best")
18 compare_parameter_vectors(pred_fns)

```

---




---

```

1 max_ls = []
2 min_ls = []
3 max_idx_ls = []
4 for i in range(len(pred_fns)):
5     max_ls.append(max(pred_fns[i]['coefs']))
6     min_ls.append(min(pred_fns[i]['coefs']))
7     max_idx_ls.append(np.argmax(pred_fns[i]['coefs']))

```

```

8
9 df_pred = pd.DataFrame(np.concatenate((max_ls, min_ls, max_idx_ls)).reshape(3,3),
10                        index = ['max', 'min', 'max index'],
11                        columns = ['target', 'no regularization', 'best regularization'])

```

---

|           | target    | no regularization | best regularization |
|-----------|-----------|-------------------|---------------------|
| max       | 2.069572  | 1.093516          | 0.222279            |
| min       | -0.717818 | -1.644220         | -0.204951           |
| max index | 25.000000 | 269.000000        | 24.000000           |

The Bayes optimal's pattern is the most sparse one. The range of its coefficients' values is the largest,  $[-0.717818, 2.069572]$ . The 25th coefficient has the largest weight. While the regularized least squared fit pattern's sparsity is between target and the regularized one. The range of its coefficients' values is  $[-1.644220, 1.093516]$ . The 269th coefficient has the largest weight. Compared to the other models, the regularized model have the most non-zero coefficients, with the range of which is  $[-0.204951, 0.222279]$ . The largest coefficient weight happens at the 24th.

## 2.3

---

```

1 any(pred_fns[2]['coefs']==0)

```

---

The chosen  $\lambda$  is 0.0316, and the coefficients are mapped into the plot showed in last question. None of the coefficients is zero.

---

```

1 from sklearn.metrics import confusion_matrix
2 threshold = [10**-6, 10**-3, 10**-2, 10**-1]
3 pred_fns[0]['coefs'][abs(pred_fns[0]['coefs']) > 0] = 1
4 target_coef = pred_fns[0]['coefs']
5 reg_coef = pred_fns[2]['coefs']
6 for i in threshold:
7     copy_val = np.copy(reg_coef)
8     copy_val[abs(reg_coef) >= i] = 1
9     copy_val[abs(reg_coef) < i] = 0
10    print(confusion_matrix(target_coef, copy_val))

```

---

The confusion matrix when threshold,  $\epsilon = 10^{-6}$ :

$$\begin{bmatrix} 5 & 385 \\ 0 & 10 \end{bmatrix}$$

The confusion matrix when  $\epsilon = 10^{-3}$ :

$$\begin{bmatrix} 9 & 381 \\ 0 & 10 \end{bmatrix}$$

The confusion matrix when  $\epsilon = 10^{-2}$ :

$$\begin{bmatrix} 61 & 329 \\ 0 & 10 \end{bmatrix}$$

The confusion matrix when  $\epsilon = 10^{-1}$ :

$$\begin{bmatrix} 376 & 14 \\ 7 & 3 \end{bmatrix}$$

## 3 Coordinate Descent for Lasso

### 3.1 Experiments with the Shooting Algorithm

#### 3.1.1

$$a_j = 2X_{.j}^T X_{.j}$$
$$c_j = 2(X_{.j}^T y - X_{.j}^T X w + w_j X_{.j}^T X_{.j})$$

#### 3.1.2

---

```
1 def soft(a, delta):
2     if abs(a) - delta > 0:
3         return np.sign(a)*(abs(a)-delta)
4     else:
5         return 0
6
7 def objective_function(X, y, w, lambd):
8     y = y.reshape(-1)
9     return np.sum((y-np.dot(X, w))**2) + lambd*sum(abs(i) for i in w)
10
11 def shooting_alg(X, y, w, lambd, iteration = 1000, tolerance = 10**-6):
12     D = X.shape[1]
13     criteria = 1
14     count = 1
15     while criteria > tolerance or count < iteration:
16         obj_prev = objective_function(X, y, w, lambd)
17         for j in range(D):
18             aj = 2*np.dot(X[:,j], X[:,j])
19             cj = 2*(np.dot(X[:,j],y) - np.dot(np.dot(X[:,j], X), w)
20                 + w[j]*np.dot(X[:,j],X[:,j]))
21             wj = soft(cj/aj, lambd/aj)
22             w[j] = wj
23             criteria = obj_prev - objective_function(X,y,w, lambd)
24             count += 1
25     return w
26 w = np.zeros(X.shape[1])
27 w_cyc = shooting_alg(X_train, y_train, w, 1)
```

---

```
1 from random import sample
2 def shooting_alg_random(X, y, w, lambd, iteration = 1000, tolerance = 10**-6):
3     D = X.shape[1]
4     criteria = 1
5     count = 1
```

---

```

6     while criteria > tolerance or count < iteration:
7         obj_prev = objective_function(X, y, w, lamdb)
8         shuffled_list = sample(list(range(D)), D)
9         for j in shuffled_list:
10             aj = 2*np.dot(X[:,j], X[:,j])
11             cj = 2*(np.dot(X[:,j],y) - np.dot(np.dot(X[:,j], X), w)
12                 + w[j]*np.dot(X[:,j],X[:,j]))
13             wj = soft(cj/aj, lamdb/aj)
14             w[j] = wj
15             criteria = obj_prev - objective_function(X,y,w, lamdb)
16             count += 1
17     return w
18
19 w_random = shooting_alg_random(X_train, y_train, w, 1)
20
21 def RMSE(X, y, w):
22     return np.sum((y-np.dot(X, w))**2)/X.shape[0]
23
24 print(RMSE(X_val, y_val, w_cyc))
25 print(RMSE(X_val, y_val, w_random))

```

---

The validation error for cyclic coordinate descent is 0.167818772197, while that of randomized coordinate descent is 0.167818825676. Those two are very much similar.

---

```

1 lambda_val = 1
2 w_ridge = np.dot(np.dot(np.linalg.inv(np.dot(X_train.transpose(), X_train)
3                                     + lambda_val*np.identity(X_train.shape[1])),
4                                     X_train.transpose()), y_train)
5 w_ridge_cyc = shooting_alg(X_train, y_train, w_ridge, 1)
6 RMSE(X_val, y_val, w_ridge_cyc)

```

---

The validation error for cyclic coordinate descent when using ridge regression solution suggested by Murphy is 0.12643800274653294, which are the smallest among all.

### 3.1.3

---

```

1 lambda_val_lst = np.arange(-6, 2, 0.5, dtype=float)
2 val_error = []
3 for val in lambda_val_lst:
4     w_ridge_cyc_val = shooting_alg(X_train, y_train, w_ridge, 10**val)
5     val_error.append(400*RMSE(X_val, y_val, w_ridge_cyc_val))
6 plt.plot(lambda_val_lst, val_error)

```

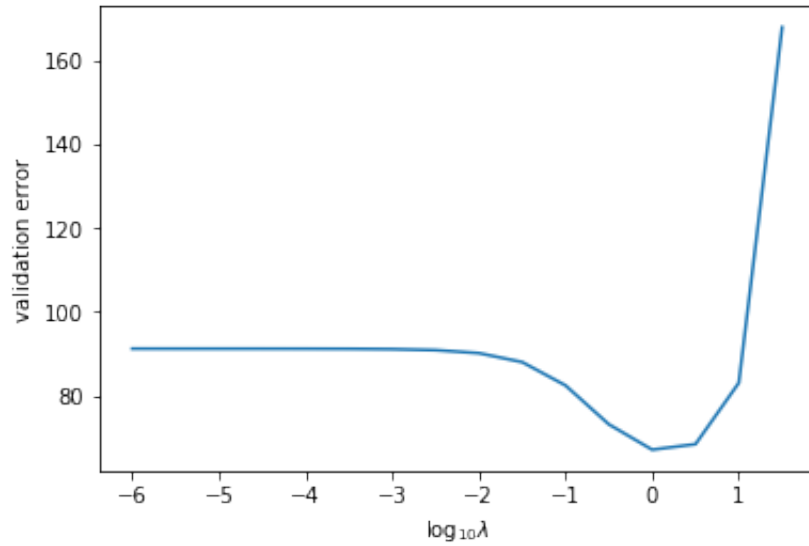


```

7 plt.xlabel('$\log_{10}\lambda$')
8 plt.ylabel("validation error")
9 plt.show()

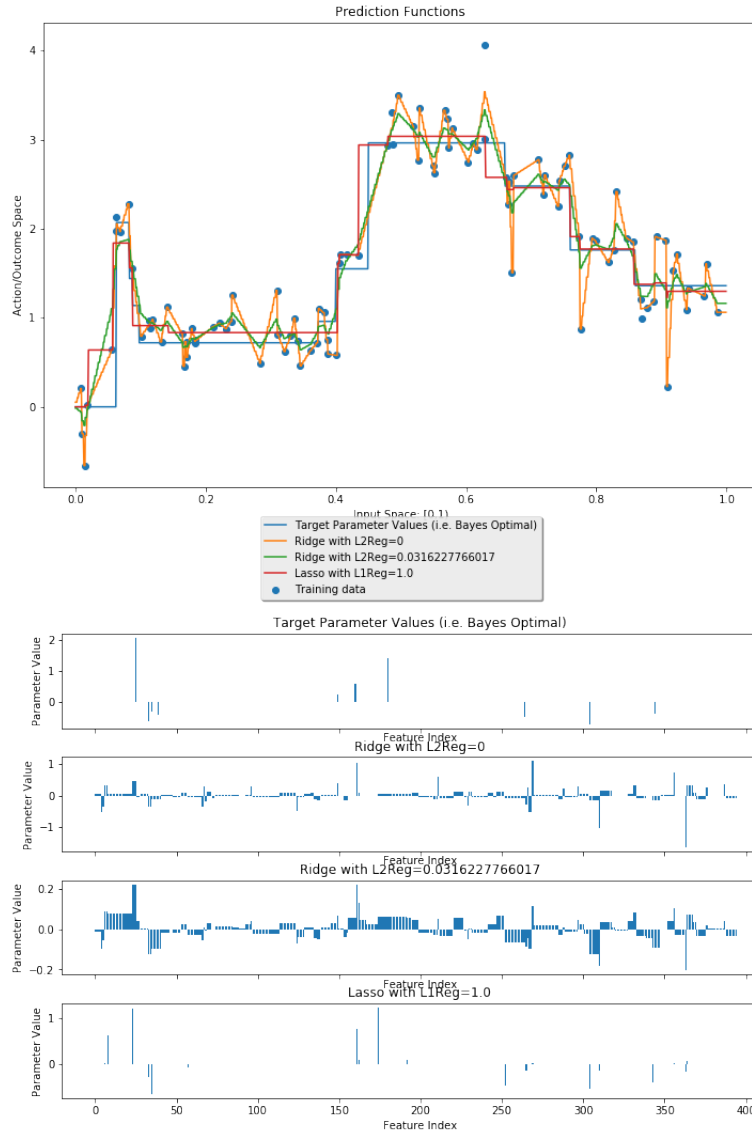
```

---



|    | lambda value (10 <sup>*</sup> ) | validation error |
|----|---------------------------------|------------------|
| 0  | -6.0                            | 91.220809        |
| 1  | -5.5                            | 91.214366        |
| 2  | -5.0                            | 91.213619        |
| 3  | -4.5                            | 91.211285        |
| 4  | -4.0                            | 91.203904        |
| 5  | -3.5                            | 91.180578        |
| 6  | -3.0                            | 91.106946        |
| 7  | -2.5                            | 90.875421        |
| 8  | -2.0                            | 90.155156        |
| 9  | -1.5                            | 88.022245        |
| 10 | -1.0                            | 82.452074        |
| 11 | -0.5                            | 73.176563        |
| 12 | 0.0                             | 67.127535        |
| 13 | 0.5                             | 68.439045        |
| 14 | 1.0                             | 83.069391        |
| 15 | 1.5                             | 167.876487       |

The  $\lambda$  that minimizes the square error on the validation set is 1 in my case. The plot of prediction functions and compare parameter vectors are showed as follow:



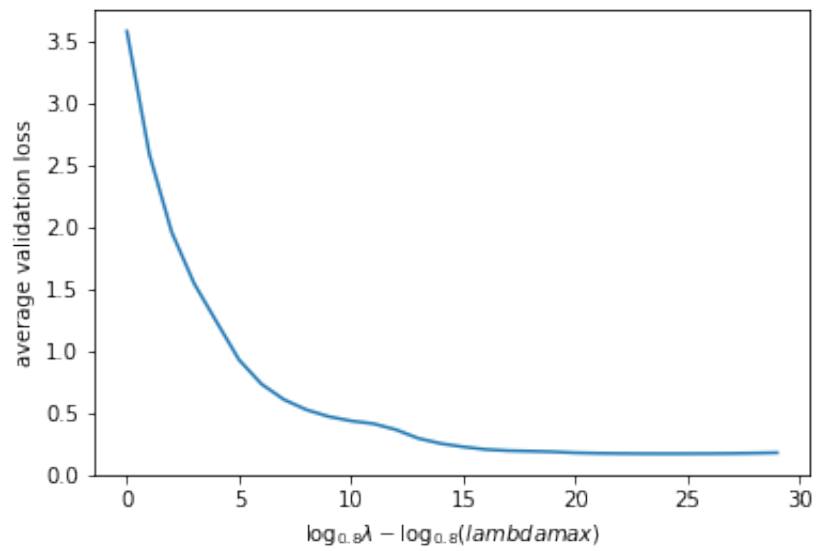
The lasso solutions' coefficients are much more sparse than the ridge ones. With less coefficients, lasso solutions are simpler and their performance are not bad compared to ridges as well. But the best model in terms of smaller validation error is the ridge regression with l2 regularization coefficient as 0.031623, and the smallest validation error is 0.139648.

### 3.1.4

---

```
1 lambda_max = 2*np.linalg.norm(np.dot(X_train.transpose(),y_train), ord = np.inf)
2 times_idx_list = np.arange(0,30, dtype=float)
3 validation_loss_homotopy = []
4 for idx in times_idx_list:
5     lambda_homotopy = lambda_max*0.8**idx
6     w = shooting_alg(X_train, y_train, w_ridge, lambda_homotopy)
7     validation_loss_homotopy.append(RMSE(X_val, y_val, w))
```

---



## 4 Lasso Properties

### 4.1 Deriving $\lambda_{Max}$

#### 4.1.1

$$\begin{aligned}
J'(0; v) &= \lim_{h \rightarrow 0} \frac{J(0 + hv) - J(0)}{h} \\
&= \lim_{h \rightarrow 0} \frac{J(hv) - J(0)}{h} \\
J(w) &= (Xw - y)^T (Xw - y) + \lambda \|w\|_1 \\
&= w^T X^T X w - 2w^T X^T y + y^T y + \lambda \|w\|_1 \\
J(hv) &= (hv)^T X^T X (hv) - 2(hv)^T X^T y + y^T y + \lambda \|hv\|_1 \\
&= h^2 v^T X^T X v - 2hv^T X^T y + y^T y + \lambda \|hv\|_1
\end{aligned}$$

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

#### 4.1.2

$$\begin{aligned}
J'(0; v) &\geq 0 \\
-2v^T X^T y + \lambda \|v\|_1 &\geq 0 \\
\lambda &\geq \frac{2v^T X^T y}{\|v\|_1}
\end{aligned}$$

The explicit expression for C is

$$\frac{2v^T X^T y}{\|v\|_1}$$

#### 4.1.3

$$\begin{aligned}
\frac{2v^T X^T y}{\|v\|_1} &= \frac{2 \sum v_i (X^T y)_i}{\|v\|_1} \leq \frac{2 \max(X^T y) \sum v_i}{\|v\|_1} \\
&= \frac{2 \|X^T y\|_\infty \sum v_i}{\|v\|_1} \leq \frac{2 \|X^T y\|_\infty \sum |v_i|}{\|v\|_1} = 2 \|X^T y\|_\infty
\end{aligned}$$

Therefore, the  $\max \frac{2v^T X^T y}{\|v\|_1} = 2 \|X^T y\|_\infty$ . If  $\lambda$  is larger than the greatest value, it means the  $J'(0; v)$  is greater than zero, which indicates J increases after  $w = 0$ .

So the minimum happens when  $w = 0$ ; If  $w = 0$  is the minimizer,  $J'(0; v) \geq 0$ , then  $\lambda \geq \frac{2v^T X^T y}{\|v\|_1}$  has to be satisfied whatever  $X, y, v$  values are given.  $\lambda \geq \max$ . So  $\lambda \geq 2\|X^T y\|_\infty$

#### 4.1.4 OPTIONAL

To take the partial derivative on the function:

$$\begin{aligned}
J_b(0, b) &= \lim_{h \rightarrow 0} \frac{((b+h)1 - y)^T((b+h)1 - y) - (b1 - y)^T(b1 - y)}{h} \\
&= \lim_{h \rightarrow 0} \frac{2bh1^T 1 + h^2 1^T 1 - 2hy^T 1}{h} \\
&= \lim_{h \rightarrow 0} 2b1^T 1 + h1^T 1 - 2y^T 1 \\
&= 2b1^T 1 - 2y^T 1
\end{aligned}$$

To minimize the partial derivative:

$$b = \frac{y^T 1}{1^T 1} = \frac{y_1 + y_2 + \dots + y_n}{n} = \bar{y}$$

$$\begin{aligned}
J(w) &= \|Xw + (\bar{y}) - y\|_2^2 + \lambda \|w\|_1 \\
&= (Xw - (y - (\bar{y})))^T (Xw - (y - (\bar{y}))) + \lambda \|w\| \\
&= w^T X^T w - 2w^T X^T (y - \bar{y}) + (y - \bar{y})^T (y - \bar{y}) + \lambda \|w\| \\
J(0; v) &= \lim_{h \rightarrow 0} \frac{J(hv) - J(0)}{h} \\
&= \lim_{h \rightarrow 0} \frac{h^2 v^T X^T X v - 2hv^T X^T (y - \bar{y}) + (y - \bar{y})^T (y - \bar{y}) + \lambda h \|v\| - (y - \bar{y})^T (y - \bar{y})}{h} \\
&= \lim_{h \rightarrow 0} \frac{h^2 v^T X^T X v - 2hv^T X^T (y - \bar{y}) + \lambda h \|v\|}{h} \\
&= \lim_{h \rightarrow 0} hv^T X^T X v - 2v^T X^T (y - \bar{y}) + \lambda \|v\| \\
&= -2v^T X^T (y - \bar{y}) + \lambda \|v\| \geq 0 \\
\lambda &\geq \frac{2v^T X^T (y - \bar{y})}{\|v\|}
\end{aligned}$$

According to previous problems:

$$\frac{2v^T X^T (y - \bar{y})}{\|v\|} \leq \frac{2\|X^T (y - \bar{y})\|_\infty \sum |v|}{\sum |v|} = 2\|X^T (y - \bar{y})\|_\infty$$

## 4.2 Feature Correlation

### 4.2.1

Suppose a and b have different signs:  $a < 0, b > 0$  or  $a > 0, b < 0$ :

$$\begin{aligned} J(\hat{\theta}) &= \|x_1(a+b) + x_r r - y\|_2^2 + \lambda|a| + \lambda|b| + \lambda\|r\|_1 \\ &= \|x_1(a+b) + x_r r - y\|_2^2 + \lambda(|a| + |b|) + \lambda\|r\|_1 \end{aligned}$$

we set  $c+d = a+b$ :

Seen from this equation, we can use  $c = \theta(a+b), d = (1-\theta)(a+b)$  to plug into the equation:

$$\begin{aligned} J(\theta^*) &= \|x_1(c+d) + x_r r - y\|_2^2 + \lambda|c| + \lambda|d| + \lambda\|r\|_1 \\ &= \|x_1(a+b) + x_r r - y\|_2^2 + \lambda\theta|a+b| + \lambda(1-\theta)|a+b| + \lambda\|r\|_1 \\ &= \|x_1(a+b) + x_r r - y\|_2^2 + \lambda|a+b| + \lambda\|r\|_1 \end{aligned}$$

since  $|a+b| \leq |a| + |b|$ .  $J(\hat{\theta})$  won't be the minimizer of  $J(\theta)$ , which is a contradiction. If  $J(\hat{\theta})$  is a minimizer, we should ensure  $|a+b| = |a| + |b|$ . In order to make the equality hold, a and b have to be the same sign, or at least one of them be zero.

The  $J(\theta^*)$  also needs to minimize, meaning  $(c, d, r^T)^T$  is the minimizer:  $J(\theta^*) = J(\hat{\theta})$ .

Hence,  $c+d = a+b, |c| + |d| = |a| + |b|$

### 4.2.2

$$\begin{aligned} J(\theta) &= \|X\theta - y\|_2^2 + \lambda\|\theta\|_2^2 \\ &= \|x_1 a + x_2 b + x_r r - y\|_2^2 + \lambda \sum \theta_i^2 \\ &= \|(a+b)x_1 + x_r r - y\|_2^2 + \lambda a^2 + \lambda b^2 + \lambda\|r\|_2^2 \end{aligned}$$

to minimize  $\lambda(a^2 + b^2)$ :

since  $a^2 + b^2 \geq 2ab$ , the minimum of  $a^2 + b^2$  happens when  $a = b$

Hence,  $a = b$  when minimizes the ridge regression objective function.