

Problem Set One

Matlab, Python and Computational Assignments

1.

• **Algorithm 1 : Least square**

- (a) Did CVX succeed?
- (b) If so, how long did it take to solve each instance?
- (c) Report the Regression error of the solution computed: $\|X\beta^* - y\|_2$ and also the Testing error: $\|X_{\text{test}}\beta - y_{\text{test}}\|_2$.

Ans:

- (a) In ps1A, there are three problems. $X_1 \in \mathbb{R}^{50 \times 500}$, $X_2 \in \mathbb{R}^{500 \times 5000}$ and $X_3 \in \mathbb{R}^{5000 \times 50000}$. CVX succeed when performing Least Square Method for first and second problems. However, when it comes to third problem, CVX is stopped because of out of memory. Maybe it succeed in a better processor. So it isn't sure whether CVX succeed for third problem or not.
- (b) In first problem, it took less than one second to solve each instance. And in the second problem, it took almost six minutes to solve each instance. However, for third problem, CVX failed. So we didn't how long it took to solve each instance.
- (c) For both first and second problems, the least square method find β^* such that the regression errors are sufficient small. However, we can find that the testing errors are large either first and second problem. So, β^* we found is still very different from the true β .

Least Square	Running Time(sec)	$\ X\beta^* - y\ _2$	$\ X_{\text{test}}\beta - y_{\text{test}}\ _2$
Problem One	0.64	6.106×10^{-11}	1.98×10^1
Problem Two	331.52	3.602×10^{-9}	2.07×10^1
Problem Three	-	-	-

• **Algorithm 2 : Sparse Recovery via an optimization-based algorithm called LASSO.**

- (a) Did CVX succeed?
- (b) If so, how long did it take to solve each instance?
- (c) Report the Regression error of the solution computed: $\|X\beta^* - y\|_2$ and also the Testing error: $\|X_{\text{test}}\beta - y_{\text{test}}\|_2$.
- (d) What is the support of β ? That is, what are the non-zero coefficients of β .

Ans:

(In this question, I set λ to be 0.1.)

- (a) In ps1A, there are three problems. $X_1 \in \mathbb{R}^{50 \times 500}$, $X_2 \in \mathbb{R}^{500 \times 5000}$ and $X_3 \in \mathbb{R}^{5000 \times 50000}$. Like least square method, CVX succeed for first and second problems. However, it failed on third problem because of out of memory. Maybe it succeed in a better processor. It isn't sure whether CVX succeed for third problem or not.
- (b) In first problem, it took almost one second to solve each instance. And in the second problem, it took several minutes to solve each instance. However, for third problem, CVX failed. So we didn't how long it took to solve each instance.
- (c) For both first and second problems, the least square method find β^* such that the regression errors are sufficient small. However, unlike Least Square Method, LASSO finds β^* closer to the true β since the testing errors are very small.

LASSO	Running Time(sec)	$\ X\beta^* - y\ _2$	$\ X_{\text{test}}\beta - y_{\text{test}}\ _2$
Problem One	1.14	1.019×10^{-9}	1.27×10^{-1}
Problem Two	104.42	2.468×10^{-8}	7.73×10^{-2}
Problem Three	-	-	-

- (d) According to the SparseRegressionData.m, we know that β is composed of five 1's and other elements are 0's. Hence, the non-zero coefficients of β is 1. On the other hand, if we see the composition of β^* we find by LASSO, we can find that there are 5 elements approach to 1, which have same index as the five 1's of β do. And the other elements of β^* are approach to zero. So, the support of β^* is exactly the same as the support of the true β

2. (OMP – Orthogonal Matching Pursuit)

- (a) What is the sparsity pattern found?
- (b) How long does the solution take?
- (c) What are the regression and testing errors?

Ans:

(a) If we take a look at each element of β^* for each problem, we can find that there are only five elements are not zero and approach to one. Hence, the sparsity pattern found by OMP is exactly the same as the original sparsity pattern for all three problems.

(b) This solution runs much faster than the above two algorithms. Even the third problem took only within a second. The detailed running time are showed on the below table.

(c) The regression errors aren't as small as the above two algorithms. But the testing errors are much smaller than other two algorithms. The detailed errors are showed on the below table.

OMP	Running Time(sec)	$\ X\beta^* - y\ _2$	$\ X_{\text{test}}\beta - y_{\text{test}}\ _2$
Problem One	0.005	6.270×10^{-2}	3.97×10^{-2}
Problem Two	0.010	2.167×10^{-1}	1.19×10^{-2}
Problem Three	0.820	6.983×10^{-1}	2.60×10^{-3}

4.

We want to run the gradient descent algorithm which iteratively computes

$$\beta^{(n+1)} = \beta^{(n)} - \gamma \nabla f(\beta^{(n)})$$

where γ is a constant step size. The initial $\beta^{(0)}$ is the all-one vector.

For each matrix, find the range of γ that the solution converges to zero and the range of γ that the algorithm diverges, and explain why. Take example values of γ to illustrate the two behaviors, convergence to zero and divergence. Plot $f(\beta^{(n)})$ over n for the two of your values.

Ans:

Because X is symmetric and positive definite matrix, there always exists an eigendecomposition of $X = U\Lambda U^T$ where $U \in \mathbb{R}^{m \times m}$ is a unitary matrix and $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ is a diagonal matrix with positive eigenvalues $\lambda_1 \geq \dots \geq \lambda_m > 0$.

We want to run the gradient descent algorithm which iteratively computes

$$\beta^{(n+1)} = \beta^{(n)} - \gamma \nabla f(\beta^{(n)}) \quad (\nabla f(\beta^{(n)}) = X\beta^{(n)})$$

$$\Rightarrow \beta^{(n+1)} = \beta^{(n)} - \gamma X\beta^{(n)} \quad (X = U\Lambda U^T)$$

$$\Rightarrow \beta^{(n+1)} = \beta^{(n)} - \gamma U\Lambda U^T \beta^{(n)}$$

$$\Rightarrow \beta^{(n+1)} = U(I - \gamma\Lambda)U^T \beta^{(n)}$$

Now we can let $\beta' = U^T \beta$

$$\Rightarrow U^T \beta^{(n+1)} = U^T U (I - \gamma\Lambda) U^T \beta^{(n)} \quad (U \text{ is unitary } \Rightarrow U^T U = I)$$

$$\Rightarrow \beta'^{(n+1)} = (I - \gamma\Lambda) \beta'^{(n)} = (I - \gamma\Lambda)^2 \beta'^{(n-1)} = \dots = (I - \gamma\Lambda)^{(n+1)} \beta'^{(0)}$$

Then if we consider for each component of β' as β'_i and $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$

$$\text{Hence, } \beta'_{i^{(n)}} = (1 - \gamma\lambda_i)^n \beta'_{i^{(0)}}$$

If we want to make sure the solution converges to zero, $|1 - \gamma\lambda_i|$ should < 1 for all $i = 1, 2, \dots, n$

$$\Rightarrow 1 - \gamma\lambda_i < 1 \text{ or } \gamma\lambda_i - 1 < 1 \Rightarrow 0 < \gamma\lambda_i \text{ or } \gamma\lambda_i < 2$$

(Because X is symmetric and positive definite matrix, $0 < \lambda_i$)

$$\Rightarrow 0 < \gamma < (2/\lambda_i)$$

For (a), all eigenvalues are one $\lambda_i=1$, that is . So it converges when $0 < \gamma < 2$ and diverges when $2 < \gamma$.

For (b), a half of the eigenvalues are one and the other half of them are very small. $0 < \gamma < (2/\lambda_i)$ must satisfy for all λ_i . Hence, it converges when $0 < \gamma < (2/\max(\lambda_i)) = 0 < \gamma < 2$ and diverges when $2 < \gamma$.

For (c), all other than a few very large eigenvalues. We can know that the large eigenvalues are 100 from the data given by the instructor. Like (b), it converges when $0 < \gamma < (2/\max(\lambda_i)) = 0 < \gamma < 2/100 = 0 < \gamma < 0.02$ and diverges when $0.02 < \gamma$.

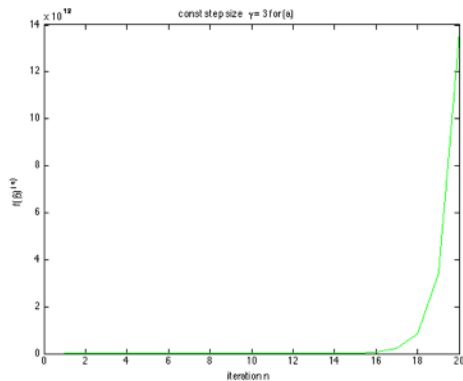
The below table are convergence and divergence conditions for (a), (b) and (c) as stated above.

	convergence	divergence
(a)	$0 < \gamma < 2$	$2 < \gamma$
(b)	$0 < \gamma < 2$	$2 < \gamma$
(c)	$0 < \gamma < 0.02$	$0.02 < \gamma$

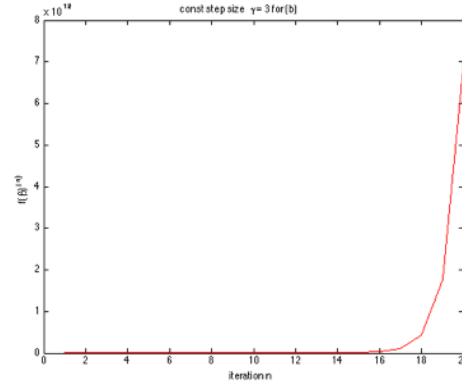
Then using Matlab to plot $f(\beta^{(n)})$ over n for two different γ :

First, $\gamma = 3$, $n = 20$

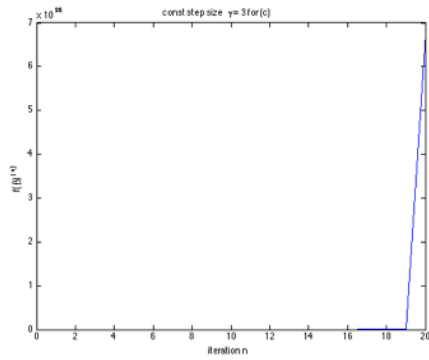
(a)



(b)



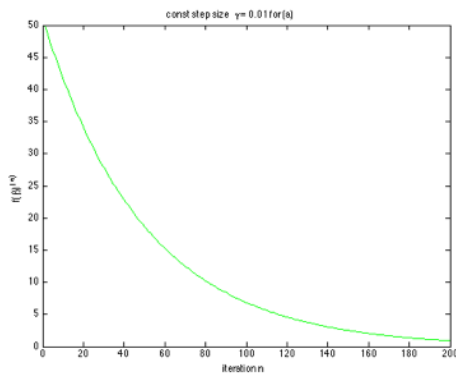
(c)



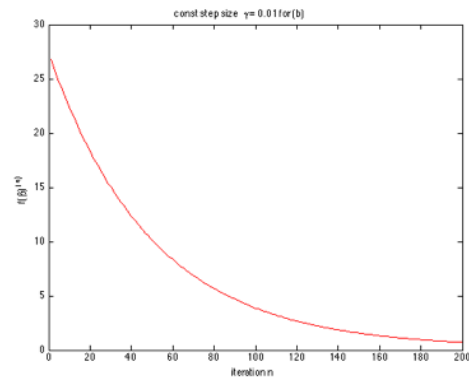
From the above three figures, it's easy to find that all three matrices diverge. And $\gamma = 3$ satisfies the divergence condition of all three matrices.

Then, $\gamma = 0.1$, $n = 20$

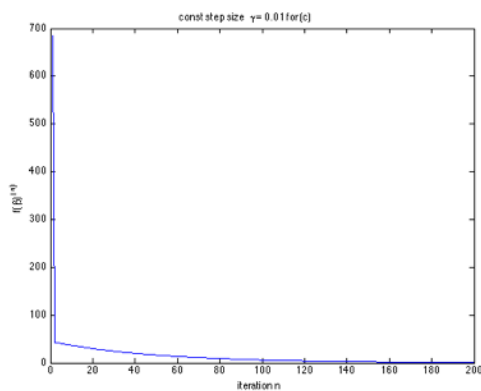
(a)



(b)



(c)

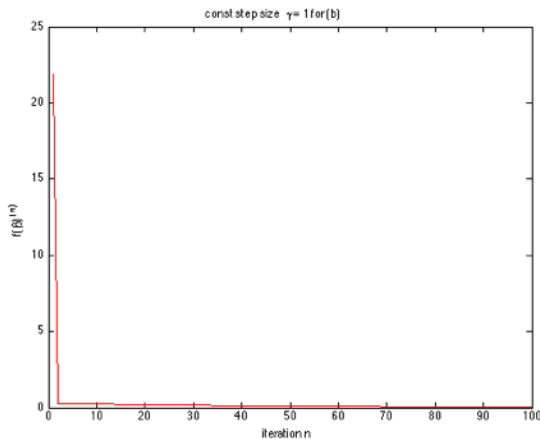


From the above three figures, it's easy to find that all three matrices converge. And $\gamma = 0.01$ satisfies the convergence condition of all three matrices.

5. Take $\gamma = 1$, and plot $f(\beta^{(n)})$ over n for the second matrix (b) of the above three. Explain the convergence behavior of the solution based on the plot.

Ans:

Take $\gamma = 1$, and plot $f(\beta^{(n)})$ over n for the second matrix (b) as below figure:



$$f(\beta) = 1/2 * \beta^T X \beta = 1/2 * \beta^T U \Lambda U^T \beta \quad (\text{If let } \beta' = U^T \beta)$$

$$\Rightarrow = 1/2 * \beta'^T \Lambda \beta' = (1/2) * \sum_{i=1}^m \lambda_i \beta_i'^2$$

As a result, we know that $f(\beta)$ depends on β' .

Besides, if we set $\gamma = 1$ for the matrix (b), because (b) a half of the eigenvalues are one and the other half of them are 0.01 (According to the data provided by the instructor.), we can know

$$\beta'_{i^{(n)}} = (1 - 1 * 1)^n \beta'_{i^{(0)}} = 0 \text{ for } i = 1, 2, \dots, 50$$

$$\beta'_{i^{(n)}} = (1 - 1 * 0.01)^n \beta'_{i^{(0)}} = 0 \text{ for } i = 51, 52, \dots, 100$$

Because of the first 50 components of β'_i equal to 0, there is a drop at the first iteration. On the other hand, because of the last 50 components of β'_i decreasing geometrically with a factor of 0.99, they contribute to the linear convergence to 0 after the first iteration.

Written Problems:

1. Over-determined

$$\textcircled{1} \|X\beta - y\|_2^2 = \beta^T X^T X \beta - \beta^T X^T y - y^T X \beta + y^T y$$

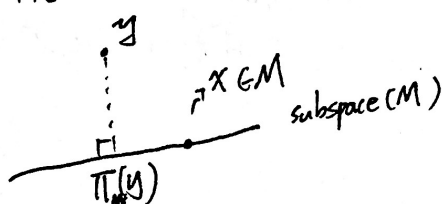
$$\nabla^2 f(\beta) = 2X^T X \quad \therefore z^T \nabla^2 f(\beta) z = 2z^T X^T X z = 2\|Xz\|^2 \geq 0, \forall z$$

$\Rightarrow f(\beta)$ is a convex function of β

$$\text{If } \nabla f(\beta_{LS}) = 0 \Rightarrow 2X^T X \beta_{LS} - X^T y - (y^T X)^T = 0$$

$$\Rightarrow \beta_{LS} = (X^T X)^{-1} X^T y$$

$\textcircled{2}$ Pic:



As the description of the question says,

$$\langle y - \pi_M(y), x \rangle = 0 \text{ for all } x \in M$$

We can say that $M = \text{Range}(X)$

$$\pi_M(y) = \beta_{LS}$$

$$\text{Hence, } \langle y - \pi_M(y), x \rangle = \langle y - X\beta_{LS}, X\beta \rangle = 0, \forall \beta$$

Because $\forall \beta$, it all satisfies

$$\Rightarrow \langle y - X\beta_{LS}, X \rangle = 0$$

$$(y - X\beta_{LS})^T \cdot X = 0$$

$$y^T X - \beta_{LS}^T X^T X = 0 \Rightarrow \beta_{LS} = (y^T X (X^T X)^{-1})^T = (X^T X)^{-1} X^T y$$

Under-determined

$$\textcircled{3} \min: \|\beta\|_2^2$$

$$\text{subject to } X\beta = y$$

Consider any other solution $\beta_1 = \beta_0 + z$

$$\|\beta_1\|^2 = \|\beta_0 + z\|^2 = \|\beta_0\|^2 + \|z\|^2 \quad (\because \beta_0 \perp z \Rightarrow \beta_0^T z = 0)$$

$$\therefore \|z\|^2 \geq 0 \Rightarrow \|\beta_1\|^2 \geq \|\beta_0\|^2$$

So, β_0 is the minimum norm solution

$$\textcircled{4} y = X\beta_0 \text{ \& } \beta_0 \perp z \text{ for any } z \in \text{Null}(X)$$

$\therefore \text{Null}(X)$ is the set of vectors perpendicular to the rows of X

\therefore The set of vectors perpendicular to $\text{Null}(X)$ must be in span of the rows of X

$\Rightarrow \beta_0$ is in the span of the rows of $X \Rightarrow \beta_0 = X^T z = \sum z_i X_i$ (X_i is the i th row of X for some vector z)

$$\begin{aligned} y &= X\beta_0 \\ &= XX^T z \\ z &= (XX^T)^{-1} y \\ \beta_0 &= X^T z \\ &= X^T (XX^T)^{-1} y \end{aligned}$$

12. $C = \{x \in \mathbb{R}^n : x^T A x + b^T x + c \leq 0\}$ where $A \in S^n$, $b \in \mathbb{R}^n$, $c \in \mathbb{R}$

(a) if $A \in S_+^n \Rightarrow x^T A x \geq 0$

Consider $\{x_1 + tv | t \in \mathbb{R}\}$ an arbitrary line

C 's intersection with $\{x_1 + tv | t \in \mathbb{R}\}$

$$\Rightarrow (x_1 + tv)^T A (x_1 + tv) + b^T (x_1 + tv) + c \leq 0$$

$$\Rightarrow v^T A v t^2 + (2x_1^T A v + b^T v) t + (b^T x_1 + c + x_1^T A x_1) \leq 0$$

It's convex because $A \in S_+^n$ s.t. $v^T A v \geq 0$

So C is convex because its intersection with an arbitrary line is convex.

(b) we consider the intersection of $C \cap H = C_1$ with an arbitrary \neq line $\{x_1 + tv | t \in \mathbb{R}\}$ Q.E.D.
Without loss of generality, we can assume $g^T x_1 + h = 0$
Hence, the intersection defined by x_1 & v is

$$\{x_1 + tv \mid v^T A v t^2 + (2x_1^T A v + b^T v) t + (b^T x_1 + c + x_1^T A x_1) \leq 0, g^T t v = 0\}$$

If $g^T v = 0$, the intersection is the singleton $\{x_1\}$

The set can be reduced to $\{x_1 + tv \mid v^T A v t^2 + (2x_1^T A v + b^T v) t + (b^T x_1 + c + x_1^T A x_1) \leq 0\}$
which is convex if $v^T A v \geq 0$

therefore C_1 is convex if $g^T v = 0 \Rightarrow v^T A v \geq 0$

Consider $\lambda g g^T \in \mathbb{R}^{n \times n}$ $\lambda \in \mathbb{R}$

$$v^T A v = v^T (A + \lambda g g^T) v \geq 0 \quad (\because \lambda g g^T v = 0)$$

In conclusion, C_1 is convex if there exists $\lambda \in \mathbb{R}$ s.t. $A + \lambda g g^T \in S_+^n$

13. $S = \{(a, b) : a^T x \leq b \forall x \in C, a^T x \geq b \forall x \in D\}$ Q.E.D. \neq

It forms a set of homogenous linear inequalities in (a, b)

That means it's the intersection of many half spaces that pass through the origin

$$\Rightarrow \because a^T x \leq b \Rightarrow a^T x - b \leq 0 \Rightarrow \begin{bmatrix} x^T & -1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \leq 0 \quad \forall x \in C \quad \& \quad a^T x \geq b \Rightarrow \begin{bmatrix} x^T & 1 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \geq 0 \quad \forall x \in D$$

Hence, it will form a convex cone.

So S is convex.

17. We can consider a hyperplane perpendicular to $(v_2 - v_1)$ & lying on $\frac{v_1 + v_2}{2}$
 Suppose this hyperplane as $(v_2 - v_1)^T x = a$

$$\Rightarrow a = (v_2 - v_1)^T \cdot \frac{v_1 + v_2}{2} = \frac{1}{2} (\|v_2\|^2 - v_1^T v_2 + v_2^T v_1 + \|v_1\|^2) \\ = \frac{1}{2} (\|v_2\|^2 - \|v_1\|^2)$$

$$\{x: \|x - v_1\| \leq \|x - v_2\|\} = \{x: \|x - v_1\|^2 \leq \|x - v_2\|^2\} \\ = \{x: \|x\|^2 - 2v_1^T x + \|v_1\|^2 \leq \|x\|^2 - 2v_2^T x + \|v_2\|^2\} \\ = \{x: 2v_2^T x - 2v_1^T x \leq \|v_2\|^2 - \|v_1\|^2\} \\ = \{x: (v_2^T - v_1^T)x \leq \frac{\|v_2\|^2 - \|v_1\|^2}{2}\} \\ c = v_2^T - v_1^T \quad d = \frac{\|v_2\|^2 - \|v_1\|^2}{2} \quad \#$$

18. $A \in \mathbb{R}^{k \times m}$ $B \in \mathbb{R}^{k \times n}$

For every $x \in \mathbb{R}^m$, $Ax = 0 \Rightarrow Bx = 0 \Rightarrow \text{Null}(A) \subseteq \text{Null}(B)$

From 11.(c), we know that $U \subseteq W \Leftrightarrow U^\perp \supseteq W^\perp$

So $\text{Null}(A) \subseteq \text{Null}(B) \Rightarrow \text{RowSpace}(A) \supseteq \text{RowSpace}(B) \\ \Rightarrow \text{Range}(A^T) \supseteq \text{Range}(B^T)$

For every $b \in \text{Range}(B^T)$, we can find a vector $c \in \mathbb{R}^n$

$$\text{s.t. } A^T c = b$$

Hence, For each column $b_i \in B^T$, we can find a $c_i \in \mathbb{R}^n$ s.t. $A^T c_i = b_i$

$$\Rightarrow A^T \begin{bmatrix} c_1 & c_2 & c_3 & \dots & c_k \end{bmatrix} = B^T = \begin{bmatrix} b_1 & b_2 & \dots & b_k \end{bmatrix}$$

$$A^T C = B^T \Rightarrow (A^T C)^T = (B^T)^T \Rightarrow C^T A = B$$

$$\because C \in \mathbb{R}^{n \times k} \quad \therefore C^T \in \mathbb{R}^{k \times n}$$

\Rightarrow There exists a $k \times n$ real matrix C
 s.t. $CA = B$

Q.E.D.

#