Q-3)

a)

It deals with the notion of Strong Convexity.

Now equation, $f_n(\beta) = \frac{1}{2N}\|y - x\beta\|_2^2$ is always Convex. But not

Strongly Convex when $N < P$.     (Statement 1)

Definition of Strong Convexity: Given a differentiable function $f : R^P \to R$, we say that it is strongly convex with parameter $\gamma > 0$ at $\theta \in R^P$ if the inequality

$$f(\theta') - f(\theta) \geq \nabla f(\theta)^T (\theta' - \theta) + \frac{\gamma}{2}\|\theta' - \theta\|_2^2$$

hold for all $\theta' \in R^P$.

In particular the function $f$ is strongly convex with parameter $\gamma$ around $\beta^* \in R^n$ if and only if the min eigen value of the Hessian matrix $\nabla^2 f(\beta)$ is atleast $\gamma$ for all vector $\beta$ in the neighbourhood of $\beta^*$.

From Statement 1

We get $\nabla^2 f(\beta) = x^T x / N$ for all $\beta \in R^P$. Thus, the least-square loss is strongly convex if and only if the eigen values of the $p \times p$ positive semidefinite matrix $x^T x$ are uniformly bounded away from zero.

Here $x^T x$ has Rand $\{N, P\}$ & hence is always Rank deficient & not Strongly Convex.

So we Relax notion of strong Convexity.

It is only necessary to impose a type of strong convexity condition for some subset $C \subset \mathbb{R}^p$ of possible perturbation vectors $v \in \mathbb{R}^p$. We say that a function $f$ satisfies __Restricted Strong Convexity__ at $\beta^*$ with respect to $C$ if there is a constant $\gamma > 0$ such that

$$\frac{v^T \nabla^2 f(\beta) v}{\|v\|_2^2} \geq \gamma \quad \text{for all nonzero } v \in C$$

and for all $\beta \in \mathbb{R}^p$ in a neighbourhood of $\beta^*$

for $f_n \to f_n(\beta) = \frac{1}{2N} \|y - X\beta\|_2^2$ (Linear Regression), this notion is equivalent to lower bounding the Restricted eigenvalues of the model matrix, requiring

$$\frac{\frac{1}{N} v^T X^T X v}{\|v\|_2^2} \geq \gamma \quad \text{for all nonzero } v \in C.$$

b) To explain $G(\hat{v}) \leq G(0)$

Definition of $f^n := \frac{1}{2N} \|y - x(\beta^* + v)\|_2^2 + \lambda_N \|\beta^* + v\|$,

$\hat{v} = \hat{\beta} - \beta^*$

Now, $G(\hat{v}) = \frac{1}{2N} \|y - x(\beta^* + \hat{v})\|_2^2 + \lambda_N \|\beta^* + \hat{v}\|$,

$$= \frac{1}{2N} \|y - x\hat{\beta}\|_2^2 + \lambda_N \|\hat{\beta}\|,$$

$$G(0) = \frac{1}{2N} \|y - x\beta^*\|_2^2 + \lambda_N \|\beta^*\|,$$

Now, $G(\hat{v}) \leq G(0)$, since

$$\frac{1}{2N} \|y - x\hat{\beta}\|_2^2 + \lambda_N \|\hat{\beta}\| \leq \frac{1}{2N} \|y - x\beta^*\|_2^2 + \lambda_N \|\beta^*\|,$$

Now Above inequality can be verified from the fact

$$\|y - x\hat{\beta}\|_2^2 \leq \|y - x\beta^*\|_2^2$$

$$\|\hat{\beta}\| \leq \|\beta^*\|,$$

c)

<u>Given</u> : $G(v) = \frac{1}{2N} \|y - X(\beta^* + v)\|_2^2 + \lambda_N \|\beta^* + v\|_1$   (11.20)

Putting $y = X\beta^* + w$ in $G(\hat{v}) \leq G(0)$

$\Rightarrow G(\hat{v}) \leq G(0) : \frac{1}{2N} \|X\beta^* + w - X\beta^* - X\hat{v}\|_2^2 + \lambda_N \|\beta^* + \hat{v}\|_1 \leq$

$\qquad\qquad\qquad \frac{1}{2N} \|X\beta^* + w - X\beta^*\|_2^2 + \lambda_N \|\beta^*\|_1$

$\Rightarrow \frac{1}{2N} \|w - X\hat{v}\|_2^2 + \lambda_N \|\beta^* + \hat{v}\|_1 \leq \frac{1}{2N} \|w\|_2^2 + \lambda_N \|\beta^*\|_1$

$\Rightarrow \frac{1}{2N}(w - X\hat{v})^T(w - X\hat{v}) - \frac{1}{2N}\|w\|_2^2 \leq \lambda_N \{\|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1\}$

$\Rightarrow \frac{1}{2N}\{\|w\|_2^2 - wX\hat{v} - (X\hat{v})^T w + \|X\hat{v}\|_2^2 - \|w\|_2^2\} \leq \lambda_N \{\|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1\}$

$\Rightarrow \frac{1}{2N}\{-2w^T X\hat{v} + \|X\hat{v}\|_2^2\} \leq \lambda_N \{\|\beta^*\|_1 - \|\beta^* - \hat{v}\|_1\}$

$\Rightarrow \frac{X\hat{v}}{2N} \leq \frac{w^T X\hat{v}}{N} + \lambda_N \{\|\beta^*\|_1 - \|\beta^* - \hat{v}\|_1\}$   (11.21)

d)

Given: $\dfrac{\|X\hat{v}\|_2^2}{2N} \leq \dfrac{w^T X\hat{v}}{N} + \lambda_N \left\{ \|\beta^*\|_1 - \|\beta^* + \hat{v}\|_1 \right\}$  (11.21)

by, $\|\beta^*\|_1 = \|\beta^*_s + \beta^*_{sc}\|_1$

$\|\beta^*\|_1 = \|\beta^*_s\|_1$,  $\because \beta^*_{sc} = 0$, and

$\|\beta^* + \hat{v}\|_1 = \|\beta^*_s + \hat{v}_s + \hat{v}_{sc}\|_1$

$= \|\beta^*_s + \hat{v}_s\|_1 + \|\hat{v}_{sc}\|_1$,  { Disjoint Joint }

$= \|\beta^*_s - (-1)\hat{v}_s\|_1 + \|\hat{v}_{sc}\|_1$,

$\geq \|\beta^*_s\|_1 - \|\hat{v}_s\|_1 + \|\hat{v}_{sc}\|_1$,  $\{\|a-b\|_1 \geqslant \|a\|_1 - \|b\|_1\}$

Substituting these relation into inequality (11.21) yields

$\dfrac{\|X\hat{v}\|_2^2}{2N} \leq \dfrac{w^T X\hat{v}}{N} + \lambda_N \left\{ \|\beta^*\|_1 - \|\beta^*_s\|_1 + \|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1 \right\}$

$\dfrac{\|X\hat{v}\|_2^2}{2N} \leq \dfrac{w^T X\hat{v}}{N} + \lambda_N \left\{ \|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1 \right\}$  ——— ①

from Holder's inequality :→

$\|f^T g\|_1 \leq \|f\|_p \|g\|_q$, where $P, q \in [1, \infty]$ with $\frac{1}{p} + \frac{1}{q} = 1$

Here, $\|w^T X\hat{v}\|_1 = \|X^T w\|_\infty \|\hat{v}\|_1$

Applying holder's in ①

$\dfrac{\|X\hat{v}\|_2^2}{2N} \leq \dfrac{w^T X\hat{v}}{N} + \lambda_N \left\{ \|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1 \right\} \leq$  ——(11.22)

$\dfrac{\|X^T w\|_\infty \|\hat{v}\|_1 + \lambda_N \left\{ \|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1 \right\}}{N}$

**e)** We know that,

$$\|\hat{v}_s\|_1 \leq \sqrt{k}\,\|\hat{v}_s\|_2 \leq \sqrt{k}\,\|\hat{v}\|_2$$

We will use this inequality **1**

**e)**

Since $\frac{1}{N}\|x^T w\|_\infty \leq \frac{\lambda N}{2}$ by assumption, eq. 11·22 becomes

$$\Rightarrow \frac{\|x\hat{v}\|_2^2}{2N} \leq \frac{\lambda N}{2}\|\hat{v}\|_1 + \lambda N\left\{\|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1\right\}$$

$$\Rightarrow \frac{\|x\hat{v}\|_2^2}{2N} \leq \frac{\lambda N}{2}\left\{\|\hat{v}_s\|_1 + \|\hat{v}_{sc}\|_1\right\} + \lambda N\left\{\|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1\right\}$$

$$\because \|v\|_1 = \|\hat{v}_s\|_1 + \|\hat{v}_{sc}\|_1 \quad \{\text{Disjoint set}\}$$

$$\Rightarrow \frac{\|x\hat{v}\|_2^2}{2N} \leq \frac{3\lambda N}{2}\|\hat{v}_s\|_1 - \frac{1}{2}\lambda N\|\hat{v}_{sc}\|_1 \leq \sqrt{k}\times\frac{3}{2}\lambda N\|\hat{v}\|_2 - \frac{1}{2}\lambda N\|\hat{v}_{sc}\|_1$$

$$\because \|\hat{v}_s\|_1 \leq \sqrt{k}\,\|\hat{v}_s\|_2 \leq \sqrt{k}\,\|\hat{v}\|_2$$

Hence,

$$\Rightarrow \frac{\|x\hat{v}\|_2^2}{2N} \leq \frac{\lambda N}{2}\left\{\|\hat{v}_s\|_1 + \hat{v}_{sc}\|_1\right\} + \lambda N\left\{\|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1\right\} \leq \frac{3}{2}\sqrt{k}\,\lambda N\|\hat{v}\|_2$$

$$(11\cdot2\underline{3})$$

f>
Lemma 11.1 allas us to apply $\gamma$-RE condition (11.10) to $\hat{v}$,

i.e. $\dfrac{\frac{1}{N} v^T X^T X v}{\|v\|_2^2} \geqslant \gamma$ OR $\dfrac{1}{N}\|X\hat{v}\|_2^2 \geqslant \gamma \|\hat{v}\|_2^2$.

Combining this with with inequality 11.23 gives the lower bound

$$\frac{\gamma \|\hat{v}\|_2^2}{2} \leq \frac{\|X\hat{v}\|_2^2}{2N} \leq \frac{3}{2}\sqrt{k}\,\lambda_N\|\hat{v}\|_2$$

$$\|v\|_2^2 \leq \frac{3}{\gamma}\sqrt{k}\,\lambda_N\|\hat{v}\|_2$$

$$\|v\|_2 \leq \frac{3}{\gamma}\sqrt{k}\,\lambda_N$$

$$\|\hat{\beta}-\beta^*\|_2 \leq \frac{3}{\gamma}\sqrt{\frac{k}{N}}\sqrt{N}\,\lambda_N$$

## 9

Given inequality: $\lambda_N \geq 2\frac{\|Xw\|_\alpha}{N}$

using this in inequality $(11.22)$ yields

$$\frac{\|X\hat{v}\|_2^2}{2N} \leq \frac{\lambda_N}{\alpha}\|\hat{v}\|_1 + \lambda_N\{\|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1\} \leq \frac{3}{\alpha}\sqrt{k}\,\lambda_N\|\hat{v}\|_2 \quad (11.23)$$

This also implies,

$\Rightarrow \quad 0 \leq \frac{\lambda_N}{\alpha}\|\hat{v}\|_1 + \lambda_N\{\|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1\}$,

$\Rightarrow \quad 0 \leq \lambda_N\{\|\hat{v}_s\|_1 + \|\hat{v}_{sc}\|_1\} + 2\lambda_N\{\|\hat{v}_s\|_1 - \|\hat{v}_{sc}\|_1\}$

$\Rightarrow \quad 0 \leq 3\lambda_N\|\hat{v}_s\|_1 - \lambda_N\|\hat{v}_{sc}\|_1$

$\Rightarrow \quad \|\hat{v}_{sc}\|_1 \leq 3\|\hat{v}_s\|_1 \qquad \{\text{proves Lemma } 11.1\}$

Hence the inequality $\lambda_N \geq 2\frac{\|Xw\|_\alpha}{N}$ used to prove lemma $11.1$

h)

From Defination of restricted eigenvalues

$$\frac{v^T \nabla^2 f(\beta) v}{\|v\|_2^2} \geq \gamma \quad \text{for all non zero } v \in C,$$

Now what constraint set C are relevant
for appropriate choices of the $l_1$-ball radius ·or equivalently, of the
regularization parameter $\lambda_N$ - it turns out that the lasso eracer
satisfies a cone cone constraint of the form

$$\|\hat{v}_{sc}\| \leq \alpha \|\hat{v}_s\|,$$

Now using this constrain we can successfully prove lemma 11·1 hence
Given a regularization parameter $\lambda_N \geq 2\|x^T \omega\|_\infty / N > 0$, any
estimate $\hat{\beta}$ from the regularized lasso (11·3) satisfies the
bound
$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{3}{\gamma} \sqrt{\frac{k}{N}} \sqrt{N} \lambda_N$$

(i) * Adv. of Theorem 11.1(b) Over Theorem 3

By Theorem 11.1(b) we know if
$\lambda_N \geq 2||X^T \omega||_\infty$ so, then $\hat{B}$ an estimate
of lasso regularized lasso statistics bound.

$$||\hat{B} - B^*||_2 \leq \frac{3}{\gamma} \sqrt{\frac{K \lambda_N \lambda_N}{N}}$$

but with theorem 3, if $\bar{A}$ is RIP
obeying matrix of $\delta_{2s} \leq 0.414$ then

$$||\hat{B} - B^*|| \leq \frac{C_0}{\sqrt{\cdot}} \cdot \frac{4\sigma \sqrt{1+\delta_{2s}}}{1 - \delta_{2s}(\sqrt{2}+1)}$$

where $B^*$ is $|S| = K$ sparse. $\hat{B}$ is
the solution for $\min ||B||_1 \; s.t \; ||Y - AB||_2 \leq$

So, by Example 11.1 we can calculate
the probable choice of $\lambda_N$ i.e

$$\lambda_N = 2\sigma \sqrt{\tau \frac{\log p}{N}}$$ where $\tau > 2$ is a
valid choice with high probability.

∴ theorem 11.1(b) becomes,

$$||\hat{B} - B^*||_2 \leq \frac{C\sigma}{\gamma} \sqrt{\frac{\tau K \log p}{N}}$$

① → which tells us that error is bounded $\frac{K}{N}$ i.e $l_2$-error decays more quickly with $\frac{K}{N}$. This error bound is

much better than error band given by theorem 3.

② → Even if we knows the support of $B^*$ then also the lasso rate (i.e 11.16 bound) is best possible w.r.t to the theorem3

→ As we decays by increase the $N$ i.e length of the vector

③ → As we increase the no. of measurement the error bounds will decrease drastically. But this is not the case withe Theorem3 because it doesnot depends on the number of measurment you take nan the spansity (when $B^*$ is $K$ spanse)

→ There Also size of $\log p$
④ → Also the effect of $\log p$ is very less with increase in the size input vector on error bound.

**\* Adv of Theorem 3 over Theorem 11.1 b**

→ Theorem 3 depends on $S_{25}$ value matrix $X$, & Theorem 11.1 b depends on $\gamma$ (restricted eigenvalue)

i.e $\underset{\wedge}{\text{Ron}}$ strong convextity

$$\frac{1}{N} \frac{v^T X^T X v}{||V||_2} \geq \gamma \text{ of all non zero } v \in C$$

where

$$C(S, \alpha) := \{ v \in \mathbb{R}^p \mid ||v_{S^c}||_1 \leq \alpha ||v_s||_1 \}$$

→ Error bounds of Theorem 3 & T1.1 b depends $S_{25}$ & $\gamma$ respectively. Therefore finding $\gamma$ w.n.t $S_{25}$ is much more difficult than that of $S_{25}$. Because if we know $X$ satisfy RIP of $S_{25} \leq 0.414$. the range of $S_{25}$ is much more limited say 0 to 0.414 an 0 to 1.

→ But in case of $\gamma$ it is much more complicated, creating $C$ set is a difficult job, also $\gamma$ can be from 0 to ∞ therefore, $\gamma$ finding is more difficult than $S_{25}$

→ Also Error bound of theorem Theorem 3 when $B^*$ is $K$ sparse is depend of $S_{25}$. So if we are able design $X$ with least $S_{25}$ then we will have least error bound. But in case Theorem 11.1 b it depends on too many parameters.