

SEEM 5380: Optimization Methods for  
High-Dimensional Statistics  
Term 2, 2017-18

Instructor: Prof. Anthony Man-Cho So  
Scribe: Zihao FU

2019.1.7

Natural of this course  
Less structured  
like a seminar  
Recent advance

方式:  
阅读论文  $\rightarrow$  讲座

What's the course about.

Integration of opt. and stat. techniques

Case study: Statistical estimation problems

- Samples:  $Z_1, Z_2, \dots, Z_n$   
n个样本,

- parameters:  $\theta_1, \theta_2, \dots, \theta_d$   
d个参数.



e.g.: linear regression

一般都先假设有一个 generative model

$$y = X \theta^* + \varepsilon \quad z = (x, y) \text{ samples}$$

$\mathbb{R}^n$                        $\mathbb{R}^d$                        $\mathbb{R}^n$  noise

$\theta^*$ : star means ground truth

Classical Setting:

$n \gg d$  (Overdetermined) 通常假设  $n$  远大于  $d$

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \|y - X\theta\|_2^2 \quad \text{Least-square estimator}$$

$\theta^*$  是 ground truth (We don't know)

$\hat{\theta}$  是 optimal solution to the opt. problem called estimator.

Note:  $\hat{\theta} \neq \theta^*$  in general

find the error bound between  $\hat{\theta}$  and  $\theta^*$

From Stat. we know:

If  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ , Then the

maximum likelihood estimator (MLE)

is given by  $\hat{\theta}$  (LSE/PF)

Issues/Observations:

① Optimization side:

how to solve for  $\hat{\theta}$

$\theta$  is convex  $\therefore$  is polynomial-time solvable

KKT 是充要.

$$\text{KKT: } \underbrace{X^T X}_{d \times d} \theta = X^T y$$

nearly full rank.

是否存在更加 lightweight method?

e.g. Gradient Descent.

Gradient  $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$

$$\theta^{k+1} \leftarrow \theta^k - \alpha_k X^T (X \theta^k - y)$$

$\alpha_k$  step size  $> 0$

那 performance 如何?

要做 convergence analysis

1) Does  $\{\theta^k\}$  converge? to where?

convex 情况下, 收敛到最优  
stably to where is important especially non-convex

2) (convergence Rate?)

$$F(\theta^k) - F(\hat{\theta}) \leq h(k)$$

current value, optimal value ← bound, 希望  $h(k) \rightarrow 0$

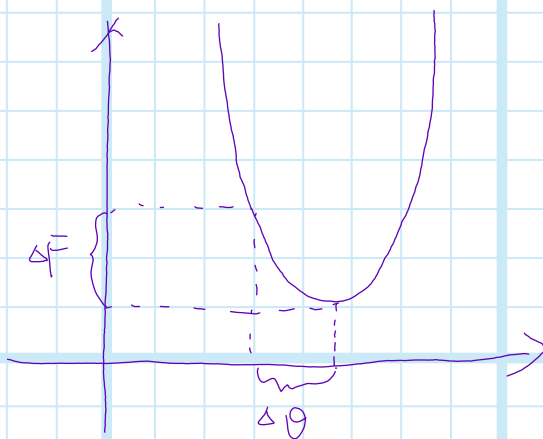
在上面的问题中  $F(\theta) = \|y - X\theta\|_2^2$

实际上  $F(\theta)$  的值并不重要, 重要的是找到  $\hat{\theta}$

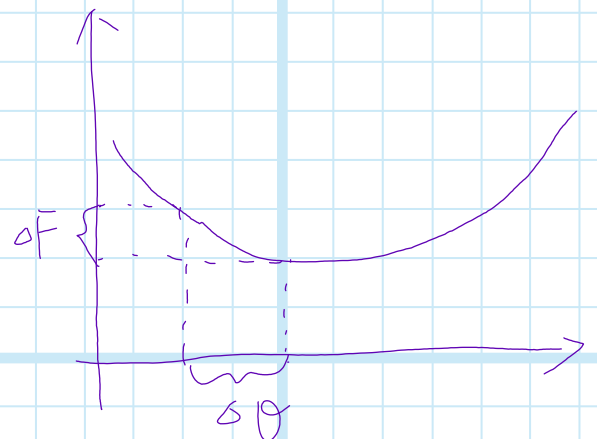
我们是否能直接找一个

$$\|\theta^k - \hat{\theta}\| \leq h_2(k) \rightarrow 0$$

因为即使  $f(\theta)$  收敛很快  $\|\theta^k - \hat{\theta}\|$  也不一定收敛快. 如下图:



High curvature



Low curvature

在 low curvature 中  $F$  快收敛了,  
 $\theta$  却还差很多. 对应高维就是 Hessian 矩阵.

(2) Statistical error:

$$\|\hat{\theta} - \theta^*\| \leq ?$$

虽然我们解出了最好的  $\hat{\theta}$ , 但是它和 ground truth 差多少?

optimational 方面只关注了  $\theta_k$  和  $\hat{\theta}$  之间的差异, 并未涉及到  $\theta^*$

Extensions to other settings

• in many applications:

$$y = X\theta^* + \varepsilon, \quad d \gg n \text{ (Underdetermined)}$$

$$n \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix} \cdot \begin{bmatrix} \theta^* \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} + \begin{bmatrix} \varepsilon \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$n \times d$

例如  $d$  是基因个数,  
 $n$  是病人样本数.

have too many parameter and few data  
can always fit it perfectly.

60~70年代, high dimensional 指的是样本很多  
现在指的是参数比样本多.

Need make more assumptions.

embed low-dimensional structure in  $\theta^*$

实际上就是假设数据具有一些结构信息，从而降维

假设  $\theta^*$  存在于某个特殊结构上。

e.g. sparsity: 虽然  $\theta$  维度很高, 但是很多维度是不起作用的。

如何加 sparsity 的信息:

Ideally: 转成 regularized problem.

$$\textcircled{1} \quad \hat{\theta} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \|y - X\theta\|_2^2 + \lambda \|\theta\|_0$$

$$\|\theta\|_0 \triangleq |\{i: \theta_i \neq 0\}|$$

$\|\theta\|_0$ :  $\theta$  中非零元的个数

$\lambda: > 0$ , regularization parameter.

同样两个问题:

① optimization aspect

② stats.

求解方法:

①  $\|\theta\|_0$  是非凸的, (都不是连续的)

0-norm 不好求, 直接用 1-norm 来近似.

$$\hat{\theta}' \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \{ \|y - X\theta\|_2^2 + \lambda \|\theta\|_1 \} \quad (\text{LASSO})$$

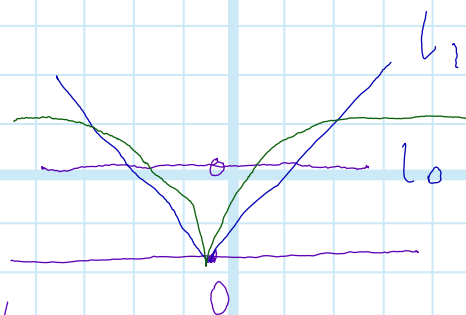
LASSO: least absolute shrinkage and

# Selection Operator.

注意到上面用的  $\hat{\theta} \in$ , 有一个位全部满足。具体用哪个要用 stats 的方法

② 先画出  $1 \times 1_0$  的图像:

$$1 \times 1_0 = \begin{cases} 1 & x \neq 0 \\ 0 & x = 0 \end{cases}$$



$l_1$  是近似  $l_0$ , 可以更进一步  
如绿线一样连续

近  $1 \times 1_0$

$$\hat{\theta}^R \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \|y - x\theta\|_2^2 + \lambda R(\theta) \right\}$$

$\downarrow$   
regularizer

好处是更逼近  $l_0$  了,

坏处是问题又变成了非凸的问题(困难)

③ 上面两种方法都是为了逼近  $l_0$ , 原因是  
① 是 intractable 的. 之所以是 intractable  
因为问题是 NP 的. 能不能直接解?

Mixed-integer optimization

问题是  $\text{practical tractable}$ , 即不在最坏情况下, 还是  $\text{tractable}$  的。所以并不能用到  $\text{deep learning}$  上。

2019.1.8

# Linear Regression

$$y = X\theta^* + w \quad \leftarrow \text{这是上节课的 } \varepsilon$$

$n$     $n \times d$     $d$     $n$

$\theta^*$ : ground-truth

$w$ : noise

$X$ : design matrix

Warm-up: classic setting  
 $n \gg d$

Least-squares estimate:

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|y - X\theta\|_2^2$$

Statistics Error:  $\|\hat{\theta} - \theta^*\|_2$

$$\textcircled{1} \quad \frac{1}{2} \|y - X\hat{\theta}\|_2^2 \leq \frac{1}{2} \|y - X\theta^*\|_2^2$$

optimality of  $\hat{\theta}$   $\therefore \hat{\theta}$  是最优解

$\textcircled{2}$  Define  $\hat{\Delta} = \hat{\theta} - \theta^*$   $\downarrow$  展开

$$\frac{1}{2} (y^T y - 2y^T X \hat{\theta} + \hat{\theta}^T X^T X \hat{\theta}) \leq \frac{1}{2} (y^T y - 2y^T X \theta^* + \theta^{*T} X^T X \theta^*)$$



利用  $\hat{\Delta}$  的定义, 以及  $X\theta^* = y - w$

展开后:  $\hat{\theta}^T X^T X \hat{\theta} - \hat{\theta}^T X^T X \theta^* \leq y^T X \hat{\Delta}$   
令  $\lambda = X\theta^* = w$ :  
 $\hat{\theta}^T X^T X \hat{\theta} + \hat{\theta}^T X^T X \theta^* - 2\hat{\theta}^T X^T X \hat{\theta} \leq 2w^T X \hat{\Delta}$

$$\|X\hat{\Delta}\|_2^2 \leq 2\hat{\Delta}^T X^T w$$

③ Assume  $X$  has full column rank (即  $n \geq d$ )  
 $\Rightarrow X^T X$  is invertible  $\lambda_{\min}(X^T X) > 0$

By the Courant-Fischer Theorem

$$\left( \lambda_{\min}(A) = \min_{\|x\|_2=1} x^T A x \quad A \in S^n \right)$$

$$\|X\hat{\Delta}\|_2^2 \geq \lambda_{\min}(X^T X) \cdot \|\hat{\Delta}\|_2^2$$

$$\begin{aligned} \text{④ } \lambda_{\min}(X^T X) \cdot \|\hat{\Delta}\|_2^2 &\leq \|X\hat{\Delta}\|_2^2 \leq 2\hat{\Delta}^T X^T w \\ &\leq 2\|\hat{\Delta}\|_2 \cdot \|X^T w\|_2 \quad (\text{Cauchy-Swartz}) \end{aligned}$$

$$\Rightarrow \|\hat{\Delta}\|_2 \leq \frac{2}{\lambda_{\min}(X^T X)} \|X^T w\|_2 \quad (*)$$

提出了一个 bound

Rmk:

① The bound (\*) is a deterministic bound with various statistical models for  $X, w$ , then the bound can be probabilistic

e.g.

①  $X_{ij} \sim \mathcal{N}(0, 1)$   $w$  bounded

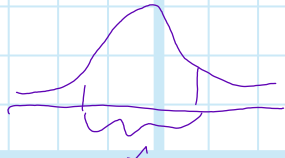
这里  $X$  确定, 这是一个 deterministic bound  
如果假设了  $X$  的矩阵分布, 有一个 high probabilistic bound.

$$X^T w = \sum_i w_i x_i \quad x_i: i\text{th col of } X^T$$

高斯变量的和仍是高斯变量.

②  $X$  has bounded col,  $w \in \mathcal{N}(0, \sigma^2 I)$

$\Rightarrow$  concentration of measure to get high-prob bounds on  $\lambda_{\min}(X^T X)$  or  $\|X^T w\|_2$

measure 还是 mass, concentration 是说  
是说  很短的  $x$  集中了很  
多的 mass.

③ We use Cauchy-Schwarz at the last step. 也可以用 Hölder-inequality

(C-S 不等式的推广), 因此:

$$2 \Delta^T X^T w \leq 2 \|\Delta\|_2 \|X^T w\|_2 \text{ 变成了:}$$

$$\leq 2 \|\Delta\|_p \|X^T w\|_q \quad \frac{1}{p} + \frac{1}{q} = 1$$

例如  $p=1, q=\infty$ , 在 sparsity 场合中更好.

$$\|u\|_2 \leq \|u\|_1 \leq \sqrt{n} \|u\|_2$$

$$\|u\|_1 = \sum_i |u_i| \leq \left(\sum_i |u_i|^2\right)^{1/2} \left(\sum_i 1^2\right)^{1/2}$$

如果  $u$  是 sparse vector,  $u = \begin{bmatrix} a \\ 0 \end{bmatrix}_{s \times n}$

把 2-norm 换成 1-norm 不会差太多 ( $\sqrt{s}$ )

Recall  $X$  is invertible  $X^T X > 0$  PD

Incidentally,  $X^T X$  is Hessian of the

loss function:  $L(\theta) = \frac{1}{2} \|y - X\theta\|_2^2$

$\Rightarrow L$  is strongly convex ( $\because$  Hessian PD)

Definition/Proposition We say  $f: \mathbb{R}^d \rightarrow \mathbb{R}$

is strongly convex with modulus  $c > 0$

(aka:  $c$ -strongly convex) if any of the

following equivalent conditions holds:

①  $\forall x, y \in \mathbb{R}^d, \alpha \in [0, 1]$

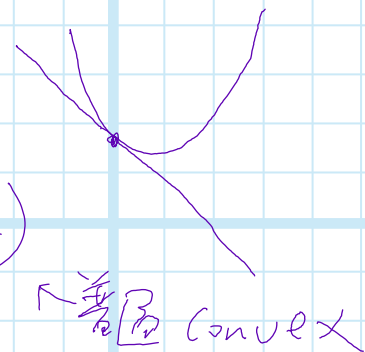
$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \underbrace{\frac{1}{2} c \alpha(1-\alpha) \|y-x\|_2^2}_{\text{减掉一些仍然 convex}}$$

② The function  $f(x) - \frac{1}{2} c \|x\|_2^2$  is convex  
(weakly convex 只是加一些东西后 convex)

③ In the presence of differentiability

$$\forall x, y \in \mathbb{R}^d$$

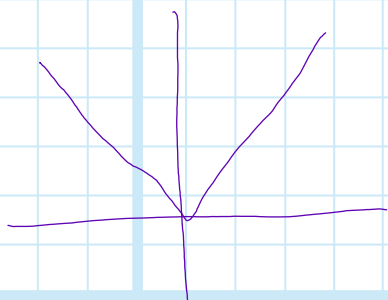
$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$



$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} c \|y - x\|_2^2$$

← strongly convex

lower bound  $L$  - 线性函数变成二次函数。



是 convex, 但不是 strongly convex,  
∴ 找不到一个二次函数的 LB

④ (In the presence of 2<sup>nd</sup> order diff)

$$\forall x \in \mathbb{R}^d, \quad v^T \nabla^2 f(x) v \geq c \cdot \|v\|_2^2 \quad \forall v \in \mathbb{R}^d$$

普遍 convex 的 Hessian 是 PSD, 而 strongly convex 的 Hessian 是 PD

下面看 strongly convex 如何帮助 optimization.

Optimization Aspect

Def: Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  be continuously diff.,

We say  $f$  has  $L$ -Lipschitz continuous gradient if  $\forall x, y \in \mathbb{R}^d$

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L \cdot \|x - y\|_2$$

原始的 Lipschitz 条件就是  $f(x) - f(y)$ , 现在扩展到梯度上。

e.g.  $\mathcal{L}(\theta) = \frac{1}{2} \|y - X\theta\|_2^2$ ,  $\nabla \mathcal{L}(\theta) = X^T(X\theta - y)$

$$\begin{aligned} \|\nabla \mathcal{L}(\theta_1) - \nabla \mathcal{L}(\theta_2)\|_2 &= \|X^T X(\theta_1 - \theta_2)\|_2 \\ &\leq \|X^T X\| \cdot \|\theta_1 - \theta_2\|_2 \\ &\quad \uparrow \text{还是那个 } L \end{aligned}$$

$\therefore$  我们之前那个  $\mathcal{L}(\theta)$  是 gradient-Lipschitz 连续的。

注意  $\mathcal{L}(\theta)$  不是 Lipschitz 连续。

Prop: Let  $f$  be  $C$ -strongly convex and have  $L$ -Lipschitz gradient. Then,  $\forall x, y \in \mathbb{R}^d$

$$(\nabla f(x) - \nabla f(y))^T (x - y) \geq \frac{CL}{C+L} \|x - y\|_2^2 + \frac{1}{C+L} \|\nabla f(x) - \nabla f(y)\|_2^2$$

Assume this prop, we study the convergence behavior of the gradient method for solving:

(P)  $\min_{\theta \in \mathbb{R}^d} f(\theta)$   $f$  is as in the prop (strongly convex, Lipschitz-grad)

Gradient method:  $\theta^{k+1} \leftarrow \theta^k - \alpha_k \nabla f(\theta^k)$

Thm: Let  $f$  be as in the prop, Suppose  $\alpha_k = \alpha$ ,  $\alpha \in (0, \frac{2}{c+L})$ , Then,

$$\|\theta^k - \hat{\theta}\|_2^2 \leq \left(1 - \frac{2\alpha cL}{c+L}\right)^k \|\theta^0 - \hat{\theta}\|_2^2 \quad (\Delta)$$

$\hat{\theta}$  is opt soln to (P)

observation:

就是说  $\theta^k$  收敛到  $\hat{\theta}$

$\left(1 - \frac{2\alpha cL}{c+L}\right)^k < 1$   $\therefore$  每一步都在减小.

Pf:  $\|\theta^{k+1} - \hat{\theta}\|_2^2 = \|\theta^k - \alpha \nabla f(\theta^k) - \hat{\theta}\|_2^2$

$$= \|\theta^k - \hat{\theta}\|_2^2 - 2\alpha \nabla f(\theta^k)^T (\theta^k - \hat{\theta}) + \alpha^2 \|\nabla f(\theta^k)\|_2^2$$

又注意到 gradient Lipschitz 有两项梯度, 而这里

只有一项, 注意到  $\nabla f(\hat{\theta}) = 0$  ( $\hat{\theta}$  is minimizer)

$$\therefore \nabla f(\theta^k)^T (\theta^k - \hat{\theta}) = (\nabla f(\theta^k) - \nabla f(\hat{\theta}))^T (\theta^k - \hat{\theta})$$

$$\geq \frac{cL}{c+L} \|\theta^k - \hat{\theta}\|_2^2 + \frac{1}{c+L} \|\nabla f(\theta^k)\|_2^2 \quad \text{代入得.}$$

= SCX prop 的结论.

$$\leq \left(1 - \frac{2\alpha cL}{c+L}\right) \|\theta^k - \hat{\theta}\|_2^2 + \underbrace{\alpha \left(\alpha - \frac{2}{c+L}\right)}_{\leq 0} \|\nabla f(\theta^k)\|_2^2$$

$\leq 0$  我们选的  $\alpha \in (0, \frac{2}{c+L})$

$$\leq \left(1 - \frac{22cL}{c+L}\right) \|\theta^k - \hat{\theta}\|_2^2$$

Q, E, D

对  $(\Delta)$  取对数:  $\ln(\|\theta^k - \hat{\theta}\|_2) \leq k \cdot c_0 + c_1$

$< 0$

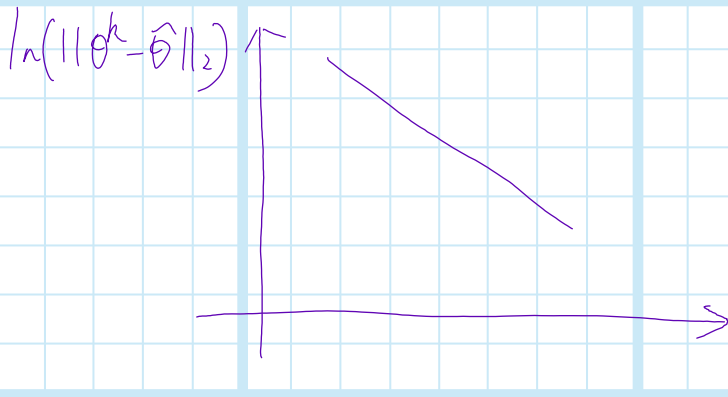
↓

Gradient method converges linearly for strongly

convex

线性收敛.

是指取对数后是linear



When moving to more general setting:

$$\gamma \|\hat{\Delta}\|_2^2 \leq \|X \hat{\Delta}\|_2^2 \leq \lambda \hat{\Delta}^T X^T W$$

如果  $n < d$ , 这  $\gamma$  就不成立了.

但能否让  $\gamma$  对当前  $\hat{\Delta}$  存在?

Restricted eigenvalues (RE)

↔ restricted strong convexity (RSC)

Loss function need not be convex everywhere, but in small region.

2019.1.14

Recall:

Standard linear model

$$y = X\theta^* + w$$

$n \quad n \times d \quad d \quad n$

consider  $n \gg d$

Least square estimator:

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \|y - X\theta\|_2^2$$

Interested:  $\|\hat{\theta} - \theta^*\|_2$

在之前的假设中,  
有  $n \gg d$ :

Crucial step:

Eigen value condition:

$$\|X\Delta\|_2^2 \geq \lambda_{\min}(X^T X) \|\Delta\|_2^2 \quad \forall \Delta \in \mathbb{R}^d$$

$\lambda_{\min}(X^T X) > 0$   
但如果  $n \ll d$ ,  $\lambda_{\min} = 0$   
就不能用上节的方法,  
需要假设一些结构.

If  $X$  has full  $rd$  rank, then  $\lambda_{\min}(X^T X) > 0$

Take  $\Delta = \hat{\theta} - \theta^*$  we can bound the error

本章研究更 general 的 case.

Today  $n \ll d$  underdetermined

$\rightarrow$  need additional constraints on the model.



→ typically, it is assumed  $\theta^*$  has some low-dim structure.  
e.g. sparsity, low-rankness, etc.

Regularization:

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \mathcal{L}(\theta, \{z\}_{z=1}^n) + \lambda R(\theta) \right\}$$

\*  $z_1, z_2, \dots, z_n$ : samples from space  $Z$   
e.g. in the linear model,

$$z_i = (x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$$

\*  $\mathcal{L}: \mathbb{R}^d \times Z^n \rightarrow \mathbb{R}$ : smooth convex loss function

\*  $R: \mathbb{R}^d \rightarrow \mathbb{R}_+$ : norm regularizer  
就是之前那个 norm.

\*  $\lambda > 0$ : regularization parameter

在线性回归中,  $\mathcal{L}(\theta, \{z_i\}) = \|y - X\theta\|_2^2$ ,

$R(\theta) = \|\theta\|_1$  且  $R(\theta)$  是 convex 的.

在  $n \gg d$  时, 之前那个不等式  $\|X\Delta\|_2^2 \geq |X^T X| \|\Delta\|_2^2$  成立,  
条件是  $X$  列满秩,  $\lambda_{\min}(X^T X) > 0$ .

当  $n \ll d$  时, 能否有一个弱一点的结论, 例如:

$$\|X\Delta\|_2^2 \geq c \|\Delta\|_2^2 \quad \forall \Delta \in \mathcal{S}$$

即在某个子域里面成立

之前是要求对  $\forall \Delta$  成立, 现在  
放松要求, 能否对部分  $\Delta$  成立  
加 regularizer 来实现.

下面看为什么加了 regularizer 就能构造出来这个新不等式.

## Decomposability of $R$

\* understand effect of  $R$  in inducing the desired structure.

\* let  $m \subseteq \bar{m} \subseteq \mathbb{R}^d$  be subspaces of  $\mathbb{R}^d$   
↓  
不是细胞, 只是大范围的子空间 (0 的空间)

•  $m$ : model subspace, intends to capture the low-dim structure in  $\mathcal{Q}^*$

•  $\bar{m}^\perp$ :  $\bar{m}^\perp = \{v \in \mathbb{R}^d : u^T v = 0 \ \forall u \in \bar{m}\}$

$\bar{m}^\perp$  叫 perturbation space, capture deviation from the model subspace

如果参数落到  $\bar{m}^\perp$  中, 则参数被惩罚

• for simplicity: take  $m = \bar{m}$

下面定义 regularizer 的 decomposability

Def:  $R$  is decomposable wrt  $(m, \bar{m}^\perp)$  if

$$R(\theta+r) = R(\theta) + R(r)$$

$$\forall \theta \in m, \forall r \in \bar{m}^\perp$$

注:  $R(\theta+r) \leq R(\theta) + R(r)$  是恒成立的, 因为  $R$  是 norm,  
 $r$  是我们不想要的, 所以想惩罚到最大,  
那么取等号的时候, 惩罚最大.  $r$  是 perturbation  
 $\bar{m}^\perp$  是 perturbation space

Rmk: The pair  $(m, \bar{m})$  can be chosen.

e.g.  $m = \mathbb{R}^d, \bar{m}^\perp = \{0\}$

此时  $r \equiv 0, R(\theta+r) = R(\theta) + R(r)$  显然成立

但这个 pair 用处不大,  
what would be a useful choice?

有用的  $(m, \bar{m})$  需满足两点:

(1)  $\theta^* \in m / \Pi_m(\theta^*) \approx \theta^*$   $\theta^*$  在  $m$  上的投影很接近自己

(2)  $m$  is "small" 不是想要的都会被惩罚

e.g. sparse linear model

$$y = X\theta^* + w \quad \theta^*: s\text{-sparse (has } s \text{ non-0 entries)}$$

$$R(\theta) = \|\theta\|_1$$

For a subset  $S \subseteq \{1, \dots, d\}$  of cardinality  $s$ ,

$$m(S) \triangleq \{\theta \in \mathbb{R}^d : \theta_j = 0 \forall j \notin S\}$$

只要不在  $S$  中的维度都为 0

if  $\theta^*$  is supported on  $S$ ,  $\Pi_{m(S)}(\theta^*) = \theta^*$

Take  $\bar{m} = m(S)$ , Then:

$\bar{m}^\perp$  是在  $S$  中的维度为 0

$$\bar{m}^\perp = \{\theta \in \mathbb{R}^d : \theta_j = 0, \forall j \in S\}$$

下面看 decomposability

(1)  $R(\theta) = \|\theta\|_1$

$$\theta \in m: \begin{array}{|c|} \hline \text{/// } S \\ \hline 0 \text{ } S^c \\ \hline \end{array}$$

$$\gamma \in \bar{m}^\perp: \begin{array}{|c|} \hline 0 \text{ } S \\ \hline \text{/// } S^c \\ \hline \end{array}$$

$$R(\theta + \gamma) = R(\theta) + R(\gamma) \text{ 成立}$$

(2) 另一个例子,  $x$  变成了一个矩阵.  $R$  用 nuclear norm

$$y_i = \langle x_i, \theta^* \rangle + w_i \quad \text{即: } y_i = \text{tr}(x_i^T \theta^*) + w_i$$

$$R(\theta) = \|\theta\|_X = \sum_{i=1}^{\min(m,d)} \sigma_i(\theta)$$

$m \times d \quad m \times d$   
rank  $r \leq \min(m,d)$

$$\sigma_i(\theta) = \sqrt{\lambda_i(\theta^T \theta)}$$

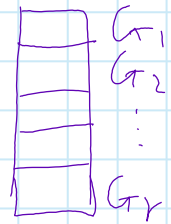
vector  $\theta$  的 non-zero entry 类似于矩阵 rank

(3) group sparsity

把  $\theta^*$  切成若干组相交的组, 希望大多数为 0

组是 sparse 的, 而组内不是 sparse 无所谓

$$R(\theta) = \sum_{i=1}^r \|\theta_{G_i}\|_2$$



1-norm 在 group, group 的 norm 无所谓

下面如何利用 decomposability 来构造  $\|\cdot\|_2 \rightarrow \|\cdot\|_2^2$   
 $\forall \Delta \in \mathcal{C}$  的形式.

### Consequence of Decomposability

Prop: Suppose  $\mathcal{L}$  is smooth convex.

$$\lambda \geq 2R^*(\nabla \mathcal{L}(\theta^*); \{z_i\}_{i=1}^n) \quad (\lambda \text{ 要足够大才有效})$$

$$R^*(v) = \text{dual norm of } R \quad \text{即: } \sup_{R(u) \leq 1} u^T v$$

$$\text{例如: } R(v) = \|\cdot\|_2, \text{ 则 } R^*(v) = \|\cdot\|_2 \quad \|\cdot\|_1 \leftrightarrow \|\cdot\|_\infty$$

dual is to minimize a linear function over a unit ball of original norm.

Then for any  $(m, \bar{m}^\perp)$  over which  $R$  is decomposable

$$\hat{\Delta} \triangleq \hat{\theta} - \theta^* \in \mathcal{C} \triangleq \left\{ \Delta \in \mathbb{R}^d : R(\Delta_{\bar{m}^\perp}) \leq 3R(\Delta_{\bar{m}}) + 4R(\theta_{\bar{m}^\perp}^*) \right\}$$

where  $\Delta_{\bar{m}} = \Pi_{\bar{m}}(\Delta)$  and so on.

可以观察到它现在不是整个空间了, 下面看每一项:

$4R(\theta_{m^\perp}^*)$ : misspecification of model

· 正常情况下  $\theta^*$  往 perturbation space  $m^\perp$  的投影应该为 0.

下面用 linear model 举例说明:

取  $d=3$ ,  $\Delta = (\Delta_1, \Delta_2, \Delta_3)$

$R(\Delta) = \|\Delta\|_1$ ,  $S = \{3\}$  假设参数只在第 3 维.

$m(S) = \{\Delta : \Delta_1 = \Delta_2 = 0\}$

$m^\perp(S) = \{\Delta : \Delta_3 = 0\}$

① 假设  $\theta_1^* = \theta_2^* = 0$  (model is exact)

$\theta^*$  就在  $m$  中, 那个  $\Theta$  变成:

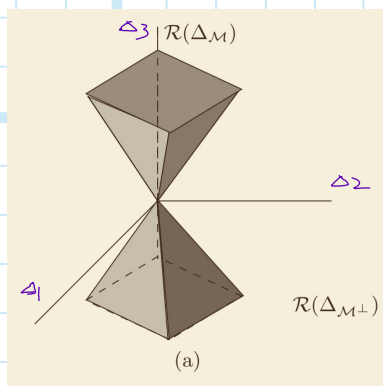
$\Theta = \{\Delta : \underbrace{|\Delta_1| + |\Delta_2|} \leq \underbrace{3|\Delta_3|}\}$

$R(\Delta_{m^\perp})$  把  $\Delta$  往  $m^\perp(S) = \{\Delta : \Delta_3 = 0\}$  上投影.

把  $\Delta$  往  $m(S) = \{\Delta : \Delta_1 = \Delta_2 = 0\}$  上投影.

画出  $\Theta$  的图形:

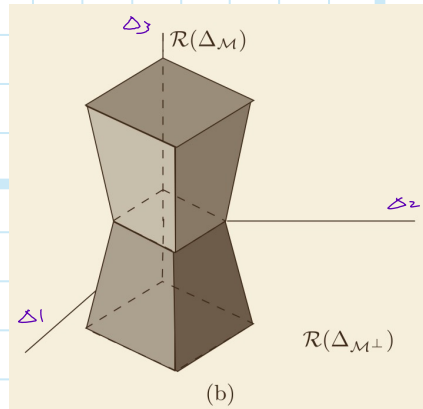
参数空间变成了一个 cone.



② 假设 model is not exact ( $R(\theta_{m^+}^*) \neq 0$ )

真实的参数  $\theta^*$  在 perturbation 的空间  $m^+$  上  
仍有分量,  $\therefore$  当  $|\Delta_3|$  取 0 时, 右边仍  $> 0$

$$\mathcal{C} = \left\{ \Delta: |\Delta_1| + |\Delta_2| \leq 3|\Delta_3| + 4R(\theta_{m^+}^*) \right\}$$



总结:

$\hookrightarrow$  0-norm?

nuclear norm 就是向量的 1-norm 向矩阵 norm 的推广,  
 $\therefore$  矩阵 1-norm 表征非 0 元素, 而矩阵的非 0 元素用 rank  
来表示, 可以这么理解: 向量非 0 指的是往这个向量上  
投影, 能剩下多少维, 而对应到矩阵, 则是其 rank 的大  
小, rank 大, 矩阵变换后剩的多。

Decomposable regularizer 的思路是: 对于  $R(\theta)$ , 首先希望  $R(\theta^*) \approx 0$ ,  
而如果在  $\theta^*$  上加一个扰动  $\gamma$ , 则希望  $R$  的反应非常剧烈, 即原来的  $R(\theta)$   
加扰动  $\gamma$  变成  $R(\theta + \gamma)$  由三角不等式  $R(\theta + \gamma) - R(\theta) \leq R(\gamma)$  即这个  $R$  的  
变化  $R(\theta + \gamma) - R(\theta)$  有一个上界  $R(\gamma)$ , 即然希望越剧烈越好, 就直接取 =  
放到最大。

$R$  是 Decomposable 的一个直接好处是 把  $\Delta$  限定在一个子空间中, 而非全空间。  
即是一个 low-dim structure。

2019.1.15

Recamp:

$\theta^*$ : ground truth

$z_1, \dots, z_n$ : Samples

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \{ \mathcal{L}(\theta; \{z_i\}_{i=1}^n) + \lambda R(\theta) \}$$

Goal: Bound stat. errors:  $\hat{\Delta} = \hat{\theta} - \theta^*$

$R$ : decomposable wrt  $(m, m^\perp)$

decomposability 是关于一对空间的.

$$R(\theta + \gamma) = R(\theta) + R(\gamma) \quad \forall \theta \in m, \gamma \in m^\perp$$

$m$ : model subspace

$m^\perp$ : perturbation subspace

就是说如果 Regularizer 合适的话, 不会被限制到一个子域中, 而是整个空间

prop: Suppose  $\mathcal{L}$  is smooth, convex,

$$\lambda \geq 2R^*(\nabla \mathcal{L}(\theta^*)) \quad R^*: \text{dual norm}$$

$R$  is decomposable wrt  $(m, m^\perp)$

Then,

$$\hat{\Delta} \in \mathcal{C}_\lambda \triangleq \{ \Delta \in \mathbb{R}^d \mid R(\Delta_{m^\perp}) \leq 3R(\Delta_m) + 4R(\theta_{m^\perp}^*) \}$$

注: 这个结论是 deterministic 的, 即与 noise 无关, 也不是什么分布.



下面证明那个奇怪的结论 Q

Pf: Define

$$D(\delta) = \mathcal{L}(\theta^* + \delta) - \mathcal{L}(\theta^*) + \lambda (R(\theta^* + \delta) - R(\theta^*))$$

(在  $\theta^*$  上加一个扰动  $\delta$ )

Obs:  $D(\hat{\theta}) \leq 0$  ( $\hat{\theta}$  是最优解,  $\mathcal{L}(\hat{\theta}) + R(\hat{\theta})$  最小)  
先证两个 Claim, 看整体的证明框架.  
Claim 1:

$$R(\theta^* + \delta) - R(\theta^*) \geq R(\hat{\Delta}_{\bar{m}^+}) - R(\hat{\Delta}_{\bar{m}}) - 2R(\theta_{m^+}^*)$$

Claim 2: Under the assumptions of Prop:

$$\mathcal{L}(\theta^* + \delta) - \mathcal{L}(\theta^*) \geq -\frac{\lambda}{2} [R(\hat{\Delta}_{\bar{m}}) + R(\hat{\Delta}_{\bar{m}^+})]$$

代入 claim 1, claim 2 得:

$$0 \geq D(\hat{\theta})$$

把最值的不等式拆成两块  
分别 bound

$$\geq \lambda [R(\hat{\Delta}_{\bar{m}^+}) - R(\hat{\Delta}_{\bar{m}}) - 2R(\theta_{m^+}^*)]$$

$$-\frac{\lambda}{2} [R(\hat{\Delta}_{\bar{m}}) + R(\hat{\Delta}_{\bar{m}^+})]$$

$$= \frac{\lambda}{2} [R(\hat{\Delta}_{\bar{m}^+}) - 3R(\hat{\Delta}_{\bar{m}}) - 4R(\theta_{m^+}^*)]$$

下面来证那两个 claim.

Pf (Claim 1)

$$R(\theta^* + \Delta) = R(\theta_m^* + \theta_{m^\perp}^* + \Delta_{\bar{m}} + \Delta_{\bar{m}^\perp})$$

注意  $\theta$  投影在  $m$  和  $m^\perp$ , 而  $\Delta$  投影在  $\bar{m}$  和  $\bar{m}^\perp$  上.

$$\geq R(\theta_m^* + \Delta_{\bar{m}^\perp}) - R(\theta_{m^\perp}^* + \Delta_{\bar{m}}) \quad \text{三角不等式}$$

$$\geq R(\theta_m^* + \Delta_{\bar{m}^\perp}) - R(\theta_{m^\perp}^*) - R(\Delta_{\bar{m}})$$

$$= R(\theta_m^*) + R(\Delta_{\bar{m}^\perp}) - R(\theta_{m^\perp}^*) - R(\Delta_{\bar{m}}) \quad (1)$$

$$\text{注意到: } R(\theta^*) \leq R(\theta_m^*) + R(\theta_{m^\perp}^*) \quad (2)$$

$$(1) - (2) \Rightarrow$$

$$R(\theta^* + \Delta) - R(\theta^*) \geq R(\Delta_{\bar{m}^\perp}) - R(\Delta_{\bar{m}}) - 2R(\theta_{m^\perp}^*)$$

Q. E. D.

Pf (Claim 2)

观察到左边是  $L(\theta^* + \Delta) - L(\theta^*)$  我们需要

对一个差来找一下界, 自然想到 gradient inequality

By convexity + smoothness,

$$\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) \geq \nabla \mathcal{L}(\theta^*)^T \Delta \geq -|\nabla \mathcal{L}(\theta^*)^T \Delta|$$

$$|\nabla \mathcal{L}(\theta^*)^T \Delta| \leq R^* (\nabla \mathcal{L}(\theta^*)) \cdot R(\Delta)$$

(generalized Cauchy-Schwarz)

当  $R$  是 2-norm 时, 就是普通的 C-S

当  $R$  是 1-norm,  $R^*$  是  $\infty$ -norm, 得出 Hölder

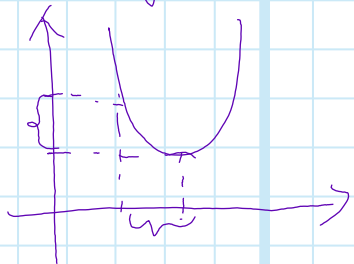
$$\leq \frac{\lambda}{2} R(\Delta)$$

$$\leq \frac{\lambda}{2} [R(\Delta_{\bar{m}}) + R(\Delta_{\bar{m}^c})] \quad (\text{三角不等式})$$

Q.E.D.

下面看, 对于  $\Delta$  的空间  $\Delta \in \mathcal{C}$ , 如果这个空间有 strong convexity, 我们就能够得到 error bound

strong convexity 直观理解: enough curved.



如果是 strong convex, 就一定能够找到一个整体的二次曲线,

Definition:  $\mathcal{L}$  satisfies the restricted strong convex (RSC) property if  $\exists \mu > 0$ , a function  $\tau(\cdot)$  s.t.:

$\wedge$  就是上面那个  $\mathcal{C}$

$$L(\theta^* + \Delta) \geq L(\theta^*) + \nabla L(\theta^*)^T \Delta + K \|\Delta\|_2^2 - \underbrace{\tau^2(\theta^*)}_{\text{tolerance}}$$


$$\forall \Delta \in \mathcal{C}$$

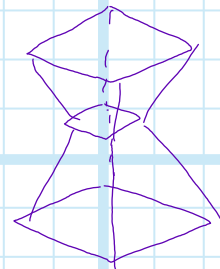
上周讲的 strong convexity  $\Delta$  是属于全空间, 而且  $\tau^2(\theta^*) = 0$ , 因此这里的这个式子要求弱一些  
同时, 我们这里的  $\theta^*$  是固定的, 只能动  $\Delta$

对比正常的 strong convexity:

$$L(v) \geq L(u) + \nabla L(u)^T (v-u) + K \|v-u\|_2^2$$

$$\forall u, v \in \mathbb{R}^d$$

$\tau^2(\theta^*)$  是 tolerance, 把  改成



直观上讲, 我们只需要  $\mathcal{C}$  curve enough, 减掉一点不影响  $\Rightarrow$

Thm: Under the assumption of Prop + RSC of  $L$

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \frac{9\lambda^2}{4K^2} \psi^2(\bar{m}) + \frac{1}{K} [2\tau^2(\theta^*) + 4\lambda R(\theta_{m^+}^*)]$$

where  $\psi(m) = \sup_{u \in m \setminus \{\emptyset\}} \frac{R(u)}{\|u\|_2}$   $\psi$  是关于空间  $m$  的 -const.

$\psi$  is the "Lipschitz const" of  $R$  restricted on  $m$

注: 这个 bounds 是一个 bounds 族,  $m \neq k$ ,  $\psi \neq k$

Rmk:

① Tradeoff: Error has 2 parts:

$$(i) \frac{9\lambda^2}{4k^2} \psi^2(\bar{m}) \quad \text{and} \quad (ii) \frac{4\lambda}{k} R(\theta_{m^*}^*) \quad (\text{assuming } \tau=0)$$

estimation error  
模型对  $\theta^*$  估计所致的误差

approx error.  
由对模型加入各种 specification 所引入的 error.  
例如假设是 sparse, 但实际上不是 sparse 的.

"bigger" the  $m$ ; lower approx error  
higher estimation error.  
如果对模型的假设, 那么模型空间就很大, estimation error 就大. (假设引入额外的信息) 此时 approx error 就小. (假设估计的假设)

对  $m$  限制小,

$k$  越大, 收敛越好.

② For concrete applications, need to compute/bound the parameters

Quick example: Sparse linear regression

$\theta^*$ :  $s$ -sparse.

$S \subseteq \{1, \dots, d\}$  of cardinality  $s$ .

$m(s) = \{\theta : \theta_j = 0 \forall j \notin S\}$

$R(\theta) = \|\theta\|_1$

$$\psi(m(s)) = \sup_{\theta \in m(s) \setminus \{0\}} \frac{\|\theta\|_1}{\|\theta\|_2} \stackrel{C-S}{=} \sqrt{s}$$

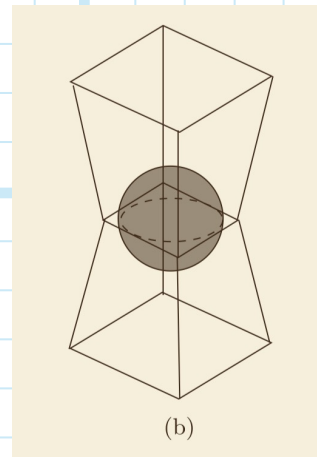
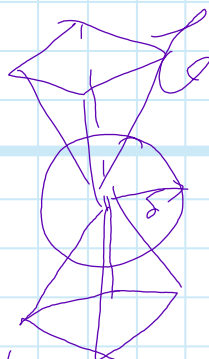
下面证明这个 Theorem

Pf:

$$D(\Delta) = \mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) + \lambda [R(\theta^* + \Delta) - R(\theta^*)]$$

$$D(\Delta) \leq 0$$

Let  $\delta > 0$ , define:  $K(\delta) = \mathcal{C} \cap \{\Delta : \|\Delta\|_2 = \delta\}$



先给一个 claim 看证明框架.

Claim: if  $D(\Delta) > 0 \forall \Delta \in K(\delta)$ , then  $\|\Delta\|_2 \leq \delta$

就是说一个  $\delta$  球上的点都不是 optimal

的, 那么, optimal 的  $\Delta$  一定在  $\delta$  球中.

下面来证 claim 中的那个条件  $D(\Delta) > 0$

Goal: Show that  $D(\Delta) > 0 \forall \Delta \in K(\delta)$  hold for some choice of  $\delta$

For  $\Delta \in K(\delta)$

$$D(\Delta) \geq \nabla \mathcal{L}(\theta^*)^T \Delta + K \|\Delta\|_2^2 - \tau^2(\theta^*) + \lambda [R(\theta^* + \Delta) - R(\theta^*)] \quad (\text{RSC})$$

(claim 1)

$$\geq \nabla \mathcal{L}(\theta^*)^T \Delta + K \|\Delta\|_2^2 - \tau^2(\theta^*) +$$

$$\lambda [R(\Delta_{\bar{m}^\perp}) - R(\Delta_{\bar{m}}) - 2R(\theta_{\bar{m}^\perp}^*)]$$

G-C-S

$$\begin{aligned} \geq & K \|\Delta\|_2^2 - \tau^2(\theta^*) + \lambda [R(\Delta_{\bar{m}^\perp}) - R(\Delta_{\bar{m}}) - 2R(\theta_{\bar{m}^\perp}^*)] \\ & - R^*(\nabla \mathcal{L}(\theta^*)) \cdot R(\Delta) \end{aligned}$$

注意到:  $R^*(\nabla \mathcal{L}(\theta^*)) \leq \frac{\lambda}{2}$

代入:  $R(\Delta) \leq R(\Delta_{\bar{m}}) + R(\Delta_{\bar{m}^\perp})$

$$\begin{aligned} \geq & K \|\Delta\|_2^2 - \tau^2(\theta^*) + \frac{\lambda}{2} [R(\Delta_{\bar{m}^\perp}) - 3R(\Delta_{\bar{m}}) - 4R(\theta_{\bar{m}^\perp}^*)] \\ & \underbrace{\qquad\qquad\qquad}_{\geq 0, \text{ 非负}} \end{aligned}$$

$$\geq K \|\Delta\|_2^2 - \tau^2(\theta^*) - \frac{\lambda}{2} [3R(\Delta_{\bar{m}}) + 4R(\theta_{\bar{m}^\perp}^*)]$$

因为  $R(\Delta_{\bar{m}}) \leq \psi(\bar{m}) \cdot \|\Delta_{\bar{m}}\|_2$

$$= \psi(\bar{m}) \cdot \|\pi_{\bar{m}}(\Delta) - \pi_{\bar{m}}(0)\|_2$$

$$\leq \psi(\bar{m}) \|\Delta - 0\|_2 \quad (\text{non-expansiveness of projection})$$

$$= \psi(\bar{m}) \|\Delta\|_2$$

综上有:

$$D(\Delta) \geq K \|\Delta\|_2^2 - \tau^2(\theta^*) - \frac{\lambda}{2} [3\psi(\bar{m}) \|\Delta\|_2 + 4R(\theta_{\bar{m}^\perp}^*)]$$

现在要凑那个 claim 中的条件  $D(\Delta) > 0$

即得到右也是  $\Delta$  的二次函数  $K > 0$ ,  $\therefore$  二次

函数开口向上, 只要最小点  $> 0$  即可.  $-\frac{b}{2a}$  代入

大和就能得出那个 Thm.

2019.1.21

Recamp:

Regularized loss minimization

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \{ \mathcal{L}(\theta) + \lambda R(\theta) \}$$

Thm: Suppose  $\mathcal{L}$  is smooth and convex,
 $\lambda \geq 2R^*(\nabla \mathcal{L}(\theta^*))$   $R$  is decomposable wrt  $(m, \bar{m}^\perp)$ . Then,

$$\hat{\Delta} = \hat{\theta} - \theta^* \in \mathcal{C} \{ \Delta : R(\Delta_{\bar{m}^\perp}) \leq 3R(\Delta_{\bar{m}}) + 4R(\theta_{\bar{m}^\perp}^*) \}$$

Note: if  $\mathcal{L}$  is RSC, i.e.

$$\mathcal{L}(\theta^* + \Delta) \geq \mathcal{L}(\theta^*) + \nabla \mathcal{L}(\theta^*)^\top \Delta + K \|\Delta\|_2^2 - \tau^2(\theta^*) \quad \forall \Delta \in \mathcal{B}$$

(只用在  $\theta^*$  或  $\bar{z}$ , 不一定对所有  $\theta$  或  $\bar{z}$ ) tolerance.

Then

$$\|\hat{\theta} - \theta^*\|_2^2 \leq \frac{9\lambda^2}{4K^2} \psi^2(\bar{m}) + \frac{2}{R} [\tau^2(\theta^*) + 2\lambda R(\theta_{\bar{m}^\perp}^*)]$$

where  $\psi(\bar{m}) = \sup_{u \in \bar{m} \setminus \{0\}} \frac{R(u)}{\|u\|_2}$



# Example 1: LASSO w/ exactly sparse models

$$y = X\theta^* + w$$

assume:  $\theta^*$  is  $s$ -sparse

$$\text{LASSO: } \hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\text{argmin}} \left\{ \underbrace{\frac{1}{2n} \|y - X\theta\|_2^2}_{L(\theta)} + \underbrace{\lambda \|\theta\|_1}_{R(\theta)} \right\}$$

is  $\hat{\theta}$  a good estimator for  $\theta^*$ ?

Take  $S = \text{Supp}(\theta^*)$ , Define the model subspace as:

$$m \equiv m(S) - \bar{m}(S) = \left\{ \theta : \theta_j = 0 \ \forall j \notin S \right\}$$

$$\Rightarrow \theta^* \in m \Rightarrow \theta_{m^\perp}^* = 0 \quad \bar{m}^\perp = \left\{ \theta : \theta_j = 0, \forall j \in S \right\}$$

$$\mathcal{L} = \left\{ \Delta : \underbrace{\|\Delta_S\|_1}_{R(\Delta_{\bar{m}^\perp})} \leq 3 \|\Delta_S\|_1 \right\}$$

下面验证 RSC:

$$\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) - \nabla \mathcal{L}(\theta^*)^\top \Delta = \frac{1}{2n} \|X\Delta\|_2^2$$

这里是把  $\mathcal{L}(\theta^*)$  在  $\theta^*$  处展开, 去掉常数项和一次项, 剩下二次项.

$\Rightarrow$  to establish RSC, suffices to show:

$$(*) \quad \frac{\|X\Delta\|_2^2}{2n} \geq K \|\Delta\|_2^2 \quad \forall \Delta \in \mathcal{L}$$

(如果成立的话, 这里取  $\tau(\cdot) \equiv 0$ )

这里将 RSC 那个条件转化成了一个类似于特征值不等式的一个式子, 只要 (\*) 成立, 就是 RSC.

下面看 (\*) 的应用.

Prop: Suppose

(i)  $X$  satisfies (\*)

(ii) (normalization)  $\frac{\|x_j\|_2}{\sqrt{n}} \leq 1 \forall j, x_j: j^{\text{th}} \text{ col of } X$

(iii)  $w$  is sub-Gaussian w/ parameter  $\sigma > 0$   
for a fixed  $\|v\|=1$   $w$  is mean zero and

$$P_x[|v^T w| \geq t] \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right) \forall t.$$

sub-Gaussian:

普通高斯分布的长尾以一个指数函数为上界。只要长尾以指数分布为上界的就是 sub-Gaussian



e.g.  $w \sim \mathcal{N}(0, I)$   $w$  是高斯, 那么  $w$  也是 sub-Gaussian

Then  $g = v^T w = \sum_i v_i w_i \sim \mathcal{N}(0, \|v\|_2^2)$

if  $\|v\|_2 = 1$ ,  $g$  is standard normal

$$P_x[|g| \geq t] = 2 \int_t^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \leq 2 \exp\left(-\frac{t^2}{2}\right)$$

Then setting  $\lambda = 4\sigma \sqrt{\frac{\log d}{n}}$ , Then

$$\|\theta - \theta^*\|_2^2 \leq \frac{36\sigma^2}{\lambda^2} \frac{s \log d}{n} \quad w/ \text{prob} \geq 1 - \frac{2}{d}$$

$\sigma$ : power of noise,  $\sigma$  越小, 效果越好.

$$\textcircled{1} \psi^2(\bar{m}) = \sup_{u \in \bar{m} \setminus \{0\}} \frac{\|u\|_1^2}{\|u\|_2^2} = s$$

$$\textcircled{2} R^*(v) = \|v\|_\infty; \quad \nabla \mathcal{L}(\theta^*) = \frac{1}{n} X^T (X\theta^* - y) = -\frac{1}{n} X^T w$$

开始推:

$$\text{Need: } \lambda \geq \frac{2}{n} \|X^T w\|_\infty = \frac{1}{n} \max_i |(X^T w)_i| = \frac{2}{n} \max_i |X_i^T w|$$

$$\text{Note: } \Pr\left[\frac{|X_i^T w|}{n} \geq t\right] \leq \Pr\left[\left|\frac{X_i}{\sqrt{n}}\right|^T w \geq \sqrt{n} t\right] \leq 2 \exp\left(-\frac{nt^2}{2\sigma^2}\right)$$

之前推导的是  $\|u\| = 1$ , 这里  $v = \frac{X_i}{\sqrt{n}}$

$\|u\| \leq 1$   $\therefore$  类似, 可以这么看:

$$\Pr\left[\left|\frac{\left(\frac{X_i}{\sqrt{n}}\right)^T w}{a}\right| \geq \frac{\sqrt{nt}}{a}\right] \leq 2 \exp\left(-\frac{nt^2}{a^2 2\sigma^2}\right) \leq 2 \exp\left(-\frac{nt^2}{2\sigma^2}\right) \quad (a \leq 1)$$

$$\Rightarrow \Pr\left[\frac{\|X^T w\|_\infty}{n} \geq t\right] \leq 2d \exp\left(-\frac{nt^2}{2\sigma^2}\right)$$

每个都有 bound, 全部加起来就是  $d$  倍的 bound.

$$\text{Set } t^2 = \frac{4\sigma^2 \log d}{n}$$

$$\Rightarrow \text{RHS} = 2d \exp(-2 \log d) \\ = 2/d$$

$$\text{w/ prob} \geq 1 - \frac{2}{d} \quad \frac{\|X^T w\|_\infty}{n} \leq 2\sigma^2 \sqrt{\frac{\log d}{n}}$$

$$\text{代入 } \lambda \geq \frac{2}{n} \|X^T w\|_\infty \text{ 得到 } \lambda = 4\sigma \sqrt{\frac{\log d}{n}} \text{ 为最小特征值}$$

下面看如何证(\*)

Thm: Suppose the rows of  $X$  are i.i.d  $\sim (0, \Sigma)$

$\Sigma > 0$ , Then,

$$\frac{\|X \Delta\|_2}{\sqrt{n}} \geq \frac{1}{4} \|\Sigma^{1/2} \Delta\|_2 - \rho(\Sigma) \sqrt{\frac{\log d}{n}} \|\Delta\|_1$$

$\forall \Delta \in \mathbb{R}^d$  with high probability

$$\text{where } \rho(\Sigma) = \max_{(i,j) \in d} \Sigma_{ij}$$

下面从这4thm 引(\*)

for  $\Delta \in \mathcal{C}$

↓ 由估计

$$\|\Delta\|_1 = \|\Delta_S\|_1 + \|\Delta_{S^c}\|_1 \leq 4\|\Delta_S\|_1 \leq 4\sqrt{5}\|\Delta_S\|_2$$

For simplicity, take  $\Sigma = I$ , by Thm,

$$\frac{\|X\Delta\|_2}{\sqrt{n}} \geq \frac{1}{4}\|\Delta\|_2 - 9\sqrt{\frac{\log d}{n}} (4\sqrt{5}\|\Delta\|_2)$$

$$= \underbrace{\left(\frac{1}{4} - 36\sqrt{\frac{\log d}{n}}\right)}_{K} \|\Delta\|_2$$

至此, 得到那么, then, example (基本证完)  
 这个Thm 改天证明, 一些观察

Fix  $\Delta$ :

$$E[\|X\Delta\|_2^2] = E[\Delta^T \underbrace{X^T X}_{\text{Random quadratic form}} \Delta] \quad X \sim N(0, I)$$

$$P_r\left[\Delta^T X^T X \Delta - E[\Delta^T X^T X \Delta] \geq t\right] \leq \dots$$

这个 concentration 不等式和Thm 很像, 但不能直接得出Thm,  $\because$  不能取+号在  
 $\Delta \in \mathcal{C}$  为  $\frac{1}{2}$  起, 因为  $\Delta$  为  $\Delta$  的  $\Delta$  很大.

Example 2: LASSO w/ weakly sparse models.

$$y = X\theta^* + w$$

并不是严格 sparse, 而是很多接近于 0.

Assume:  $\theta^* \in B_q(R_q) = \{\theta : \|\theta\|_q \leq R_q\} \quad q \in [0, 1]$

$$\|\theta\|_q^q \leq \sum_i |\theta_i|^q$$

$$\sum_i |\theta_i|^q \in \mathbb{R}_q \quad \text{impose decay rate on } \theta$$

如是  $\theta_i$  大, 后面的就小.

下面看用 Thm 在 Example 1 的问题.

$$\mathcal{C} \text{ 定义为: } \mathcal{C} = \left\{ \Delta : \|\Delta_S\|_1 \leq \beta \|\Delta\|_1 + 4 \underbrace{\|\theta_{m+1}^*\|_1}_{\text{多了这一项}} \right\}$$

是一个 ball,

而 (X)  $\frac{\|X\Delta\|_2^2}{2n} \geq K \|\Delta\|_2^2$  是二次的, 如果

$\Delta$  在一个 ball 中成立, 它在整个空间都成立, 这显然是不行的. 因此需要借助  $\mathcal{C}^2(\theta^*)$  这一项, 让它不  $\equiv 0$ .

2019.1.22

上节课讲了 LASSO 的例子, 其中关键的一步:

Need  $\exists K > 0$ :

$$\frac{\|X\Delta\|_2^2}{n} \geq K \|\Delta\|_2^2 \quad \forall \Delta \in \{\Delta: \|\Delta_{S^c}\|_1 \leq 3\|\Delta_S\|_1\}$$

Thm: Suppose the rows of  $X$  are iid  $\mathcal{N}(0, I)$

Then whp (with high probability)

$$\frac{\|X\Delta\|_2}{\sqrt{n}} \geq \frac{1}{4}\|\Delta\|_2 - 9\sqrt{\frac{\log d}{n}} \|\Delta\|_1, \quad \forall \Delta \in \mathbb{R}^d$$

上节课讲用这个 thm 证明上面那个 Need.

本节课证明这个 thm.

首先给出框架:

WLOG, assume  $\|\Delta\|_2 = 1$  (因为这个不等式是齐次的, 所以如果对  $\|\Delta\|_2 = 1$  成立, 则对所有  $\Delta$  成立)

Step 1: Want LB on

$$\inf_{\|\Delta\|_2=1} \frac{\|X\Delta\|_2}{\sqrt{n}}$$

这里  $X$  是随机变量  $\therefore$  整个式子是随机变量

For given  $r > 0$  define

$$V(r) \triangleq \{\Delta \in \mathbb{R}^d: \|\Delta\|_2 = 1, \|\Delta\|_1 \leq r\}$$

proposition 1:

$$\mathbb{E} \left[ \inf_{\Delta \in V(r)} \frac{\|X \Delta\|_2}{\sqrt{n}} \right] \geq 3 \left[ \frac{1}{4} - \sqrt{\frac{\log d}{n}} r \right] \text{ whenever } V(r) \neq \emptyset$$

∵ 对不同的  $x$ ,  $V(r)$  不同, ∴  $\mathbb{E}$  不能直接得到  $\inf$  的

Step 2: Concentration property of

$$Q(r, x) = \inf_{\Delta \in V(r)} \frac{\|X \Delta\|_2}{\sqrt{n}}$$

proposition 2: Let  $r > 0$ , be  $s, \{ \Delta \in V(r) \neq \emptyset$ , Then,

$$\Pr \left[ |Q(r, x) - \mathbb{E}[Q(r, x)]| \geq \frac{1}{2} t(r) \right] \leq 2 \exp\left(-\frac{nt^2(r)}{8}\right)$$

$$\text{where, } t(r) = \frac{1}{4} + 3 \sqrt{\frac{\log d}{n}} r$$

这个  $\uparrow$  是  $Q(r, x)$  偏离其中心的概率由

一个指数函数控制

Step 3: Prop 1+2 implies w/ prob  $\geq 1 - 2 \exp(-\frac{nt^2(r)}{8})$

$$Q(r, x) = \inf_{\Delta \in V(r)} \frac{\|X \Delta\|_2}{\sqrt{n}} \geq \mathbb{E}[Q(r, x)] - \frac{1}{2} t(r)$$

$$\geq 3 \left[ \frac{1}{4} - \sqrt{\frac{\log d}{n}} r \right] - \frac{1}{2} \left( \frac{1}{4} + 3 \sqrt{\frac{\log d}{n}} r \right)$$

$$= \frac{5}{8} - \frac{9}{2} \sqrt{\frac{\log d}{n}} r$$



# Bound for all levels of $r$

下面将要用到的工具

Tools

1) Comparison inequality of Gaussian processes,

随机过程只是随机变量 + 编号

2) Concentration inequality for Lipschitz functions of Gaussian RVs.

在高斯分布上加一个 Lipschitz 函数,

高斯分布的 concentration 性质依然保留

3) "Peeling argument" from empirical process theorem.

今天 focus 在 proposition 1.

Pf (prop 1)

$$\tilde{Q}(r, X) \triangleq \inf_{\Delta \in V(r)} \|\mathbf{X} \Delta\|_2 = \inf_{\Delta \in V(r)} \sup_{\substack{u: \|u\|_2=1 \\ u \in S^{n-1}}} u^T \mathbf{X} \Delta$$

↙ C-S 不等式

Note: for each  $(u, \Delta) \in S^{n-1} \times V(r)$   
把每个  $(u, \Delta)$  pair 看成 index set

$$Y_{u,\Delta} = u^T X \Delta \quad (\text{取一个 index } (u, \Delta) \text{ 并有一个随机变量 } Y_{u,\Delta})$$

$$= \sum_{i,j} u_i \Delta_j X_{ij}$$

is a mean-zero Gaussian RV.

want: LB on  $E[\tilde{Q}(r, X)]$

$$\Leftrightarrow \text{UB on } E[-\tilde{Q}(r, X)]$$

$$\Leftrightarrow \text{UB on } E \left[ \sup_{\Delta \in V(r)} \inf_{u \in S^{n-1}} (-u^T X \Delta) \right]$$

注意到  $u \in S^{n-1} \therefore -u \in S^{n-1}$  可以去掉负号

$$= \text{UB on } E \left[ \sup_{\Delta \in V(r)} \inf_{u \in S^{n-1}} u^T X \Delta \right] \leq ?$$

$\parallel$   
 $Y_{u,\Delta}$

$\{Y_{u,\Delta}\}_{u,\Delta}$  is a GP (这个集合每个元素都是高斯, 均值为0, 方差各不相同)

Idea: construct another GP  $\{Z_{u,\Delta}\}_{u,\Delta}$

s.t. (1)  $E \left[ \sup_{\Delta \in V(r)} \inf_{u \in S^{n-1}} Z_{u,\Delta} \right]$  is "easy" to compute

(2) related to  $E \left[ \sup_{\Delta \in V(r)} \inf_{u \in S^{n-1}} Y_{u,\Delta} \right]$

Fact: (Gordon's inequality)

Let  $U, V$  be arbitrary index sets, consider  $\{Y_{u,v}\}, \{Z_{u,v}\}$  families of zero-mean Gaussian RVs. Suppose

$$\sigma(Y_{u,v} - Y_{u',v'}) \leq \sigma(Z_{u,v} - Z_{u',v'})$$

$$\forall (u,v), (u',v') \in U \times V$$

$$\sigma(Y_{u,v} - Y_{u,v'}) = \sigma(Z_{u,v} - Z_{u,v'})$$

$$\forall u \in U, v, v' \in V$$

Then,

$$E \left[ \sup_{u \in U} \inf_{v \in V} Y_{u,v} \right] \leq E \left[ \sup_{u \in U} \inf_{v \in V} Z_{u,v} \right]$$

$Z$  更 spread out,  $E$  更大  $\frac{1}{Z}$  更大

下面看如何来构造  $Z$ .

In our setting,

$$\sigma^2(Y_{u,\delta} - Y_{u',\delta'}) = E \left[ (Y_{u,\delta} - Y_{u',\delta'})^2 \right]$$

$$= E \left[ (u^T X \delta - u'^T X \delta')^2 \right]$$

$$= E \left[ \left( \sum_{i,j} x_{ij} (u_i \delta_j - u'_i \delta'_j) \right)^2 \right] \rightarrow \text{展开} \sum_{i,j} \sum_{k,l} x_{ij} x_{kl} (u_i \delta_j - u'_i \delta'_j) (u_k \delta_k - u'_k \delta'_k)$$

$$= \sum_{i,j} \sum_{k,l} E(x_{ij} x_{kl}) ( \dots ) \quad \begin{matrix} \text{这里若 } (i,j) \neq (k,l) \\ \text{则 } E = 0, \text{ 否则 } = 1 \end{matrix}$$

$$= \sum_{i,j} (u_i \Delta_j - u'_i \Delta'_j)^2$$

$$= \|u \Delta^T - u' \Delta'^T\|_F^2 \quad (\text{写成矩阵形式})$$

$$= \|(u-u') \Delta^T + u'(\Delta-\Delta')^T\|_F^2$$

$$= \|\Delta\|_2^2 \|u-u'\|_2^2 + \|u'\|_2^2 \|\Delta-\Delta'\|_2^2$$

$$+ 2(u^T u' - \|u'\|_2^2)(\|\Delta\|_2^2 - \Delta^T \Delta')$$

$$\leq \|u-u'\|_2^2 + \|\Delta-\Delta'\|_2^2$$

$$\|\Delta\|_2^2 = 1, \|u'\|_2^2 = 1$$

$$\Delta^T \Delta' \leq 1 \quad (C-S)$$

equality holds when  $u=u'$  or  $\Delta=\Delta'$   $\therefore u^T u' - \|u'\|_2^2 \leq 0$

$$\|\Delta\|_2^2 - \Delta^T \Delta' \geq 0.$$

下面根据这个不等式来凑 Gordon inequality 的那个东西。  
希望这个 UB 是整个 RV 的 variance.

This suggests

$$Z_{u,a} = g_1^T u + g_2^T \Delta$$

where  $g_1, g_2 \sim \mathcal{N}(0, I)$  and independent.

check:

$$g_1^T (u - u') \sim \mathcal{N}(0, \|u - u'\|_2^2)$$

$$\sum_i g_{1,i} (u_i - u'_i)$$

$$g_2^T (\Delta - \Delta') \sim \mathcal{N}(0, \|\Delta - \Delta'\|_2^2)$$

$$Z_{u, \Delta} - Z_{u', \Delta'} = g_1^T (u - u') + g_2^T (\Delta - \Delta')$$

这样以来  $\|u - u'\|_2^2 + \|\Delta - \Delta'\|_2^2 = \sigma^2 (Z_{u, \Delta} - Z_{u', \Delta'})^2$

By Gordon's ineq.

$$\mathbb{E} \left[ \sup_{\Delta \in V(r)} \inf_{u \in S^{n-1}} \underbrace{(g_1^T u + g_2^T \Delta)}_{Z_{u, \Delta}} \right]$$

这个子的期望是可以直接拆开的.

$$= \mathbb{E} \left[ \underbrace{\inf_{u \in S^{n-1}} g_1^T u}_{\downarrow C-S} \right] + \mathbb{E} \left[ \sup_{\Delta \in V(r)} g_2^T \Delta \right]$$

$$\leq \mathbb{E}[-\|g_1\|_2] + \mathbb{E} \left[ \sup_{\Delta \in V(r)} \|\Delta\|_1 \cdot \|g_2\|_\infty \right]$$

↓ Hölder

$$= \mathbb{E}[-\|g_1\|_2] + r \left[ \|g_2\|_\infty \right]$$

由  $V(r)$  中  $\|\Delta\|_1 = r$

Claim:  $\mathbb{E}[\|g_2\|_\infty] \leq 3\sqrt{\log d}$

$$\mathbb{E}[\|g_2\|_2] \geq \frac{3}{4}\sqrt{n} \quad \text{for } n \geq 10 \text{ say.}$$

$$\Rightarrow \mathbb{E}[-\tilde{Q}(r, X)] \leq 3r\sqrt{\log d} - \frac{3}{4}\sqrt{n}$$

$$\Rightarrow \mathbb{E}[\tilde{Q}(r, X)] \geq \frac{3}{4}\sqrt{n} - 3r\sqrt{\log d}$$

[对称性记录]  
proposition 1

下面看那 (Claim)

$$g_2 \sim N(0, I) \quad E[\|g_2\|_\infty] = \int_0^\infty \Pr[\|g_2\| \geq t] dt$$

$$(E(x) = \int_0^\infty \Pr[x > t] dt \quad \text{for } x > 0.)$$

$$\Pr[\|g_2\|_\infty > t] \Rightarrow \Pr[|g| > t] \quad g \sim N(0, I)$$

存在  $-t/2 \leq g \leq t/2$

$$\Pr[|g| > t] = 2 \int_t^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

下面看 Gordon 引理 等式成立

$$(u-u') \Delta^T \cdot (u-u') \Delta^T$$

$$= \text{tr}((u-u') \Delta^T \Delta (u-u')^T)$$

$$A \cdot B = \text{tr}(A B^T) = \text{tr}(A^T B)$$

$$= \|\Delta\|_2^2 - \|u-u'\|_2^2$$

2019.1.28

回顾:

Thm: Suppose  $X \in \mathbb{R}^{n \times d}$  has iid  $\mathcal{N}(0,1)$  entries  
 $n \ll d$

Then, whp

$$\frac{\|X \Delta\|_2}{\sqrt{n}} \geq \frac{1}{4} \|\Delta\|_2 - \underbrace{9 \sqrt{\frac{\log d}{n}} \|\Delta\|_1}_{\because n \ll d \therefore \text{需要这项来修正}}$$

Pf: Define

$$V(r) = \left\{ \Delta \in \mathbb{R}^d : \|\Delta\|_2 = r, \|\Delta\|_1 \leq r \right\}$$

$\because$  那个不是完备的  
 $\swarrow$   
 $\leftarrow$  为了控制  $Q(r, x)$   
 看最多能得到多少。

$$Q(r, x) = \inf_{\Delta \in V(r)} \frac{\|X \Delta\|_2}{\sqrt{n}}$$

Step 1:  $E[Q(r, x)] \geq 3 \left[ \frac{1}{4} - \sqrt{\frac{\log d}{n}} \right] r$

以上是上节课, 证明了 mean 有一个下界, 本节课讲  $Q$  大概率集中在其 mean 的周围

Step 2: Concentration

prop: for  $r > 0$   $\swarrow$  deviation of  $Q$  from its mean

$$P_r \left[ |Q(r, x) - E[Q(r, x)]| \geq \frac{1}{2} t(r) \right] \leq 2 \exp(-n t^2(r)/8)$$

where  $t(r) = \frac{1}{4} + 3 \sqrt{\frac{\log d}{n}} r$

就是说  $Q(v|x)$  离  $E[Q(v|x)]$  的距离  
指数衰减。

Pf of prop:

Fact: (concentration of measure for Lipschitz  
Fns of Gaussians)

Let  $g \sim \mathcal{N}(0, I_m)$  and  $F: \mathbb{R}^m \rightarrow \mathbb{R}$  be  
an  $L$ -Lipschitz function i.e.

$$|F(x) - F(y)| \leq L \|x - y\|_2 \quad \forall x, y \in \mathbb{R}^m$$

Then,  $\forall t \geq 0$ ,

$$\Pr \left[ |F(g) - E[F(g)]| \geq t \right] \leq 2 \exp\left(-\frac{t^2}{2L^2}\right)$$

$L$  越小, 连续性越好, 则这个概率的UB越紧

例:  $F(g) = g$

$F$  不能取二次函数,  $\because$  二次函数不是  $L$ -Lipschitz

可以猜到, 这里  $Q$  应该是  $L$ -Lipschitz 的, 我们  
再找到那个  $L$ , 就好办了。那我们直接验证  
 $Q$  是不是  $L$ -Lipschitz 就好了。



Goal: verify  $x \rightarrow Q(r, x)$  is Lipschitz.

$$\sqrt{n} [Q(r, x) - Q(r, Y)] = \inf_{\Delta \in V(r)} \|x\Delta\|_2 - \inf_{\Delta \in V(r)} \|Y\Delta\|_2$$

$$= \inf_{\Delta \in V(r)} \|x\Delta\|_2 - \|Y\hat{\Delta}\|_2 \quad \text{where } \hat{\Delta} \in \arg \min_{\Delta \in V(r)} \|Y\Delta\|_2$$

(为什么这个 minimizer  $\hat{\Delta}$  存在?  $\because$  是  
最小化一个连续函数在紧集上,

由 Weierstrass, 最小值存在)

$$\leq \|x\hat{\Delta}\|_2 - \|Y\hat{\Delta}\|_2 \quad (\hat{\Delta} \text{ 比 } \inf_{\Delta \in V(r)} \|x\Delta\|_2 \text{ 的值大})$$

$$\leq \sup_{\Delta \in V(r)} \|(x-Y)\Delta\|_2$$

$$\leq \|x-Y\| \cdot \sup_{\Delta \in V(r)} \|\Delta\|_2$$

$\left\{ \begin{array}{l} \|A u\|_2 \leq \|A\| \cdot \|u\|_2 \\ \|A\| = \sup_{\|u\|_2=1} \|A u\|_2 \\ \uparrow \\ \text{spectral norm of } A. \end{array} \right.$

$$= \|x-Y\| \cdot 1 \leftarrow \text{由 } V(r) \text{ 定义.}$$

$$\leq \|x-Y\|_F$$

$\left\{ \begin{array}{l} \text{2-norm of vector} \\ \text{is similar to } F\text{-norm} \\ \text{of a matrix.} \end{array} \right.$

综上:  $\sqrt{n}|Q(r, X) - Q(r, Y)| \leq \|X - Y\|_F$  is similar to  $\|u\|_\infty$   
 $\Rightarrow Q: \frac{1}{\sqrt{n}}$  Lipschitz.

spectral norm of vector  
 $\|u\|_\infty \leq \|u\|_2$   
 $\|A\| \leq \|A\|_F$

$\Rightarrow$  by fact, set  $t \equiv \frac{1}{2}t(r)$   
 $F(x) = Q(r, x)$

这里  $L = \frac{1}{\sqrt{n}}$  代入即得:

$$\Pr[|Q(r, X) - E[Q(r, X)]| \geq \frac{1}{2}t(r)] \leq 2\exp(-nt^2(r)/8)$$

$L^2 \leq \text{证 } \square A \downarrow$  step 2,  $T \rightarrow \square$  证 step 3.

Step 3: Bound uniformly over  $r$

For a fixed  $r > 0$ , whp,

$$\begin{aligned}
 Q(r, X) &\geq E[Q(r, X)] - \frac{1}{2}t(r) \\
 &\geq \frac{5}{8} - \frac{9}{2} \sqrt{\frac{\log d}{n}} r
 \end{aligned}$$

Pf of prop:

Lemma 1: Let  $A \subseteq \mathbb{R}$  be non-empty  $f: \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}$   
 $h: \mathbb{R}^d \rightarrow \mathbb{R}_+$  are given. Consider

r.v.  $\left\{ \begin{array}{l} \sup_{v \in A} f(v, X) \\ h(v) \leq r \end{array} \right\}$   $X$  is a random vector.  
 通过控制  $r$  来控制这个期望。

Suppose  $\exists g: \mathbb{R} \rightarrow \mathbb{R}$  be non-negative and strictly increasing, s.t. (i)  $g(r) \geq \mu \forall r \geq 0$  and

(ii)  $\exists c > 0, \forall r > 0,$   $\rightarrow$  这里  $c$  是常数,  $a$  是个函数。

$$(*) \Pr \left[ \sup_{\substack{v \in A \\ h(v) \leq r}} f(v, X) \geq g(r) \right] \leq 2 \exp(-c \cdot a \cdot g^2(r))$$

$$\text{Then, } \Pr[\mathcal{E}] \leq \frac{2 \exp(-4 \cdot c \cdot a \cdot \mu^2)}{1 - \exp(-4 \cdot c \cdot a \cdot \mu^2)}$$

where  $\mathcal{E} = \left\{ \exists v \in A: f(v, X) \geq 2g(h(v)) \right\}$

Then 后面的结论去掉  $r$ ,  $r \rightarrow r$  无关。  
 总是之前对所有  $r$  的 level 都成立, 现在是对所有的, 根号都很小。

下面看这个 Lemma 如何用到 Step 3 上:

In our setting,

$$A = \{s: \|s\|_2 = 1\}$$

$$h(v) = \|v\|,$$

就是把那个  $v(v)$  拆开。

$$f(\Delta, X) = 1 - \frac{\|\Delta\|_2}{\sqrt{n}}$$

下面看  $g$  是什么.

$$t(r) = \frac{1}{4} + 3\sqrt{\frac{\log d}{n}} r$$

$$\begin{aligned} Q(r, X) &\geq \frac{5}{8} - \frac{9}{2n} \sqrt{\frac{\log d}{n}} r \\ &= 1 - \frac{3}{2} t(r). \end{aligned}$$

下面渣形式来求  $g(r)$

(\*) 变式了:

$$\Pr \left[ \sup_{\Delta \in V(r)} \left( 1 - \frac{\|\Delta\|_2}{\sqrt{n}} \right) \geq \frac{3}{2} t(r) \right]$$

$$= \Pr \left[ - \inf_{\Delta \in V(r)} \frac{\|\Delta\|_2}{\sqrt{n}} \geq \frac{3}{2} t(r) \right]$$

$$= \Pr \left[ Q(r, X) \leq 1 - \frac{3}{2} t(r) \right] \leq 2 \exp(-c \cdot a \cdot \dots)$$

$$\therefore g(r) = \frac{3}{2} t(r) \geq \frac{3}{8} = \mu.$$

有了  $g$ , 下面套入 step 2:

$$\begin{aligned} \Pr \left[ Q(r, X) \leq 1 - \frac{3}{2} t(r) \right] &\leq 2 \exp(-n t^2(r) / 8) \\ &= 2 \exp(-n g^2(r) / 18) \end{aligned}$$

由 lemma 的结论有:

whp:

$$\forall \|\Delta\|_2 = 1$$

$$1 - \frac{\|x\Delta\|_2}{\sqrt{n}} \leq 2 \cdot \frac{3}{2} t(\|\Delta\|_1)$$

$$= \frac{3}{4} + 9 \sqrt{\frac{\log d}{n}} \|\Delta\|_1$$

$$\Rightarrow \frac{\|x\Delta\|_2}{\sqrt{n}} \geq \frac{1}{4} - 9 \sqrt{\frac{\log d}{n}} \|\Delta\|_1$$

这样就得出了那个 Thm.

下面看看那个 Lemma 如何证。

Pf of Lemma 1.

Since  $g(r) \geq \mu \quad \forall r \geq 0$ , define, for  $m=1, 2, \dots$

$$A_m = \{v \in A : 2^{m-1} \mu \leq g(h(v)) \leq 2^m \mu\}$$

证明中的常用技巧: divide-and-conquer.

Fix  $x$ , If  $v \in A$  is s.t.  $f(v, x) \geq 2g(h(v))$

Then  $v \in A_m$  for some  $m$ .

$$\Rightarrow \Pr[\mathcal{E}] \leq \Pr \left[ \bigcup_{m \geq 1} \left\{ \exists v \in A_m \text{ s.t. } f(v, x) \geq 2g(h(v)) \right\} \right]$$

union bound

$$\leq \sum \Pr \left[ \exists v \in A_m, \text{ s.t. } f(v, x) \geq 2g(h(v)) \right]$$

Now,  $\forall v \in A_m$  and  $f(v, x) \geq 2g(h(v)) \Rightarrow$

$$(1) \quad f(v, x) \geq 2 \cdot 2^{m-1} \mu = 2^m \mu$$

$$(2) \quad g(h(v)) \leq 2^m \mu \Rightarrow h(v) \leq g^{-1}(2^m \mu)$$

$$\Rightarrow P_r[\varepsilon] \leq \sum_{m \geq 1} P_r \left[ \sup_{\substack{h(v) \leq g^{-1}(2^m \mu) \\ v \in A}} f(v, x) \geq 2^m \mu \right]$$

$$\leq 2 \sum_{m \geq 1} \exp(-c \cdot a \cdot \underbrace{(g^{-1}(2^m \mu))^2}_r)$$

$$= 2 \sum_{m \geq 1} \exp(-c \cdot a \cdot 2^{2m} \mu^2)$$

用几个序列来做UB, 然后几个序列求和即可

下节课讲

$\hat{\theta} \in \arg \min_{\theta \in \mathcal{R}^d} \{ \mathcal{L}(\theta) + \lambda R(\theta) \}$  如何高效计算.

2019.1.29

本节讲如何求  $\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \{L(\theta) + \lambda R(\theta)\}$  (P)

Assumptions:

(A1)  $L$  is smooth and convex

(A2)  $R$  is a norm with  $\lambda=1$

( $\lambda \neq 1 \rightarrow$  fix scale  $R$  with  $\lambda=1$ )

Obs: (P) is in general non-smooth convex optimization problem.

First-order optimality condition =

$$0 \in \nabla L(\theta) + \partial R(\theta) \quad (\text{necessary \& sufficient})$$

where  $\partial R(\theta) = \{s \in \mathbb{R}^d : R(y) \geq R(\theta) + s^T(y - \theta) \forall y\}$

( $\partial R(\theta)$  是 sub-gradient 是一个集合,

$\because R(\theta)$  不一定是光滑, 但保证 convex,  
 $\therefore$  FOC 可以用 sub-gradient 来写)

$\because$  是 convex  $\therefore$  FOC 是充要

Aside: Exercise: when  $R$  is a norm,

$$\partial R(\theta) = \{s \in \mathbb{R}^d : R^*(s) \leq 1, s^T \theta = R(\theta)\}$$

$R^*$  is the dual norm of  $R$ .

( $R$  是 norm 的话,  $\partial R(\theta)$  可以用其 dual norm 更简单地表示.)

注意到这里 FOC 是一个集合包含 0, 比等式难解, 下面看怎么弄.

Define the proximal map associated with  $R$  as!

$$\mathbb{R}^d \ni \text{prox}_R(\theta) = \underset{y \in \mathbb{R}^d}{\text{argmin}} \left\{ R(y) + \frac{1}{2} \|\theta - y\|_2^2 \right\}$$

↑  
非光滑部分

proximity term

希望  $y$  离当前  $\theta$  不要太远.

希望找一个  $y$  使得  $R(y)$  足够小; 同时离  $\theta$  不太远.

Note:  $\text{prox}_R(\theta)$  is well-defined (exists & unique)

due to the strongly convexity of

$$y \mapsto R(y) + \frac{1}{2} \|\theta - y\|_2^2$$

一个 strong convex + convex 还是 strong convex  
strong convex 不在一个平面内吗?



claim:  $\hat{\theta}$  is optimal for (P) iff

$$\hat{\theta} = \text{prox}_R(\hat{\theta} - \nabla \mathcal{L}(\hat{\theta})) \quad (\text{fixed-point equation})$$

于是通过这个 map 将优化问题转化成了不动点问题

Pf: Consider first-order opt cond of  $\text{prox}_R(\cdot)$   
写 FOC 得:

$$0 \in \partial R(\theta) + \gamma - \theta$$

$$0 \in \partial R(\text{prox}_R(\theta)) + \text{prox}_R(\theta) - \theta \quad \text{for any } \theta$$

Hence,  $\hat{\theta}$  is optimal for (P)

$$\Leftrightarrow 0 \in \nabla \mathcal{L}(\hat{\theta}) + \partial R(\hat{\theta})$$

$$\Leftrightarrow 0 \in (\nabla \mathcal{L}(\hat{\theta}) + \hat{\theta}) + \hat{\theta} + \partial R(\hat{\theta})$$

$$\hat{\theta} = \text{prox}_R(\hat{\theta} - \nabla \mathcal{L}(\hat{\theta})) \quad \checkmark \text{ 代入验证:}$$

验证:  
把  $\hat{\theta}$  和  $\partial R(\hat{\theta})$  中的  $\hat{\theta}$  替换  
 $0 \in (\hat{\theta} - \nabla \mathcal{L}(\hat{\theta})) + \hat{\theta} + \partial R(\hat{\theta})$   
 $0 \in (\hat{\theta} - \nabla \mathcal{L}(\hat{\theta})) + \text{prox}_R(\hat{\theta} - \nabla \mathcal{L}(\hat{\theta})) + \partial R(\text{prox}_R(\hat{\theta} - \nabla \mathcal{L}(\hat{\theta})))$   
满足上面那个 FOC 条件

以上证明了那个 claim. 下面举个例子.

$$\text{Example: } R(\gamma) = \mathbb{1}_C(\gamma) = \begin{cases} 0 & \text{if } \gamma \in C \\ +\infty & \text{o/w} \end{cases}$$

$C$  is closed convex.

$$\text{prox}_R(\theta) = \arg \min_{\gamma \in \mathbb{R}^d} \left\{ \mathbb{1}_C(\gamma) + \frac{1}{2} \|\theta - \gamma\|_2^2 \right\}$$

$\gamma$  肯定是  $\gamma \in C$   $\therefore$

$$= \arg \min_{\gamma \in C} \left\{ \frac{1}{2} \|\theta - \gamma\|_2^2 \right\} = \Pi_C(\theta)$$

可见 proximal map 是 projection 的一个扩展

下面看这个 fixed-point equation 如何算法求解

Algorithm:

$$\text{Idea: } H(\hat{\theta}) = \hat{\theta}$$

就是直接暴力迭代求不动点

$$H(\theta^0) \stackrel{?}{=} \theta^0 \begin{array}{l} \nearrow \text{dive} \\ \searrow \theta^1 = H(\theta^0) \end{array}$$

问题: (1) 能否收敛

(2) 能否收敛到 fixed point.

Proximal gradient Method (PGM)

$$\theta^{k+1} \leftarrow \text{prox}_{\alpha_k R} \left( \theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) \right) \quad \alpha_k > 0$$

Another interpretation of PGM. 代入定义得:

$$\theta^{k+1} = \arg \min_{\gamma \in \mathbb{R}^d} \left\{ \alpha_k R(\gamma) + \frac{1}{2} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \gamma\|_2^2 \right\}$$

$$= \arg \min_{\gamma \in \mathbb{R}^d} \left\{ R(\gamma) + \frac{1}{2\alpha_k} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \gamma\|_2^2 \right\}$$

$$= \arg \min_{\gamma \in \mathbb{R}^d} \left\{ \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\gamma - \theta^k) + \frac{1}{2\alpha_k} \|\gamma - \theta^k\|_2^2 + R(\gamma) + \underbrace{\alpha_k^2 \|\nabla \mathcal{L}(\theta^k)\|_2^2}_{\text{常数, 省略}} \right\}$$

前三项其实只是  $\mathcal{L}(\gamma)$  的前三阶泰勒展开

$$\text{即: } \gamma \mapsto \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\gamma - \theta^k) + \frac{1}{2\alpha_k} \|\gamma - \theta^k\|_2^2$$

is a quadratic approx of  $\mathcal{L}$  at  $\theta^k$ .

∴ 也可以这么解释, 每一步迭代都是用一个二次函数去估计 smooth 的项.

扩展: 也可以用真正的泰勒二阶去估计, 这里只是有近似泰勒, 但还是有差别

下面来看: (1) 能否收敛 (2) 能否收敛到 fixed-point  
之后再研究收敛速度.

加入另一个 Assumption

(A3)  $\nabla \mathcal{L}$  is  $L$ -Lipschitz

下面套用经典算法分析框架

Proposition: Suppose  $\{\alpha_k\}$  satisfies  $0 < \underline{\alpha} \leq \alpha_k \leq \bar{\alpha} < \frac{1}{L}$

Then:

(a) (Sufficient Descent)  $\exists K_1 > 0$ :

$$F(\theta) = \mathcal{L}(\theta) + \lambda R(\theta)$$

$$F(\theta^k) - F(\theta^{k+1}) \geq K_1 \|\theta^k - \theta^{k+1}\|_2^2$$

(b) (Safe guard)  $\exists K_2 > 0$ :

$$\|E(\theta^k)\|_2 \leq K_2 \|\theta^k - \theta^{k+1}\|_2$$

where  $E(\theta) = \underset{R}{\text{prox}}(\theta - \nabla \mathcal{L}(\theta)) - \theta$  is the residual function

(Rmk:  $E(\theta) = 0 \Leftrightarrow \theta$  is optimal for (P))

下面看 proposition 的应用.

(a)  $\Rightarrow \{F(\theta^k)\}_{k \geq 0}$  monotonically decreasing

since  $F(\theta^k) \geq \hat{V} = F(\hat{\theta}) \forall k$ .

$\Rightarrow F(\theta^k) \rightarrow V$  (单调有界序列收敛)  
i.e.  $\{F(\theta^k)\}$  converges.

$\Rightarrow \theta^k - \theta^{k-1} \rightarrow 0$  (by (a) again)

(b)  $\Rightarrow \|E(\theta^k)\|_2 \rightarrow 0$

(1)  $\Rightarrow V = \hat{V}$

(2) Every accumulation point of  $\{\theta^k\}_{k \geq 0}$   
is optimal for (P)

$\{\theta^k\}$  可能不收敛, 但其中的一些  
子序列可能收敛, 这些收敛的  
点就是 (P) 的最优解.

(Subsequential convergence)

下面证 BA (a)

Pf (a)

$$\theta^{k+1} = \arg \min \left\{ \frac{1}{2} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \gamma\|_2^2 + \alpha_k R(\gamma) \right\}$$

$$\Rightarrow \frac{1}{2} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \theta^{k+1}\|_2^2 + \alpha_k R(\theta^{k+1})$$

$$\leq \frac{1}{2} \alpha_k^2 \|\nabla \mathcal{L}(\theta^k)\|_2^2 + \alpha_k R(\theta^k) \quad \because \theta^{k+1} \text{ 是 optimizer}$$

Rearrange:

$$\textcircled{a.1} \quad R(\theta^{k+1}) + \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \theta^k) + \frac{1}{2\alpha_k} \|\theta^{k+1} - \theta^k\|_2^2 \leq R(\theta^k)$$

For  $\mathcal{L}$ :

$$\text{define } g(t) = \mathcal{L}(\theta^k + t(\theta^{k+1} - \theta^k))$$

$$\mathcal{L}(\theta^{k+1}) - \mathcal{L}(\theta^k) = g(1) - g(0)$$

$$= \int_0^1 g'(t) dt$$

$$= \int_0^1 \nabla \mathcal{L}(\theta^k + t(\theta^{k+1} - \theta^k))^T (\theta^{k+1} - \theta^k) dt$$

应用 Lipschitz

$$= \int_0^1 \left[ \nabla \mathcal{L}(\theta^k + t(\theta^{k+1} - \theta^k)) - \nabla \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k) \right]^T (\theta^{k+1} - \theta^k) dt$$

$$\leq \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \theta^k) + \int_0^1 t L \|\theta^{k+1} - \theta^k\|_2^2 dt$$

先用一道 Cauchy 不等式变形:

$$[\nabla \mathcal{L}(\theta^k + t(\theta^{k+1} - \theta^k)) - \nabla \mathcal{L}(\theta^k)]^T (\theta^{k+1} - \theta^k)$$

$$\leq \|\nabla \mathcal{L}(\theta^k + t(\theta^{k+1} - \theta^k)) - \nabla \mathcal{L}(\theta^k)\| \|\theta^{k+1} - \theta^k\|$$

再用 Lipschitz:

$$\leq L \|\theta^{k+1} - \theta^k\| \|\theta^{k+1} - \theta^k\|$$

$$\textcircled{a.2} \quad = \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \theta^k) + \frac{L}{2} \|\theta^{k+1} - \theta^k\|_2^2$$

$$\text{Hence, } F(\theta) - F(\theta^k) = \mathcal{L}(\theta^{k+1}) + R(\theta^{k+1}) - \mathcal{L}(\theta^k) - R(\theta^k)$$

$$\leq \underbrace{-\frac{1}{2} \left( \frac{1}{2\alpha_k} - L \right)}_{\kappa_1} \|\theta^{k+1} - \theta^k\|_2^2$$

$\textcircled{a.1} + \textcircled{a.2}$  即得

证毕 (a)

2019.2.11

(a) + (b)  $\Rightarrow$

①  $F(\theta^k) \rightarrow \hat{V}$

② Every accumulation point of  $\{\theta^k\}$

is an optimal solution

(subsequential convergence)

每个聚点都是一个最优解.

以上是 Rockaford.

Pf (Prop 1(b))

Lemma: Let  $\theta, \gamma$  be arbitrary. Then,

$0 < \alpha \mapsto g_1(\alpha) \triangleq \frac{1}{\alpha} \|\text{prox}_{\mathcal{R}}(\theta - \alpha\gamma) - \theta\|_2$   
is decreasing

( $g_1$  是关于  $\alpha$  的 单变量函数, 关于  $\alpha \downarrow$ )

and

$0 < \alpha \mapsto g_2(\alpha) \triangleq \|\text{prox}_{\mathcal{R}}(\theta - \alpha\gamma) - \theta\|_2$   
is increasing

接着怎么用这4 Lemma,

Assume the lemma:

$$\|\theta^{k+1} - \theta^k\|_2 = \|\text{prox}_{\alpha R}(\theta^k - \alpha \nabla \mathcal{L}(\theta^k)) - \theta^k\|_2$$

(根据PGM的定义)  
这4步其实就是  $g_2(\alpha)$

之前的假设  $0 < \underline{\alpha} \leq \alpha_k \leq \bar{\alpha} < \frac{1}{L}$

$$\geq \|\text{prox}_{\underline{\alpha} R}(\theta^k - \underline{\alpha} \nabla \mathcal{L}(\theta^k)) - \theta^k\|_2$$

下面需要把这4LB和

$$\|\text{prox}_R(\theta^k - \nabla \mathcal{L}(\theta^k)) - \theta^k\|_2 \text{ 联系起来,}$$

当  $\alpha \geq 1$  时直接放

当  $\alpha < 1$  时, 由  $g_1(\alpha) \downarrow \therefore$

$$\frac{1}{\alpha} \|\text{prox}_{\underline{\alpha} R}(\theta^k - \underline{\alpha} \nabla \mathcal{L}(\theta^k)) - \theta^k\|_2 \geq$$

$$\|\text{prox}_R(\theta^k - \nabla \mathcal{L}(\theta^k)) - \theta^k\|_2$$

$$\geq \min\{1, \underline{\alpha}\} \cdot \|\text{prox}_R(\theta^k - \nabla \mathcal{L}(\theta^k)) - \theta^k\|_2$$

$E(\theta^k)$  的定义.



以上已经证明了 (b) (safeguard).

下面看如何证明这个 Lemma.

Pf of Lemma: Define  $h: \mathbb{R}_{++} \times \mathbb{R}^d \rightarrow \mathbb{R}$

$$\text{by } h(\alpha, \omega) = \gamma^T(\omega - \theta) + \frac{1}{2\alpha} \|\omega - \theta\|_2^2 + R(\omega)$$

( $\gamma, \theta$  由前面固定,  $\alpha$  是正数,  $\omega$  是向量)

观察到  $h(\alpha, \omega)$  后两项恰好就是那个 minimizer, 只是多了前面一项.

and define Moreau envelope of  $h$  by

$$(*) \quad H(\alpha) = \inf_{\omega} h(\alpha, \omega)$$

Moreau envelope 正在研究 partial optimization 的一个工具, How minimizer change its behaviour. minimizer 由一个参数  $\alpha$  控制

(claim: The opt. soln to (\*) given by

$$\omega^* = \text{prox}_{\alpha R}(\theta - \alpha \gamma)$$

Pf: 观察(\*)每一项关于  $w$  都是 convex 的,

$\therefore$  写 FO condition for (\*):

$$0 \in \gamma + \frac{1}{\alpha}(w - \theta) + \alpha R(w)$$

$w^* \triangleq \text{prox}_{\alpha R}(\theta - \alpha\gamma)$  is opt soln to

$$\frac{1}{2} \underbrace{\|\theta - \alpha\gamma - w\|_2^2}_{\theta} + \alpha R(w)$$

$\alpha\gamma \in \text{prox}_{\alpha R}(\theta, \gamma)$

$$\left( \text{prox}_R(\theta) = \underset{r}{\text{argmin}} \left\{ \frac{1}{2} \|\theta - r\|_2^2 + R(r) \right\} \right)$$

FO condition:

$$0 \in w - (\theta - \alpha\gamma) + \alpha \cdot \partial R(w)$$

$$= \alpha\gamma + w - \theta + \alpha \cdot \partial R(w)$$

$$\Leftrightarrow 0 \in \gamma + \frac{(w - \theta)}{\alpha} + \alpha R(w)$$

This is unique ( $\because$   $\frac{1}{2}\|\cdot\|_2^2$  strong convex  $\therefore$  没有平点,  $\therefore$   $(w^* = \text{prox}_{\alpha R}(\theta - \alpha\gamma))$  有唯一  $\frac{1}{2}\|\cdot\|_2^2$  的 minimizer)

$$\Rightarrow H(\alpha) = h(\alpha, w^*)$$

Here, heuristically,

$$H'(\alpha) = \frac{\partial h(\alpha, w^*)}{\partial \alpha}$$

$$= -\frac{1}{2\alpha^2} \|w^* - \theta\|_2^2$$

$$= -\frac{1}{2} g_1(\alpha)^2$$

这样就吧  $H(\alpha)$  和  $w^*$  联系在一起

这里没有  $\frac{\partial w^*}{\partial \alpha}$ , 只是一个 heuristic 的结论.

同理,

$\therefore w^*$  也和  $\alpha$  相关.

$$\text{Set } \tilde{H}(\alpha) = H\left(\frac{1}{\alpha}\right)$$

$$\tilde{H}'(\alpha) = -\frac{1}{\alpha^2} H'\left(\frac{1}{\alpha}\right)$$

$$= -\frac{1}{\alpha^2} \left(-\frac{\alpha^2}{2}\right) \left\| \text{prox}_{\frac{1}{2}R} \left(\theta - \frac{1}{\alpha}\theta\right) - \theta \right\|_2^2$$

$$= \frac{1}{2} g_2\left(\frac{1}{\alpha}\right)^2$$

$$\text{Note: } \tilde{H}(\alpha) = \inf_w h\left(\frac{1}{\alpha}, w\right)$$

$$= \inf_w \left\{ \theta^T (w - \theta) + \frac{\alpha}{2} \|w - \theta\|_2^2 + R(w) \right\}$$

This is a pointwise minimize of affine function of  $\alpha$ ,  $\therefore \tilde{H}(\alpha)$  is concave function



pointwise of concave is concave

$\Rightarrow H'$  is decreasing

(这里的意思是 concave 前面的系数  $\alpha$ )

$\Rightarrow g_2(\alpha) \uparrow$

下面来看  $g_1$ .

Envelope theorem

\* deal differentiability of Mordieu envelope

\* Belongs to Danskin theorem.

(directional differentiability)

大概就是说  $\frac{dw^*}{d\alpha}$  可以省掉?

下面看那些 accumulation point 能否是同一个  $\alpha$

Prop 2: Let  $\Theta$  be the opt. soln set of our RLM problem, assumed to be non-empty under previous setting.

(c) (lost-to-go Estimate)

$$F(\theta^{k+1}) - \hat{V} \leq K_3 \left[ \text{dist}(\theta^k, \Theta)^2 + \|\theta^k - \theta^{k+1}\|_2^2 \right]$$

where  $\text{dist}(\theta, \Theta) = \inf_{r \in \Theta} \|\theta - r\|_2$  这是投影

Rmk:  $\text{dist}$  is well defined

$\because \Theta$  is close, convex set under our assumptions

To get convergence rate, we need estimate of  $\text{dist}(\theta^k, \Theta)$

Towards that end, consider

<sup>Assumption</sup>  
(A4) For any  $V \geq \hat{V}$ ,  $\exists \mu, \varepsilon$ , s.t.

$\text{dist}(\theta, \Theta) \leq \mu \|E(\theta)\|_2$  for any  $\theta$  satisfying

$$F(\theta) \leq V \text{ and } \|E(\theta)\|_2 \leq \varepsilon$$

(Local error bound condition)

就是说这个不等式在一个小邻域中成立 (故称局部条件)

with (a) (b) (c) + (A4)

$$F(\theta^{k+1}) - \hat{V} \stackrel{(c, A4)}{\leq} K_4 \left[ \|E(\theta)\|_2^2 + \|\theta^k - \theta^{k+1}\|_2^2 \right]$$

$$\stackrel{(b)}{\leq} K_5 \|\theta^k - \theta^{k+1}\|_2^2$$

$$\stackrel{(a)}{\leq} K_6 [F(\theta^k) - F(\theta^{k+1})]$$

$$= K_6 [F(\theta^k) - \hat{V} - (F(\theta^{k+1}) - \hat{V})] \stackrel{\text{recurrence}}{+ \dots}$$

$$\Rightarrow F(\theta^{k+1}) - \hat{V} \leq \underbrace{\frac{R_6}{1+R_6}}_{< 1} (F(\theta^k) - \hat{V})$$

Moreover: For the iterate:

by (d).

$$\|\theta^k - \theta^{k+1}\|_2^2 \leq \frac{1}{R_1} (F(\theta^k) - \hat{V}) \quad \because F(\theta^{k+1}) \geq \hat{V}$$

by (A4), (b):

$$\begin{aligned} \text{dis}(\theta^k, \Theta) &\leq R \|\theta^k - \theta^{k+1}\|_2 \\ &\leq R' \sqrt{F(\theta^k) - \hat{V}} \end{aligned}$$

这就关联了 value 的收敛速度和  $\theta$  的收敛速度。

至此证明收敛以及速度, 还差 (c) 和 (A4) 没证明。

(a) (b) (c) 都是 PGM 的性质, 都依赖于 PGM, 而 (A4) 是与算法无关的。

Recap.

$$(*) \quad \hat{v} = \min_{\theta} \{ F(\theta) \triangleq \mathcal{L}(\theta) + R(\theta) \}$$

$$\text{PGM: } \theta^{k+1} \leftarrow \text{prox}_{\alpha_k R} \left( \theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) \right)$$

$$\text{where } \text{prox}_R(\theta) = \arg \min_{\gamma} \left\{ \frac{1}{2} \|\theta - \gamma\|_2^2 + R(\gamma) \right\}$$

prop: Suppose  $0 < \underline{\alpha} \leq \alpha_k \leq \bar{\alpha} < \frac{1}{L}$

(a) (Sufficient Decrease)

$$F(\theta^k) - F(\theta^{k+1}) \geq K_1 \|\theta^k - \theta^{k+1}\|_2^2$$

make it making progress

(b) (Safe guard)

$$\|E(\theta^k)\|_2 \leq K_2 \|\theta^k - \theta^{k+1}\|_2$$

where

$$E(\theta) = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) - \theta$$

(c) (Cost-to-Go Estimate)

$$F(\theta^{k+1}) - \hat{v} \leq K_3 \left[ \text{dist}(\theta^k, \Theta)^2 + \|\theta^k - \theta^{k+1}\|_2^2 \right]$$

where  $\Theta$  is the opt soln set of (\*)

How far it is from the optimal value

下面先证明 (c)

Pf of (c)

close, convex  
set

Let  $\bar{\theta}^k$  be the projection of  $\theta^k$  onto  $\Theta$

Then,  $\text{dist}(\theta^k, \Theta) = \|\theta^k - \bar{\theta}^k\|_2$

We compute:

$\because \bar{\theta}^k \in \Theta$  is opt soln

$$F(\theta^{k+1}) - \hat{V} = F(\theta^{k+1}) - F(\bar{\theta}^k)$$

$$= \mathcal{L}(\theta^{k+1}) - \mathcal{L}(\bar{\theta}^k) + R(\theta^{k+1}) - R(\bar{\theta}^k)$$

$$= \nabla \mathcal{L}(\hat{\theta}^k)^T (\theta^{k+1} - \bar{\theta}^k) + R(\theta^{k+1}) - R(\bar{\theta}^k)$$

( $\because$  mean-value theorem,  $\hat{\theta}^k \in [\theta^{k+1}, \bar{\theta}^k]$ )

这就是某个中值定理

$$= \underbrace{\nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \bar{\theta}^k) + R(\theta^{k+1}) - R(\bar{\theta}^k)}_I + \underbrace{\nabla \mathcal{L}(\hat{\theta}^k) - \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \bar{\theta}^k)}_{II}$$

下面分别来 bound I 和 II.

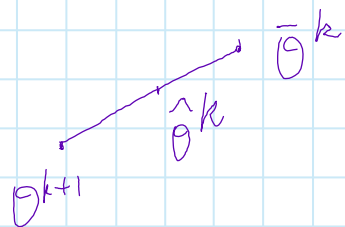
$$(I): I \leq \|\nabla \mathcal{L}(\hat{\theta}^k) - \nabla \mathcal{L}(\theta^k)\|_2 \cdot \|\theta^{k+1} - \bar{\theta}^k\|_2 \quad (C-5)$$

$$\leq L \cdot \|\hat{\theta}^k - \theta^k\|_2 \cdot \|\theta^{k+1} - \bar{\theta}^k\|_2$$

Note:  $\|\theta^{k+1} - \bar{\theta}^k\|_2 \leq \|\theta^{k+1} - \theta^k\|_2 + \|\theta^k - \bar{\theta}^k\|_2 \leftarrow \text{dist}(\theta^k, \Theta)$   
= 三角不等式.



下式  $\|\hat{\theta}^k - \theta^k\|_2$



$$\|\hat{\theta}^k - \theta^k\|_2 \leq \|\theta^{k+1} - \theta^k\|_2 + \|\theta^{k+1} - \bar{\theta}^k\|_2 \quad (\text{三角不等式})$$

$$\leq \|\theta^{k+1} - \theta^k\|_2 + \|\theta^{k+1} - \bar{\theta}^k\|_2 \quad (\text{由图中的距离关系})$$

$$\leq L \left[ 2\|\theta^{k+1} - \theta^k\|_2 + \text{dist}(\theta^k, \Theta) \right] \cdot \left[ \|\theta^{k+1} - \theta^k\|_2 + \text{dist}(\theta^k, \Theta) \right]$$

$$\leq 2L \left( \|\theta^{k+1} - \theta^k\|_2 + \text{dist}(\theta^k, \Theta) \right)^2$$

$$\leq 4L \left[ \|\theta^{k+1} - \theta^k\|_2^2 + \text{dist}^2(\theta^k, \Theta) \right]$$

$$\because (a+b)^2 \leq 2a^2 + 2b^2$$

(I): By definition of  $\theta^{k+1}$ ;

$$\frac{1}{2} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \theta^{k+1}\|_2 + \alpha_k R(\theta^{k+1})$$

$$\left( \text{P.P.}: \theta^{k+1} \text{ minimize } \frac{1}{2} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \theta^{k+1}\|_2^2 + \alpha_k R(\theta^{k+1}) \right)$$

$$\leq \frac{1}{2} \|\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k) - \bar{\theta}^k\|_2^2 + \alpha_k R(\bar{\theta}^k)$$

( $\because \theta^{k+1}$  是 optimal  
把  $\theta^{k+1}$  换成  $\bar{\theta}^k$  会变大)

化简后得:

$$R(\theta^{k+1}) - R(\bar{\theta}^k) + \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \bar{\theta}^k)$$

$$\leq \frac{1}{2\alpha_k} \|\bar{\theta}^k - \theta^k\|_2^2 \leq \frac{1}{2\alpha_k} \text{dist}^2(\theta^k, \Theta)$$

$\therefore$  (I) 和 (II) 都小于  $L[\text{dist}^2(\theta^k, \Theta) + \|\theta^k - \theta^{k+1}\|_2^2]$

下面来看上节课的 Assumption (A4)

To get convergence rate, we make the following assumption:

(A4) (Local) Error Bound:

For any  $v \geq \hat{v}$ ,  $\exists \mu, \epsilon > 0$ , s.t.

$$\text{dist}(\theta, \Theta) \leq \mu \|E(\theta)\|_2$$

← 这个就是说我们知道  $E(\theta) = 0$  那么

for any  $\theta$ :

$$F(\theta) \leq v, \|E(\theta)\|_2 \leq \epsilon$$

$\theta \in \Theta$  它是一个 opt. soln.  
那么如果  $E(\theta) \neq 0$ , 但很小, 我们问 opt soln set 是否与之接近.

This relates to the geometry of the optimization problem.

在这里是成立的, 但不意味着对所有问题成立

和算法性质 (例如 (a) (b) (c) 无关).

下面看成立的情况

Scenario I:

$L$  strongly convex,  $\nabla L$  Lipschitz cont.

strongly convex:  $L(x) \geq L(\theta) + \nabla L(\theta)^T (x - \theta) + \frac{\mu}{2} \|x - \theta\|_2^2$   
can find a quadratic support rather than just a linear support.

Fact:

$$(\nabla L(x) - \nabla L(\theta))^T (x - \theta) \geq \mu \|x - \theta\|_2^2$$

wikipedia 上对 strongly convex 的 第二种定义

Let  $\hat{\theta}$  be the optimal solution to  $(*)$ . Then,   
 只有-个.

for any  $\theta$ :

$$K \cdot \text{dist}(\theta, \Theta)^2 = K \|\theta - \hat{\theta}\|_2^2$$

$$\leq (\nabla \mathcal{L}(\theta) - \nabla \mathcal{L}(\hat{\theta}))^T (\theta - \hat{\theta}) \quad (\text{strong convex 的 fact})$$

下面需要和  $E(\theta)$  联系起来.

By the FO condition,  $0 \in \nabla \mathcal{L}(\hat{\theta}) + \partial R(\hat{\theta})$

$$\Rightarrow \boxed{-\nabla \mathcal{L}(\hat{\theta}) \in \partial R(\hat{\theta})} \quad (1)$$

On the other hand.

$$E(\theta) = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) - \theta$$

$$\Rightarrow \theta + E(\theta) = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) = \arg \min_y \|\theta - \nabla \mathcal{L}(\theta) - y\|_2^2 + R(y)$$

using the FO condition for the  $\text{prox}_R$

$$0 \in \theta + E(\theta) - \theta + \nabla \mathcal{L}(\theta) + \partial R(\theta + E(\theta))$$

(代入 prox 的定义, 这里  $y = \theta + E(\theta)$ )

$$\Rightarrow \boxed{-[\nabla \mathcal{L}(\theta) + E(\theta)] \in \partial R(\theta + E(\theta))} \quad (2)$$

接下来看如何把 sub-differential 转换成不等式.

By def of subdifferential

$$(1) \Rightarrow R(\theta + E(\theta)) \geq R(\hat{\theta}) - \nabla \mathcal{L}(\hat{\theta})^T (\theta + E(\theta) - \hat{\theta})$$

$$(2) \Rightarrow R(\hat{\theta}) \geq R(\theta + E(\theta)) - [\nabla \mathcal{L}(\theta) + E(\theta)]^T (\hat{\theta} - (\theta + E(\theta)))$$

加起来  $\Rightarrow$

$$0 \geq (\underbrace{\nabla f(\theta) + E(\theta)} - \underbrace{\nabla f(\hat{\theta})})^T (\theta + E(\theta) - \hat{\theta})$$

$$(\nabla f(\theta) - \nabla f(\hat{\theta}))^T (\theta - \hat{\theta}) + \|E(\theta)\|_2^2 \leq$$

$$E(\theta)^T (\nabla f(\hat{\theta}) - \nabla f(\theta) + \hat{\theta} - \theta)$$

$$\leq \|E(\theta)\|_2 \left[ \|\nabla f(\hat{\theta}) - \nabla f(\theta)\|_2 + \text{dist}(\theta, \Theta) \right]$$

L-S + 三角不等式,

$$\leq \|E(\theta)\|_2 \left( \underbrace{L \cdot \|\hat{\theta} - \theta\|_2}_{\text{Lipschitz 条件}} + \text{dist}(\theta, \Theta) \right)$$

Lipschitz 条件

$$\leq (L+1) \cdot \text{dist}(\theta, \Theta) \cdot \|E(\theta)\|_2$$

$$\Rightarrow \text{dist}(\theta, \Theta) \leq \frac{L+1}{K} \|E(\theta)\|_2$$

注:  $\frac{L+1}{K}$  - 一般表征 condition number, of the problem.

L 衡量 smoothness, K 衡量曲率.

L 越小, K 越大越好.

Scenario 2:

$\mathcal{L}$  takes the form  $\mathcal{L}(\theta) = h(A\theta)$

for some  $A \in \mathbb{R}^{m \times d}$   $h$  is strongly convex

on compact sets, smooth, w/ Lipschitz continuous

grad.  $R$  has polyhedral epigraph.

① strongly convex on compact  $C \ni K = K(C) > 0$   

$$h(x) \geq h(\theta) + \nabla h^T(\theta)(x - \theta) + \frac{K}{2} \|x - \theta\|_2^2 \quad \forall x, \theta \in C$$

这时候  $L$  不一定是 strongly convex 了。因为如果  $A$  存在 null space, 则  $L$  上有一个平面等值。

②  $\text{epi}(R) = \{(\theta, t) \in \mathbb{R}^d \times \mathbb{R} : R(\theta) \leq t\}$

e.g.  $R(\theta) = \|\theta\|_1$

polyhedral epigraph =  $\text{epi}(R)$  is a polyhedron  
 即 epigraph 是有限个半平面相交而成。



e.g.  $R(\theta) = \|\theta\|_1$   
 $R(\theta) = \|\theta\|_\infty$

Application under Scenario 2:

$$R(\theta) = \lambda \|\theta\|_1 \quad \lambda > 0.$$

(1) least squares.

$$L(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2 = h(X\theta)$$

$$h(u) = \frac{1}{2n} \|y - u\|_2^2$$

$$\nabla^2 h = \frac{1}{2n} I \quad \therefore \text{globally convex}$$

(2) logistic

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i a_i^T \theta))$$

$$= h(A\theta)$$

$$h(u) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i u_i))$$

(not convex over the entire space)

Hessian 矩阵恒等于0,  $\therefore$  没有一个全局的 LB, 但如果限制在一个 compact set 上就能有一个 LB.

2019.2.18

Recall

(Local) Error Bound

$$\hat{v} = \min_{\theta} \{ \mathcal{L}(\theta) + R(\theta) \} \quad (P)$$

$\mathcal{H}$ : optimal solution set

(EB) For any  $v \geq \hat{v}$ ,  $\exists \mu, \varepsilon > 0$

$$\text{dist}(\theta, \mathcal{H}) \leq \mu \|E(\theta)\|_2 \quad \text{右边比左边好计算}$$

for any  $\theta: F(\theta) \leq v, \|E(\theta)\|_2 \leq \varepsilon$ .

Recall:  $E(\theta) = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) - \theta$

$E(\theta)$  is the residual measure.

(S1):  $\mathcal{L}$  strongly convex,  $\nabla \mathcal{L}$  is Lipschitz

continuous  $\Rightarrow$  (EB) hold.

(S2):  $\mathcal{L}(\theta) = h(A\theta)$ ,  $h$ : strongly convex  
on compact sets,  $A$  is a linear operator  
 $h$ : smooth, Lipschitz continuous  $\nabla h$ .

$R$  has polyhedral epigraph.  $\mathcal{H}$  is compact.  
i.e.  $\text{epi}(R) = \{(x, t) : R(x) \leq t\}$  is a polyhedron  
(intersection of limited half-space)

e.g. ① LASSO

$$\mathcal{L}(\theta) = \frac{1}{2} \|y - A\theta\|_2^2 \quad R(\theta) = \|\theta\|_1$$

② Logistic regression

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_i \log(1 + \exp(-b_i a_i^T \theta))$$

$$R(\theta) = \|\theta\|_1$$

本节主要记录 (S2):

先给出大概步骤:

Step 1. characterize  $\Theta$  先看 opt. soln set 长啥样.

$$\begin{aligned} \text{Idea 1: } \Theta &= \{\theta : 0 \in \nabla \mathcal{L}(\theta) + \partial R(\theta)\} \quad (\text{FOL}) \\ &= \{\theta : \theta = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta))\} \end{aligned}$$

This is true but not very useful  
We should explore the structure by  
using the assumptions.

Idea 2:

$$\begin{aligned} \text{prop: } \exists \bar{y} \text{ s.t. } \forall \hat{\theta} \in \Theta, A\hat{\theta} = \bar{y} \text{ and} \\ \nabla \mathcal{L}(\hat{\theta}) = A^T \nabla h(\bar{y}) \triangleq \bar{g} \end{aligned}$$

直观想法: 在之前 strongly convex 的 case 中,  
有唯一性保证, 但是这里即使  $h(\cdot)$  是  
strongly convex, 但也不能保证  $h(Ax)$  是 strongly



convex 的。这样  $\theta$  可能就在一个 flat part 上。那这个 prop 的作用就在于即使  $\theta$  不是唯一的，但  $A\hat{\theta} = \bar{y}$  的  $\bar{y}$  却是唯一的。

$$\mathcal{L}(\theta) = h(A\theta) \quad \nabla \mathcal{L}(\theta) = A^T \nabla h(A\theta)$$

$$\text{if } \hat{\theta} \in \Theta: \nabla \mathcal{L}(\hat{\theta}) = A^T \nabla h(A\hat{\theta}) = A^T \nabla h(\bar{y}) \\ = \bar{y} \text{ 是唯一的!}$$

In particular,

$$0 \in \Theta \Rightarrow 0 \in \nabla \mathcal{L}(\hat{\theta}) + \partial R(\hat{\theta}) \\ = \bar{y} + \partial R(\hat{\theta})$$

$$\Leftrightarrow -\bar{y} \in \partial R(\hat{\theta})$$

$\therefore$  可以写  $\Theta$ : Ideal, Ideal  
FOC

$$\Theta = \{ \theta : A\theta = \bar{y}, -\bar{y} \in \partial R(\theta) \}$$

以下来证明  $\bar{y}$  的唯一性。

Pf: Let  $\theta_1, \theta_2 \in \Theta$ , Set  $\bar{y}_1 = A\theta_1, \bar{y}_2 = A\theta_2$

By strong convexity of  $h$ ,

$$h\left(\frac{\bar{y}_1 + \bar{y}_2}{2}\right) < \frac{1}{2}h(\bar{y}_1) + \frac{1}{2}h(\bar{y}_2) \quad (\because \text{strong convex})$$

$$h\left(\frac{A\theta_1 + A\theta_2}{2}\right)$$

代入得:

$$\textcircled{1} \mathcal{L}\left(\frac{\theta_1 + \theta_2}{2}\right) < \frac{1}{2} \mathcal{L}(\theta_1) + \frac{1}{2} \mathcal{L}(\theta_2)$$

Also, by convexity of  $R$ ,

$$\textcircled{2} R\left(\frac{\theta_1 + \theta_2}{2}\right) \leq \frac{1}{2} R(\theta_1) + \frac{1}{2} R(\theta_2)$$

$\textcircled{1} + \textcircled{2} \Rightarrow$

$$F\left(\frac{\theta_1 + \theta_2}{2}\right) < \frac{1}{2} F(\theta_1) + \frac{1}{2} F(\theta_2) = \hat{V}$$

$$\parallel \hat{V} \quad \parallel \frac{1}{2} \hat{V} \quad \parallel \frac{1}{2} \hat{V}$$

$\hat{V}$  ( $\because \theta_1, \theta_2$  都是 optimal)

$$\therefore \hat{V} < \hat{V} \text{ 矛盾, } \therefore \bar{\eta}_1 = \bar{\eta}_2 \quad \square$$

In particular,

$$\Theta = \Theta_L \cap \Theta_R \text{ where}$$

$$\Theta_L = \{\theta : A\theta = \bar{y}\} \quad \Theta_R = \{\theta : -\bar{q} \in \partial R(\theta)\}$$

这样就把  $\Theta$  写成了两个集合的交集

下面看  $\text{dist}(\theta, \Theta)$  能不能拆, 如果能拆开就简单了.

$$\text{dist}(\theta, \Theta) = \text{dist}(\theta, \Theta_L \cap \Theta_R)$$

$$\left\{ \begin{array}{l} \text{dist}(\theta, \Theta_L), \text{dist}(\theta, \Theta_R) \end{array} \right.$$

Question:

dist to intersection  $\stackrel{?}{\approx}$  dist to individual component

Step 2: Relationship between  $\text{dist}(\theta, \Theta_L \cap \Theta_R)$   
and  $\text{dist}(\theta, \Theta_L)$ ,  $\text{dist}(\theta, \Theta_R)$

先做笔记:

Plan: (1) show  $\Theta_L \cap \Theta_R$  are polyhedral

$$(2) \text{dist}(\theta, \Theta_L \cap \Theta_R) \leq c [\text{dist}(\theta, \Theta_L) + \text{dist}(\theta, \Theta_R)]$$

Note:

$\Theta_L$  is obviously polyhedral since it's solution to linear system.

$\Theta_R$ , consider the following facts:

(1) if  $R$  has polyhedral epigraph, so does its conjugate  $\tilde{R}$ ,  $\tilde{R}(y) = \sup_{\theta} \{\theta y - R(\theta)\}$

(2)  $\partial \tilde{R} = (\partial R)^{-1}$ , i.e.

$$\partial \tilde{R}(y) = (\partial R)^{-1}(y) \triangleq \{\theta : y \in \partial R(\theta)\}$$

(3). If  $R$  has polyhedral epigraph and  $R(\theta)$  is finite then  $\partial R(\theta)$  is a polyhedron.

可以在 convex Wiki: RockFeller 书

下面来看  $\Theta_R$ :

$$\textcircled{H} R = (\partial R)^{-1}(-\bar{q}) \stackrel{\text{by def'n}}{=} \partial \tilde{R}(-\bar{q}) \quad (2)$$

下面的检查是否  $\tilde{R}(-\bar{q}) < \infty$  只需要证明可以  
以用(3)来证明  $\partial \tilde{R}(-\bar{q})$  是 polyhedral.

Note:

$$\textcircled{H} R \neq \emptyset \quad \because \hat{\theta} \in \textcircled{H} \text{ will belong to } \textcircled{H} R$$

Fact: If  $R$  is a norm, then

$$\tilde{R} = \begin{cases} 0 & \text{if } R^*(x) \leq 1 \\ +\infty & \text{o/w} \end{cases}$$

$$\therefore \tilde{R}(-\bar{q}) < \infty$$

故  $\textcircled{H} R$  和  $\textcircled{H} R$  都是 polyhedral.

Next: Estimate point-to-polyhedron dist.

Fact: (Hoffman Error Bound)

Let  $P = \{z \mid Cz \leq d\}$  be a non-empty polyhedron

Then,  $\exists c > 0$  (which depends only on  $C$ ) s.t.

$$\text{dist}(x, P) \leq c \| (Cx - d)^+ \|_2 \quad \forall x$$

with this fact, we first prove the following

Corollary: Let  $\{P_1, \dots, P_m\}$  be a finite collection of polyhedron. s.t.  $P = \bigcap_{i=1}^m P_i \neq \emptyset$ . Then,  $\exists \alpha > 0$ :

$$\text{dist}(x, P)^2 \leq 2 \sum_{i=1}^M \text{dist}(x, P_i)^2 \quad \forall x.$$

i.e.  $\{P_1, \dots, P_M\}$  is linearly regular

regular means it's somehow nice,

since we can bound the distance by individual dist.

Apply the corollary to  $\{\Theta_L, \Theta_R\}$ :

$$\text{dist}(\theta, \Theta) \leq C \cdot [\text{dist}(\theta, \Theta_L) + \text{dist}(\theta, \Theta_R)]$$

Also,  $\text{dist}(\theta, \Theta_L) \leq C' \|A\theta - \bar{y}\|_2$  (By Hoffman BD)

目前的问题是在  $\text{dist}(\theta, \Theta_L) \leq C' \|A\theta - \bar{y}\|_2$  和  $\text{dist}(\theta, \Theta_R) \leq C'' \|\theta\|_2$  形式不太一样: ②  $\text{dist}(\theta, \Theta_R)$  怎么用?

Key result to prove the corollary:

Let  $H = \{z : c^T z \leq \delta\}$  Then,

$$\text{dist}(x, H) = \frac{(c^T x - \delta)^+}{\|c\|_2}$$

Hint:  $\text{dist}(x, H)^2 = \min \{\|x - z\|_2^2 : c^T z \leq \delta\}$   
write down KKT.

2019.2.19

Recamp:

$$\hat{V} = \min_{\theta} \{ F(\theta) \stackrel{\Delta}{=} \underbrace{\mathcal{L}(\theta) + R(\theta)}_{h(A\theta)} \}$$

$h$ : str. convex on compact sets

$\nabla h$ : Lipschitz cont.

$R$ : norm, polyhedral epigraph

$\mathbb{H}$ : opt. set, assumed to be compact.

Claim:  $\mathbb{H} = \mathbb{H}_L \cap \mathbb{H}_R$  where

$$\mathbb{H}_L = \{ \theta : A\theta = \bar{y} \} \text{ for some } \bar{y}$$

$$\mathbb{H}_R = \{ \theta : -\bar{g} \in \partial R(\theta) \}; \quad \bar{g} = A^T \nabla h(\bar{y})$$

Claim:  $\mathbb{H}_L, \mathbb{H}_R$  are polyhedral

Corollary:  $\text{dist}(\theta, \mathbb{H}) = \text{dist}(\theta, \mathbb{H}_L \cap \mathbb{H}_R)$

$$\leq c [\text{dist}(\theta, \mathbb{H}_L) + \text{dist}(\theta, \mathbb{H}_R)] \quad \forall \theta$$

(linear regularity of  $\{\mathbb{H}_L, \mathbb{H}_R\}$ )

这里  $\mathbb{H}_L, \mathbb{H}_R$  是 polyhedral

当它们不是 polyhedron 时, 如何构造反例?

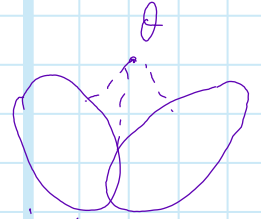
答: 构造一个  $\theta$  的序列

$$\text{dist}(\theta^k, \Theta_L) \rightarrow 0, \text{dist}(\theta^k, \Theta_R) \rightarrow 0$$

$$\text{dist}(\theta^k, \Theta) \rightarrow 0$$

同时让右边两个趋于 0 的速度比左边快,  
这样就没法找到那个  $c$ .

## EXERCISE



这个不行,  $\theta$  只有一个数.

By Hoffman Bound

$$\text{dist}(\theta, \Theta_L) \leq c' \|A\theta - \bar{y}\|_2$$

$$\therefore \text{dist}(\theta, \Theta) \leq c' [\|A\theta - \bar{y}\|_2 + \text{dist}(\theta, \Theta_R)] \leq \mu \|E(\theta)\|_2$$

今天的目标.



$$E(\theta) = \text{prox}_{\Theta}(\theta - \nabla L(\theta)) - \theta$$

Observe:  $\Theta_R = (\partial R)^{-1}(-\bar{q})$

Idea: If we view  $\theta \in (\partial R)^{-1}(q)$  for some  $q$

$$\text{then } \text{dist}(\theta, \Theta_R) \approx \text{dist}((\partial R)^{-1}(q), (\partial R)^{-1}(-\bar{q}))$$

$(\partial R)^{-1}$  是一个映射, 把  $q$  映射成一个集合.

如果它有某种 Lipschitz 性质, 那么这个集合  
就可以用  $q$  来 bound.

$$(\partial R)^{-1}(q) \triangleq \{\theta: q \in \partial R(\theta)\} \text{ 把 } q \text{ 映射成集合.}$$

下面来定义  $(\partial R)^{-1}$  的连续性. OLC

Def: Outer Lipschitz continuity of  $(\partial R)^{-1}$

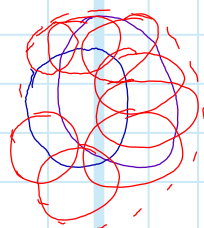
We say  $(\partial R)^{-1}$  is outer Lipschitz continuous if  $\exists \beta > 0$ , s.t., for any  $g'$ ,  $\exists$  neighborhood  $V_{g'}$  of  $g'$  s.t.,  $(\partial R)^{-1}(g'') \subseteq (\partial R)^{-1}(g') + \beta \|g' - g''\|_2 B(0, 1) \quad \forall g'' \in V_{g'}$

集合的 Lipschitz 连续:

改变输入, 会得到一个新集合, 新集合比原来的集合只改变了一点。

如何刻画它?

把原来的集合的每点扩大成一个圆, 能盖住新的集合。



Outer 的意思就是外面把它包住。

e.g.  $R(\theta) = |\theta|$

$$\partial R(\theta) = \begin{cases} \{1\} & \theta > 0 \\ [-1, 1] & \theta = 0 \\ \{-1\} & \theta < 0 \end{cases}$$

For  $g' = 1$ :

$$(\partial R)^{-1}(1) = \{\theta : 1 \in \partial R(\theta)\} = [0, \infty)$$

For  $g' = -1$ :

$$(\partial R)^{-1}(-1) = \{\theta : -1 \in \partial R(\theta)\} = (-\infty, 0]$$



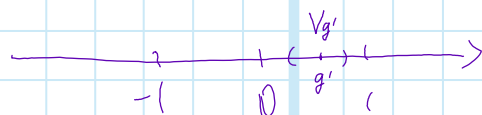
for  $g' \in (-1, 1) =$

$$(\partial R)^{-1}(g') = \{0\}$$

下面考虑 Outer Lipschitz.

take  $g' \in (-1, 1)$

$$\{0\} \subseteq \{0\} + \beta \|g' - g''\|_2 B(0, 1) \quad \forall \beta > 0$$



take  $g' = 1$

$$\cdot \subseteq \mathbb{R}_+ + \beta \|g' - g''\|_2 \cdot B(0, 1)$$

当  $g'$  取到  $-1$  时  $\mathbb{R}_- \not\subseteq \mathbb{R}_+ + \beta \|g' - g''\|_2 B(0, 1)$

Fact: If  $R$  has polyhedral epigraph then,  $(\partial R)^{-1}$  is OLL

$$\Rightarrow (\partial R)^{-1}(-g) \subseteq (\partial R)^{-1}(-\bar{g}) + \beta \|g - \bar{g}\|_2 B(0, 1) \quad \forall g \in V_{-\bar{g}}$$

↑  
这仍然是  $\Theta_R$

$$\Rightarrow \text{if } \theta \in (\partial R)^{-1}(-g) \text{ then } \text{dist}(\theta, \Theta_R) \leq \beta \|g - \bar{g}\|_2$$

for  $\theta \in (\partial R)^{-1}(-g), -g \in V_{-\bar{g}}$

这样我们就得到了  $\text{dist}(\theta, \Theta_R)$  的 bound,

∴ 之前那个不等式可以继续写为:

$$\leq c' [\|A\theta - \bar{y}\|_2 + \text{dist}(\theta, \Theta_R)] \quad \forall \theta$$

$$\leq c'' [\|A\theta - \bar{y}\|_2 + \|g - \bar{g}\|_2] \text{ for } \theta \in (\partial h)^{-1}(-g), -g \in V_{-\bar{g}}$$

下面再研究如何合并这两项. 注意到  $\theta$  和  $g$  有对应关系  $\theta \in (\partial h)^{-1}(-g), -g \in V_{-\bar{g}}$ .

Now, consider,  $-(\underbrace{\nabla \mathcal{L}(\theta)}_{-g} + \underbrace{E(\theta)}_{-\theta'}) \in \partial R(\theta + E(\theta))$

$$\text{b/c } \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) = \underset{\gamma}{\text{argmin}} \left\{ \frac{1}{2} \|\theta - \nabla \mathcal{L}(\theta) - \gamma\|_2^2 + R(\gamma) \right\}$$

写的 FOL 就是上面那个

$$0 \in \underbrace{\theta - \nabla \mathcal{L}(\theta)}_{E(\theta)} + \partial R(\gamma)$$

$$\gamma = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta))$$

$$\Rightarrow \theta \in (\partial R)^{-1}(-g')$$

provided  $-(\nabla \mathcal{L}(\theta) + E(\theta)) \in V_{-g'}$ ,  $g'$  和  $g$  代入:

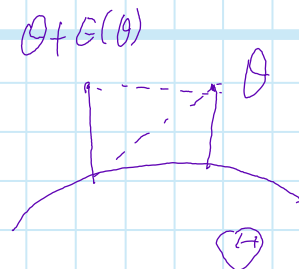
$$\text{dist}(\theta + E(\theta), \Theta) \leq c'' [\|A(\theta + E(\theta)) - \bar{y}\|_2 + \|\nabla \mathcal{L}(\theta) + E(\theta) - \bar{g}\|_2]$$

$$\leq c'' [\|A\theta - \bar{y}\|_2 + \|\nabla \mathcal{L}(\theta) - \bar{g}\|_2 + (\|A\| + 1) \|E(\theta)\|_2]$$

$$\leq c'' [\|A\theta - \bar{y}\|_2 + \|E(\theta)\|_2]$$

$$\begin{aligned} \|\nabla \mathcal{L}(\theta) - \bar{g}\|_2 &= \|A^T \nabla h(A\theta) - A^T \nabla h(\bar{y})\|_2 \\ &\leq \|A\| \cdot \|\nabla h(A\theta) - \nabla h(\bar{y})\|_2 \\ &\leq L \cdot \|A\| \cdot \|A\theta - \bar{y}\|_2 \end{aligned}$$

$$\begin{aligned} \therefore \text{dist}(\theta, \Theta) &\leq \text{dist}(\theta + E(\theta), \Theta) + \|E(\theta)\|_2 \\ &\leq c [\|A\theta - \bar{y}\|_2 + \|E(\theta)\|_2] \end{aligned}$$



Now, by strong convex of  $h$ ,

$$\bar{y} = A\hat{\theta}$$

$$\downarrow \theta = \pi_{\Theta}(\theta)$$

$$K \|A\theta - \bar{y}\|_2^2 \leq (\nabla h(A\theta) - \nabla h(\bar{y}))^T (A\theta - \bar{y})$$

$$\left[ \begin{array}{l} \text{由: } h(y) \geq h(x) + \nabla h(x)^T (y-x) + \frac{K}{2} \|y-x\|_2^2 \quad (\text{for convex}) \\ h(x) \geq h(y) + \nabla h(y)^T (x-y) + \frac{K}{2} \|x-y\|_2^2 \\ \text{加起来} \end{array} \right] \text{ 得到上面那个.}$$

$$= (A^T (\nabla h(A\theta) - \nabla h(\bar{y})))^T (\theta - \hat{\theta})$$

$$= (\nabla \mathcal{L}(\theta) - \bar{g})^T (\theta - \hat{\theta})$$

Hence:

$$(a+b)^2 \leq a^2 + b^2$$

$$\text{dist}^2(\theta, \hat{\theta}) \leq \tau' [ (\nabla \mathcal{L}(\theta) - \bar{g})^T (\theta - \hat{\theta}) + \|E(\theta)\|_2^2 ]$$

$$\stackrel{\text{claim}}{\leq} \tau'' \text{dist}(\theta, \hat{\theta}) \cdot \|E(\theta)\|_2$$

$$\text{dist}(\theta, \hat{\theta}) \leq \tau'' \|E(\theta)\|_2$$

□ 证毕.

下面来 verify 这个 claim

$$- \nabla \mathcal{L}(\hat{\theta}) \in \partial R(\hat{\theta})$$

$$- (\nabla \mathcal{L}(\theta) + E(\theta)) \in \partial R(\theta + E(\theta))$$

$$\Leftrightarrow R(\theta + E(\theta)) \geq R(\hat{\theta}) - \nabla \mathcal{L}(\hat{\theta})^T (\theta + E(\theta) - \hat{\theta})$$

$$R(\hat{\theta}) \geq R(\theta + E(\theta)) - (\nabla \mathcal{L}(\theta) + E(\theta))^T (\hat{\theta} - \theta - E(\theta))$$

以上并证完 Scenario 2.

再看这种 Error Bound 分析能否推广到  
second order 的算法中.

2019.2.25

回顾: 之前讲了各种一阶方法.

$$\min_{\theta \in \mathbb{R}^d} \{ F(\theta) = \mathcal{L}(\theta) + R(\theta) \} \quad (P)$$

$$\text{PGM: } \theta^{k+1} \leftarrow \text{prox}_{\alpha_k R} (\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k))$$

$$\text{prox}_R(\theta) = \underset{\gamma}{\text{argmin}} \left\{ \frac{1}{2} \|\theta - \gamma\|_2^2 + R(\gamma) \right\}$$

We only use the first order information

这里  $\text{prox}_R(\theta)$  也是一个优化问题, 如果有解析解则能加快计算, 如果没有: 要用两重循环.

2 loop:

outer: prox-update

inner: computing the prox.

今天来讲 second order methods.

## Second order method

o Assume  $f$  is smooth consider

$$\inf_x f(x)$$

classic Newton's method: 回顾牛顿法.

In each iteration, solve a local quadratic model of  $f$ .

i.e. At  $x^k$ : 用泰勒展开.

$$d^k = \arg \min_d \left\{ f(x^k) + \nabla f(x^k)^T d + \frac{1}{2} d^T \nabla^2 f(x^k) d \right\}$$

and set next iterate to be

$$x^{k+1} \leftarrow x^k + \alpha_k d^k, \quad \alpha_k > 0 \text{ step size}$$

Need  $d^k$  to be well-defined.

If  $f$  is strongly convex, then  $d^k$  is well-defined in each iteration (unique)

if strongly convex,  $\nabla^2 f(x^k)$  是 PD: 最佳值唯一

o SOM for (P) SOM 会有如下几个问题.

-  $F$  is not twice differentiable. (因为存在  $\mathbb{R}$ )

能不能设计一个 sub-gradient 的二阶扩展?

- Is SOM reasonable?

= 阶的优势在哪?

收敛更快, 之前一阶的收敛速度是  $c^k$ ,

= 阶是  $2^k$  : 可能外层循环更少但  
内部循环更多.

. fewer outer iterations  
more expensive inner iterations.

Issue: Non-differentiability of  $F$

Recall: PGD,

PGD 的两种解释:

① 是一个求不动点的过程.

$$\textcircled{2} \theta^{k+1} \leftarrow \underset{r \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ L(\theta^k) + \nabla L(\theta^k)^T (r - \theta^k) + \frac{1}{2\alpha_k} \|r - \theta^k\|^2 + R(r) \right\}$$

从②中可见 PGD 将其分为两块, 前面可以计算, 算不出的  $R$  直接写为加

To Generalize, consider

$$\textcircled{2} \theta^{k+1} \leftarrow \underset{r \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ L(\theta^k) + \nabla L(\theta^k)^T (r - \theta^k) + \frac{1}{2\alpha_k} (r - \theta^k)^T H_k (r - \theta^k) + R(r) \right\}$$

where  $H_k$  is an (approximate) Hessian at  $\theta^k$

Example of  $H_k$ :

①  $H_k = \underline{I}$  (PGD) 当  $k_k = \underline{I}$ , 退化为 PGD

②  $H_k = \nabla^2 L(\theta^k)$

③  $H_k =$  quasi-Newton strategies.

例如对  $\nabla^2 L(\theta^k)$  的 low-rank approx  
也可以选 subsample of  $\nabla^2 L(\theta^k)$ .

BFGS

以上通过把 PGM 的形式进行扩展, 使用  $\theta^k$  的 Hessian 的近似

For simplicity, consider  $H_k = \nabla^2 L(\theta^k)$

How to solve (\*)?

• Do we need exact soln?

if no, how exact we need?

• Is  $\theta^{k+1}$  well-defined.

例如  $L(\theta) = \|y - A\theta\|_2^2$ ,  $\therefore$  Hessian  $A^T A$

一般都是 singular.  $\therefore$  不可逆, 不可求逆.

$\therefore$  不是 well-defined.  $\leftarrow$  我们引入 regularizer

To fix this, use regularization.  $\leftarrow$  之前那想为引入 structure

use  $H_k = \underbrace{\nabla^2 L(\theta^k)}_{\geq 0} + \mu_k I, \mu_k > 0$

$\underbrace{\qquad\qquad\qquad}_{> 0}$



observe,  $H_k > 0$ , 这样, 原问题是强凸

子问题的近似:

$$\theta^{k+1} \leftarrow \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\theta - \theta^k) + \frac{1}{2\lambda_k} (\theta - \theta^k)^T H_k (\theta - \theta^k) + R(\theta) \right\}$$

where  $H_k = \nabla^2 \mathcal{L}(\theta^k) + \mu_k I$ ,  $\mu_k > 0$ .

Issues to address:

- Choice of  $\lambda_k, \mu_k$

- Inexactness criterion for the sub-problem.

prop: Define

$$l_k(\theta) = \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\theta - \theta^k) + R(\theta)$$

$$q_k(\theta) = l_k(\theta) + \frac{1}{2} (\theta - \theta^k)^T H_k (\theta - \theta^k)$$

Suppose  $\theta^k$  is not optimal. Then,

$$q_k(\theta^k + \Delta^k) - q_k(\theta^k) \leq \frac{1}{2} (l_k(\theta^k + \Delta^k) - l_k(\theta^k))$$

where  $\Delta^k = \underset{\Delta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T \Delta + \frac{1}{2} \Delta^T H_k \Delta + R(\theta^k + \Delta) \right\}$

这个 prop 的结论 = 代价的差和 - 代价的导数

prop 2: For any  $k \geq 0$  and  $\beta \in (0, 1)$ , we have

$$F(\theta^k) - F(\theta^k + \alpha \Delta^k) \geq \beta [l_k(\theta^k) - l_k(\theta^k + \alpha \Delta^k)]$$

whenever  $\alpha \in (0, \bar{\alpha}]$  for some  $\bar{\alpha} > 0$ .

Next task: ① argue  $\{F(\theta^k)\}$  converges.

② 由 prop 2,  $l_k(\theta^k) - l_k(\theta^k + \alpha_k \Delta^k) \rightarrow 0$ .

③ use prop 1 and def'n of  $H_k$  to argue  $\epsilon(\theta^k) \rightarrow 0$ , for some choice of  $\mu_k$ .

(由 prop 1, 证明  $l_k$  和  $q_k$  的差  $\rightarrow 0$ .)

下面先证两个 prop.

Pf of Prop 1:

Observe,  $\theta^k + \Delta^k$  minimize  $q_k$ .  $\Rightarrow$

$$0 \in \partial q_k(\theta^k + \Delta^k)$$

$$\partial q_k(\theta) = \nabla l(\theta^k) + H_k(\theta - \theta^k) + \partial k(\theta)$$

$$= \nabla l(\theta^k) + H_k(\theta - \theta^k)$$

$$0 \in \partial q_k(\theta^k + \Delta^k) \Leftrightarrow$$

$$-H_k \Delta^k \in \partial \ell_k(\theta^k + \Delta^k)$$

$\Rightarrow \Delta^k \neq 0$ , for o/w  $\theta^k$  is optimal,  
(check the FOC) 和/或设置矛盾.

Now, the sub-differential of  $\ell_k$  is:

$$\ell_k(\eta) \geq \ell_k(\theta^k + \Delta^k) + (-H_k \Delta^k)^T (\eta - \theta^k - \Delta^k)$$

$$\text{let } \eta = \theta^k \Rightarrow$$

$$\ell_k(\theta^k) \geq \ell_k(\theta^k + \Delta^k) + (-H_k \Delta^k)^T (-\Delta^k)$$

$$= \ell_k(\theta^k + \Delta^k) + (\Delta^k)^T H_k \Delta^k$$

$$> \ell_k(\theta^k + \Delta^k) \quad (\because H_k \succ 0)$$

$$\Rightarrow \ell_k(\theta^k + \Delta^k) - \ell_k(\theta^k) = \ell_k(\theta^k + \Delta^k) + \frac{1}{2} (\Delta^k)^T H_k \Delta^k - \ell_k(\theta^k)$$

$$\leq \frac{1}{2} [\ell_k(\theta^k + \Delta^k) - \ell_k(\theta^k)]$$

2019.2.26

Recall

$$\min_{\theta \in \mathbb{R}^d} \{ F(\theta) \stackrel{\Delta}{=} \mathcal{L}(\theta) + R(\theta) \}$$

SOM:

$$\textcircled{1} \text{ Find } \Delta^k = \underset{\Delta}{\operatorname{argmin}} \left\{ \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T \Delta + \frac{1}{2} \Delta^T H_k \Delta + R(\theta^k + \Delta) \right\}$$

$$\text{where } H_k = \nabla^2 \mathcal{L}(\theta^k) + \mu_k I, \mu_k \nearrow$$

② Find step size  $\alpha$  s.t.

$$(*) \quad F(\theta^k) - F(\theta^k + \alpha \Delta^k) \geq \beta \left[ \ell_k(\theta^k) - \ell_k(\theta^k + \alpha \Delta^k) \right]$$

$$\textcircled{3} \text{ Update } \theta^{k+1} = \theta^k + \alpha \Delta^k$$

$$\ell_k(\theta) = \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\theta - \theta^k) + R(\theta)$$

$$q_k(\theta) = \ell_k(\theta) + \frac{1}{2} (\theta - \theta^k)^T H_k (\theta - \theta^k) \quad \begin{array}{l} \partial \ell_k(\theta^k + \alpha \Delta^k) + H_k \Delta^k \\ \parallel \end{array}$$

$$\Rightarrow \theta^k + \Delta^k \text{ minimize } q_k \Leftrightarrow 0 \in \partial q_k(\theta^k + \Delta^k)$$

# Preliminary Results

$$\textcircled{1} \quad l_k(\theta^k) \geq l_k(\theta^k + \Delta^k) + (\Delta^k)^T H_k(\Delta^k)$$

$$\textcircled{2} \quad q_k(\theta^k + \Delta^k) - q_k(\theta^k) \leq \frac{1}{2} [l_k(\theta^k + \Delta^k) - l_k(\theta^k)]$$

prop 2: for any  $\beta \in (0, 1)$ , (\*) holds whenever  $\alpha \in (0, \bar{\alpha}]$ ,  $\bar{\alpha} = \min \left\{ 1, \frac{2\mu_k(1-\beta)}{L_G} \right\}$

$L_G$ : Lipschitz const. for  $\nabla \mathcal{L}$

Pf: By prelim result (1):

$$l_k(\theta^k) - l_k(\theta^k + \Delta^k) \geq (\Delta^k)^T H_k \Delta^k \geq \mu_k \|\Delta^k\|_2^2 \quad (\Delta)$$

drop the  $\nabla^2$

on the other hand,

$$l_k(\theta^k + \alpha \Delta^k) = l_k((1-\alpha)\theta^k + \alpha(\theta^k + \Delta^k)) \quad \alpha \in [0, 1]$$

利用  $l_k$  的 convexity.

$$\leq (1-\alpha)l_k(\theta^k) + \alpha l_k(\theta^k + \Delta^k)$$

$$\Rightarrow l_k(\theta^k) - l_k(\theta^k + \alpha \Delta^k) \geq \alpha [l_k(\theta^k) - l_k(\theta^k + \Delta^k)]$$

$$\geq \alpha \mu_k \|\Delta^k\|_2^2 \quad (\text{由 } (\Delta))$$

利用  $\nabla \mathcal{L}$  Lipschitz cont.

$$\Rightarrow \mathcal{L}(\theta^k + \alpha \Delta^k) - \mathcal{L}(\theta^k) \leq \alpha \nabla \mathcal{L}(\theta^k)^T \Delta^k + \frac{\alpha^2}{2} L_G \|\Delta^k\|_2^2$$

(这个结论不要求  $\mathcal{L}$  是 convex)

$$\begin{aligned} \text{证: } g(\alpha) &= f(\theta^k + \alpha \Delta^k) \\ \therefore f(\theta^k + \alpha \Delta^k) - f(\theta^k) &= g(\alpha) - g(0) \\ &= \int_0^\alpha g'(t) dt. \end{aligned}$$

Hence,

$$\begin{aligned} & \underbrace{f(\theta^k)}_{\downarrow} - \underbrace{f(\theta^k + \alpha \Delta^k)}_{\rightarrow} \\ & l_k(\theta^k) + R(\theta^k) - (l_k(\theta^k + \alpha \Delta^k) + R(\theta^k + \alpha \Delta^k)) \\ & \geq l_k(\theta^k) - l_k(\theta^k + \alpha \Delta^k) - \frac{\alpha^2 L_G}{2} \|\Delta^k\|_2^2 \end{aligned}$$

$\Rightarrow$  for any  $\beta \in (0, 1)$ :

$$\begin{aligned} & f(\theta^k) - f(\theta^k + \alpha \Delta^k) - \beta [l_k(\theta^k) - l_k(\theta^k + \alpha \Delta^k)] \\ & \geq (1-\beta) [l_k(\theta^k) - l_k(\theta^k + \alpha \Delta^k)] - \frac{\alpha^2 L_G}{2} \|\Delta^k\|_2^2 \\ & \geq [(1-\beta)\mu_k - \frac{\alpha^2 L_G}{2}] \|\Delta^k\|_2^2 \end{aligned}$$

我们希望这项  $\geq 0$  才能解出那个  $\alpha$

Convergence.

以上证明了  $F$  会单调下降, 但是能否收敛到最优值?

①  $E(\theta) = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) - \theta$   
for the orig. problem.

② Note: 在我们新的近似问题中:

$$q_k(\theta) = \underbrace{\mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\theta - \theta^k) + \frac{1}{2} (\theta - \theta^k)^T H_k (\theta - \theta^k)}_{\tilde{\mathcal{L}}(\theta)} + R(\theta)$$

$\tilde{\mathcal{L}}(\theta)$  strongly convex

$$\Rightarrow E_k(\theta) = \text{prox}_R^X(\theta - \nabla \tilde{\mathcal{L}}(\theta)) - \theta$$

$$= \text{prox}_R^X(\theta - \nabla \mathcal{L}(\theta^k) - H_k(\theta - \theta^k)) - \theta$$

$$= \text{prox}_R^X((I - H_k)\theta - (\nabla \mathcal{L}(\theta^k) - H_k \theta^k)) - \theta$$

$$\text{FACT: } E(\theta^k) = E_k(\theta^k)$$

(\*)  $\Rightarrow \{F(\theta^k)\}$  monotonic decreasing  $\Rightarrow F(\theta^k) \rightarrow \bar{v}$

$$\text{Recall: } \mathcal{L}_k(\theta^k) - \mathcal{L}_k(\theta^k + \alpha_k \Delta^k) \geq \alpha_k \mu_k \|\Delta^k\|_2^2 \geq 0$$

$$\Rightarrow \lim_{k \rightarrow \infty} \alpha_k \mu_k \|\Delta^k\|_2^2 = 0$$

我们希望能由此得出  $\Delta^k = 0$ , 这样  $\theta^k$  就是 opt.

observe:

$$\|E(\theta^k)\|_2 = \|E_k(\theta^k) - E_k(\theta^k + \Delta^k)\|_2$$

$\leftarrow \because \theta^k + \Delta^k$  是当前问题的 opt  
 $\therefore$  找相当于减了个 0

$$\leq (L_G + \mu_k + 2) \|\Delta^k\|_2 \quad (\text{由 } E_k \text{ 的 Lipschitz})$$

$$\left\{ \begin{array}{l} \| \text{prox}_R^X(x) - \text{prox}_R^X(y) \|_2 \leq \|x - y\|_2 \end{array} \right.$$

$\therefore \text{prox}$  是 non-expansive operator.

类似的, projection 也有类似的中性, 经过 projection 后不会变长.

non-expansiveness of  $\text{prox}_R(\cdot)$ .

$$\Rightarrow 0 \leq \frac{2\mu_k^2(1-\beta)\|E(\theta^k)\|_2^2}{L_G(L_G + \mu_k + 2)^2} \leq \alpha_k \mu_k \|\Delta^k\|_2^2$$

$$\Rightarrow \lim_{k \rightarrow \infty} \frac{\mu_k^2 \|E(\theta^k)\|_2^2}{(L_G + \mu_k + 2)^2} = 0 \quad \text{两边夹逼.}$$

$\mu_k$  是我们选的, 有很多选法, 如果选  $\mu_k$  是个常数, 那么  $\|E(\theta^k)\|_2^2 = 0$ , 这意味着, 尽管我们一直用一个有误差的  $H_k$ , 还是可以?

也可以 choose  $\mu_k = C \cdot \|E(\theta^k)\|_2^p, p \geq 0$

$\Rightarrow \|E(\theta^k)\|_2^2$  goes to 0

Set  $\mu_k = C \cdot \|E(\theta^k)\|_2^p, p \geq 0, C > 0$ .

以上证明了 convergence, 下面来看看 convergence rate 是不是把  $p$  取得越大, 收敛越快?



prop3: Under the Error Bound assumption, the above setting  $\mu_k$ , then for large  $k$ , we can take  $\delta_k =$

$$(EB): \text{dist}(\theta, \Theta) \leq \mu \|E(\theta)\|_2 \quad \forall \theta: \|E(\theta)\|_2 \leq \varepsilon, \dots$$

$$\mu_k = \mu_k = C \cdot \|E(\theta^k)\|_2^p, \quad p \geq 0, C > 0$$

Sketch:

$$\text{Key inequality: } F(\theta^k + \delta^k) - F(\theta^k) \leq \frac{L_H}{6} \|\delta^k\|_2^3 - \frac{\mu_k}{2} \|\delta^k\|_2^2 + \frac{1}{2} [l_k(\theta^k + \delta^k) - l_k(\theta^k)]$$

$L_H$ : Lip const. of  $\nabla^2 \mathcal{L}$

$$\leq \beta [l_k(\theta^k + \delta^k) - l_k(\theta^k)] + \left(\frac{1}{2} - \beta\right) [l_k(\theta^k + \delta^k) - l_k(\theta^k)] + \frac{L_H}{6} \|\delta^k\|_2^3 - \frac{\mu_k}{2} \|\delta^k\|_2^2$$

$\leq -\mu_k \|\delta^k\|_2^2$

$$\leq \beta [l_k(\theta^k + \delta^k) - l_k(\theta^k)] + \underbrace{\frac{L_H}{6} \|\delta^k\|_2^3 - \frac{\mu_k}{2} (2 - 2\beta) \|\delta^k\|_2^2}_{\text{provided } \beta \in (0, \frac{1}{2})}$$

$$\leq \beta [l_k(\theta^k + \delta^k) - l_k(\theta^k)] \quad \text{provided } \downarrow \leq 0$$

$$\Leftrightarrow \|\delta^k\|_2 \leq \frac{6\mu_k(1-\beta)}{L_H}$$

↙ 先找一个UB, 如果这个UB对于右也,  $k \rightarrow \infty$

$$\text{Lemma: } \forall k: \|\delta^k\|_2 \leq \left[ \frac{L_H}{2\mu_k} \text{dist}(\theta^k, \Theta) + 2 \right] \text{dist}(\theta^k, \Theta)$$

$$\text{Then, } \|\delta^k\|_2 \leq \frac{L_H}{2C\|E(\theta^k)\|_2^p} \cdot [\mu\|E(\theta^k)\|_2 + 2] \cdot \mu\|E(\theta^k)\|_2$$

$$= \frac{\mu^2 L_H}{2C} \|E(\theta^k)\|_2^{2-p} + 2\mu\|E(\theta^k)\|_2 \stackrel{\text{Want}}{\leq} \frac{6(1-\beta)}{L_H} \cdot C \cdot \|E(\theta^k)\|_2^p$$

分类讨论  $\begin{cases} \rho \in [0, 1) & \text{成立} \\ \rho = 1 & , \text{成立. 左右都是一次方, 可以调 } c \text{ 成立.} \\ \rho > 1 & , \text{不成立, 无法成为 UBS.} \end{cases}$

. 'In order to guarantee prop 3,  $\rho$  should be  $\in [0, 1]$

2019.3.4

$$(P) \min_{\theta} \{ \mathcal{L}(\theta) + R(\theta) \}$$

Second Order Method

① Solve

$$\Delta^k = \arg \min_{\Delta} \left\{ \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T \Delta + \frac{1}{2} \Delta^T H_k \Delta + R(\theta^k + \Delta) \right\}$$

where  $H_k = \nabla^2 \mathcal{L}(\theta^k) + \mu_k I$ ,  $\mu_k = c \|E(\theta^k)\|_2^p$

$$E(\theta) = \text{prox}_R(\theta - \nabla \mathcal{L}(\theta)) - \theta \quad c > 0, p \in [0, 1]$$

② Update:  $\theta^{k+1} = \theta^k + \alpha_k \Delta^k$ 

prop: Under the error bound assumption, i.e.,

$$\text{dist}(\theta, \Theta) \leq \mu \cdot \|E(\theta)\|_2 \quad \text{for } \theta: \|E(\theta)\|_2 \leq \varepsilon,$$

then for sufficiently large  $k$ ,  $\alpha_k = 1$ , and descent condition can be satisfied.

在这种情况下,  $\alpha_k$  直接取 1 就行啦。

$$\text{Lemma: } \forall k, \quad \|\Delta^k\|_2 \leq \frac{L_H}{2\mu_k} \text{dist}(\theta^k, \Theta) \text{dist}(\theta^k, \Theta)$$

以上是引理。

以上已知其可以 converge, 下面看 convergence rate.

## Local Convergence Rate Analysis.

Based on EB.

说 local: 要  $k$  足够大.

For large  $k$ ,  $\Delta^k = \theta^{k+1} - \theta^k$  ( $k$  足够大  $\mu > \frac{1}{2}$ ,  $d_k = 1$ )  
 $E(\theta^k) \rightarrow \theta$

$$\|\Delta^k\|_2 = \|\theta^{k+1} - \theta^k\|_2 \leq \left[ \frac{L_H}{2C\|E(\theta^k)\|_2^p \text{dist}(\theta^k, \theta) + 2} \right] \text{dist}(\theta^k, \theta)$$

这两项有不等式关系.

by EB  $\sim \text{dist}(\theta^k, \theta)^p$

$$\sim \text{dist}(\theta^k, \theta)^{2-p} \dots \text{dist}(\theta^0, \theta)^1$$

$$\leq \mathcal{O}(\text{dist}(\theta^0, \theta))$$

$$\text{dist}(\theta^{k+1}, \theta) \stackrel{\text{EB}}{\leq} \mu \|E(\theta^{k+1})\|_2 = \mu \|E(\theta^{k+1}) - E_k(\theta^{k+1})\|_2$$

$$E_k(\theta) = \text{prox}_r(\theta - \nabla \tilde{\mathcal{L}}(\theta)) - \theta.$$

$$\tilde{\mathcal{L}}_k(\theta + \Delta) = \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T \Delta + \frac{1}{2} \Delta^T H_k \Delta.$$

$$\begin{cases} \theta^{k+1} \min \left\{ \mathcal{L}(\theta^k) + \nabla \mathcal{L}(\theta^k)^T (\theta - \theta^k) + \frac{1}{2} (\theta - \theta^k)^T H_k (\theta - \theta^k) + r(\theta) \right\} \\ \Rightarrow \bar{E}_k(\theta^{k+1}) = 0 \end{cases}$$

$$\|E(\theta^{k+1}) - E(\theta^k)\|_2$$

(by def)

$$= \left\| \text{prox}_R(\theta^{k+1} - \nabla \mathcal{L}(\theta^{k+1})) - \theta^{k+1} - \left[ \text{prox}_R(\theta^{k+1} - \nabla \mathcal{L}(\theta^k) - H_k(\theta^{k+1} - \theta^k)) - \theta^{k+1} \right] \right\|_2$$

$\nabla \mathcal{L}(\theta^{k+1})$

(non-expansiveness of prox)

$$\leq \| \nabla \mathcal{L}(\theta^{k+1}) - \nabla \mathcal{L}(\theta^k) - H_k(\theta^{k+1} - \theta^k) \|_2$$

(by def of  $H_k$  und  $\Delta$ -ineq)

$$\leq \| \nabla \mathcal{L}(\theta^{k+1}) - \nabla \mathcal{L}(\theta^k) - \nabla^2 \mathcal{L}(\theta^k)(\theta^{k+1} - \theta^k) \|_2 + \mu_k \| \theta^{k+1} - \theta^k \|_2$$

(Lipschitz cond of  $\nabla^2 \mathcal{L}$ )

$$\leq \frac{L_H}{2} \| \theta^{k+1} - \theta^k \|_2^2 + C \cdot \|E(\theta^k)\|_2^p \cdot \| \theta^{k+1} - \theta^k \|_2$$

← 下面分析这一项看是什么级别的

$$\|E(\theta^k)\|_2 = \|E(\theta^k - \bar{\theta}^k)\|_2 \quad \bar{\theta}^k = \Pi_{\Theta}(\theta^k)$$

← 其子集=0, 由定义

$$= \left\| \text{prox}_R(\theta^k - \nabla \mathcal{L}(\theta^k)) - \theta^k - \left[ \text{prox}_R(\bar{\theta}^k - \nabla \mathcal{L}(\bar{\theta}^k)) - \bar{\theta}^k \right] \right\|_2$$

non-expansiveness

$$\leq \| \theta^k - \bar{\theta}^k \|_2 + \| \theta^k - \bar{\theta}^k \|_2 + \| \nabla \mathcal{L}(\theta^k) - \nabla \mathcal{L}(\bar{\theta}^k) \|_2$$

$$\leq (L_G + 2) \text{dist}(\theta^k, \Theta)$$

$$\text{Hence, } \text{dist}(\theta^{k+1}, \Theta) \leq O(\text{dist}(\theta^k, \Theta)^2) + O(\text{dist}(\theta^k, \Theta)^{1+p})$$

$$= O(\text{dist}(\theta^k, \Theta)^{1+p})$$

superlinear when  $p \in (0, 1)$   
quadratic if  $p=1$

linear if  $\rho=0$ , provided  $\epsilon > 0$  is chosen carefully

那么我们能直接选  $\rho=1$ ?

不行! 注意到之前的那个 sub-problem,

$$L(\theta^k) + \nabla L(\theta^k)^T \Delta + \frac{1}{2} \Delta^T H_k \Delta + r(\|\theta^k + \Delta\|_2)$$

$$H_k = \nabla^2 L(\theta^k) + \mu_k I \quad \mu_k = (\|\nabla L(\theta^k)\|_2)^\rho$$

当  $\rho=0$  时,  $\mu_k$  是常数,  $\therefore$  always guaranteed

$\rho$  太大时,  $\mu_k \downarrow$ ,  $H_k$  变得 ill-conditioned.

$\therefore$  要做一些权衡.

big outer vs. small outer  
small inner vs. big inner,

下面看下一个问题.

$$(P) \quad \min_{\theta} \{ F(\theta) = L(\theta) + \lambda R(\theta) \}$$

Non-convex Instances of (P).

$L, R$ : both can be non-convex.

Example: Linear Model. w/ additively corrupted covariates

$$\text{Recall } y = X\theta^* + \varepsilon$$

$$X = \begin{bmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{bmatrix} \quad \begin{array}{l} X: \text{Covariate vector} \\ \Sigma: \text{mean } 0, \text{ iid} \end{array}$$

$$(E) \hat{\theta} \in \arg \min_{\theta} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda \|\theta\|_1 \right\}$$

Assumption:  $x_i$  is iid. w/ covariance  $\Sigma_x > 0$  (known)

P.g:  $x_i \sim \mathcal{N}(0, \Sigma_x)$

Then, an idealized version of the est. problem is:

$$\min_{\theta} \left\{ \frac{1}{2} \underbrace{\theta^T \Sigma_x \theta}_F - \underbrace{\theta^T \Sigma_x \theta^*}_y \right\}$$

这个问题的 FOC:  $\Sigma_x \theta - \Sigma_x \theta^* = 0 \Rightarrow \theta = \theta^*$

如果知道  $\theta^*$  就直接解出来了, 这里可以把  $\Sigma_x$  和  $\Sigma_x \theta^*$  当成已知参数。

How is this related to (E)?

$$\frac{1}{2n} \|y - X\theta\|_2^2 = \frac{1}{2n} [y^T y - 2 \theta^T X^T y + \theta^T X^T X \theta]$$

拆成三项。

$$\sim \frac{1}{2n} \left[ \underbrace{\theta^T X^T X \theta}_F - 2 \underbrace{\theta^T X^T y}_y \right]$$

consider  $\hat{f} = \frac{1}{n} x^T x = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$

$$E_{(x, \epsilon)} [\hat{f}] = \frac{1}{n} \sum_{i=1}^n E[x_i x_i^T] = \Sigma_x$$

$\therefore \hat{f}$  is a sample average of  $\Sigma_x$   
unbiased est. of  $\Sigma_x$

consider  $\hat{y} = \frac{1}{n} x^T y$        $y = x \theta^* + \epsilon$

$$E_{(x, \epsilon)} [\hat{y}] = E_x [E_{\epsilon} [\frac{1}{n} x^T y]]$$

$$= \frac{1}{n} E_x [x^T E_{\epsilon} [y]]$$

$$= \frac{1}{n} E_x [x^T x \theta^*]$$

$$= \Sigma_x \theta^*$$

下面看如何根据这个性质来建模.

Consider: we do not observe  $x_i$  directly,

but through

现在  $x_i$  不能直接观察到,  
也是有噪声,

$$z_i = x_i + w_i$$

$w_i$ : = random, mean 0, covariance  $\Sigma_w$  (assume to be known)



$$\frac{1}{n} E_w [z^T z] = \frac{1}{n} X^T X + \Sigma_w$$

$$\frac{1}{n} E_w [z^T y] = \frac{1}{n} X^T y$$

In this setting: We can use:

$$\hat{\beta} = \frac{1}{n} z^T z - \Sigma_w$$

$$\hat{\gamma} = \frac{1}{n} z^T y$$

Problem:  $\leftarrow$  两个 PSD - 话, 并不知道是不是 PSD 了.

$$\frac{1}{2} \left[ \theta^T \left[ \frac{1}{n} z^T z - \Sigma_w \right] \theta - \frac{1}{n} \theta^T z^T y \right]$$

$$z: n \times d \quad \Sigma_w: d \times d$$

$$\text{rank} \leq n$$

$\hookrightarrow L = \frac{1}{2} \theta^T \theta$  - 一个 loss function  $\nabla_{\theta}^2 L$  convex 的 例子

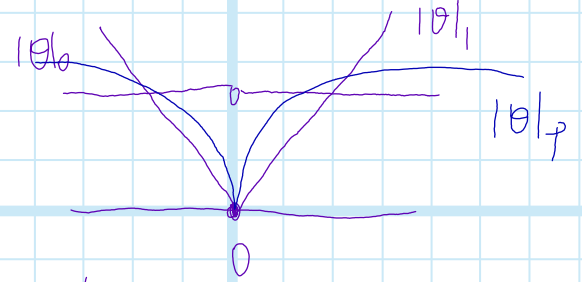
2019.3.5.

# Non-convex Regularized loss

$$(P) \quad \min_{\theta} \{ F(\theta) \triangleq \mathcal{L}(\theta) + \lambda R(\theta) \}$$

$R(\theta) = \|\theta\|_1$  convex approx. of  $\|\theta\|_0$

how many  
non-zeros  
↓



$$\|\theta\|_p^p = \sum_i |\theta_i|^p \quad 0 < p < 1$$

$\ell_p$ -quasi-norm. (bridge-penalty)

Observe:

$$\lim_{\theta \rightarrow 0} (|\theta|^p)' = +\infty$$

$\Rightarrow \hat{\theta} = 0$  is a local min to (P)

$\Rightarrow$  cannot bound  $\|\hat{\theta} - \theta^*\|_2$  for all local  
min in general

$\therefore \hat{\theta} = 0$  是一最优化解, 但并不是  $\theta^*$ .

Consider a family of regularizers 形式化定义下 regularizer.

$R_\lambda$ :

① separable:  $R_\lambda(\theta) = \sum_{i=1}^d R_\lambda(\theta_i)$  可以拆成每个维度

$$\textcircled{1} R_\lambda(0) = 0, R_\lambda(t) = R_\lambda(-t) \quad \forall t \in \mathbb{R}$$

$\textcircled{2} R_\lambda$  is non-decreasing on  $\mathbb{R}_+$

$\textcircled{3}$  For  $t > 0$ ,  $t \rightarrow \frac{R_\lambda(t)}{t}$  is non-increasing in  $t$

$\therefore$  为门接的  $\| \cdot \|_0$

$\textcircled{4} R_\lambda$  is differentiable  $\forall t \neq 0$ , and subdifferentiable at  $t=0$ , with  $\lim_{t \downarrow 0} R'_\lambda(t) = \lambda L$  for some  $L > 0$

希望在 0 点处 导数不为  $+\infty$

$\textcircled{5} \exists \mu > 0$ , s.t.,  $t \mapsto R_{\lambda, \mu}(t) = R_\lambda(t) + \frac{\mu}{2} t^2$  is convex

(weak-convexity of  $R_\lambda$ )

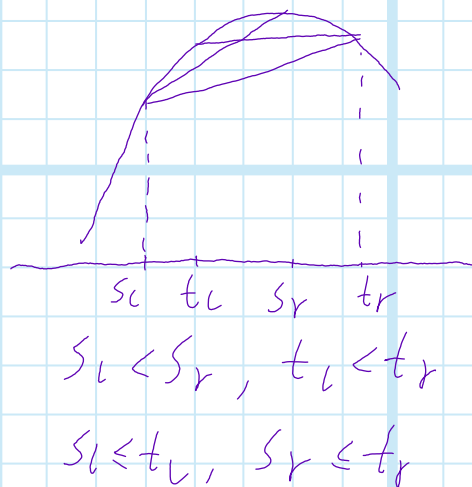
和 strong-convex 相反.

注意:  $\textcircled{3}$ : If  $R_\lambda$  is concave, Then it satisfies  $\textcircled{3}$

由图看出, 不成立:

$$\frac{R_\lambda(s_r) - R_\lambda(s_l)}{s_r - s_l} \geq \frac{R_\lambda(t_r) - R_\lambda(s_l)}{t_r - s_l}$$

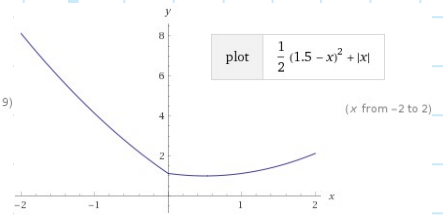
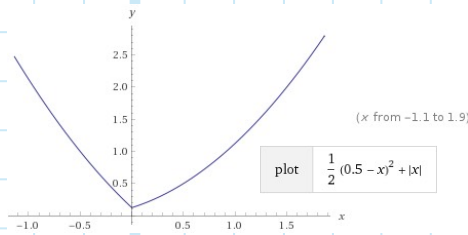
$$\geq \frac{R_\lambda(t_r) - R_\lambda(t_l)}{t_r - t_l}$$



This family implements a thresholding rule.

Consider:

$$\operatorname{argmin}_t \left\{ \frac{(z-t)^2}{2} + R_\lambda(t) \right\}$$



e.g.  $R_\lambda(t) = \lambda|t|$

给定一个  $z$ , 求出一个最佳的  $t$ .



则  $0 \in \operatorname{argmin}_t \left\{ \frac{(z-t)^2}{2} + R_\lambda(t) \right\}$  iff  $|z| \leq \lambda$

More generally, define

$$\lambda^* = \inf_{t > 0} \left\{ \frac{t}{z} + \frac{R_\lambda(t)}{t} \right\}$$

e.g.  $R_\lambda(t) = \lambda|t|$       $\lambda^* = \lambda$  (if  $z > 0$ )

Claim:  $0 = \operatorname{argmin}_t \left\{ \frac{(z-t)^2}{2} + R_\lambda(t) \right\}$  iff  $|z| \leq \lambda^*$

sketch:

$$\frac{(z-t)^2}{2} + R_\lambda(t) - \frac{z^2}{2} = t \left( \frac{t}{z} + \frac{R_\lambda(t)}{t} - z \right)$$

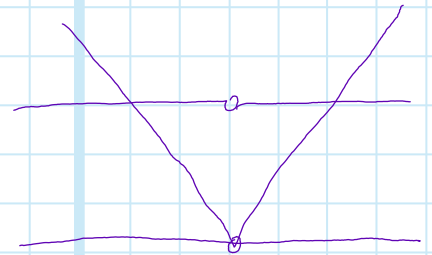
下面看  $\frac{t}{z} + \frac{R_\lambda(t)}{t}$  regularizer 的导数。

Examples:

①  $R_\lambda(t) = \lambda|t|$

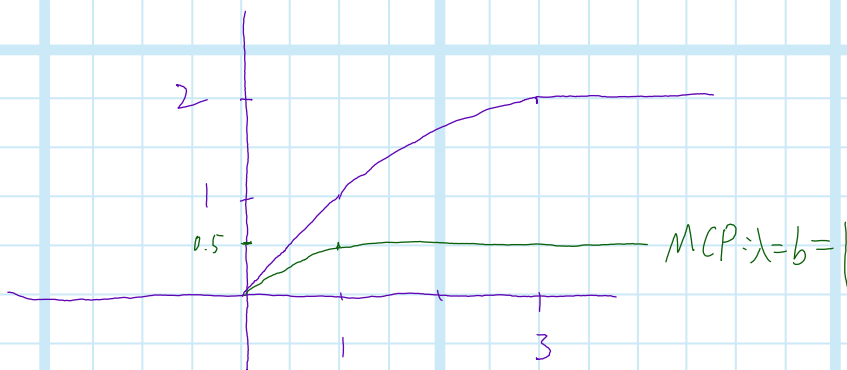
② Smoothly Clipped Absolute Deviation (SCAD)

$$R_\lambda(t) = \begin{cases} \lambda|t| & \text{for } |t| \leq \lambda \\ \frac{t^2 - 2a\lambda|t| + \lambda^2}{2(a-1)} & \text{for } \lambda < |t| \leq a\lambda \\ \frac{(a+1)\lambda^2}{2} & \text{for } |t| > a\lambda \end{cases}$$



here,  $a > 2$  is fixed

SCAD:  $\lambda=1, a=3$



以下验证 condition ④:

$$R'_\lambda(t) = \begin{cases} \text{sgn}(t) \cdot \lambda & 0 < t < \lambda \\ \frac{2t - 2a\lambda}{2(a-1)} = \frac{a\lambda - t}{a-1} & \lambda < t < a\lambda \\ 0 & t > a\lambda \end{cases}$$

t 左右对称

写成 compact 形式:

$$= \text{sgn}(t) \left[ \lambda \mathbb{1}_{\{|t| \leq \lambda\}} + \frac{(a\lambda - t)_+}{a-1} \mathbb{1}_{\{|t| > a\lambda\}} \right]$$

for  $t \neq 0$ , For  $t = 0$ ,

$$\partial R_\lambda(0) = [-\lambda, \lambda] \Rightarrow \lim_{t \rightarrow 0} R'(t) = \lambda \quad (L=1)$$

以下验证 condition ⑤

Take  $\mu = \frac{1}{a-1}$ . Then,  $R_{\lambda, \mu}(t) = R_\lambda(t) + \frac{\mu}{2} t^2 = \begin{cases} \lambda|t| + \frac{\mu}{2} t^2 \\ \frac{a\lambda}{a-1}|t| + \frac{\lambda^2}{2(a-1)} \\ \frac{(a+1)\lambda^2}{2} + \frac{\mu}{2} t^2 \end{cases}$

### ③ Minimax Contrace Penalty (MCP)

$$R_\lambda(t) = \lambda \int_0^{|t|} \left(1 - \frac{z}{\lambda b}\right)_+ dz, \quad b > 0 \text{ is fixed}$$

下面画图.

For  $t > 0$

$$R_\lambda(t) = \lambda \int_0^t \left(1 - \frac{z}{\lambda b}\right)_+ dz$$

$$= \lambda \int_0^{\min\{t, \lambda b\}} \left(1 - \frac{z}{\lambda b}\right) dz$$

$$= \lambda \left[ \min\{t, \lambda b\} - \frac{1}{2\lambda b} (\min\{t, \lambda b\})^2 \right]$$

画图: MCP:  $\lambda = b = 1$  见上一页

下面来验证 MCP 是否满足条件.

Condition ④:

$$R'_\lambda(t) = \text{sgn}(t) \cdot \lambda \cdot \left(1 - \frac{|t|}{\lambda b}\right)_+ \quad \text{for } t \neq 0.$$

$$\text{For } t=0, \quad \lim_{t \rightarrow 0} R'_\lambda(t) = \lambda \quad (L=1)$$

Condition ⑤

$$\text{Take } \mu = \frac{1}{b}$$

下面继续讨论

① stat error bound  $\|\hat{\theta} - \theta^*\|_2$  现在不是 convex 的 ↓

② alg, conv analysis.

2019.3.11

Reramp.

感觉方法还是为了做non-convex的分析, 去掉了Regularizer的conv性质, 需引入regularizer的其他性质以便分析

Non-Convex Regularizers

$$R_\lambda: \mathbb{R} \rightarrow \mathbb{R}$$

$$(1) R_\lambda(0) = 0, R_\lambda(t) = R_\lambda(-t) \quad \forall t.$$

(2)  $R_\lambda$  is non-decreasing on  $\mathbb{R}_+$

(3)  $\forall t, t \rightarrow \frac{R_\lambda(t)}{t}$  is non-increasing

(4)  $R_\lambda$  differentiable  $\forall t \neq 0$ , sub-differentiable

at  $t=0$ , with  $\lim_{t \downarrow 0} R'_\lambda(t) = \lambda$  for some  $\lambda$

(5)  $\exists \mu > 0, \forall t, t \mapsto R_{\lambda, \mu}(t) \triangleq R_\lambda(t) + \frac{\mu}{2} t^2$  is convex  
不要非凸得太严重.

上节课讲的两个例子都是 piece-wise quadratic

$\therefore$  加上一个二次项就能 convex.

$$\theta \in \min_{\|\theta\|_1 \leq R} \{F(\theta) = \mathcal{L}(\theta) + R_\lambda(\theta)\} \quad (P)$$

$R$ : radius is chosen s.t.  $\|\theta^*\|_1 \leq R$ .

Goal: Estimation/statistical error

$$\|\hat{\theta} - \theta^*\|_2$$

for any first-order point  $\hat{\theta}$ !

满足 FOC 的点, 也可能是最大点

先写 FOC:

$\tilde{\theta}$  is a first order point if:

$$\tilde{\theta} = \Pi_B(\tilde{\theta} - \nabla F(\tilde{\theta}))$$

$$B = \{\theta : \|\theta\|_1 \leq R\}$$

或者写为不等式: (projection 不好往下写, 写不等式)

$$(\tilde{\theta} - \nabla F(\tilde{\theta}) - \tilde{\theta})(\theta - \tilde{\theta}) \leq 0 \quad \forall \theta \in B$$

$$\Leftrightarrow (\nabla L(\tilde{\theta}) + \nabla \lambda_n(\tilde{\theta}))^T (\theta - \tilde{\theta}) \geq 0 \quad \forall \theta \in B$$

下面先定义一下性质.

(restricted strong convex)

Def:  $L$  satisfies the following RSC:

$$\begin{aligned} & (\nabla L(\theta^* + \Delta) - \nabla L(\theta^*))^T \Delta \\ & \geq \begin{cases} \alpha_1 \|\Delta\|_2^2 - \tau_1 \frac{\log d}{n} \|\Delta\|_1^2 & \text{if } \|\Delta\|_2 \leq 1 \\ \alpha_2 \|\Delta\|_2 - \tau_2 \sqrt{\frac{\log d}{n}} \|\Delta\|_1 & \text{if } \|\Delta\|_2 \geq 1 \end{cases} \end{aligned}$$

$n$  是样本数  
 $d$  是维度.



下面看如何从 strong convexity 事得出.

$$g(y) \geq g(x) + \nabla g(x)^T (y-x) + \frac{\alpha}{2} \|x-y\|_2^2$$

$$\underline{+)} \quad g(x) \geq g(y) + \nabla g(y)^T (x-y) + \frac{\alpha}{2} \|x-y\|_2^2$$

$$0 \geq \underbrace{(\nabla g(x) - \nabla g(y))}_{-\Delta} (y-x) + \underbrace{\alpha}_{\Delta} \|x-y\|_2^2$$

下面给出 Thm.

Thm: Under the above setting ( $L$  满足 RSC,  $R$  满足之前那五条), with:

$$\alpha_1 \geq \frac{\mu}{2}, \theta^* \text{ is feasible for (P)} \quad \theta^* \in \mathcal{B}$$

$$\frac{4}{L} \max \left\{ \|\nabla \mathcal{L}(\theta^*)\|_\infty, \alpha_2 \sqrt{\frac{\log d}{n}} \right\} \leq \lambda \leq \frac{\alpha_2}{6LR}$$

希望  $\lambda$  尽量大, 这样权重更符合 sparsity.

and  $n \geq \frac{16R^2 \max\{\tau_1^2, \tau_2^2\}}{\alpha_2^2} \log d$  every first-order point

$\tilde{\theta}$  satisfies

$$\|\tilde{\theta} - \theta^*\|_2 \leq \frac{3\lambda L \sqrt{K}}{2\alpha_1 - \mu} \quad \text{where } K = \|\theta^*\|_0$$

Pf: Define  $\tilde{\delta} = \tilde{\theta} - \theta^*$

Claim:  $\|\tilde{\delta}\|_2 \leq$

注意到 RSC 是分两段的.

By RSC,

$$(\nabla \mathcal{L}(\tilde{\theta}) - \nabla \mathcal{L}(\theta^*))^T \tilde{\delta} \geq$$

$\geq \alpha_1 \|\tilde{\Delta}\|_2^2 - \tau_1 \frac{\log d}{n} \|\tilde{\Delta}\|_1^2$  下面把不一致的也往  $\|\tilde{\Delta}\|_2^2$  凑。

Since  $R_{\lambda, \mu}(t) + \frac{\mu}{2} t^2$  is convex, (Regularizer ⑤)

$$R_{\lambda, \mu}(\theta^*) - R_{\lambda, \mu}(\tilde{\theta}) \geq \nabla R_{\lambda, \mu}(\tilde{\theta})^T (\theta^* - \tilde{\theta})$$

$$= (\nabla R_{\lambda}(\tilde{\theta}) + \mu \tilde{\theta})^T (-\tilde{\Delta})$$

$$\Rightarrow \nabla R_{\lambda}(\tilde{\theta})^T (\theta^* - \tilde{\theta}) \leq R_{\lambda}(\theta^*) - R_{\lambda}(\tilde{\theta}) + \frac{\mu}{2} \|\tilde{\theta} - \theta^*\|_2^2$$

$$\therefore \alpha_1 \|\tilde{\Delta}\|_2^2 - \tau_1 \frac{\log d}{n} \|\tilde{\Delta}\|_1^2$$

$$\leq \underbrace{\nabla \mathcal{L}(\tilde{\theta})^T \tilde{\Delta}}_{\text{FOC}} - \nabla \mathcal{L}(\theta^*)^T \tilde{\Delta} \quad (\text{def of RSC})$$

$$\leq -\nabla R_{\lambda}(\tilde{\theta})^T \tilde{\Delta} - \nabla \mathcal{L}(\theta^*)^T \tilde{\Delta}$$

$$\leq -\nabla \mathcal{L}(\theta^*)^T \tilde{\Delta} + R_{\lambda}(\theta^*) - R_{\lambda}(\tilde{\theta}) + \frac{\mu}{2} \|\tilde{\Delta}\|_2^2$$

$$\leq \underbrace{\|\nabla \mathcal{L}(\theta^*)\|_{\infty}}_{\leq \frac{\lambda L}{4}} \|\tilde{\Delta}\|_1 + \underbrace{R_{\lambda}(\theta^*) - R_{\lambda}(\tilde{\theta})}_{\text{需要和 } \|\tilde{\Delta}\|_1 \text{ 产生联系}} + \frac{\mu}{2} \|\tilde{\Delta}\|_2^2$$

为了让  $R_{\lambda}$  和  $\|\tilde{\Delta}\|_1$  产生联系, 引入如下:

Claim 2 Let  $\theta \in \mathbb{R}^d$  and  $S$  be index set of the

$k$  largest entries (in magnitude) of  $\theta$ . Then

$R_{\lambda}(\theta_S) - R_{\lambda}(\theta_{S^c})$ , Then,

$$R_\lambda(\theta_s) - R_\lambda(\theta_{sc}) \leq \lambda L (\|\theta_s\|_1 - \|\theta_{sc}\|_1)$$

Also, if  $\theta^*$  is  $k$ -sparse. then,  $\forall \theta$ :

$$R(\theta^*) - R_\lambda(\theta) \leq \lambda L (\|\Delta_s\|_1 - \|\Delta_{sc}\|_1),$$

$\Delta = \theta - \theta^*$ ,  $s$ : index set of  $k$  largest entries of  $\Delta$ .

$$\leq \frac{\lambda L}{4} \underbrace{\|\tilde{\Delta}\|_1}_{\|\tilde{\Delta}_s\|_1 + \|\tilde{\Delta}_{s^c}\|_1} + \lambda L (\|\tilde{\Delta}_s\|_1 - \|\tilde{\Delta}_{sc}\|_1) + \frac{\mu}{2} \|\tilde{\Delta}\|_2^2$$

下面合并下同类项.

$$(\alpha_1 - \frac{\mu}{2}) \|\tilde{\Delta}\|_2^2 \leq \frac{5\lambda L}{4} \|\tilde{\Delta}_s\|_1 - \frac{3\lambda L}{4} \|\tilde{\Delta}_{sc}\|_1 + \tau_1 \frac{\log d}{n} \|\tilde{\Delta}\|_1^2$$

右边的  $\|\cdot\|_1$  有平方, 有的没有. 取用那个 Ball 来打平方.

$$\tau_1 \frac{\log d}{n} \|\tilde{\Delta}\|_1^2 \leq 2R \tau_1 \frac{\log d}{n} \|\tilde{\Delta}\|_1 \quad (\text{三角不等式})$$

$$\leq \frac{5\lambda L}{4} \|\tilde{\Delta}_s\|_1 - \frac{3\lambda L}{4} \|\tilde{\Delta}_{sc}\|_1 + 2R \tau_1 \frac{\log d}{n} \|\tilde{\Delta}\|_1$$

$$\leq \dots \dots \dots + \alpha_2 \sqrt{\frac{\log d}{n}} \|\tilde{\Delta}\|_1$$

这里是用  $n \geq \frac{16n^2 \max\{\tau_1^2, \tau_2^2\}}{\alpha_2^2} \log d$  那个假设

$$\leq \dots \dots \dots + \frac{\lambda L}{4} \|\tilde{\Delta}\|_1$$

(利用前面的 Assumption.)

注意 2)  $\|\tilde{\Delta}\|_1 = \|\tilde{\Delta}_S\|_1 + \|\tilde{\Delta}_{S^c}\|_1$

$$\leq \frac{3\lambda L}{2} \|\tilde{\Delta}_S\|_1 - \frac{\lambda L}{2} \|\tilde{\Delta}_{S^c}\|_1$$

$\Rightarrow$

$$2\left(\alpha_1 - \frac{\mu}{2}\right) \|\tilde{\Delta}\|_2^2 \leq 3\lambda L \|\tilde{\Delta}_S\|_1 \leq 3\lambda L \sqrt{K} \|\tilde{\Delta}_S\|_2$$

$$\leq 3\lambda L \sqrt{K} \|\tilde{\Delta}\|_2$$

下面来证明 Claim 2.

Pf of Claim 2: Define  $f(t) = \frac{t}{R_\lambda(t)}$  for  $t > 0$ ,

$$\|\theta_{S^c}\|_1 = \sum_{j \in S^c} R_\lambda(\theta_j) \underbrace{\frac{|\theta_j|}{R_\lambda(|\theta_j|)}}_{f(|\theta_j|)} \stackrel{\text{由 } R_\lambda \text{ 的 (3)}}{\leq} \sum_{j \in S^c} R_\lambda(\theta_j) f(\|\theta_{S^c}\|_\infty)$$

$$= R_\lambda(\theta_{S^c}) f(\|\theta_{S^c}\|_\infty)$$

Similarly

$$R_\lambda(\theta_S) \cdot f(\|\theta_{S^c}\|_\infty) \leq \|\theta_S\|_1 \quad \underline{\text{EXERCISE}}$$

注意  $\theta_{S^c}$  中的每一维都小于  $\theta_S$  中的每一维

$$R_\lambda(\theta_s) - R_\lambda(\theta_{s_0}) \leq \frac{1}{f(\|\theta_{s_0}\|_\infty)} (\|\theta_s\|_1 - \|\theta_{s_0}\|_1)$$

A4)

For  $t \geq s > 0$

$$f(t) \geq f(s) = \frac{s-0}{R_\lambda(s) - R_\lambda(0)} \quad \text{取个 lim 也成}$$

$$\Rightarrow f(t) \geq \lim_{s \downarrow 0} f(s) = \frac{1}{\lambda L}$$

代回去即证完了 Claim 2.

下节课讲如何计算 first order point.

2019.3.12

本节主要讲算法.

Recamp:

$$(P) \min \{L(\theta) + R_\lambda(\theta)\}$$

$L$ : smooth

non-convex

$R_\lambda$ : Regularizer

non-smooth

$\rightarrow$  (P) is non-convex non-smooth.

Algorithmily,

Find descent direction at  $x$

For smooth problems, e.g.

$$\min f(x)$$

a descent direction  $-\nabla f(x)$

For non-convex, non-smooth:

NP-hard.

下面来证 NP-hard.

Consider  $c \in \mathbb{Z}_+^n$  s.t.  $\gamma = \sum_i c_i \geq 1$  (不全为0)

Define  $f(x) = (1 - \frac{1}{\gamma}) \max_i |x_i| - \min_i |x_i| + |c^T x|$

Note:  $f(0) = 0$

(Claim (Nesterov)) It is NP-hard to decide if  $\exists x$ , s.t.  $f(x) < 0$

Pf: This is equivalent to deciding if

$\exists \sigma \in \{\pm 1\}^n$  s.t.  $c^T \sigma = 0$  (就是把集合分成两个和相等的问题)

(Partition Problem, NP-hard)

下面证这两个问题是等价.

( $\Leftarrow$ ) Suppose  $\exists \sigma \in \{\pm 1\}^n$  s.t.  $c^T \sigma = 0$

Then  $f(\sigma) = -\frac{1}{\gamma} < 0$

( $\Rightarrow$ ) Suppose  $\exists x$ , s.t.  $f(x) < 0$ .

Normalize:  $\max |x_i| = 1$

Set  $\delta = |c^T x|$ , Then

$0 > f(x) = 1 - \frac{1}{\gamma} - \min_i |x_i| + \delta$

$\Leftrightarrow |x_i| > 1 - \frac{1}{\gamma} + \delta \quad \forall i$ .

Set  $\sigma_i = \text{sgn}(x_i)$ . Then,

$$|x_i| = \sigma_i x_i > \frac{1}{r} - \delta$$

$$\Rightarrow |\sigma_i - x_i| = |\sigma_i| \cdot |\sigma_i - x_i| \\ = 1 - \sigma_i x_i < \frac{1}{r} - \delta$$

$$\Rightarrow |c^T \sigma| = |c^T x| + |c^T (\sigma - x)| \quad (\text{triangle inequality})$$

$$\leq \delta + \|c\|_1 \cdot \|\sigma - x\|_\infty$$

$$< \delta + r \left( \frac{1}{r} - \delta \right) \quad (-S)$$

$$= 1 + (1-r)\delta \leq 1$$

$$\Rightarrow c^T \sigma = 0 \quad \text{b/c } c \in \mathbb{Z}_+^n, \sigma \in \{\pm 1\}^n$$

如果  $|c^T \sigma| < 0$  它只能  $= 0$ .

∴ 我们找到了, 只要取  $\sigma_i = \text{sgn}(x_i)$  那么就能直接  $\Rightarrow c^T \sigma = 0$ .

根据  $(\Leftarrow)$  和  $(\Rightarrow)$  推出了两个问题是可以在多项式时间内相互转换, 因而是等价的复杂度.

Use structure of (P):

$$\mathcal{L}(\theta) + R_\lambda(\theta)$$

利用  $R_\lambda(\theta)$  是 weakly convex.

$$= \underbrace{\left( \mathcal{L}(\theta) - \frac{\mu}{2} \|\theta\|_2^2 \right)}_{\text{Smooth non-convex}} + \underbrace{\left( R_\lambda(\theta) + \frac{\mu}{2} \|\theta\|_2^2 \right)}_{\text{convex non-smooth}}$$



Nesterov: composite gradient.

"converge to a point with no descent direction".

下面来看算法:

Want to solve:

$$\min_{\bar{R}_{\lambda, \mu}(\theta) \leq R} \left\{ \underbrace{\left( \mathcal{L}(\theta) - \frac{\mu}{2} \|\theta\|_2^2 \right)}_{\bar{\mathcal{L}}(\theta)} + \lambda \bar{R}_{\lambda, \mu}(\theta) \right\}$$

$$\text{where } \bar{R}_{\lambda, \mu}(\theta) = \frac{1}{\lambda} R_{\lambda, \mu}(\theta) = \frac{1}{\lambda} \left[ R_{\lambda}(\theta) + \frac{\mu}{2} \|\theta\|_2^2 \right]$$

Alg:

$$(A) \quad \theta^{t+1} = \underset{\bar{R}_{\lambda, \mu}(\theta) \leq R}{\operatorname{argmin}} \left\{ \frac{1}{2} \left\| \theta - \left( \theta^t - \frac{1}{\eta} \nabla \bar{\mathcal{L}}(\theta^t) \right) \right\|_2^2 + \frac{\lambda}{\eta} \bar{R}_{\lambda, \mu}(\theta) \right\}$$

是  $\uparrow$  strongly convex problem, w.r.t convex region

compare:

$$\operatorname{prox}_{\frac{1}{\eta}(\lambda \bar{R}_{\lambda, \mu})} \left( \theta^t - \frac{1}{\eta} \nabla \bar{\mathcal{L}}(\theta^t) \right) \triangleq \underset{\theta}{\operatorname{argmin}}(\dots)$$

也可以用两个算子:

2-stage approach

Step 1 = Compute

$$\theta^{t+1} = \operatorname{prox}_{\frac{1}{\eta}(\lambda \bar{R}_{\lambda, \mu})} \left( \theta^t - \frac{1}{\eta} \nabla \bar{\mathcal{L}}(\theta^t) \right)$$

Step 2 =

If  $\bar{R}_{\lambda, \mu}(\tilde{\theta}^{t+1}) \subseteq R$  then

$$\theta^{t+1} \leftarrow \tilde{\theta}^{t+1}$$

Else, set

$$\theta^{t+1} = \operatorname{argmin}_{\bar{R}_{\lambda, \mu}(\theta) \subseteq R} \left\{ \frac{1}{2} \left\| \theta - \left( \theta^t - \frac{\nabla \tilde{L}(\theta^t)}{\eta} \right) \right\|_2^2 \right\}$$

$$= \Pi_{\bar{R}_{\lambda, \mu}(\theta) \subseteq R} \left[ \theta^t - \frac{\nabla \tilde{L}(\theta^t)}{\eta} \right]$$

下面来 verify 用两步的是等价的  
Correctness:

If  $\bar{R}_{\lambda, \mu}(\tilde{\theta}^{t+1}) \subseteq R$ , Then  $\checkmark$

Else, observe:

有两个优化问题, 一个有约束, 一个没有  
 $x^* = \operatorname{argmin} f(x)$ ,  $\hat{x} = \operatorname{argmin}_{x \in C} f(x)$   
若  $x^* \notin C$ ,  $\Rightarrow \hat{x} \in \text{bd}(C)$   
如果不在  $C$  内, 那么  $\hat{x}$  应该在  $C$  的边界上。

$$\bar{R}_{\lambda, \mu}(\theta^{t+1}) = R \quad (\text{在边界上})$$

Claim:  $\theta^{t+1}$  solves step 2.

Pf: Suppose not. Then,  $\exists \bar{\theta}$ , s.t.

$$\bar{R}_{\lambda, \mu}(\bar{\theta}) \leq R \text{ and}$$

$$\begin{aligned} & \frac{1}{2} \|\bar{\theta} - (\theta^t - \frac{\nabla \bar{J}(\theta^t)}{\eta})\|_2^2 + \frac{\lambda}{\eta} \bar{R}_{\lambda, \mu}(\bar{\theta}) \\ & < \frac{1}{2} \|\theta^{t+1} - (\theta^t - \frac{\nabla \bar{J}(\theta^t)}{\eta})\|_2^2 + \frac{\lambda}{\eta} \bar{R}_{\lambda, \mu}(\theta^{t+1}) \end{aligned}$$

↖ 下面那项  
↓ 更大.

Contradicts optimality of  $\theta^{t+1}$  wrt (A)

Sometimes, step 1 can be computed in closed form.

Recall in step 1:

$$\operatorname{argmin}_{\theta} \left\{ \frac{1}{2} \|\theta - (\theta^t - \frac{\nabla \bar{J}(\theta^t)}{\eta})\|_2^2 + \frac{1}{\eta} R_{\lambda}(\theta) + \frac{\mu}{2\eta} \|\theta\|_2^2 \right\}$$

$$= \operatorname{argmin}_{\theta} \left\{ \frac{1}{2} (1 + \frac{\mu}{\eta}) \|\theta\|_2^2 - \theta^T (\theta^t - \frac{\nabla \bar{J}(\theta^t)}{\eta}) + \frac{1}{\eta} R_{\lambda}(\theta) \right\}$$

$$= \operatorname{argmin}_{\theta} \left\{ \frac{1}{2} \|\theta - \frac{1}{1 + \frac{\mu}{\eta}} (\theta^t - \frac{\nabla \bar{J}(\theta^t)}{\eta})\|_2^2 + \frac{1/\eta}{1 + \mu/\eta} R_{\lambda}(\theta) \right\}$$

The last expression takes the form

$$\frac{1}{2} (x - c)^2 + \nu R_{\lambda}(x)$$

in each coordinate

e.g. SCAD.

$$SCAD: R_\lambda(t) = \begin{cases} \lambda|t| & |t| \leq \lambda \\ -\frac{t^2 - 2a\lambda|t| + \lambda^2}{2(a-1)} & \lambda < |t| \leq a\lambda \\ \frac{(a+1)t^2}{2} & |t| > a\lambda \end{cases}$$

FOC:

$$0 \in \partial \left[ \frac{1}{2}(x-c)^2 + \nu R_\lambda(x) \right]$$

$$= x - c + \nu \partial R_\lambda(x)$$

$$= \begin{cases} -c + \nu[-\lambda, \lambda] & x=0 \\ x - c + \nu\lambda & 0 < x \leq \lambda \\ x - c + \frac{\nu(a\lambda - x)}{(a-1)} & \lambda < x \leq a\lambda \\ x - c & x > a\lambda \end{cases}$$

..... 关于轴对称

$$\hat{x} = \begin{cases} 0 & \text{if } |c| \leq \nu\lambda \\ c - \text{sgn}(c)\nu\lambda & \text{if } \nu\lambda \leq |c| \leq (\nu+1)\lambda \\ \left(1 - \frac{\nu}{a-1}\right)^{-1} \left(c - \text{sgn}(c) \frac{a\nu\lambda}{a-1}\right) & \text{if } (\nu+1)\lambda \leq |c| \leq a\lambda \\ c & \text{if } |c| \geq a\lambda \end{cases}$$

以上是 SCAD 的估计, 来计算 step 1.

本讲讲了算法(A) 和其等价算法.

方法在于转化为能求解和 convex 的部分.

Recall

2019.3.18

$$\begin{aligned} & \min_{\theta} \{ \mathcal{L}(\theta) + R_{\lambda}(\theta) \} \\ & = \min_{\theta} \left\{ \underbrace{\left( \mathcal{L}(\theta) - \frac{\mu}{2} \|\theta\|_2^2 \right)}_{\substack{\text{smooth} \\ \text{convex}}} + \underbrace{\left( R_{\lambda}(\theta) + \frac{\mu}{2} \|\theta\|_2^2 \right)}_{\substack{\text{non-smooth} \\ \text{convex}}} \right\} \end{aligned}$$

非 smooth, non-convex 拆成了两个.  
Algorithm:

$$\theta^{t+1} = \underset{\bar{R}_{\lambda, \mu}(\theta)}{\operatorname{argmin}} \left\{ \frac{1}{2} \left\| \theta - \left( \theta^t - \frac{\nabla \bar{\mathcal{L}}(\theta^t)}{\eta} \right) \right\|_2^2 + \frac{\lambda}{\eta} \bar{R}_{\lambda, \mu}(\theta) \right\}$$

$$\bar{R}_{\lambda, \mu} = \frac{1}{\lambda} \left( R_{\lambda}(\theta) + \frac{\mu}{2} \|\theta\|_2^2 \right)$$

Def:  $\mathcal{L}$  satisfies RSC if

$$\mathcal{L}(\theta_1) - \mathcal{L}(\theta_2) - \nabla \mathcal{L}(\theta_2)^T (\theta_1 - \theta_2) \geq$$

$$\begin{cases} \alpha_1 \|\theta_1 - \theta_2\|_2^2 - \tau_1 \frac{\log d}{n} \|\theta_1 - \theta_2\|_2^2 & \text{if } \|\theta_1 - \theta_2\|_2 \leq C \\ \alpha_2 \|\theta_1 - \theta_2\|_2^2 - \tau_2 \sqrt{\frac{\log d}{n}} \|\theta_1 - \theta_2\|_2 & \text{if } \|\theta_1 - \theta_2\|_2 \geq C \end{cases}$$

Def:  $\mathcal{L}$  satisfies restricted smoothness (RS) if

$$\mathcal{L}(\theta_1) - \mathcal{L}(\theta_2) - \nabla \mathcal{L}(\theta_2)^T (\theta_1 - \theta_2) \leq$$

$$\alpha_3 \|\theta_1 - \theta_2\|_2^2 + \tau_3 \frac{\log d}{n} \|\theta_1 - \theta_2\|_1^2$$

如果右边没有后面那个修正项, 就变成了普通的 Lipschitz continuous gradient.

这两条提供了局部收敛速度的保证.

Thm (Informal)

Suppose  $\mathcal{L}$  satisfies RSC & RS, Under appropriate choice of parameters, for sufficiently large  $t$ ,

$$\max \left\{ \|\theta^t - \hat{\theta}\|_2, \|\hat{\theta} - \theta^*\|_2 \right\} = O\left(\sqrt{\frac{k \log d}{n}}\right)$$

where  $k = \|\theta^*\|_0$  代表 sparsity

it establish a bound on  $\|\theta^t - \hat{\theta}\|_2$ , the bound does not converge to 0.

Even we found the optimality  $\hat{\theta}$ , we still have statistical error.

下面看另一类问题.

Phase Synchronization.

(有更广的一类叫 group synchronization)

Goal:

measurements: 观测到的信号

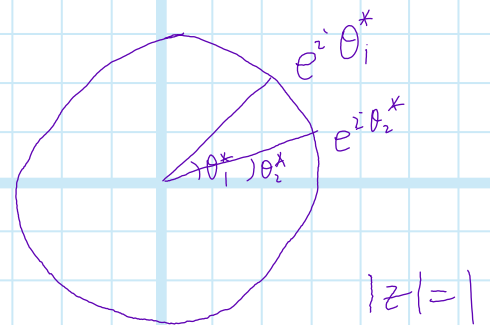
$$e^{i(\theta_j^* - \theta_k^*)} \quad \forall j, k, \quad \text{这里 } i \text{ 是虚数.}$$

recover: 希望 recover 每个  $j$  的数据.

$$e^{i\theta_j^*} \quad \forall j = 1, \dots, n$$

类似于一个 Localization 问题.

已知距离求座标.



Remark: There's a "group" of interpretation of this problem

一个群是一个代数结构, 定义了 identity operation, and inverse.

$$g_j^* = e^{i\theta_j^*} \quad g_j^* (g_k^*)^{-1}$$

↑ group element.

问题是用群来解释就是：已知两个群元素的组合，能否还原成各个元素。

★应用：有两组点云，能否组合到一起，  
 底下 group synchronization.

Let  $z^* \in \mathbb{T}^n = \{w \in \mathbb{C}^n : |w_i| = 1 \forall i\}$  be the ground truth.

$\mathbb{T}$  表示 Torus?

measurement model:

$$C_{j,l} = z_j^* \bar{z}_l^* + \Delta_{j,l} \quad 1 \leq j < l \leq n.$$

实际上  $\bar{z}_k^* = e^{-i\theta_k^*}$

$\Delta_{j,l}$  是噪声，加了噪声后，可能就不是 group element

$$z_j^* \bar{z}_l^* \Delta_{j,l} \in \mathbb{T}$$

Least-squares formulation:

$$\hat{z} \in \operatorname{argmin}_{z \in \mathbb{T}^n} \sum_{j,l} |C_{j,l} - z_j \bar{z}_l|^2$$

compact, Weierstrass



Assume:  $\Delta_{jj} = 0 \Rightarrow C_{jj} = 1$  我们知道转换为

把那个平方打开, 利用  $|z_j - \bar{z}_j|^2 = 1$  EXERCISE  
问题可转化为:

$$\hat{z} \in \arg \max_{z \in \mathbb{T}^n} z^H C z \quad (P)$$

$z^H$ : Hermitian transpose of  $z$ .

Assume  $\Delta = \Delta^H$ , then  $C = C^H$

这里如是  $C = C^H$  那么  $z^H C z \in \mathbb{R}$

证:  $(z^H C z)^H = z^H C^H z = z^H C z$

Observe: If  $\hat{z}$  is optimal, then so is

$e^{i\theta} \hat{z}$  for any  $\theta \in [0, 2\pi)$

如是  $\hat{z}$  是最优的, 把  $\hat{z}$  转一定的角度, 仍然是最优的.

可以依此来定又  $\hat{z}$  optimal 还有多远.

$\Rightarrow$  distance measure:

$$d_2(z, w) \triangleq \min_{\theta \in [0, 2\pi)} \|z - e^{i\theta} w\|_2$$

$z, w \in \mathbb{T}^n$ .

SDP

一种办法是用 semi-definite approximation

这里我们直接求解.

Q: Estimation error  $d_2(\hat{z}, z^*)$

Optimization error  $d_2(z^t, \hat{z})$

← Relaxation of  $\pi^n$

prop: Let  $z \in \mathbb{C}^n$  be s.t.  $\|z\|_2^2 = n$  and

$(z^*)^H C z^* \leq z^H C z$  Then,

$$d_2(z, z^*) \leq \frac{4\|\Delta\|}{\sqrt{n}}$$

Pf:  $d_2(z, z^*)^2 = \min_{\theta \in [0, 2\pi)} \|z - e^{i\theta} z^*\|_2^2$

$$= 2 \left[ n - \max_{\theta \in [0, 2\pi)} \operatorname{Re}(e^{i\theta} z^H z^*) \right]$$

$e^{i\theta} z^H z^* = |z^H z^*|$  时最大.

$$= 2(n - |z^H z^*|)$$

$$C_{jl} = z_j^* \bar{z}_l^* + \Delta_{jl}$$

$$C = z^* (z^*)^H + \Delta$$

$$\begin{aligned} z^H C z &= |z^H z^*|^2 + z^H \Delta z \geq (z^*)^H C z^* \\ &= n^2 + (z^*)^H \Delta z^* \end{aligned}$$

$$\Rightarrow n^2 - |z^H z^*|^2 \leq z^H \Delta z - (z^*)^H \Delta z^*$$

$$(n - |z^H z^*|)(n + |z^H z^*|) \leq z^H \Delta z - (z^*)^H \Delta z^*$$

$$n - |z^H z^*| \leq \frac{1}{n} (z^H \Delta z - (z^*)^H \Delta z^*)$$

$$\text{4.) } \|V\|_2^2 = V^H V$$

$$\|z\|_2^2 = \|z^*\|_2^2 = 1$$

加上  $z^H \Delta z^* = (z^*)^H \Delta z$  (一定是實數) 然後只保留實部。

$$= \frac{1}{n} \operatorname{Re} \left[ (z - z^*)^H \Delta (z + z^*) \right]$$

$$\leq \frac{1}{n} \|\Delta\| \cdot \|z - z^*\|_2 \cdot \|z + z^*\|_2$$

$$\|z + z^*\| \leq \|z\|_2 + \|z^*\|_2 = \sqrt{n} + \sqrt{n}$$

$$\leq \frac{2}{\sqrt{n}} \|\Delta\| \cdot \|z - z^*\|_2 = d_2(z, z^*) \leftarrow \underline{\text{EXER(ISE)}}$$

2019.3.19

# Phase Synchronization

$$C = z^* (z^*)^H + \Delta$$

$$z^* \in \mathbb{T}^n = \{w \in \mathbb{C}^n : |w_i| = 1\}$$

$$\hat{z} \in \operatorname{argmax}_{z \in \mathbb{T}^n} z^H C z$$

↑  
non-convex, ∴ 只在边界上  
而不是边界内。

Prop 1: (Est. Error) Let  $z \in \mathbb{C}^n$  be s.t.  $\|z\|_2^2 = n$

and  $(z^*)^H C z^* \leq z^H C z$ , Then,  $d_2(z, z^*) \leq \frac{4\|\Delta\|}{\sqrt{n}}$

$$d_2(z, w) = \min_{\theta \in (0, 2\pi)} \|z - e^{i\theta} w\|$$

来看看传统SDP怎么做,

SDP:

$$Z = Z^H$$

Semidefinite Relaxation

$$\max C \cdot Z$$

$$\text{s.t. } \operatorname{Diag}(Z) = e,$$

$$Z \succeq 0 \quad (w^H Z w \geq 0, \forall w \in \mathbb{C}^n)$$

我们研究另一种方法:

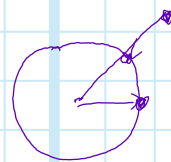
# Projected Gradient Method

$$w^k \leftarrow z^k + \frac{\alpha^k}{n} C z^k \quad \text{梯度上升}$$

$$z^{k+1} \leftarrow \frac{w^k}{\|w^k\|} \quad \text{往 feasible 集合上投影.}$$

when given  $w \in \mathbb{C}^n$

$$\left(\frac{w}{\|w\|}\right)_i = \begin{cases} \frac{w_i}{|w_i|} & \text{if } w_i \neq 0 \\ 1 & \text{if } w_i = 0 \end{cases}$$



圆心往圆上投影  
可取任一点, 这里取 1

那么关于这个算法我们可以研究:

Convergence?

initialization?

下面研究 initialization.

Consider spectral estimator:

(spectral estimator 就是用其个最大 eigen value 来估)

Let  $u \in \mathbb{C}^n$  be a leading eigenvector of  $C$

$u^H C u = \|u\|^2 \therefore$  eigen vector is real.

$V_C = \frac{u}{\|u\|} \in \mathbb{R}^n$  这算是 component wise normalization 而不是定义.

$$\text{prop 2: } d_2(v_c, z^*) \leq \frac{8 \|\Delta\|}{\sqrt{n}}$$

和 prop 1 只差了一个常数, 现在有了  $z^*$  这个点的 bound 能否下降若干步来满足 prop 1.

Thm (Est. Error of PG)

Suppose  $\|\Delta\| \leq \frac{n}{16}$ , If  $z^0 = v_c$ ,  $\alpha_n \geq 2$ , Then,

$$d_2(z^{k+1}, z^*) \leq \underbrace{\mu^{k+1} d_2(z^0, z^*)}_{\textcircled{1}} + \frac{\nu}{1-\mu} \underbrace{\frac{8 \|\Delta\|}{\sqrt{n}}}_{\textcircled{2}} \quad \forall k$$

$$\text{where, } \mu = \frac{16(\alpha \|\Delta\| + n)}{(7\alpha + 8)n} < 1 \quad \nu = \frac{2\alpha}{7\alpha + 8}$$

①: 由 optimization 的误差, 随着迭代  $\downarrow$

②: 由 noise 引入的误差

prop 3: For any  $w \in \mathbb{C}^n$  and  $z \in \Pi^n$ ,

$$\left\| \frac{w}{|w|} - z \right\|_2 \leq 2 \|w - z\|_2$$

在 convex 的情况下, 我们还有 non-expansiveness:  
 $\|\pi(x) - \pi(y)\|_2 \leq \|x - y\|_2$   
 这里没有 convex, we pay a factor 2.

下面用 prop 3 来证 prop 2.

pf (prop 2). Choose  $u$ , s.t.  $\|u\|_2^2 = n$  and

$u^H z^* = |u^H z^*|$ , (w/o g) By def'n of  $u$ ,

$$(z^*)^H C (z^*) \leq u^H C u$$

$$d_2(v_c, z^*) \leq \|v_c - z^*\|_2 \quad \text{这里固定了 } \theta = 0 \text{ 故大了.}$$

$$\leq 2 \|u - z^*\|_2 \quad (\text{prop 3})$$

$$= 2 d_2(u, z^*)$$

$$d_2^2(u, z^*) = \min_{\theta} \|u - e^{i\theta} z^*\|_2^2$$

$$= 2 \left[ n - \max_{\theta} \operatorname{Re}(e^{i\theta} u^H z^*) \right]$$

$$\leq \frac{8 \|\Delta\|}{\sqrt{n}} \quad (\text{prop 1})$$



Pf (prop 3) It suffices to prove  $\left| \frac{w_j}{|w_j|} - z_j \right| \leq 2 \left| \frac{w_j}{|w_j|} - z_j \right|$

即: 如果每一维都成立即总体也成立

WLOG, assume  $z_j = 1$ ,

If  $w_j = 0 \Rightarrow$  trivial

Consider  $w_j \neq 0$ , let  $\frac{w_j}{|w_j|} = e^{i\phi}$  for some  $\phi \in [0, 2\pi)$

Claim:  $|e^{i\phi} - 1| \leq 2 |r e^{i\phi} - 1|$  for any  $\phi \in [0, 2\pi)$

and  $r \geq 0$ .  
要证这个不等式, 只用找到  $r$  让右也最小时也成立即可.

Pf (claim):

$$|r e^{i\phi} - 1|^2 = r^2 - 2r \cos \phi + 1 \triangleq g(r)$$

Thus,

$$g'(r) = 2r - 2\cos\phi$$

$$\Rightarrow \operatorname{argmin}_{r \geq 0} |re^{i\phi} - 1|^2 = \begin{cases} 0 & \text{if } \phi \in \left[\frac{\pi}{2}, \frac{3\pi}{2}\right] \\ \cos\phi & \text{if } \phi \in [0, \frac{\pi}{2}] \cup \left(\frac{3\pi}{2}, 2\pi\right) \end{cases}$$

$$\Rightarrow \min_{r \geq 0} |re^{i\phi} - 1|^2 = \begin{cases} 1 & \text{if } \phi \in \left[\frac{\pi}{2}, \frac{3\pi}{2}\right], \\ \sin^2\phi & \text{o/w} \end{cases}$$

$$\text{For } \phi \in \left[\frac{\pi}{2}, \frac{3\pi}{2}\right]$$

$$|e^{i\phi} - 1| \leq 2 \leq 2|r e^{i\phi} - 1|$$

$$\text{o/w, } |e^{i\phi} - 1| = \sqrt{g(\phi)} = \sqrt{2(1 - \cos\phi)}$$

$$= 2|\sin\frac{\phi}{2}| \quad (\text{half-angle formula})$$

$$\leq 2|\sin\phi| \leq 2|r e^{i\phi} - 1|$$

□

Prop 4: Let  $\{z^k\}$  be generated by PG, w/  $\alpha_k = \alpha \geq 0$ .

$$\text{Define } \theta_k = \operatorname{argmin}_{\theta \in [0, 2\pi)} \|z^k - e^{i\theta} z^*\|_2$$

$$\varepsilon_k = e^{-i\theta_k} (z^k - e^{i\theta_k} z^*)$$

$$\beta_k = 1 + \alpha + \frac{\alpha}{n} (z^*)^H \varepsilon^k$$



Then for any  $r \in \mathbb{C}$ ,  $k \geq 0$ ,

$$d_2(z^{k+1}, z^*) \leq 2 \|r g^k - z^*\|_2$$

where,

$$g^k = \beta_k z^* + \left(I + \frac{\alpha}{n} \Delta\right) \varepsilon^k + \frac{\alpha}{n} \Delta z^*$$

Pf: By def'n,

$$w^k = \left(I + \frac{\alpha}{n} C\right) z^k$$

$$= \left(I + \frac{\alpha}{n} (z^* (z^*)^H + \Delta)\right) z^k$$

$$= e^{i\theta_k} \left(I + \frac{\alpha}{n} (z^* (z^*)^H + \Delta)\right) (z^* + \varepsilon^k)$$

乘上  $g^k$  EXERCISE

$$= \left[ \left(1 + \alpha + \frac{\alpha}{n} (z^*)^H \varepsilon^k\right) z^* + \left(I + \frac{\alpha}{n} \Delta\right) \varepsilon^k + \frac{\alpha}{n} \Delta z^* \right] e^{i\theta_k}$$

$$= g^k e^{i\theta_k}$$

$$\Rightarrow z^{k+1} = \frac{w^k}{|w^k|} = e^{i\theta_k} \frac{g^k}{|g^k|}, \text{ Hence, by Prop 3,}$$

$$d_2(z^{k+1}, z^*) \leq \left\| \frac{g^k}{|g^k|} - z^* \right\|_2 = \left\| \frac{r g^k}{|r g^k|} - z^* \right\|_2$$

$$\leq 2 \|r g^k - z^*\|_2$$

EXERCISE  $\uparrow$  加  $\downarrow$   $r$  对复数  
同模成立

2019.3.25

继续集:

## Phase Synchronization

The generative model:  $C = z^* (z^*)^H + \Delta$ 

$$\hat{z} \in \operatorname{argmax}_{z \in \mathbb{T}^n} z^H C z$$

$$\mathbb{T}^n = \{z \in \mathbb{C}^n; |z_i| = 1\}$$

Projected Grad.

证  $z^H C z \leq 4$ ,  $\frac{z^H C z}{2} \leq \frac{4}{2} = 2$   
 $\Rightarrow \mathbb{T}^n$

$$w^k \leftarrow z^k + \frac{\alpha^k}{n} C z^k$$

$$z^{k+1} \leftarrow \frac{w^k}{|w^k|}$$

$$\left[ \frac{w}{|w|} \right]_j = \begin{cases} \frac{w_j}{|w_j|} & \text{if } w_j \neq 0 \\ 1 & \text{if } w_j = 0 \end{cases}$$

Spectral Estimator

-  $u$ : leading eigenvector of  $C$ 

$$- v_c \stackrel{\Delta}{=} \frac{u}{|u|}$$

$$\text{Fact: } d_2(V_c, z^*) \leq \frac{8\|\Delta\|}{\sqrt{n}}$$

does not depend on the distribution of the noise.

Thm: Suppose  $\|\Delta\| \leq \frac{n}{16}$ ,  $\alpha_k = \alpha \geq 2$ , Then

$$d_2(z^{k+1}, z^*) \leq \mu^{k+1} d_2(z^0, z^*) + \frac{\nu}{1-\mu} \cdot \frac{8\|\Delta\|}{\sqrt{n}},$$

$$\mu = \frac{16(\alpha\|\Delta\| + n)}{(7\alpha + 8)n} < 1, \quad \nu = \frac{2\alpha}{7\alpha + 8}$$

$$\text{Recall: } d_2(z, w) = \min_{\theta \in [0, 2\pi)} \|z - e^{i\theta} w\|_2$$

prop 4: Let  $\{z^k\}$  be the iterates, Define

$$\theta_k = \arg \min_{\theta \in [0, 2\pi)} \|z^k - e^{i\theta} z^*\|_2$$

$$\varepsilon^k = e^{-i\theta_k} (z^k - e^{i\theta_k} z^*)$$

$$\beta_k = 1 + \alpha + \frac{\alpha}{n} (z^*)^H \varepsilon^k$$

Then, for any  $y \in \mathbb{C}$ ,  $k \geq 0$

$$d_2(z^{k+1}, z^*) \leq 2\|y g^k - z^*\|_2$$

$$\text{where, } g^k = \beta_k z^* + \left(I + \frac{\alpha}{n} \Delta\right) \varepsilon^k + \frac{\alpha}{n} \Delta z^*$$

Today:

Pf: (Thm 1)

We prove by induction, that

$$(i) \|\varepsilon^k\|_2 \leq \frac{\sqrt{n}}{2}$$

$$(ii) d_2(z^{k+1}, z^k) \leq \mu d_2(z^k, z^*) + \nu \cdot \frac{\delta \|\Delta\|}{\sqrt{n}}$$

unroll and calc the sum.

Base case:  $k=0$

$$\begin{aligned} \|\varepsilon^0\|_2 &= \|z^0 - e^{i\theta_0} z^*\|_2 && \text{by Assumptions } \|\Delta\| \leq n/16 \\ &= d_2(z^0, z^*) \leq \frac{\delta \|\Delta\|}{\sqrt{n}} \leq \frac{\sqrt{n}}{2} \\ & \quad \text{这里 } z^0 = v_c \end{aligned}$$

$$|\beta_0| \geq \left| 1 + \alpha + \frac{\alpha}{n} \operatorname{Re}(z^{*H} \varepsilon^0) \right| \quad \text{prop 4中, 去掉虚部}$$

$$= \left| 1 + \alpha + \frac{\alpha}{2n} (\|z^* + \varepsilon^0\|_2^2 - \|z^*\|_2^2 - \|\varepsilon^0\|_2^2) \right|$$

$$= \left| 1 + \alpha + \frac{\alpha}{2n} (\|z^0\|_2^2 - \|z^*\|_2^2 - \|\varepsilon^0\|_2^2) \right|$$

↑  
长度相同, 都是  $n$ .

$$\geq 1 + \frac{7\alpha}{8}$$

利用前面的  $\varepsilon_0$  的 bound  $\|\varepsilon^0\|_2 \leq \frac{\sqrt{n}}{2}$

Take  $y = \frac{1}{\beta_0}$  in prop 4:

$$\begin{aligned}
d_2(z^1, z^*) &\leq 2 \left\| \frac{1}{\beta_0} (I + \frac{\alpha}{n} \Delta) \varepsilon^0 + \frac{1}{\beta_0} \frac{\alpha}{n} \Delta z^* \right\|_2 \\
&\leq 2 \frac{1}{|\beta_0|} \left( \|(I + \frac{\alpha}{n} \Delta) \varepsilon^0\|_2 + \frac{\alpha}{n} \|\Delta z^*\|_2 \right) \\
&\leq \frac{16}{7\alpha + 8} \left( \left(1 + \frac{\alpha}{n} \|\Delta\| \right) \underbrace{\|\varepsilon^0\|_2}_{d_2(z^0, z^*)} + \frac{\alpha}{\sqrt{n}} \|\Delta\| \right) \\
&= \mu d_2(z^0, z^*) + \frac{\nu}{1-\mu} \frac{8 \|\Delta\|}{\sqrt{n}}
\end{aligned}$$

Inductive Step

$$\|\varepsilon^{k+1}\|_2 = d_2(z^{k+1}, z^*)$$

$$\leq \mu \cdot \underbrace{d_2(z^k, z^*)}_{\|\varepsilon^k\|_2 \leq \frac{\sqrt{n}}{2}} + \nu \frac{8 \|\Delta\|}{\sqrt{n}}$$

$$\leq \frac{8(\alpha \|\Delta\| + n)}{(7\alpha + 8)\sqrt{n}} + \frac{16\alpha}{7\alpha + 8} \cdot \frac{\|\Delta\|}{\sqrt{n}}$$

$$\leq \frac{\sqrt{n}}{2}$$

$$|\beta_{k+1}| \geq 1 + \frac{7\alpha}{8}$$

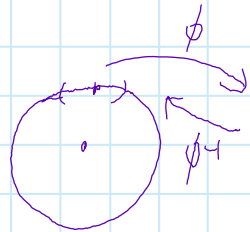
$$\Rightarrow d_2(z^{k+1}, z^*) \leq \mu d_2(z^{k+1}, z^*) + \nu \frac{8 \|\Delta\|}{\sqrt{n}}$$

by Prop 4.

# Optimization aspect

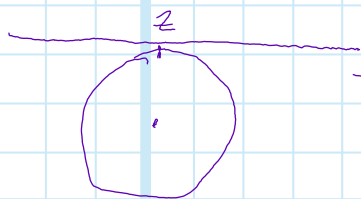
First-order optimality condition

obs:  $\mathbb{T}$  is a smooth manifold.



For every point, a neighbour can be mapped into a Euclidean-space.

It locally looks like a Euclidian space



$T_z T$  Tangent space.

$$T_z \mathbb{T} = \{w \in \mathbb{C} : \operatorname{Re}(z\bar{w}) = 0\}$$

$$\text{例: } \frac{1}{2} z = e^{i\theta}$$

$$T_z \mathbb{T} \ni w = r e^{i(\theta \pm \frac{\pi}{2})}$$

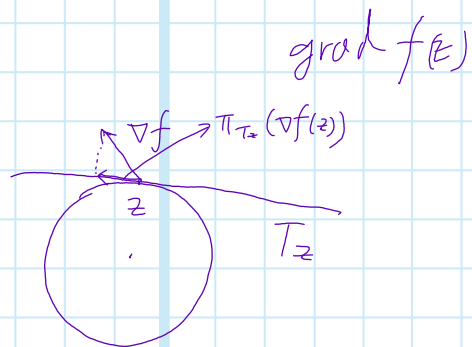
$$z\bar{w} = r e^{i(\pm \frac{\pi}{2})}$$

$$\mathbb{T}^n = \underbrace{\mathbb{T} \times \dots \times \mathbb{T}}_{n \text{ times}} \Rightarrow T_z \mathbb{T}^n = \{w \in \mathbb{C}^n : \operatorname{Re}(z_i \bar{w}_i) = 0, \forall i\}$$

This is a linear subspace.

Def:  $z$  is first order critical if

$$\Pi_{T_z}(\nabla f(z)) = 0.$$



对于无约束的 FOC:

$$\min_z f(z)$$

$$\text{FOC: } \nabla f(z) = 0,$$

现在变成投影  $\Pi_{T_z}(\nabla f(z)) = 0$ .

$\Pi_{T_z}(\nabla f(z))$  is called the Riemannian gradient of  $f$  at  $z$  and is denoted by  $\text{grad } f(z)$

$$\therefore \text{FOC: } \text{grad } f(z) = 0.$$

For  $\mathbb{T}^n$ : (EXERCISE)

$$\Pi_{T_z \mathbb{T}^n}(w) = w - \text{Diag}(\text{Re}(z_j \bar{w}_j)) z$$

$\therefore$  In original problem:

$$f(z) = -z^H c z$$

$$\Rightarrow \text{grad } f(z) = \Pi_{T_z \mathbb{T}^n}(\nabla f(z))$$

$$= 2(c z - \text{Diag}(\text{Re}(c z_j \bar{z}_j)) z)$$

Define  $S(z) = \text{Diag}(\text{Re}((z)_j \bar{z}_j)) - C$

$\Rightarrow$  1<sup>st</sup> order optimality condition is

$$S(z)z = 0.$$

上面定义了黎曼梯度, 下面来定义 Riemannian Hessian.  
对于普通的函数  $f$ , 若知道  $\nabla^2 f(z)w$ , 对所有  $w$ , 便能确定  $\nabla^2 f(z)$   
Riemannian Hessian: Projection of the directional derivative of the Riemannian grad onto the tangent space.

$$\text{Hess } f(z)(w) = \Pi_{T_z \mathbb{T}^n} (D \text{grad } f(z)[w])$$

$$= \Pi_{T_z \mathbb{T}^n} (2S(z)w)$$

这里是隐式定义的 Hessian, 而是 define its action on each direction

Def: 2<sup>nd</sup> order optimality condition

对比普通 Euclidean 的情况:

$$\min_z f(z)$$

$$\nabla f(z) = 0$$

$$\nabla^2 f(z) \succcurlyeq 0$$

manifold 情况下:



$$w^H (\text{Hess } f(z)) w$$

$$= z w^H S(z) w \geq 0 \quad \forall w \in T_z \mathbb{T}^n$$

$$\left[ \begin{aligned} & w^H \Pi_{T_z \mathbb{T}^n} (z S(z) w) \\ &= w^H P (z S(z) w) \quad \left\{ \begin{array}{l} \leftarrow \text{投影. 自伴算子.} \\ \because P \text{ is self-adjoint} \\ \therefore P = P^H \end{array} \right. \\ &= (Pw)^H (z S(z) w) \end{aligned} \right.$$

2019.3.26

$$\hat{z} \in \operatorname{argmax}_{z \in \mathbb{T}^n} \{f(z) \triangleq z^H C z\}$$

$C$  is generated by:  
 $C = z^* z^{*H} + \Delta$

$$\mathbb{T}^n = \{z \in \mathbb{C}^n : |z_i| = 1 \forall i\}$$

$$\operatorname{grad} f(z) = 2 [C - \operatorname{Diag}(\operatorname{Re}((Cz)_j; \bar{z}_j))] z$$

$$\operatorname{Hess} f(z)(w) = \mathbb{T}_{T_z \mathbb{T}^n} (2S(z)w),$$

$$\text{where } S(z) = \operatorname{Diag}(\operatorname{Re}((Cz)_j; \bar{z}_j)) - C$$

Optimality conditions:

$$1^{\text{st}} \text{ order : } S(z) z = 0$$

$$2^{\text{nd}} \text{ order : } w^H S(z) w \geq 0, \forall w \in T_z \mathbb{T}^n$$

where,

$$T_z \mathbb{T}^n = \{w \in \mathbb{C}^n, \operatorname{Re}(z_j \bar{w}_j) = 0, \forall j\}$$

Fact:  $\hat{z}$  satisfies the  $1^{\text{st}} + 2^{\text{nd}}$  order conditions

Prop: (i)  $z \in \mathbb{T}^n$  satisfies the  $1^{\text{st}}$  order condition iff

$$(Cz)_j \bar{z}_j \text{ is real } \forall j$$

(ii) If  $\text{diag}(c) \geq 0$  and  $z \in \mathbb{T}^n$  satisfies  
 1<sup>st</sup> + 2<sup>nd</sup> order condition, then

$$(Cz)_j \bar{z}_j \geq 0, \forall j$$

$$\Rightarrow (Cz)_j \bar{z}_j = |(Cz)_j|$$

Pf: (i)  $0 = S(z)z$

$$\Leftrightarrow \text{Re}((Cz)_j \bar{z}_j) z_j = (Cz)_j \quad \text{同乘 } \bar{z}_j, \text{ 用 } \|z_j\| = |z_j|$$

$$\Leftrightarrow \text{Re}((Cz)_j \bar{z}_j) = (Cz)_j \bar{z}_j$$

(ii) 观察到我们对 2<sup>nd</sup> order 定义是任取  $w$ ,  
 使得  $w^H S(z) w \geq 0$  我们可以取  $w = e_j$ ;

$$\begin{aligned} 0 &\leq e_j^H S(z) e_j \\ &= \text{Re}((Cz)_j \bar{z}_j) - C_{jj} \end{aligned}$$

$$= (Cz)_j \bar{z}_j - C_{jj}$$

还需验证  $e_j \in T_z \mathbb{T}^n$

$$\text{Verify } \text{Re}(z_j) \stackrel{?}{=} 0$$

This cannot be verified.

那么 How to fix this argument?

需满足如下条件:

$$(i) \text{Re}(z_j \bar{w}_j) = 0$$

$$(ii) w_j = \eta \mathbf{1} \Rightarrow (\eta e_j)^H S(z) (\eta e_j) \\ = (Cz)_j \bar{z}_j - C_{jj}$$

$$\Rightarrow |\eta| = 1$$

$$\eta = iz_j$$

Let  $w = (iz_j) e_j \in T_z \mathbb{T}^n \quad \forall j$

Then,

$$0 \leq w^H S(z) w = (Cz)_j \bar{z}_j - \underbrace{C_{jj}}_{\geq 0} \Rightarrow (Cz)_j \bar{z}_j \geq 0$$

Corollary:

consider  $\hat{z}$ , It satisfies both 1<sup>st</sup> + 2<sup>nd</sup> order condition. If  $\text{diag}(c) \geq 0$ , then

$$\boxed{\underbrace{(\text{Diag}(|Cz|) - C)}_{S(\hat{z})} \hat{z} = 0}$$

这个结论是在  $\text{diag}(c) \geq 0$  下成立的。A 是否 AA 可解 这个解

Obs:

$$(i) S(z) = \text{Diag}[\text{Re}((Cz)_j \bar{z}_j)] - C$$

$$C \rightarrow C + \frac{n}{2} I \triangleq \tilde{C}$$

$$(\tilde{C}z)_j = (Cz)_j + \frac{n}{2} z_j$$

$$\therefore \text{Diag}[\text{Re}((\tilde{C}z)_j \bar{z}_j)] = \text{Diag}[\text{Re}((Cz)_j \bar{z}_j + \frac{n}{2} z_j \bar{z}_j)]$$

$$= \text{Diag}[\text{Re}(Cz); \bar{z}_j] + \frac{n}{\alpha} I$$

$$\therefore \text{Diag}[\text{Re}(\tilde{C}z); \bar{z}_j] - \tilde{C}$$

$$= \text{Diag}[\text{Re}(Cz); \bar{z}_j; +\frac{n}{\alpha}] - \tilde{C}$$

$$= \text{Diag}[\text{Re}(Cz); \bar{z}_j] + \frac{n}{\alpha} I - \left( C + \frac{n}{\alpha} I \right)$$

(ii) For  $\hat{z}$ ,

$$\boxed{[\text{Diag}(|\tilde{C}\hat{z}|) - \tilde{C}] \hat{z} = 0}$$

## Convergence Analysis of PG

Alg: PG

$$w^k \leftarrow z^k + \frac{\alpha}{n} C z^k$$

$$z^{k+1} \leftarrow \frac{w^k}{|w^k|}$$

High-Level Ideas:

- ① Error Bound. (Depend only on the problem, not related to Alg)

之前14-15题。

$$\text{dist}(x, X) \leq \mu \cdot \|R(x)\|_2$$

$$\|R(x)\| \leq \text{dist}(0, \partial F(x))$$

② Alg properties.

(i) sufficient ascent

(ii) cost-to-go      how far away from the optimal value

(iii) Safeguard.      if residual is 0 we should stop.

在这个问题中, Error Bound 和  $\hat{z}$  不同.

下面来看每一句.

① Error Bound. (e.g.  $\text{dist}(x, X) \leq \mu \|R(x)\|_2$ )  
 如何确定 LHS 和 RHS      what  $\mu$  to put here.

LHS:

$$d_2(z, \hat{z})$$

RHS:

在 Corollary (ii) 中, 如果  $\hat{z}$  是最优, 那么等号, 如果等号一个很小的数, 那么我们就是否能够接近?

Candidate:

$$\Sigma(z) = \text{Diag}(\| \tilde{C} z \|) - \tilde{C}$$

$$\text{Define } \rho(z) = \| \Sigma(z) z \|_2$$

$$\text{Note: } \rho(\hat{z}) = 0$$

注意, 我们已知  $\hat{z}$  is optimal  $\Rightarrow P(\hat{z}) = 0$   
但反过来不知道.

$$\boxed{d_2(z, \hat{z}) \leq c \cdot P(z)} \leftarrow \text{我们尝试证这个.}$$

② Alg properties.

(a) (Sufficient Ascent)

$$f(z^{k+1}) - f(z^k) \geq a_0 \|z^{k+1} - z^k\|_2^2$$

(b) (Cost-to-go)

How far you are

$$f(\hat{z}) - f(z^k) \leq a_1 d_2(z^k, \hat{z})^2$$

(c) (Safeguard)

$$P(z^k) \leq a_2 \|z^{k+1} - z^k\|_2$$

Using these properties,

$$\underbrace{f(\hat{z}) - f(z^{k+1})}_{\text{(cost-to-go)}} = [f(\hat{z}) - f(z^k)] - [f(z^{k+1}) - f(z^k)]$$

$$\leq a_1 d_2(z^k, \hat{z})^2 - [f(z^{k+1}) - f(z^k)]$$

$$\text{(EB)} \leq a' \cdot P(z^k)^2 - [f(z^{k+1}) - f(z^k)]$$

$$\text{(safeguard)} \leq a' \cdot \|z^{k+1} - z^k\|_2^2 - [f(z^{k+1}) - f(z^k)]$$

$$\leq \underbrace{(a'' - 1)}_{> 0} [f(z^{k+1}) - f(\hat{z}) + f(\hat{z}) - f(z^k)]$$

Rearrange

$$a''' [f(\hat{z}) - f(z^{k+1})] \leq (a''' - 1) [f(\hat{z}) - f(z^k)]$$

$$[f(\hat{z}) - f(z^{k+1})] \leq \frac{a''' - 1}{a'''} [f(\hat{z}) - f(z^k)]$$

$\Rightarrow$  linear convergence.

$\exists \lambda \in (0, 1)$ :

$$f(\hat{z}) - f(z^{k+1}) \leq \lambda^k (f(\hat{z}) - f(z^0))$$

How about the iterates,

$$d_2(z^k, \hat{z})^2 \leq c \cdot \rho (z^k)^2 \quad (\text{EB})$$

$$\leq c' \|z^{k+1} - z^k\|_2^2 \quad (\text{Safeguard})$$

$$\leq c'' (f(z^{k+1}) - f(z^k)) \quad (\text{Sufficient ascent})$$

$$\leq c'' (f(\hat{z}) - f(z^k))$$

$$\leq \lambda^{k-1} [f(\hat{z}) - f(z^k)] - c'''$$

注意这是  $d_2(z^k, \hat{z})^2$   $\therefore d_2(z^k, \hat{z})$  的收敛速度

大约是  $\lambda^{\frac{k-1}{2}}$ , 比  $f$  的速度慢.

这个分析和之前类似, 唯一区别是

这里是非凸的.



2019.4.2

Generative model:

$$C = z^* (z^*)^H + \Delta$$

Estimator:

$$\hat{z} \in \operatorname{argmax}_{z \in \mathbb{T}^n} z^H C z$$

PG

$$w^k \leftarrow \left( I + \frac{\alpha}{n} C \right) z^k$$

$$z^{k+1} \leftarrow \frac{w^k}{|w^k|}$$

Def:  $\tilde{C} = C + \frac{n}{\alpha} I$

$$\Rightarrow w^k = \frac{\alpha}{n} \tilde{C} z^k$$

$$\Rightarrow z^{k+1} = \frac{\tilde{C} z^k}{|\tilde{C} z^k|}$$

也叫 (GPM): Generalized Power Method

Facts:

$$(1) d_2(\hat{z}, z^*) \leq \frac{4\|\Delta\|}{\sqrt{n}}$$

$$(2) d_2(V_c, z^*) \leq \frac{8 \|\Delta\|}{\sqrt{n}}$$

(3) define

$$\Sigma(z) = \text{Diag}(|\tilde{C}(z)|) - \tilde{C}$$

Then

$$\Sigma(\hat{z}) \hat{z} = 0 \quad \left( \begin{array}{l} 1^{\text{st}} + 2^{\text{nd}} \\ \text{opt. cond.} \end{array} \right)$$

Define residual measure.

$$P(z) = \|\Sigma(z)z\|_2$$

### Error Bound

Thm. For any  $z \in \mathbb{T}^n$ , satisfying  $d_2(z, z^*) \leq \frac{\sqrt{n}}{2}$   
+ additional assumptions (on  $\|\Delta\|, \alpha, \dots$ )  
and any opt soln  $\hat{z}$

$$d_2(z, \hat{z}) \leq \frac{8}{n} P(z)$$

这是一个 local error bound, 但是如果  
起始在这个 region 内, 整个序列都在.

Pf:

$$P(z) = \|\Sigma(z)z\|_2 \geq \|\Sigma(\hat{z})z\|_2 - \|(\Sigma(z) - \Sigma(\hat{z}))z\|_2$$

need LB

need UB

通常 UB 更好算

UB:

看如何把  $d_2$  关联起来。

$$\begin{aligned} & \|(\Sigma(z) - \Sigma(\hat{z}))z\|_2 \\ &= \left\| \left[ \text{Diag}(|\tilde{c}(z)|) - \text{Diag}(|\tilde{c}(\hat{z})|) \right] z \right\|_2 \end{aligned}$$

$$= \left( \sum_j \left| (|\tilde{c}(z)|_j - |\tilde{c}(\hat{z})_j|) z_j \right|^2 \right)^{1/2}$$

$z_j$  模长始终为 1, 可以扔掉

加这项, 模长不变

$$= \left\| |\tilde{c} e^{-i\hat{\theta}} z| - |\tilde{c} \hat{z}| \right\|_2$$

$$\hat{\theta} = \arg \min_{\theta \in [0, 2\pi]} \|z - e^{i\theta} \hat{z}\|_2$$

$$\leq \left\| \tilde{c} (e^{-i\hat{\theta}} z - \hat{z}) \right\|_2 \quad (\text{tri-inequality})$$

$$\leq \left\| z^* (z^*)^H (e^{-i\hat{\theta}} z - \hat{z}) \right\|_2 + \left\| \Delta (e^{-i\hat{\theta}} z - \hat{z}) \right\|_2$$

$$+ \frac{n}{\alpha} \left\| e^{-i\hat{\theta}} z - \hat{z} \right\|_2 \quad (\because \tilde{c} = c + \frac{n}{\alpha} I)$$

$$\leq \underbrace{\sqrt{n} \cdot |(z^*)^H (e^{-i\hat{\theta}} z - \hat{z})|}_{\downarrow} + (\|\Delta\| + \frac{n}{\alpha}) d_2(z, \hat{z})$$

$$|(z^*)^H (e^{-i\hat{\theta}} z - \hat{z})| \leq \left| (z^* - e^{-i\hat{\theta}^*} \hat{z})^H (e^{-i\hat{\theta}} z - \hat{z}) \right| +$$

$$\left| (e^{-i\hat{\theta}^*} \hat{z})^H (e^{-i\hat{\theta}} z - \hat{z}) \right|$$

$$\left( \hat{\theta}^* = \arg \min_{\theta} \|\hat{z} - e^{i\theta} z^*\|_2 \right)$$

$$\leq \frac{4\|\Delta\|}{\sqrt{n}} d_2(z, \hat{z}) + \left| (e^{i\theta^*} \hat{z})^H (e^{-i\theta^1} z - \hat{z}) \right|$$

Note:  $d_2(z, \hat{z})^2 = \|e^{i\theta^1} z - \hat{z}\|_2^2$

$$= 2 \left( n - |\hat{z}^H z| \right)$$

$-\max_{\theta} 2 \operatorname{Re} \left( (e^{i\theta} \hat{z})^H z \right)$   
 当  $e^{i\theta} \hat{z}$  is real, get max

$$\leq \frac{4\|\Delta\|}{\sqrt{n}} d_2(z, \hat{z}) + \left| e^{i\theta^*} (|\hat{z}^H z| - n) \right|$$

取模长可以忘掉

$$\leq \frac{4\|\Delta\|}{\sqrt{n}} d_2(z, \hat{z}) + \frac{1}{2} d_2(z, \hat{z})^2$$

$\therefore$  UB:

$$\|(\Sigma(z) - \Sigma(\hat{z}))z\|_2$$

$$\leq \left( 5\|\Delta\| + \frac{n}{2} \right) d_2(z, \hat{z}) + \frac{\sqrt{n}}{2} d_2(z, \hat{z})^2$$

下面来看看  $\|\Sigma(z)z\|_2$  的 LB,

如果  $z$  在  $\Sigma(z)$  的 null space 中, 那么  $\Sigma(z)z$  是 0, 但观察到  $\Sigma(z)$  由  $z$  决定,  $z$  是否能让其不在 Null space 中.

Want:  $\Sigma(\hat{z}) \geq 0$ ,  $\Sigma(\hat{z})$  pd on certain subspace

$\therefore \Sigma(\hat{z})\hat{z} = 0$  (opt condition)  $\therefore$  suppose  $\hat{z}$  is  $(\operatorname{span}(\hat{z}^\perp))^\perp$

Define:  $\hat{u} = \left( I - \frac{1}{n} \hat{z} \hat{z}^H \right) (e^{i\theta^1} z - \hat{z})$

be projection of  $e^{-i\theta} \hat{z} - \hat{z}$  onto  $\text{span}(\hat{z})^\perp$

Note: ①  $\hat{u}^H \hat{z} = 0 \therefore \hat{u}$  is the projector

$$\begin{aligned} \text{② } \|\Sigma(\hat{z})z\|_2 &= \|\Sigma(\hat{z})(e^{-i\theta} \hat{z} - \hat{z})\|_2 \\ &= \|\Sigma(\hat{z})\hat{u}\|_2 \leftarrow \text{看这个是在 LB by sth positive} \end{aligned}$$

注意到  $\lambda_{\min}(\Sigma(\hat{z})) = \min_{\hat{u}} \frac{\hat{u}^H \Sigma(\hat{z}) \hat{u}}{\|\hat{u}\|_2^2}$

那么  $\hat{u}^H \Sigma(\hat{z}) \hat{u} \geq c \|\hat{u}\|_2$  ?

$$\hat{u}^H \Sigma(\hat{z}) \hat{u} = \hat{u}^H [\text{Diag}(|\tilde{c}(\hat{z})|) - \tilde{c}] \hat{u}$$

$$= \hat{u}^H [\text{Diag}(|c(\hat{z})|) - c] \hat{u}$$

$$|(\tilde{c}(\hat{z}))_j| = |(\tilde{c}(\hat{z}))_j \bar{z}_j|$$

(prop 5 of prev notes)

$$= \sum_j |c(\hat{z})_j| \cdot |\hat{u}_j|^2 - |(z^*)^H \hat{u}|^2 = \hat{u}^H \Delta \hat{u}$$

$$\geq (|(z^*)^H \hat{z}| - \|\Delta \hat{z}\|_\infty) \|\hat{u}\|_2^2 \quad c(\hat{z}) = z^* (z^*)^H \hat{z} + \Delta \hat{z}$$

$$= |\hat{u}^H (z^* - e^{-i\theta} \hat{z})|^2 - \|\Delta\| \cdot \|\hat{u}\|_2^2$$

$\leftarrow \because \hat{u}^H \hat{z} = 0 \therefore$  可以随便找系数

$$\geq \left[ n - \|\Delta \hat{z}\|_\infty - \frac{3}{2} d_2(\hat{z}, z^*)^2 - \|\Delta\| \right] \|\hat{u}\|_2^2$$

$$\geq \left( n - \|\Delta \hat{z}\|_\infty - \|\Delta\| - \frac{24\|\Delta\|^2}{n} \right) \|\hat{u}\|_2^2$$

下面利用这个来分析 Bound.

$$\begin{aligned}\|\hat{u}\|_2 &\geq \|e^{-i\hat{\theta}} z - \hat{z}\|_2 - \frac{1}{n} \|\hat{z} \hat{z}^H (e^{-i\hat{\theta}} z - \hat{z})\|_2 \\ &= d_2(z, \hat{z}) - \frac{1}{\sqrt{n}} \left| \hat{z}^H (e^{-i\hat{\theta}} z - \hat{z}) \right| \\ &= d_2(z, \hat{z}) - \frac{1}{2\sqrt{n}} d_2(z, \hat{z})^2\end{aligned}$$

Hence,

$$\begin{aligned}\|\Sigma(\hat{z})z\| &= \|\Sigma(\hat{z})\hat{u}\|_2 \geq \lambda_{\min}(\Sigma(\hat{z})) \|\hat{u}\|_2 \\ &\geq \left[ n - \|\Delta \hat{z}\|_{\infty} - \|\Delta\| - \frac{24\|\Delta\|^2}{n} \right] \cdot \left( d_2(z, \hat{z}) - \frac{1}{2\sqrt{n}} d_2(z, \hat{z})^2 \right)\end{aligned}$$

Note:  $d_2(z, \hat{z}) \leq d_2(z, z^*) + d_2(\hat{z}, z^*)$

$$\leq \frac{\sqrt{n}}{2} + \frac{4\|\Delta\|}{\sqrt{n}}$$

$$\Rightarrow d_2(z, \hat{z})^2 \leq \left( \frac{\sqrt{n}}{2} + \frac{4\|\Delta\|}{\sqrt{n}} \right) \cdot d_2(z, \hat{z})$$

Note:  $\|\Delta \hat{z}\|_{\infty} \leq \|\Delta z^*\|_{\infty} + \|\Delta (e^{-i\theta^*} \hat{z} - z^*)\|_{\infty}$

$$\leq \|\Delta z^*\|_{\infty} + \|\Delta\| \cdot d_2(\hat{z}, z^*)$$

$$\leq \|\Delta z^*\|_{\infty} + \frac{4\|\Delta\|^2}{\sqrt{n}}$$

$\therefore$  Assumptions:

$$\|\Delta\| = \mathcal{O}(n^{-3/4})$$

加上 Assumptions 来保证结论.

$$\|\Delta z^x\|_\infty = O(n)$$

2019.4.8.

今天换一个题目.

信号处理问题.

## MIMO Detection

$$y = Hx^* + v \quad (\text{Generative model})$$

$$H \in \mathbb{C}^{m \times n}$$

$m$ : # of output

$n$ : # of input

Channel matrix

$x^* \in \mathbb{C}^n$ : vector of transmitted symbols

$y \in \mathbb{C}^m$ : received vector

$v \in \mathbb{C}^m$ : additive noise

Typically,

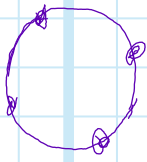
each,  $x_i^*$  draw from a discrete  
discrete constellation  $\mathcal{S}$

Two examples:

↑ 就是集合的意思 ↘



(1) M-ary Phase shift keying (MPSK)

$M=4$    $\mathcal{S}_M = \left\{ \exp\left(\frac{2\pi i k}{M}\right), k=0, 1, \dots, M-1 \right\}$

(2)  $(4M^2)$ -Quadrature Amplitude Modulation (QAM)

$$\mathcal{Q}_M = \left\{ z \in \mathbb{C} : \operatorname{Re}(z), \operatorname{Im}(z) = \pm 1, \pm 3, \dots, \pm(z-1) \right\}$$

$V_i$ : Complex Gaussian  $\mathcal{CN}(0, \sigma_i^2)$  circular symmetry

Problem:

Given  $y$  and  $H$ , goal is recover  $x^*$

ML:  $\hat{x} \in \arg \min_{x \in \mathcal{S}^n} \|y - Hx\|_2^2$  与信道回归相似。  
x 有了离散约束

和 Phase Sync 相同, 有两种方法:

Approach 1: SDR

在 (1) 中, 直接用  $\|x\|=1$  然后 rounding

(2) 模长不一样怎么做?

Approach 2: PG projected gradient

这里我们不但不是投影到  $\mathcal{S}$  集合上,  
而是投影到离散集合上.

PG:

$$\tilde{w}^k \leftarrow 2H^H (HX^k - y)$$

$$x^{k+1} \leftarrow \Pi_{\mathcal{S}^n} \left( X^k - \frac{\alpha_k}{m} \tilde{w}^k \right)$$

我们为什么能这么做:

Rmk: ① Projection. step is easy

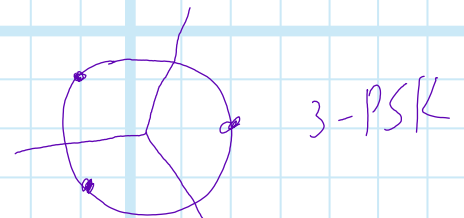
- projection can be done coordinate-wise

比如, 只有四个点, 那就计算每个距离, 然后取最小的那个.

Then, ① brute-force computation

② Voronoi diagram (算最近邻点)

例:



② 在连续的问题中, convergence 很重要, 可以无限.  
但在离散问题中, 不存在这种情况. (不能无限)

want: Finite convergence  
preferably, less than  $|\mathcal{S}^n|$  iterations.  
至少比暴力搜索简单些.

Thm: Let  $c = \frac{4}{\min_{s \neq s' \in \mathcal{S}} |s - s'|} < \infty$

任两个点不能无限接近

Suppose:

(i)  $\frac{2\alpha_k}{m} \|H^H v\|_\infty < \frac{1}{c}$

control noise

(ii)  $\|I - \frac{2\alpha_k}{m} H^H H\| \leq \beta < \frac{1}{4}$

conditioning of the channel

Then,

$$\|x^{k+1} - x^*\|_2 \leq 4\beta \|x^k - x^*\|_2 \leq \dots \leq (4\beta)^{k+1} \|x^0 - x^*\|_2 \leq \frac{c}{2}$$

这里没有限制起始点, 只要控制了

noise 和 conditioning, 就有 linear convergence

注意, 这也意味着迭代步数是有限的,

imply finite convergence. 离散空间是有距离的, 不能两个点距离太小。

In particular, after at most

$$k^* = \left\lceil \log_{4\beta} \left( \frac{2}{c \|x^0 - x^*\|_2} \right) \right\rceil \text{ iterations,}$$

$$x^k = x^* \quad \forall k \geq k^*$$

Proof: By defin of PG

$$z^k \triangleq x^k - \frac{\alpha_k}{m} w^k \quad (\text{what we are going to project})$$

$$= x^k - \frac{2\alpha_k}{m} H^H (H x^k - y)$$

$$= x^* + \underbrace{\left( I - \frac{2\alpha_k}{m} H^H H \right)}_{\text{Assumption 1 bound}} (x^k - x^*) + \underbrace{\frac{2\alpha_k}{m} H^H v}_{\text{Assumption 2 bound}}$$

$$\begin{cases} H^H y = H^H (H x^* + v) \end{cases}$$

and  $x^{k+1} = \Pi_{\mathcal{S}}(z^k)$

Let  $w^k = \left( I - \frac{2\alpha_k}{m} H^H H \right) (x^k - x^*)$

$$J_k = \left\{ j : |w_j^k| \geq \frac{1}{c} \right\} \quad \leftarrow \text{key tricks}$$

这里  $J_k$  让  $w_j^k$  足够大; 因此  $\{z^k\}$  就变成了  $x^*$  加两个小项. - 投影后就是  $x^*$  了.

Note:

$$z_l^k = x_l^* + w_l^k + \frac{2\alpha_k}{m} (H^H v)_l \quad \text{Coordinate-wise}$$

$$\Rightarrow |z_l^k - x_l^*| \leq |w_l^k| + \frac{2\alpha_k}{m} |(H^H v)_l|$$

$$\left( \leq \frac{1}{c} + \frac{1}{c} = \frac{2}{c} \quad \text{if } l \notin J_k \right)$$

$$\Rightarrow x_i^* = \prod_{\mathcal{S}}(z_i^k) = x_i \dots$$

Fact:

$$x_i^{k+1} = x_i^* \quad \forall i \notin J_k \quad \text{只要不在 } J_k \text{ 中, } x_i \text{ 不变这一项.}$$

prop: let  $z \in \mathbb{C}^n$  and  $x \in \mathcal{S}^n$  be given,  
where  $\mathcal{S} = \mathcal{S}_m$  or  $\mathbb{C}_u$ .

$$\text{Then, } \|\Pi_{\mathcal{S}^n}(z) - x\|_2 \leq 2\|z - x\|_2$$

先假设成 3 -

Assume this, ... 只要不在  $J_k$  中的值, 取 0 相等

$$\|x^{k+1} - x^*\|_2 = \|x_{J_k}^{k+1} - x_{J_k}^*\|_2 \leq 2\|z_{J_k}^* - x_{J_k}^*\|_2$$

由那个 prop

$$= 2\left\|w_{J_k}^k + \left(\frac{2\alpha k}{m} H^H v\right)_{J_k}\right\|_2.$$

$$\left[ \left\| \frac{2\alpha k}{m} H^H v \right\|_{\infty} < \frac{1}{c} \leq |w_j^k| \quad \text{for } j \in J_k \right]$$

Assumption (i)

$$\leq 4\|w_{J_k}^k\|_2 \leq 4\|w^k\|_2$$

$$\leq 4\beta\|x^k - x^*\|_2 \quad (\text{Assumption 2})$$

以此可以证明明了 Then, 下面求着 prop.

Pf of Prop.

Focus on  $\mathcal{D}_m$

$$\textcircled{1} \min_{r \geq 0} |re^{i\phi} - 1|^2 = \begin{cases} 1 & \text{if } \phi \in [\frac{\pi}{2}, \frac{3\pi}{2}] \\ \sin^2 \phi & \text{if } \phi \in [0, \frac{\pi}{2}) \cup (\frac{3\pi}{2}, 2\pi) \end{cases}$$

$$\textcircled{2} \text{ Let } z = re^{i\phi} \text{ if } \phi \in [\frac{\pi}{2}, \frac{3\pi}{2}].$$

$$|\pi_{\mathcal{D}_m}(z) - 1| \leq 2 \leq 2|re^{i\phi} - 1|$$

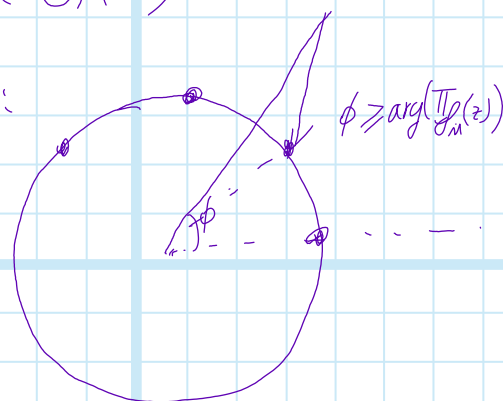
$$\textcircled{3} \text{ Say } \phi \in [0, \frac{\pi}{2}] \text{ (a) If } \phi \geq \arg(\pi_{\mathcal{D}_m}(z))$$

$$|\pi_{\mathcal{D}_m}(z) - 1| \leq |e^{i\phi} - 1| \quad \text{by (1):}$$

$$\leq 2|\sin \phi|$$

↳ Half-angle formula.

$$\leq 2|re^{i\phi} - 1|$$



$$\text{(b) } \phi < \arg(\pi_{\mathcal{D}_m}(z))$$

和 (a) 一样, 我们需要  $\phi$  的 LB.  
obs  $\phi \geq \frac{1}{2} \arg(\pi_{\mathcal{D}_m}(z))$

$$|\pi_{\mathcal{D}_m}(z) - 1| = 2 \left| \sin \frac{\arg(\pi_{\mathcal{D}_m}(z))}{2} \right|$$

$$\stackrel{\parallel}{=} |e^{i \arg(\pi_{\mathcal{D}_m}(z))}| \leq 2|\sin \phi| \leq \dots$$

下面来验证那四个 Assumption 能不能同时满足。

Thm 2: Suppose

(i) entries of  $H$  iid standard complex

Gaussian (i.e.  $H_{ij} \sim \mathcal{CN}(0, 1)$ )

$$g_{ij}^R + i g_{ij}^I \quad \leftarrow \text{complex normal}$$

$$g_{ij}^R, g_{ij}^I \sim \mathcal{N}(0, 1/2)$$

(ii) entries of  $V$  have variance

$$\sigma_V^2 \leq \frac{m}{4c^2 \log n}$$

(iii) aspect ratio:  $\gamma \triangleq \frac{m}{n} \geq \frac{20}{\beta^2} > 1$

(因此输出比输入大  $\frac{20}{\beta^2}$ )

(iv) step size  $\alpha_k = \frac{1}{2}$

Then, whp, (i)+(ii) in Thm hold  
↑ with high probability

Comments:

$\gamma \leq 1$  是否可能存在, open problem.

下节将主要看 (ii).

2019.4.9

上节课的 Assumption 2

$$\|I - \frac{1}{m} H^H H\| \leq \beta < \frac{1}{4}$$

spectral norm

eigen values:

$$\max_i \left| 1 - \frac{1}{m} \lambda_i(H^H H) \right|$$

可以取:

$$\sigma_{\max}(H) \left| 1 - \frac{1}{m} \lambda_{\max}(H^H H) \right|$$

$$\sigma_{\min}(H) \left| 1 - \frac{1}{m} \lambda_{\min}(H^H H) \right|$$

下面来研究随机矩阵的奇异值。

Largest singular value of Gaussian

Random Matrix

$$A \in \mathbb{R}^{m \times n} \quad m \geq n$$

$$A_{ij} \sim \mathcal{N}(0, 1)$$

$$\sigma_{\max}(A) = \sup_{\substack{\|u\|_2=1 \\ \|v\|_2=1}} u^T A v \quad (\text{Courant-Fischer})$$

Idea: Fix  $u, v$ , consider

$$u^T A v = \sum_{i,j} u_i v_j A_{ij}$$



$$= \sum_i u_i \left( \sum_j v_j A_{ij} \right) \rightarrow \sim \mathcal{N}(0, \sum_j v_j^2)$$

$$\sim \sum_i u_i g_i \quad g_i \sim \mathcal{N}(0, 1) \quad = \mathcal{N}(0, 1)$$

$$\sim \mathcal{N}(0, \sum_i u_i^2) = \mathcal{N}(0, 1)$$

$\therefore u^T A v \sim \mathcal{N}(0, 1) \quad \therefore$  可以 bound

Fact:  $\Pr(|u^T A v| \geq t) \leq \frac{1}{\sqrt{2\pi} t} \exp(-t^2/2)$

这里不能用 union bound,  $\because u, v$  是无限个,  
不能取完  $u, v$  来取最小。

obs: We cannot take union bound over  
 $S^{m-1} \times S^{n-1}$  ( $S^{m-1} = \{x \in \mathbb{R}^m; \|x\|_2 = 1\}$ )

Idea: Discretize  $S^{m-1} \times S^{n-1}$

(Find representative points and take union bound over this finite subset)

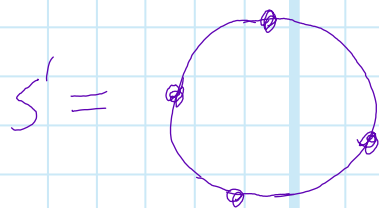
We should control the error of the discretization

$\rightarrow \epsilon$ -net

Def A set  $N \subseteq S^{n-1}$  is an  $\epsilon$ -net if  $\forall u \in S^{n-1}$   
 $\exists v \in N, \text{ s.t. } \|u - v\|_2 \leq \epsilon$

就是集合中任一点到  $\varepsilon$ -net 的距离都小于  $\varepsilon$

The main question is how large  $|N|$  will be.



Prop: Let  $\varepsilon > 0$  be given, There exists an  $\varepsilon$ -Net of size  $|N| \leq 2n(1 + 2/\varepsilon)^{n-1}$

看这个怎么用在前面的那个命题上.

Thm: Let  $M$  be  $\delta$ -net of  $S^{m-1}$  and  $N$  be an  $\varepsilon$ -net of  $S^{n-1}$ , Then,

$$\sigma_{\max}(A) \leq \frac{1}{(1-\varepsilon)(1-\delta)} \sup_{u \in M, v \in N} |u^T A v|$$

Here is similar to Lovasz-Fischer, Difference is  $M, N$  is finite.

Pf Take  $z \in S^{n-1}$ , We can write  $z = v + h$ , where  $v \in N$ ,  $\|h\|_2 \leq \varepsilon$

Note:  $\sigma_{\max}(A) = \sup_{\substack{u \in S^{m-1} \\ v \in S^{n-1}}} u^T A v = \sup_{v \in S^{n-1}} \|A v\|_2$   $\leftarrow \sigma_{\max}(A)$

$$\therefore \|A z\|_2 \leq \|A v\|_2 + \|A h\|_2 \leq \|A v\|_2 + \varepsilon \|A\|$$

$\therefore$  两边同取 sup

$$\sigma_{\max}(A) \triangleq \sup_{z \in S^{n-1}} \|A z\|_2 \leq \sup_{v \in N} \|A v\|_2 + \varepsilon \sigma_{\max}(A)$$

$$\sigma_{\max}(A) \leq \frac{1}{1-\epsilon} \sup_{v \in N} \|Av\|_2$$

Next  $\|Av\|_2 = \sup_{w \in S^{n-1}} |w^T Av|$   $w = u + q, u \in M, \|q\|_2 \leq \delta$

$$\leq \sup_{u \in M} |u^T Av| + \delta \|Av\|_2$$

□

Consequently, since for each  $u, v$ :

$$\Pr(|u^T Av| \geq t) \leq \frac{1}{\sqrt{2\pi}t} \exp(-t^2/2)$$

So, use union bound

$$\Pr \left[ |u^T Av| \geq t \text{ for some } u \in M, v \in N \right] \leq |M| \cdot |N| \cdot \frac{1}{\sqrt{2\pi}t} \exp(-t^2/2)$$

$$O(m^2 \cdot C^{2m}) \leftarrow \text{由前面那个 } |N| \leq 2n \left(1 + \frac{1}{\epsilon}\right)^{n-1}$$

choose  $t = O(\sqrt{m})$ , then exp small prob.

以上就得到了奇异值的 UB

下面补充 Prop 的证明.

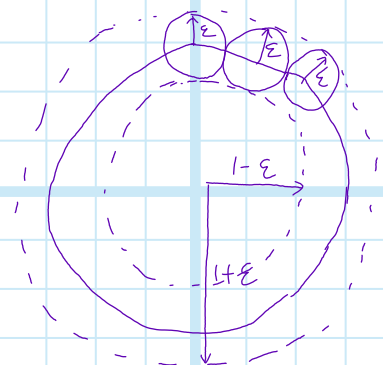
Pf of Prop

↙ a useful trick

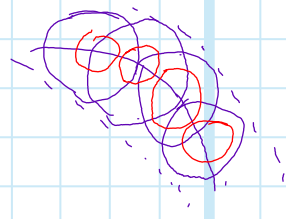
By volume argument

数有多少球, 直接用那个环的体积除以每个小球的体积.

Let  $N$  be a maximal cardinality subset of  $S^{n-1}$ , s.t.



$$u, v \in N \Rightarrow \|u - v\|_2 > \varepsilon$$



$\Rightarrow$  by construction,  $N$  is  $\varepsilon$ -net

obs:  $u, v \in N$ :

$$B(u, \frac{\varepsilon}{2}) \cap B(v, \frac{\varepsilon}{2}) = \emptyset$$

$$\bigcup_{u \in N} B(u, \frac{\varepsilon}{2}) \subseteq B(0, 1 + \frac{\varepsilon}{2}) \setminus B(0, 1 - \frac{\varepsilon}{2})$$

$$|N| \cdot \text{Vol}(B(u, \frac{\varepsilon}{2}))$$

$$\leq \text{Vol}(B(0, 1 + \frac{\varepsilon}{2})) - \text{Vol}(B(0, 1 - \frac{\varepsilon}{2}))$$

Fact:

$$\because \text{Vol}(B(0, r)) = r^n \cdot \text{Vol}(B(0, 1))$$

$\Rightarrow$  divide by  $\text{Vol}(B(0, 1))$

$$|N| \cdot \left(\frac{\varepsilon}{2}\right)^n \leq \left(1 + \frac{\varepsilon}{2}\right)^n - \left(1 - \frac{\varepsilon}{2}\right)^n$$

$$\left[ \text{Ineq: } (1+x)^l - (1-x)^l \leq 2lx(1+x)^{l-1} \right.$$

Q.E.D

注意到  $H$  是复矩阵,

$$u^H Q a$$

$$\equiv \begin{pmatrix} u^R & u^I \end{pmatrix} \begin{bmatrix} \text{Re}(Q) & \text{Im}(Q) \\ \text{Im}(Q) & \text{Re}(Q) \end{bmatrix} \begin{bmatrix} a^R \\ a^I \end{bmatrix}$$

看另一种方法, 利用 Lipschitz 性质.

$\sigma_{\max}(A)$ : 1-Lipschitz

$$\left( \left| \sigma_{\max}(A) - \sigma_{\max}(B) \right| \leq \|A - B\|_2 \right)$$

$$\Pr[\sigma_{\max}(A) \geq \mathbb{E}[\sigma_{\max}(A)] + t] \leq O(e^{-t^2})$$

(concentration, ineq for Lipschitz fns of Gaussian RVs)

$$\mathbb{E}[\sigma_{\max}(A)] = \mathbb{E} \left[ \sup_{\substack{\|u\|_2=1 \\ \|v\|_2=1}} u^T A v \right]$$

Slepian lemma

- 一个由  $(u, v)$  index 的  
高斯过程可以由另一个  
高斯过程来 bound.

2019.4.15

新问题 LR:

## Community Detection

- Goal: identify groups of densely connected nodes in a network.

- Stochastic block model (SBM)

Data:

undirected network of  $n$  nodes,

- specified by  $n \times n$  adjacent matrix  $A$

$$A_{ij} = \begin{cases} 0 & \text{if } (i, j) \text{ is not an edge} \\ 1 & \text{o/w} \end{cases}$$

-  $K$  communities in network, each node belong to one community

→ membership vector  $z_i$  for node  $i$

$z_{ik} = 1$  if node  $i$  belongs to the community  $k$

→  $z_i \in \{0, 1\}^K$ , not observed

- SBM is parametrized by  $\Psi$  ( $K \times K$  symmetric)

$\Psi_{kr}$  = prob of an edge forming between a pair of nodes from community  $k$  and  $r$

For simplicity, assume

$$n = mK, \quad m \geq 1 \text{ Integer.}$$

这里并没说每个 community 的 member 相等.

- network generation

Given  $\{z_i\}$ ,  $\Psi$ .  $A_{ij}$  independent Bernoulli RVs  $i < j$  (对称)

$$\text{s.t. } \mathbb{E}[A_{ij} | z_i, z_j] = z_i^T \Psi z_j$$

属于某类的概率乘以某类的概率.

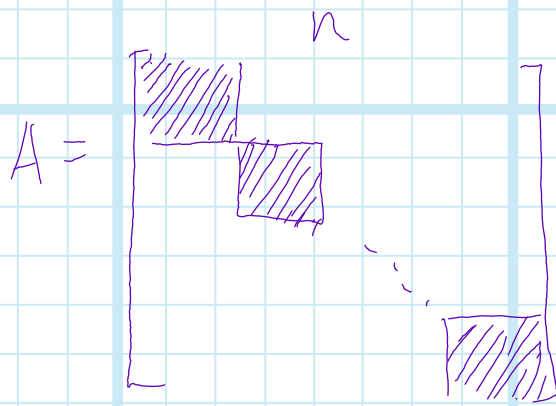
$$\left( = \Psi_{kr} \text{ if } \begin{matrix} z_{ik}=1 \\ z_{jr}=1 \end{matrix} \right) \leftarrow \text{类别}$$

More compactly

$$M_z = \mathbb{E}[A | z] = Z \Psi Z^T$$

$Z$ :  $i^{\text{th}}$  row is  $z_i^T$

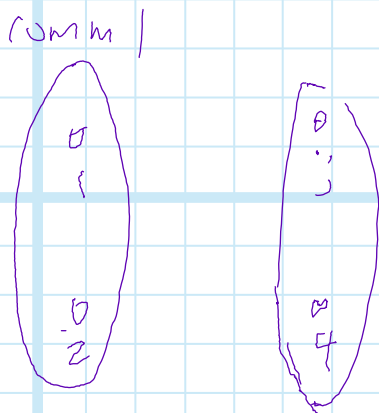
Intuitively, if  $\Psi_{kr} \approx 0$  whenever  $k \neq r$   
 (Between communities have no talking)



Diagonal block will be very dense, otherwise, sparse.

我們要找的基是  $A$  的 permute.

e.g. 4 nodes, 2 communities



$$n=4, k=2,$$

$$z_{11} = z_{21} = z_{32} = z_{42} = 1$$

$$E[A|z] = z \Psi z^T$$

$$= \begin{bmatrix} \Psi_{11} & \Psi_{11} & \Psi_{12} & \Psi_{12} \\ \Psi_{11} & \Psi_{11} & \Psi_{12} & \Psi_{12} \\ \Psi_{12} & \Psi_{12} & \Psi_{22} & \Psi_{22} \\ \Psi_{12} & \Psi_{12} & \Psi_{22} & \Psi_{22} \end{bmatrix}$$

$$\text{rank} = k = 2$$

generally, rank-k matrix.



# Models of $\Psi$

(实际上  $\Psi$  可以是任意对称阵)

① planted partition  $PP(p, q)$

$$\Psi = \begin{bmatrix} p & q \\ q & \dots & p \end{bmatrix}$$

$$p > q$$

↑  
community 中连接  
的权重比 community 间的大

$$= qE_q + (p - q)I_k$$

② Balance  $pp$ ,  $PP_b(p, q)$

every community has the same size  $m = \frac{n}{k}$

问题描述:

Given  $A$  generated from SBM, estimate  $Z$

log-likelihood function:  $\rightarrow \max$

$$\ell(Z, \Psi) = \sum_{i,j} A_{ij} \log(M_z)_{ij} + (1 - A_{ij}) \log(1 - (M_z)_{ij})$$

$$= \sum_{i,j} A_{ij} \left( f_0(M_z)_{ij} + (g_0(M_z)_{ij}) \right)$$

↑ apply  $f$  component-wise

$$f(x) = \log \frac{x}{1-x} \quad g(x) = \log(1-x)$$

Obs:  $f_0 M_z = f_0(Z\Psi Z^T) = Z(f_0\Psi)Z^T$

Under the  $pp$  model;  $PP(p, q)$

$$f_0\Psi = f(q)E_k + (f(p) - f(q))I_k$$

$$\Rightarrow \ell(z, \Psi) = \sum_{i < j} A_{ij} \left[ f(q) z E_k z^T + (f(p) - f(q)) z z^T \right]_{ij} \\ + \sum_{i < j} \left[ g(q) z E_k z^T + (g(p) - g(q)) z z^T \right]_{ij}$$

$$z E_k z^T = \sum_{k < n} z E_{k \times n} = E_n \quad \text{看起來是依賴于 } z, \text{ 但實際不是}$$

$$= \sum_{i < j} \left[ (f(p) - f(q)) A_{ij} (z z^T)_{ij} + g(p) - g(q) (z z^T)_{ij} \right] \\ + \text{constant}$$

obs:  $[z z^T]_{ii} = 1$

$$2 \ell(z, \Psi) = (f(p) - f(q)) \langle A, z z^T \rangle + \\ (g(p) - g(q)) \langle E_n, z z^T \rangle + \\ \text{const.}$$

Divide  $f(p) - f(q)$ : we obtain:

優化問題:

$$\hat{z} \in \arg \max_{z \in \mathcal{Z}} \langle A, z z^T \rangle - \lambda \langle E_n, z z^T \rangle \quad (P)$$

$\mathcal{Z}$ : set of admissible membership matrices.

constraint is discrete

non-convex

也可視為 SDP relaxation 求解.

consider  $X = z z^T \in \{0, 1\}^{n \times n}$

Note:  $X_{ij} = \begin{cases} 1 & \text{if } i, j \text{ in same community} \\ 0 & \text{o/w} \end{cases}$

→

$$(Q) \hat{X} \in \arg \max_{X \in \mathcal{X}} \langle A, X \rangle - \lambda \langle E_n, X \rangle$$

$\mathcal{X}$  = set of admissible clustering matrices

Under the  $p_b$  model, <sup>↓ 假设类 # / member 个数相同</sup> all  $Z$  take the form

up to permutation,

$$Z_0 = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ n & & & \ddots \end{bmatrix}$$

相同类的节点在一起

$$Z_0 = I_k \otimes \mathbb{1}_m$$

$$\Rightarrow \mathcal{Z} = \{ P Z_0 Q^T : P, Q : \text{perm matrices} \}$$

(correspondingly

$$\mathcal{X} = \{ P X_0 P^T, P : \text{perm matrices} \}$$

$$X_0 = Z_0 Z_0^T = I_k \otimes E_m$$

e.g.:

$$\begin{bmatrix} 1 \\ 0 \\ 26 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 3 \\ 40 \end{bmatrix}$$

$$z_0 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$x_0 = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ & & & & \ddots & \\ & & & & & 1 \end{bmatrix}$$

以上只是问题的 formulation

Idea: Semidefinite relaxation of (Q) SDR

obs:

$\underline{1}^0$ : For any  $X \in \mathcal{X}$ ,  $X \succeq 0$  ( $\because X = z z^T$ ),  
 $0 \leq x_{ij} \leq 1$

$x_{ii} = 1$  ( $i$  和  $i$  是在一个 community, 显然存在)

$\underline{z}^0$ :  $\langle E_n, X \rangle = \mathbf{1}_n^T X \mathbf{1}_n$  ( $\because E_n = \mathbf{1}_n \mathbf{1}_n^T$ )  
 $\uparrow$  全为 1

$$X \mathbf{1}_n = P X_0 P^T \mathbf{1}_n = P X_0 \mathbf{1}_n = m P \mathbf{1}_n = m \mathbf{1}_n$$

$$\Rightarrow \langle E_n, X \rangle = m n \text{ is const. !}$$

$\therefore$  只需要优化第一项了.

(SDR) of (Q):

$$\max \langle A, X \rangle$$

$$\text{s.t. } X \mathbb{1}_n = m \mathbb{1}_n$$

$$\text{diag}(X) = \mathbb{1}_n$$

$$X \succeq 0 \quad (\text{PSD})$$

$$X \succeq 0$$

$\because X \succeq 0 \therefore \begin{bmatrix} 1 & x_{ij} \\ x_{ij} & 1 \end{bmatrix} \succeq 0 \therefore$  不用写  $x \leq 1$ ,  
自动蕴含。

Question: Under what conditions would (SDP) give  $X_0$  as the optimal soln?

Idea: Construct a primal-dual certificate.  
我们希望  $X_0$  是 opt. soln. 能否建一个 dual?

Derivation of the dual

Obs:  $2(X \mathbb{1}_n)_i = \langle X, \Phi_i \rangle$

$$\Phi_i = e_i \mathbb{1}_n^T + \mathbb{1}_n e_i^T$$

$$\left( \langle X, \mathbb{1}_n e_i^T \rangle = e_i^T X \mathbb{1}_n \right)$$

$$X \mathbb{1}_n = m \mathbb{1}_n \Leftrightarrow \mathcal{L}(X) = 2m \mathbb{1}_n,$$

where  $L(x) = (\langle x, \Phi_1 \rangle, \dots, \langle x, \Phi_n \rangle)$

下面求其对偶.

(原问题(DP))  $\max \langle A, x \rangle$

$$\text{s.t. } x \mathbb{1}_n = m \mathbb{1}_n \quad (\mu)$$

$$\text{diag}(x) = \mathbb{1}_n, \quad (\nu)$$

$$x \neq 0 \quad x \geq 0 \quad (\Gamma)$$

(A)

$\rightarrow$  (SDD)  $\min 2m \mathbb{1}_n^T \mu + \mathbb{1}_n^T \nu$

$$\text{s.t. } \Lambda \triangleq L^*(\mu) + \text{diag}^*(\nu) - A - \Gamma \succeq 0$$

$$\Gamma \succeq 0.$$

$L^*$ : adjoint operator of  $L$ , defined by

$$L(x) = (\langle x, \Phi_1 \rangle, \dots, \langle x, \Phi_n \rangle)$$

$$\mu^T L(x) = \langle x, L^*(\mu) \rangle \quad \text{类似矩阵转置}$$

$$\mu^T L(x) = \sum_i \mu_i \langle x, \Phi_i \rangle = \langle x, \underbrace{\sum_i \mu_i \Phi_i}_{L^*(\mu)} \rangle$$

2019.4.16

SBM:

$$E[A|z] = z \Psi z^T$$

$A: n \times n$ ,  
membership matrix

$$z: n \times k$$

$$\Psi: k \times k$$

$$m = n/k$$

balanced PP<sub>b</sub>(p, q)

$$\bar{\Psi} = \begin{bmatrix} p & & q \\ & \ddots & \\ q & & p \end{bmatrix}$$

$$= q \cdot E_k + (p - q) \bar{L}_k$$

(SDP) max  $\langle A, X \rangle$ 

$$\text{s.t. } X \mathbb{1}_n = m \mathbb{1}_n$$

$$\text{diag}(X) = \mathbb{1}_n$$

$$X \succeq 0 \quad X \succeq 0$$

$$\mathcal{L}(X) = (\langle X, \Phi_1 \rangle, \dots, \langle X, \Phi_n \rangle)$$

$$\Phi_i = e_i \mathbb{1}^T + \mathbb{1}_n e_i^T$$

dual:  
(SDP)

$$\min 2m\mathbb{1}_h^T \mu + \mathbb{1}_n^T \nu$$

$$\text{s.t. } \Lambda \triangleq L^*(\mu) + \text{diag}^*(\nu) - A - T \succeq 0$$

$$T \succeq 0$$

WLOG:  $X_0 = Z_0 Z_0^T = \mathbb{I}_k \otimes E_m$  is the ground truth

$$Z_0 = \begin{bmatrix} | \\ | \\ | \\ \vdots \\ | \\ | \\ | \end{bmatrix} \quad X_0 = \begin{array}{c} \begin{array}{c} s_1 \quad s_2 \quad \dots \quad s_k \\ \hline E_m \end{array} \\ \begin{array}{c} \vdots \\ E_m \\ \vdots \end{array} \\ \begin{array}{c} \vdots \\ E_m \\ \vdots \end{array} \end{array} \begin{array}{l} s_1 \\ s_2 \\ \vdots \\ s_k \end{array}$$

$\leftarrow \sum_i X_{S_i}$   
 $S_i$ : set of indices of nodes in community  $i$

$\sum_{i=1}^k \mathbb{1}_{S_i}^T X_0 \frac{\mathbb{1}_{S_i}}{|S_i|}$  unique solution

Idea: Find a dual solution  $(\mu, \nu, T)$  s.t.

$(X_0, (\mu, \nu, T))$  form a primal-dual opt. pair.

$$(PF) L(X) = 2m\mathbb{1}, \text{diag}(X) = \mathbb{1}_n, X \succeq 0, X \geq 0$$

$$(DF) \Lambda \succeq 0, T \succeq 0$$

$$(CS) \langle \Lambda, X \rangle = 0, T_{ij} X_{ij} = 0 \quad \forall i, j$$

观察: (SDP)  $\mathcal{F}$ , region is compact, opt a continuous function on a compact set, it has soln.

下面求解这个 KKT:



From (CS), ①  $(X_0)_{s_k, s_k} = E_m \Rightarrow \Gamma_{s_k, s_k} = \Lambda$

$$\textcircled{2} \quad \Lambda X_0 = 0 \iff \langle \Lambda, X_0 \rangle = 0$$

(spectral decomposition)

$$\Rightarrow \ker(\Lambda) \supseteq \text{range}(X_0) = \{X_0 u\} \quad \forall u$$

Observe:  $\underline{I}_{s_k} \in \mathbb{R}^n$  is an eigenvector of  $X_0 \forall k$ , where

$$\underline{I}_{s_k} = \begin{bmatrix} 0 \\ \vdots \\ \mathbb{I}_m \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} s_1 \\ \vdots \\ s_k \\ \vdots \\ s_K \end{matrix}$$

$$\Rightarrow \text{range}(X_0) = \text{span}\{\underline{I}_{s_1}, \dots, \underline{I}_{s_K}\}$$

$$\therefore \text{span}\{\underline{I}_{s_1}, \dots, \underline{I}_{s_K}\} \subseteq \ker(\Lambda)$$

Prop 1: Let  $\mu, \nu, \Gamma$  be s.t.  $\Gamma \geq 0$ , Suppose

$$\left. \begin{array}{l} \text{(A1)} \quad \ker(\Lambda) = \text{span}\{\underline{I}_{s_1}, \dots, \underline{I}_{s_K}\}, \quad \Lambda \succ 0 \\ \text{(A2)} \quad \Gamma_{s_k, s_k} = 0 \quad \forall k, \end{array} \right\} \text{viz opt.}$$

(A3) For  $k \neq \ell$ ,  $\Gamma_{s_k, s_k}$  has at least 1 non-zero element  
 $\leftarrow$  viz uniqueness

Then,  $x_0$  is the unique opt soln to (SDR)

$(\mu, \nu, \Gamma)$  is opt for (SDID)

For (A1): It holds if

$$(A1-1) \quad \Lambda \mathbb{I}_{S_k} = 0 \quad \forall k$$

$$(A1-2) \quad u^T \Lambda u \geq \varepsilon \|u\|_2^2 \quad \forall u \in \text{span}\{\mathbb{I}_{S_1}, \dots, \mathbb{I}_{S_k}\}^\perp$$

Consider (A1-1):

$$0 = \Lambda \mathbb{I}_{S_k} = \begin{bmatrix} \square & & & \\ & \square & & \\ & & \dots & \\ & & & \square \end{bmatrix} \begin{bmatrix} \square \\ \square \\ \square \\ \square \end{bmatrix} \mathbb{I}_m$$

$$\Lambda_{S_k, S_k} \mathbb{1}_m = 0 \quad \Lambda_{S_k^c, S_k} \mathbb{1}_m = 0$$

Use the def'n of  $\Lambda$ :

$$L^*(\mu) = \sum_i \mu_i \Phi_i = \mu \mathbb{1}_n^T + \mathbb{1}_n \mu^T$$

$$\text{diag}^*(v) = \text{Diag}(v)$$

$$\text{diag}(x) = (\langle x, e_1 e_1^T \rangle, \dots, \langle x, e_n e_n^T \rangle)$$

$$\text{diag}^*(v) = \sum_i v_i e_i e_i^T = \text{Diag}(v)$$

diag: take the diag of matrix as a vector

Diag: take vector and form a matrix with it as diagonal

$$(\mu \mathbb{1}_n^T)_{S_k, S_k} = \mu_{S_k} \mathbb{1}_m^T$$

$$(\mathbb{1}_n \mu^T)_{S_k, S_k} = \mathbb{1}_m \mu_{S_k}^T$$

$$0 = \Lambda_{S_k, S_k} \mathbb{1}_m$$

$$= (\mu_{S_k} \mathbb{1}_m^T + \mathbb{1}_m \mu_{S_k}^T + \text{Diag}(v_{S_k}) - A_{S_k, S_k} - \overset{0}{\Gamma_{S_k, S_k}}) \mathbb{1}_m$$

$$= \underbrace{m \mu_{S_k} + (\mu_{S_k}^T \mathbb{1}_m) \mathbb{1}_m}_{\text{如果 } \mu_{S_k} = \frac{1}{2} \phi_k \mathbb{1}_m \text{ 则可以合并}} + v_{S_k} - \underbrace{A_{S_k, S_k} \mathbb{1}_m}_{\parallel}$$

如果  $\mu_{S_k} = \frac{1}{2} \phi_k \mathbb{1}_m$   
则可以合并

$$[d(S_k)]_{S_k}$$

$$\text{所以 } v_{S_k} = [d(S_k)]_{S_k} - m \phi_k \mathbb{1}_m$$

$$d(S_k) = A \mathbb{1}_{S_k}$$

同时  $\Lambda_{S_l, S_k} \mathbb{1}_m = 0$  即条件。

$$\Lambda_{S_l, S_k} \mathbb{1}_m = 0 \quad \forall k \neq l \quad \mathcal{L}^*(\mu) = \mu \mathbb{1}_n^T + \mathbb{1}_n \mu^T$$

$$\Rightarrow [\mu_{S_l} \mathbb{1}_m^T + \mathbb{1}_m \mu_{S_k}^T - (A + \Gamma)_{S_l, S_k}] \mathbb{1}_m = 0$$

代入  $\mu_{S_k}$

$$\Rightarrow m \left( \frac{1}{2} \phi_l + \frac{1}{2} \phi_k \right) \mathbb{1}_m = (A + \Gamma)_{S_l, S_k} \mathbb{1}_m$$

$\Rightarrow \mathbb{1}_m$  is an eigenvector of  $(A + \Gamma)_{S_l, S_k}$

$$(A+T)_{s_k, s_k} = \frac{1}{2}(\phi_k + \psi_k) E_m + B_{s_k, s_k}$$

where  $B_{s_k, s_k}$  acts on  $\text{span}\{\mathbb{1}_m\}^\perp$

下面来验证 (A1-2)

Note:  $\text{span}\{\mathbb{I}_{s_1}, \dots, \mathbb{I}_{s_k}\}^\perp$

$$= \left\{ u = \begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix} = \sum_k e_k \otimes u_k : \mathbb{1}_m^T u_k = 0 \forall k \right\}$$

Take  $u = (u_1, \dots, u_k) \in \text{span}\{\mathbb{I}_{s_1}, \dots, \mathbb{I}_{s_k}\}^\perp$

$$u^T \Lambda u = \sum_{k, \ell} u_k^T \Lambda_{s_k, s_\ell} u_\ell$$

$$= \sum_k u_k^T \Lambda_{s_k, s_k} u_k + \sum_{k \neq \ell} u_k^T \Lambda_{s_k, s_\ell} u_\ell$$

$$\Lambda_{s_k, s_k} = U_{s_k} \mathbb{1}^T + \mathbb{1}_m U_{s_k}^T + \text{diag}^*(\nu_{s_k}) A_{s_k, s_k}$$

$$\Rightarrow u_k^T \Lambda_{s_k, s_k} u_k$$

$$= u_k^T \left[ \text{diag}^*(\nu_{s_k}) - A_{s_k, s_k} \right] u_k$$

$$\left[ \Delta \triangleq A - \mathbb{E}[A|z] \right]$$

$$\Rightarrow \Delta_{s_k, s_k} = A_{s_k, s_k} - P E_m$$

$$= u_k^T \left[ \text{diag}^*(\nu_{s_k}) - P E_m - \Delta_{s_k, s_k} \right] u_k$$

can be drop:  $E_m = \mathbb{1}_m \mathbb{1}_m^T$

$\mathbb{1}_m^T u_k = 0$

$$= U_k^T \left[ \text{diag}^* \left( [d(s_k)]_{s_k} \right) - \phi_k m I_m - \Delta_{s_k, s_k} \right] U_k$$

$$A_{s_k, s_k} = U_{s_k} \mathbb{1}_m^T + \mathbb{1}_m U_{s_k}^T - (A + \Gamma)_{s_k, s_k}$$

$$= \frac{1}{2} (\phi_k + \phi_e) E_m - (A + \Gamma)_{s_k, s_k} = -B_{s_k, s_k}$$

$$\Rightarrow U^T \Lambda U = \sum_k U_k^T \left[ \text{diag}^* \left( [d(s_k)]_{s_k} \right) - \phi_k m I_m - \Delta_{s_k, s_k} \right] U_k$$

$$- \sum_{k \neq l} U_k^T B_{s_k, s_l} U_l$$

Set  $B_{s_k, s_l} = P_{\mathbb{1}_m^+} A_{s_k, s_l} P_{\mathbb{1}_m^+}$

obs:  $\Gamma_{s_k, s_l}$  is determined

$$U_k^T B_{s_k, s_l} U_l = U_k^T P_{\mathbb{1}_m^+} A_{s_k, s_l} P_{\mathbb{1}_m^+} U_l$$

$$\left[ \Delta_{s_k, s_l} = A_{s_k, s_l} - q E_m \right]$$

$$= U_k^T P_{\mathbb{1}_m^+} \left( \underbrace{q E_m}_{\text{projection} = 0} + \Delta_{s_k, s_l} \right) P_{\mathbb{1}_m^+} U_l$$

$$= U_k^T \Delta_{s_k, s_l} U_k$$

$$= \sum_k U_k^T \left[ \text{diag}^* \left( [d(s_k)]_{s_k} \right) - \Delta_{s_k, s_k} \right] U_k -$$

$$\sum_k m \phi_k \|u_k\|_2^2 - \sum_{k \neq l} a_k^T \Delta_{s_k, s_l} u_l$$

If  $\text{diag}^* [d(s_k)]_{s_k} \succcurlyeq m \rho_k I_m$  for some  $\rho_k$ , then

$$u^T \Lambda u \geq \sum_k \left[ (\rho_k - \phi_k) m - \|\Delta_{s_k, s_k}\| \right] \|u_k\|_2^2 - \sum_{k \neq l} a_k^T \Delta_{s_k, s_l} u_l$$

$$\geq \underbrace{\min_k \left[ (\rho_k - \phi_k) m - \|\Delta_{s_k, s_k}\| \right]}_I \|u\|_2^2 - \underbrace{\|E_{s_0} \circ \Delta\|}_{II} \|u\|_2^2$$

as long as  $I > II$ , done.

↓ need to use prob model.  
use matrix concentration inequality

LV  $\leq \text{fl}_{\frac{D}{2}}$  SDR approach, from practical, you cannot run in large scale, is it possible to develop some light weight, with theory guarantee.

Sum:

How prob & opt. get tied together.

Several algo ideas















