

Regression (1D)

Given $\{(x_i, y_i)\}_{i=1}^n$, "Training Data"

$x_i \in \mathbb{X}$: "input" features, representations, covariates, var

$y_i \in \mathbb{Y}$: "output" targets, labels, dependent variable.

$$\{(x_i, y_i)\}_{i=1}^n \quad f(x) = \beta x + c$$

ind.

OLS (Ordinary Least Squares) $(\hat{\beta}, \hat{c}) = \operatorname{argmin} \sum_{i=1}^n (y_i - \beta x_i - c)^2$

$$\hat{\beta} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{\operatorname{cov}(x_i, y_i)}{\operatorname{var}(x_i)}, \quad \hat{c} = \bar{y} - \hat{\beta} \bar{x}$$

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \operatorname{var}(x_i)$$

Goal: Find a map $f: \mathbb{X} \rightarrow \mathbb{Y}$

{ if $\mathbb{X} = \{\text{null}\}$ "unsupervised" }

Generalization Error = $\int \Psi(\hat{f}(x), y) d\mu(x, y)$ $\Psi: \mathbb{Y} \times \mathbb{Y} \rightarrow \mathbb{R}$
(lower is better)

High Dimensional

$$\mathbb{X} = \mathbb{R}^{d_n} \quad \mathbb{Y} = \mathbb{R}^{c_n}$$

Analysis $d_n, c_n \uparrow n$ e.g. $d_n = \beta n$ "proportional regime"
 $\beta \in (0, \infty)$

Operator Norm

$$\|M\|_{op} := \sup_{\vec{v} \neq 0} \frac{\|\vec{M}\vec{v}\|_2}{\|\vec{v}\|_2}$$

$$\stackrel{*}{=} \left[\max_i |\lambda_i(M)| \right] \text{ if } M \text{ is symmetric.}$$

Problem

$$y_1, y_2, \dots, y_n \stackrel{iid}{\sim} N(0, \Sigma_{d_n}) \quad \vec{y} \in \mathbb{R}^{d_n} \text{ Jt. Gaussian}$$

$$\begin{aligned} \text{Estimate } \Sigma_{d_n} \\ \hat{\Sigma}_n = \langle y_i y_i^T \rangle = \frac{1}{n} \sum_{i=1}^n y_i y_i^T \end{aligned}$$

Classical Regime: $d_n = d_0$ fixed. $\hat{\Sigma}_n \rightarrow \Sigma_{d_0}$ a.s.

if $d_n = \beta n$ $\hat{\Sigma}_n$ is random

Special Case if $\Sigma_{d_n} = I_{d_n}$

$$\text{then } \lim_{n \rightarrow \infty} \|\hat{\Sigma}_n - \Sigma_{d_n}\|_{op} = \beta + 2\sqrt{\beta} \text{ a.s. for } \beta \in (0, 1)$$

For a symmetric positive semidefinite matrix $A \in \mathbb{R}^{d \times d}$

$$F_A(t) = \frac{1}{d} \sum_{i=1}^d \mathbb{1}\{\lambda_i(A) \leq t\}$$

$$X = \begin{bmatrix} \vec{x}_1 & \vec{x}_2 & \cdots & \vec{x}_n \end{bmatrix} \in \mathbb{R}^{n \times d} \quad (\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n) \sim N(0, I_{d \times d})$$

$$\text{Then } F_n \xrightarrow{d} F \quad d = \beta n$$

PDF of location of EV
of \sum as $n \uparrow$

$$\frac{dF(s)}{ds} = \frac{\sqrt{(b-s)(s-a)}}{2\pi\beta t} \quad b = (1+\sqrt{\beta})^2 \quad a = (1-\sqrt{\beta})^2 \quad 0 < \beta < 1$$

Pre/Post multiplying by diagonal matrix scales rows/columns

High Dimensions $f(x) = \beta^T x + c \quad \beta \in \mathbb{R}^d$

$$\hat{\beta} = \bar{y} - \hat{\beta} \bar{x} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n \vec{x}_i$$

$$\hat{\beta}_* = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{y}$$

$$\bar{X}_{n \times d} = X_{n \times d} - \bar{I}_{1 \times d} = \begin{bmatrix} (\vec{x}_1 - \bar{x})^T \\ \vdots \\ (\vec{x}_n - \bar{x})^T \end{bmatrix}_{n \times d}$$

can also be written as $(\hat{\beta}, \hat{c}) = \operatorname{argmin} \|\bar{X}\beta - \bar{y}\|_2^2$

$\hat{\beta}_*$ is not the only solⁿ but least norm

$$\hat{\beta} = \hat{\beta}_* + \alpha \beta_{\text{null}}, \quad \beta_{\text{null}} \in \text{Null}(\bar{X})$$

$$\text{SVD } X_{n \times d} = U_{n \times n} S_{n \times d} V_{d \times d}^T \quad \begin{array}{l} U, V: \text{unitary} \\ S: \text{diagonal with non-ve values} \\ (S^+)^{ii} = \begin{cases} \frac{1}{S_{ii}} & \text{if } S_{ii} > 0 \\ 0 & \text{if } S_{ii} = 0 \end{cases} \end{array}$$

Probability Measure: countable additive map from a σ -field of subsets of Ω to real numbers which assigns value 0 to the empty set.

Notation $P(X \in A) = P(\{w \in \Omega : X(w) \in A\})$

'Probability that X is in A'

$X: (\Omega, \mathcal{F}) \mapsto (S, \mathcal{B}(S)), \quad \bar{X}(A) \in \mathcal{F} \quad \forall A \in \mathcal{B}(S)$

is called a R.V. (\approx a measurable map $(\Omega, \mathcal{F}) \mapsto (S, \mathcal{B}(S))$)

$$\bar{X}(A) = \{w \in \Omega \mid X(w) \in A\} \in \mathcal{F}$$

$$P(X \in A) := P(\bar{X}(A))$$

$$\text{ID CDF } F_X(t) = P(X \leq t)$$

c) Non-Decreasing b) $\lim_{t \rightarrow -\infty} F_X(t) = 0$ & $\lim_{t \rightarrow \infty} F_X(t) = 1$ c) right continuous

However, left limit exists: $P(X=t) = F_X(t) - \lim_{t' \rightarrow t^-} F_X(t')$

$$\text{Right limit } F_X(t) - F_X(t-) = P(X \in (t, t+\delta)) \xrightarrow{\delta \downarrow 0} 0$$

However in higher dimensions, it must also satisfy:

For every finite hypercube $A = [a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_d, b_d]$

$$P(A) = \sum_{v \in V(A)} \operatorname{sgn}(v) F_X(v) \geq 0 \quad \text{where}$$

$V(A) = \{a_1, b_1\} \times \{a_2, b_2\} \times \cdots \times \{a_d, b_d\}$ are vertices of A .

$$\operatorname{sgn}(v) = (-1)^{\# \text{ of } a_i's \text{ in } v}$$

Convergence of sequence of R.V. X_1, X_2, \dots, X_n

dealing with

For every $f: \mathbb{R}^d \rightarrow \mathbb{R}$ that satisfies these 4 conditions, \exists a unique \mathbb{R} -V. Furthermore, that R.V. satisfies

$$P(X \in A) = \sum_{v \in V(A)} \text{sgn}(v) F_X(v)$$

$$\text{Median}(m) \quad P(X \leq m) \geq 1/2 \quad \& \quad P(X \geq m) \geq 1/2, \lim_{x \rightarrow m^-} F(x) \leq \frac{1}{2} \leq F(x)$$

Quantile Function

$F_X(t)$ may not be ONTO

$$F_X^{-1}: [0, 1] \rightarrow \mathbb{R} \quad F_X^{-1}(q) = \inf_{t \in \mathbb{R}} \{t : F_X(t) > q\}$$

$F_X^{-1}(1/4) = 1^{\text{st}}$ quantile, $F_X^{-1}(3/4) = 3^{\text{rd}}$ quantile.

$$f_{X|Y}(x|y) = f_{X,Y}(x,y) \quad \text{for support}(y) = \{y : f_Y(y) > 0\}$$

$$f_Y(y)$$

$$\int f_{X|Y}(x|y) dx = \int_0^1 \mathbf{1}_{\{y \in \text{supp}(Y)\}} dx$$

Expectation of R.V.

Jensen's Inequality

$$E[X] = \int x dF_X(x) = \int x f_X(x) dx \quad (\text{if density exists})$$

If ϕ is convex: $\phi(\alpha x_1 + (1-\alpha)x_2) \leq \alpha \phi(x_1) + (1-\alpha)\phi(x_2)$

$$E[\phi(X)] \geq \phi(E[X])$$

If $p > q \geq 1$ Proof uses Holder: $\|X\|_p \geq \|X\|_q$

Holder: if $\frac{1}{p} + \frac{1}{q} = 1$ $|E[XY]| \leq \|X\|_p \|Y\|_q$

k^{th} moment exists $E|X|^k < \infty$ means smaller moments also exist

$$\text{Var}(X) = \min_{c \in \mathbb{R}} E((X-c)^2) = E[X^2] - E^2[X]$$

if $X \in \mathbb{R}^d$, $E[X] \in \mathbb{R}^d$

Covariance Matrix / Variance Covariance Matrix $E[XX^T] - E[X]E^T[X] = E(X - E[X])(X - E[X])^T$

Markov If $X \geq 0$ & $E[X] < \infty$ $P(X \geq t) \leq \frac{E[X]}{t}$

Stronger Version Random walk $U \perp\!\!\!\perp X$ & U uniform in $[0, 1]$. then

$$P(X \geq tU) \leq E[X]/t$$

MGF does not always exist

Characteristic F^n always exists

$$M_X(t) = E(e^{tX})$$

$$\phi_X(t) = E[e^{itX}]$$

$$\phi_X(0) = 1, |\phi_X(t)| \leq 1$$

$$M_X^{(k)}(0) = E[X^k]$$

$$\phi'_X(t) = j E[X e^{itX}]$$

$$\text{if } E|X|^2 < \infty \quad P(|X - E[X]| > t) \leq \frac{\text{Var}(X)}{t^2}$$

Chebychev

Chernoff

$$P(X \geq t) = P(uX \geq ut) \quad \text{if } u > 0$$

$$= P(e^{uX} \geq e^{ut}) \leq \frac{E[e^{uX}]}{e^{ut}}$$

$$\leq \inf_{u>0} \bar{e}^{ut} M_X(u)$$

$$\text{② } X_n \rightarrow X \text{ i.p. } X_n \xrightarrow{P} X$$

$$\text{if } \forall \epsilon > 0 \quad \lim_{n \rightarrow \infty} P(\lim_{n \rightarrow \infty} |X_n - X| > \epsilon) = 0 \quad \text{joint}(X_n, X)$$

$$a.s. \Rightarrow i.p.$$

In statistics $\hat{\theta}_n$: estimate from n samples of θ^*

$$\text{if } \hat{\theta}_n \xrightarrow{P} \theta^* \quad \text{"consistency"}$$

↳ allowed to be random.

$$\text{WLLN } \{X_i\}_{i=1}^{\infty} \text{ i.i.d. } F_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n F_{X_i}(x_i) \text{ independence}$$

$$= \prod_{i=1}^n F_{X_i}(x_i) \text{ identically distributed}$$

$$\left| \begin{array}{l} \text{with } E[X_i] = \mu \\ \rightarrow \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mu \\ \text{random const} \end{array} \right.$$

↳ stronger version: For pairwise independence also holds.

$y = X \oplus Z$, $X, Z \sim \text{Ber}(\frac{1}{2})$. Pairwise but not jointly independent

SLLN $\{X_n\}$ pairwise independence. Let x_1, x_2, \dots, x_n

i.i.d. R.V. & assume $E|X| < \infty$. Let $\mu = E[X]$

$$\sum_{i=1}^n X_i/n \rightarrow \mu \text{ a.s.}$$

(Weak*) Convergence in distribution/Law of large numbers $\{X_i\} \xrightarrow{d} X$

if $F_{X_i}(t) \rightarrow F_X(t) \Leftrightarrow \phi_{X_i}(t) \rightarrow \phi_X(t) \Leftrightarrow E[f(X_i)] \rightarrow E[f(X)]$

for all continuity points of F_X Pointwise convergence

CLT $\{X_i\}_{i=1}^{\infty}$ i.i.d. $E|X_i|^2 < \infty$ $\mu = E[X_i]$ $\sigma^2 = \text{Var}(X_i)$

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{X_i - \mu}{\sigma} \xrightarrow{d} N(0, 1)$$

↳ if it was n it would converge to 0 by SLLN

"8-pa"

$$M_X^{(k)}(0) = E[X^k]$$

$$\phi'_X(t) = j E[X e^{itX}]$$

$$\text{if } E|X|^2 < \infty \quad P(|X - E[X]| > t) \leq \frac{\text{Var}(X)}{t^2}$$

Chebychev

Chernoff

$$P(X \geq t) = P(uX \geq ut) \quad \text{if } u > 0$$

$$= P(e^{uX} \geq e^{ut}) \leq \frac{E[e^{uX}]}{e^{ut}}$$

$$\leq \inf_{u>0} \bar{e}^{ut} M_X(u)$$

$$\text{if } E|X|^2 < \infty \quad P(|X - E[X]| > t) \leq \frac{\text{Var}(X)}{t^2}$$

Chebychev

Chernoff

$$P(X \geq t) = P(uX \geq ut) \quad \text{if } u > 0$$

$$= P(e^{uX} \geq e^{ut}) \leq \frac{E[e^{uX}]}{e^{ut}}$$

$$\leq \inf_{u>0} \bar{e}^{ut} M_X(u)$$

$$\text{if } E|X|^2 < \infty \quad P(|X - E[X]| > t) \leq \frac{\text{Var}(X)}{t^2}$$

Chebychev

Chernoff

The standard inner product on $\mathbb{R}^{m \times n}$
is given by
 $\langle x, y \rangle = \text{tr}(x^T y) = \sum_{i=1}^m \sum_{j=1}^n x_{ij} y_{ij}$
 inner product of associated vectors in \mathbb{R}^m

classmate

Date _____

Page _____

→ log Partition Function

Statistics

data: y_1, y_2, \dots, y_n i.i.d. observations of R.V. $\gamma \sim P$ distribution

$A(\theta)$ is a convex fⁿ of θ

P_θ : Parametrised family of distribution

Proof 1] $\nabla_\theta^2 A(\theta) \geq 0$ Proof 2] Using Hölder's Inequality

$\hat{\theta}_n$: "estimator of θ " $\hat{\theta}_n : \mathbb{Y}^n \rightarrow \Theta$,

$$\int f(y)g(y)dy \leq \left(\int |f(y)|^p dy \right)^{1/p} \left(\int |g(y)|^q dy \right)^{1/q}$$

where $\frac{1}{p} + \frac{1}{q} = 1, p \geq 1$

$\theta \in \Theta$ "parameter space"

$\hat{\theta}_n = \hat{\theta}_n(y_1, \dots, y_n)$ R.V., $E[\hat{\theta}_n]$: constant

Likelihood Function $L_n(\theta) = \prod_{i=1}^n f(x_i; \theta)$

Bias of $\hat{\theta}_n$

$$E[\hat{\theta}_n] - \theta^* = \hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n y_i, \hat{\sigma}_n^2 = \frac{1}{n} (y_i - \hat{\mu}_n)^2$$

Eg. Exponential Family i.i.d. data $y_1, y_2, \dots, y_n \stackrel{iid}{\sim} f_{Y,\theta}(y)$

Consider

$$L(\theta) = \prod_{i=1}^n e^{\langle \theta, T(y_i) \rangle - A(\theta)} h(y_i)$$

For all distributions whose 1st & 2nd moment exists

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} e^{\langle \theta, \sum T(y_i) \rangle - n A(\theta)}$$

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} A(\theta) - \langle \theta, \frac{1}{n} \sum_{i=1}^n T(y_i) \rangle$$

Convex + Affine = Convex

Then $\hat{\mu}_n$ is unbiased estimator of μ

(convex Opt problem when Θ is convex set,

$$\text{Bias}(\hat{\sigma}_n^2/n) = -\frac{\sigma^2}{n}, \text{Bias}(\frac{n}{n-1} \hat{\sigma}_n^2) = 0$$

$$\nabla A(\theta) - \frac{1}{n} \sum_{i=1}^n T(y_i) = 0 \quad \Theta \subseteq \int_{y \sim f_{Y,\theta}(y)} T(y) = \frac{1}{n} \sum_{i=1}^n T(y_i)$$

$$\text{MSE}(\hat{\theta}_n) = E[|\hat{\theta}_n - \theta^*|^2] = E[(\hat{\theta}_n - E\hat{\theta}_n + E\hat{\theta}_n - \theta^*)^2]$$

[Unique Solⁿ] if $\nabla^2 A(\theta) > 0 \iff T_1(y), \dots, T_S(y)$ are linearly independent

$$(not \text{ random}) = E[|\hat{\theta}_n - E\hat{\theta}_n|^2] + \cancel{E[|\theta^* - E\hat{\theta}_n|^2]}$$

can drop E

If sufficient statistics are L.I.

$$= \text{Variance}(\hat{\theta}_n) + \text{Bias}^2(\hat{\theta}_n)$$

Exponential Families are called minimal.

Theorem

If $\text{bias} \rightarrow 0$ & $\text{Var}(\hat{\theta}_n) \rightarrow 0$, then $\hat{\theta}_n$ is consistent

∴ Conv. in Mean Square Sense \Rightarrow Convergence i.p.

Conjugate of a Function:

$$f^*(y) = \sup_{x \in \text{dom } f} (y^T x - f(x))$$

↳ Convex El if f not convex

$(f^*)^*$ = Biconjugate of f = $\begin{cases} f & \text{if f is convex, lower semi continuous} \\ \text{Dual Norm} & \end{cases}$

Exponential Family Distributions

$$f_y(y) = \underbrace{e^{\langle \theta, T(y) \rangle - A(\theta)}}_{Z(\theta)} h(y)$$

Let $\|\cdot\|$ be a norm on \mathbb{R}^n , Associated dual norm $\|\cdot\|_*$

$$\|z\|_* = \sup \{ z^T x \mid \|x\| \leq 1 \}$$

$$z^T x \leq \|x\| \|z\|_*$$

Dual of l_p norm is l_q norm $\frac{1}{p} + \frac{1}{q} = 1$

$$f(x) = \|x\|_p^p \quad f^*(y) = \|y\|_2^2$$

$$f(x) = \|x\| \quad f^*(y) = \begin{cases} 0 & \|y\|_* \leq 1 \\ \infty & \text{otherwise.} \end{cases}$$

θ : natural param, $h(y)$: Base measure/ind. of θ

$$s\text{-param family} \quad \langle \theta, T(y) \rangle = \sum_{i=1}^s \theta_i T_i(y) \quad T_i : \mathbb{Y} \rightarrow \mathbb{R}, T : \mathbb{Y} \rightarrow \mathbb{R}^s$$

T: sufficient statistic fⁿ

Z: Partition Fⁿ / A: log Partition Fⁿ

$$Z \Leftrightarrow \int f_y(y) dy = 1 \Rightarrow Z(\theta) = \int h(y) e^{\langle \theta, T(y) \rangle} dy$$

$\Theta = \{\theta : Z(\theta) < \infty\}$ Natural param space.

In $N(\mu, \Sigma)$ Σ : Covar Matrix, Σ^{-1} : Precision

$$y^T \Sigma^{-1} y = \langle \text{flatten}(\Sigma^{-1}), \text{flatten}(yy^T) \rangle$$

Properties $A(\theta) = \ln Z(\theta)$

$$1: \nabla_\theta A(\theta) = E_\theta T(y)$$

2: $\nabla_\theta^2 A(\theta)$: Covariance Matrix of sufficient statistics

$$(\nabla_\theta^2 A(\theta))_{ij} = \frac{\partial}{\partial \theta_i} \left(\frac{\partial A(\theta)}{\partial \theta_j} \right) = \text{cov}(T_i(y), T_j(y))$$

$\{\lambda_i\}, \{v_i\}$ EV of $A \in \mathbb{R}^{n \times n}$

Properties of Symmetric Matrix

- ① $A = A^T$
- ② All Eigenvalues Real
- ③ λ_i, \vec{v}_i & λ_j, \vec{v}_j EV, EV pair

$$\text{④ } \lambda_i \neq \lambda_j \Rightarrow \vec{v}_i^T \vec{v}_j = 0$$

b) $\lambda_i = \lambda_j$, Eigenspace of $\lambda_i = \text{span}(v_i, v_j)$

WLOG $\{\vec{v}_i\}_{i=1}^n$ orthonormal $\vec{v}_i^T \vec{v}_j = \delta_{ij}$

+ $\vec{x} \in \mathbb{R}^n, \vec{x} = \sum_{i=1}^n \alpha_i \vec{v}_i$ where $\alpha_i = \vec{v}_i^T \vec{x}$

$$A \vec{x} = \sum_{i=1}^n \lambda_i \vec{x} \vec{v}_i^T \vec{v}_i = \left(\sum_{i=1}^n \lambda_i \vec{v}_i \vec{v}_i^T \right) \vec{x}$$

Eigen Decomposition of A

$$A = \sum_{i=1}^n \lambda_i \vec{v}_i \vec{v}_i^T$$

$$V = [\vec{v}_1 \vec{v}_2 \dots \vec{v}_n] \quad \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$$

$$r = \text{rank}(A), \dim(\text{null}(A)) = n-r$$

$$V_{n \times n} = \begin{bmatrix} V_r \\ V_{n-r} \end{bmatrix} \quad \Lambda = \begin{bmatrix} \Lambda_r \\ 0_{n-r} \end{bmatrix}$$

$\{\vec{v}_i : \lambda_i \neq 0\}$ $\{\vec{v}_i : \lambda_i = 0\}$

$$B = V \Lambda V^T = V_r \Lambda_r V_r^T = \sum_{i=1}^r \lambda_i \vec{v}_i \vec{v}_i^T$$

Non-zero EV

($A^T A$) $_{n \times n}$

a) Symmetric

$$A^T A \vec{v} = \lambda \vec{v}$$

$$\lambda_i = \|A\vec{v}_i\|^2 / \|\vec{v}_i\|^2$$

b) all EV are non-negative

$$(C) \text{ If } (\lambda_i, \vec{v}_i) \text{ is EP of } A^T A, \text{ then } (\lambda_i, A\vec{v}_i) \text{ is EP of } A A^T$$

$$(A A^T) A \vec{v}_i = A(A^T A \vec{v}_i) = \lambda_i A \vec{v}_i, \quad \vec{u}_i = \frac{A \vec{v}_i}{\sqrt{\lambda_i}}, \quad \|\vec{u}_i\| = 1$$

$$\|A\vec{v}_i\| = \sqrt{\lambda_i}$$

$$A^T A = V \Lambda_r V^T \quad A A^T = U \Lambda_m U^T$$

$$\Lambda_r = \begin{bmatrix} \Lambda_r \\ 0_{n-r} \end{bmatrix}_{n \times n} \quad \Lambda_m = \begin{bmatrix} \Lambda_r \\ 0_{m-r} \end{bmatrix}_{m \times m}$$

$$A^T A = V_r \Lambda_r V_r^T \quad A A^T = U_r \Lambda_r' U_r^T$$

$$\text{span}(V_r) = \text{range}(A^T A) = \text{range}(A A^T)$$

$$\text{span}(V_{n-r}) = \text{null}(A)$$

$$\text{span}(U_r) = \text{range}(A)$$

$$\text{span}(U_{n-r}) = \text{null}(A^T)$$

$$x \in \mathbb{R}^n: X_{n(A)} + X_{n(A)} = \sum_{i=1}^n \alpha_i \vec{v}_i + \sum_{i=r+1}^n \alpha_i \vec{v}_i$$

$$A \vec{x} = \sum_{i=1}^r \alpha_i A \vec{v}_i = \sum_{i=1}^r \alpha_i \vec{u}_i \sqrt{\lambda_i}$$

$$AV = A[\vec{v}_1 \vec{v}_2 \dots \vec{v}_n]$$

$$= [\vec{u}_1 \sqrt{\lambda_1}, \vec{u}_2 \sqrt{\lambda_2}, \dots, \vec{u}_n \sqrt{\lambda_n}] = U_r \Lambda_r^{\frac{1}{2}}$$

$$A = U_n \Lambda_n^{\frac{1}{2}} V^T$$

$$A = \sum_{i=1}^r \sqrt{\lambda_i} u_i v_i^T = U_r \Lambda_r^{\frac{1}{2}} V_r^T$$

assume $m > n$ U_n has n columns of A

Economy version of SVD

$\sqrt{\lambda_i}$: singular values of $A \geq 0$

Diagonal from right (left) scales column (rows)

if $m < n$ work with A^T $A^T = V_m \Lambda_m^{\frac{1}{2}} U^T$

? Projection onto $\text{span}(A^T)$ $V_r V_r^T = A^T A$

$\dot{x}(t) = -\Lambda x(t)$ $x(t) \in \mathbb{R}^{n \times 1}$ $A \in \mathbb{R}^{n \times n}$

$x(t) = e^{-At} x(0)$, $A = V \Lambda V^T$

$\dot{x} = -V \Lambda V^T x \Rightarrow \dot{u} = -\Lambda u$ $u := V^T x$

$u_i(t) = e^{-\lambda_i t} u_i(0)$ wlog $\lambda_1 > \lambda_2 > \dots > \lambda_n$

slowest decaying $u_n(t) = e^{-\lambda_n t} u_n(0)$

$u^T u = x^T V V^T x = x^T x \quad \|x(t)\|_2 = \|u(t)\|_2 \leq e^{-\lambda_n t} \|u_n(0)\|_2$

$e^A = I + A + \frac{A^2}{2!} + \dots = V \Lambda V^T$ $x(t) = e^{-At} x(0)$

$x(t) = V e^{-At} V^T x(0)$

0 has multiplicity $n-r$, $\lambda_{r+1}, \dots, \lambda_n$ $u = V^T x = [V_r^T \quad V_{n-r}^T]^T x$

$\lim_{t \rightarrow \infty} \|x(t)\|^2 = \sum_{k=r+1}^n u_k^2(0) = \|V_{n-r}^T x(0)\|^2$

\hookrightarrow Squared length of projection into $N(A)$

$x(t) = V e^{-At} V^T x(0) = [V_r \quad V_{n-r}] \begin{bmatrix} e^{-\lambda_1 t} 0 & \\ 0 & I \end{bmatrix} \begin{bmatrix} V_r^T \\ V_{n-r}^T \end{bmatrix} x(0)$

$= V_r e^{-\lambda_1 t} V_r^T x(0) + V_{n-r} I V_{n-r}^T x(0)$

$I = V V^T = [V_r \quad V_{n-r}] \begin{bmatrix} V_r^T \\ V_{n-r}^T \end{bmatrix} = V_r V_r^T + V_{n-r} V_{n-r}^T$

Digression Maximal Inequalities.

Forward Euler Discretization

At stepsize/tearing rate/resolution

$$\frac{x_{t+\Delta t} - x_t}{\Delta t} = -A \mathbf{I}_d \quad t = k \Delta t$$

$$x_{k+1} = (\mathbf{I} - \Delta t A) x_k$$

$$B = W \sum_i W^T, u_k := W^T \mathbf{I}_d$$

$$\vec{u}_{k+1} = \sum_i \vec{u}_k \quad \Sigma = \text{diag}(\sigma_i), \sigma_i = 1 - \Delta t \lambda_i$$

$$|\lambda_i| < 1 \Rightarrow 0 < \Delta t < \frac{2}{\lambda_1} \quad \lambda_1 > \lambda_2 > \dots$$

Max step size for forward Euler to converge ODE

$$x_k = (V V^T - \Delta t V \Lambda V^T)^k x_0 = V (\mathbf{I} - \Delta t \Lambda)^k V^T x_0$$

$$= [V_n \ V_{n-r}] \begin{bmatrix} (\mathbf{I} - \Delta t \Lambda_n)^k & 0 \\ 0 & \mathbf{I}_{n-r} \end{bmatrix} \begin{bmatrix} V_n^T \\ V_{n-r}^T \end{bmatrix} x_0$$

$$x_k = V_n (\mathbf{I} - \Delta t \Lambda_n)^k V_n^T x_0 + V_{n-r} V_{n-r}^T x_0$$

$$\|V_n (\mathbf{I} - \Delta t \Lambda_n)^k V_n^T x_0\|_2 \leq \epsilon$$

$$\Rightarrow \max(|1 - \Delta t \lambda_1|, |1 - \Delta t \lambda_r|)^k \|u_0\| \leq \epsilon$$

Let X_1, \dots, X_N be N R.V. such that $X_i \sim \text{subG}(\sigma^2)$

$$\text{Then } E[\max_{1 \leq i \leq N} X_i] \leq \sigma \sqrt{2 \log N} \quad E[\max_{1 \leq i \leq N} |X_i|] \leq \sigma \sqrt{2 \log 2N}$$

$$P(\max_{1 \leq i \leq N} X_i > t) \leq N e^{-t/\sigma^2}$$

Note that Random Variables need not be independent:

$$\text{Proof: } E[\max_i X_i] = \frac{1}{\delta} \mathbb{E}[\log e^{\delta \max_i X_i}] \leq \frac{1}{\delta} \log \mathbb{E}[e^{\delta \max_i X_i}] \quad \text{Jensen}$$

$$\leq \frac{1}{\delta} \log \sum_i E[e^{\delta X_i}] \leq \frac{1}{\delta} \log N e^{\frac{\sigma^2}{2}} \quad B = \sqrt{\frac{2 \log N}{\sigma^2}}$$

$$\max_{1 \leq i \leq N} |X_i| = \max_{1 \leq i \leq N} X_i \quad X_{N+i} = X_i \text{ for } i=1, \dots, N$$

\rightarrow $\max_{1 \leq i \leq N} |X_i| = \max_{1 \leq i \leq N} X_i$

Addⁿ Structural Assumption: Sparsity, deterministic

$$\bar{Y}_n \sim N(\theta^*, \frac{I_d}{n}), \theta^* \text{ 1-sparse wlog } \theta_i^* \neq 0$$

$$i^{\text{th}} \text{ component } \bar{Y}_{ni} \sim N(\theta_i^*, \frac{1}{n})$$

$$i^* = \arg \max_{i \in [d]} |\bar{Y}_{ni}|, \hat{\theta} = \begin{cases} \hat{\theta}_{i^*} = \bar{Y}_{ni^*} \\ \theta_i^* = 0 \quad \forall i \neq i^* \end{cases}$$

Claim

$$E\|\hat{\theta} - \theta^*\|^2 \leq \frac{C \log d}{n}, \quad d_n = o(n^2) \text{ & have } \hat{\theta} \xrightarrow{L^2} \theta^*$$

S1] Show $P(\|\hat{\theta} - \theta^*\|_2 > t) \leq 2d e^{-\frac{nt^2}{4}}$

S2] Use for $Z \geq 0$, $E[Z] = \int_0^\infty P(Z > t) dt$

$$S1] \|\hat{\theta} - \theta^*\| = \sqrt{(\theta_1^* - \hat{\theta}_1)^2 \mathbf{1}\{i^*=1\}^2 + \sum_{i=2}^d (\theta_i^* - \hat{\theta}_i)^2 \mathbf{1}\{i^*=i\}^2}$$

$$E\|\hat{\theta} - \theta^*\|^2 \geq t \quad \hat{\theta}_i \mathbf{1}\{i^*=i\} = Y_{ni} \mathbf{1}\{i^*=i\}$$

$$Q \Rightarrow \max((\theta_1^* - \bar{Y}_{n1})^2, \bar{Y}_{n2}^2, \dots, \bar{Y}_{nd}^2) \geq t^2/2 \quad E \Rightarrow Q$$

$$P(E) \leq P(Q) \leq P(\|\theta^* - \bar{Y}_{n1}\| > t) + \sum_{i=2}^d P(\|\bar{Y}_{ni}\| > \frac{t}{\sqrt{2}})$$

$$\leq 1 \cdot 2 e^{-\frac{nt^2}{4}} + (d-1) 2 e^{-\frac{nt^2}{4}}$$

$$= 2d e^{-\frac{nt^2}{4}}$$

S2] Having a Tail Bound Helps to get a bound in expectation

$$E[\|\hat{\theta} - \theta^*\|^2] = \int_0^\infty P(\|\hat{\theta} - \theta^*\|^2 \geq t) dt = \int_0^\infty 2d e^{-\frac{nt^2}{4}} dt$$

if a factor of d comes up split the integral,

$$= \int_0^\infty P(\quad) dt + \int_u^\infty 2d e^{-\frac{nt^2}{4}} dt$$

$$\leq u + \frac{8d}{n} \bar{e}^{-nu/4} \quad \text{minimized at } u = \frac{4}{n} \ln \frac{d}{2}$$

$$d_n = o(n)$$

$$\lim_{n \rightarrow \infty} Q\left(\frac{dn}{2}, \frac{\sigma n}{2}\right) = \frac{1}{2}$$

$$\lim_{n \rightarrow \infty} P(\|\bar{Y}_n - \theta^*\| > \sqrt{\sigma}) = \frac{1}{2}$$

Def'

Convergence in Probability in High Dimensions

$$X_n \in \mathbb{R}^d, X_n \xrightarrow{P} X$$

$$\lim_{n \rightarrow \infty} P(\sup_i |X_{ni} - X_i| > t) = 0$$

$$\lim_{n \rightarrow \infty} P(\|X_n - X\|_\infty > t) = 0$$

(Approximately) Sparse Mean Estimation.

 $y_i \sim N(\theta^*, I_d)$ θ^* is approximately S -sparse.

$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$ (can be other estimator also)

 S/\bar{S} : top s magnitude coordinates of $\hat{\theta}/\theta^*$ Deterministic
bound.

$S1 \quad \|\hat{\theta} - \theta^*\| \geq 3 \sqrt{\frac{2s \log 2d}{n}} + \sqrt{2} \|\theta_{-S^*}\|_1$

$\Rightarrow \exists i \in [d] \text{ s.t. } |\bar{Y}_{ni} - \theta_i^*| > \sqrt{\frac{2s \log 2d}{n}}$

$\Leftrightarrow \forall i \in [d] \quad |\bar{Y}_{ni} - \theta_i^*| \leq \sqrt{\frac{2s \log 2d}{n}} \text{ (ct)} \Rightarrow \|\hat{\theta} - \theta^*\| \leq t$

$S2 \quad P(\|\hat{\theta} - \theta^*\| \geq t) \leq P(\exists i \in [d] \text{ s.t. } |\bar{Y}_{ni} - \theta_i^*| > ct)$

$\leq 3\sqrt{ct} \cdot \frac{1}{\sqrt{n}} \leq 2d e^{-\frac{nct^2}{2}} \in \delta \quad t_c = \sqrt{\frac{2s \log 2d}{n}}$

begin with
this

$i \in S \quad \text{then} \quad \hat{\theta}_i = \bar{Y}_{ni} - \theta_i^* \leq ct$

$\hat{\theta}_i = \bar{Y}_{ni} - \theta_i^* \leq \bar{Y}_{ni} + |\bar{Y}_{ni}|$

$\leq ct + |\bar{Y}_{nj}| \text{ for some } j \in \bar{S}$

$\leq ct + |\bar{Y}_{nj}| + |\theta_j^*|$

$\leq 2ct + |\theta_j^*|$

$\hat{\theta}_i = (\theta_i^*)$

$$\begin{aligned} \|\hat{\theta} - \theta^*\|^2 &= \sum_{i \in \bar{S} \cup \bar{S}^*} |\hat{\theta}_i - \theta_i^*|^2 \leq 8ct^2 \\ &\quad + \sum_{i \in \bar{S}^*} |\hat{\theta}_i - \theta_i^*|^2 \left(\leq \sum_{j \in \bar{S}} (2ct + |\theta_j^*|)^2 \leq \sum_{j \in \bar{S}} 2(8c^2t^2 + |\theta_j^*|^2) \right) \\ &\quad + \sum_{i \in \bar{S}^*} |\hat{\theta}_i - \theta_i^*|^2 \left(= \sum_{i \in \bar{S}^*} |\theta_i^*|^2 \right) \\ &\leq 98ct^2 + 2\|\theta_{-S^*}\|^2 \left(= \sum_{j \in \bar{S} \cup \bar{S}^*} |\theta_j^*|^2 \right) \end{aligned}$$

consider

$A \quad y_i \stackrel{iid}{\sim} N(\theta^*, \frac{1}{n} \Sigma) \quad \Sigma \text{ known} \quad \theta^* \text{ sparse & unknown}$

$B \quad y_i \stackrel{iid}{\sim} N(A\theta^*, I_m) \quad A \text{ known}, \quad \begin{matrix} \downarrow \\ \text{① can be reduced} \\ \text{to ② if } \Sigma \text{ full rank} \end{matrix}$

$y_i = A\theta^* + \xi_i \quad \xi_i \sim N(0, I_m)$

$\text{Sample Average} \quad \bar{Y}_n = A\theta^* + \xi \quad \xi \sim N(0, \frac{I_m}{n})$

Consider noiseless case $n \rightarrow \infty$ so $\xi = 0$

$\text{Opt problem} \quad \min_{\theta \in R^d} \|\theta\|_1 \text{ such that } X\theta = \bar{y} \quad \begin{matrix} \text{Cost f'' non-differentiable} \\ \downarrow \text{Convex Relaxation} \end{matrix}$

$\min_{\theta \in R^d} \|\theta\|_1 \text{ such that } X\theta = y$

Exact Recovery & Restricted Nullspace.

When is solving BP equivalent to solving L0 problem?

$T(\theta^*) = \{ \Delta \in R^d \mid \|\theta^* + t\Delta\|_1 \leq \|\theta^*\|_1, \text{ for some } t \in \mathbb{R} \}$

 $\theta^* + \text{null}(X)$ passes through Tangent cone

Tangent Cone

Unfavourable
Case

$\text{null}(X) := \{ \Delta \in R^d \mid X\Delta = 0 \}$

$C(S) := \{ \Delta \in R^d \mid \|\Delta_S\|_1 \leq \|\Delta\|_1 \}$

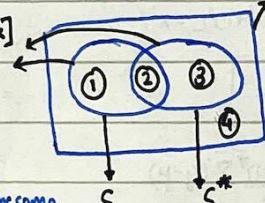
• The matrix X satisfies the Restricted Nullspace Propertywrt S if $C(S) \cap \text{null}(X) = \{0\}$ $C(S)$ captures behaviour of tangent cone $T(\theta^*)$ ind. of θ^*
For any S -sparse vector θ^* , tangent cone $T(\theta^*)$ is contained in union of such tangent cones.

The following two properties are equivalent

a) For any vector $\theta^* \in R^d$ with support S , BP applied with $y = X\theta^*$ has unique soln $\hat{\theta} = \theta^*$ b) Matrix X satisfies restricted nullspace property wrt S .Restricted Isometry Property (RIP) For given set $\{1, \dots, d\}$ we say that $X \in R^{n \times d}$ satisfies RIP of order s with const $\delta_s(X) > 0$ if

$$\left\| \frac{X_S^T X_S - I_S}{n} \right\|_2 \leq \delta_s(X) \quad \begin{matrix} \text{for all subsets } S \text{ of} \\ \text{size at most } s. \end{matrix}$$

$$\text{if } A \text{ satisfies RIP} \quad (1 - \delta_s) \|\theta\|_2^2 \leq \frac{\|A\theta\|_2^2}{n} \leq (1 + \delta_s) \|\theta\|_2^2$$

 θ : s sparse

$$\left\| \frac{X_S^T X_S - I_S}{n} \right\|_2 \leq \delta_s(X)$$

$$\text{if } A \text{ satisfies RIP} \quad (1 - \delta_s) \|\theta\|_2^2 \leq \frac{\|A\theta\|_2^2}{n} \leq (1 + \delta_s) \|\theta\|_2^2$$

 θ : s sparse

Prediction $f(x) = \theta^T x$

 $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_i (y_i - \hat{x}_i^T \theta)^2 + \lambda \|\theta\|^2 = (\hat{X}^T \hat{X} + \lambda I) \hat{X}^T y$

"Robust Regression" $\hat{\theta}_{\text{robust}} = \underset{\theta}{\operatorname{argmin}} \sum_i \underbrace{|x_i^T \theta - y_i|}_{\text{CVL}} + \lambda \|\theta\|^2$

Logistic Regression $P(y_i = +1 | x_i) = P_{\text{cl}} = \frac{e^{x_i^T \theta}}{1 + e^{x_i^T \theta}}$

$\hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | x_i) = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \sum_i \log(1 + e^{-x_i^T \theta})$

Gradient of loss lies in span of data.

Convexity \Rightarrow GD converges

$\text{prox}_f(x) = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left(f(\theta) + \frac{1}{2} \|\theta - x\|^2 \right)$

$\nabla \cdot \text{prox}_f = \sum_i \frac{\partial}{\partial \theta_i} \text{prox}_f, \quad \nabla \cdot F = \text{div } F = \sum_{i=1}^d \frac{\partial F}{\partial x_i}$

$R = X + N(0, \sigma^2)$

observed $P_x = f_x(t) \propto e^{-S(t)} = \frac{-(R-x)^2}{2\sigma^2}$

$P(R|x)$ likelihood $\propto e^{-\frac{(R-x)^2}{2\sigma^2}}$

$P(X) : \text{prior}, \quad P(X|R) = \frac{P(X|R)}{P(R)} = \frac{P(X)}{P(R)}$

posterior $P(X|R) \propto \frac{-(R-x)^2}{2\sigma^2} + S(x)$

MAP $\Rightarrow \hat{X}(R) = \underset{x}{\operatorname{argmax}} P(x|R) = \text{prox}_S(R)$

$\hat{X}(R) = E[X|R] = \int x t \frac{e^{-\frac{(R-t)^2}{2\sigma^2} + S(t)}}{\int e^{-\frac{(R-t)^2}{2\sigma^2} + S(t)} dt} dt$

$E[X|R] = \underset{f}{\operatorname{argmin}} E[(X - f(R))^2] = E[E[X|f(R)]|R]$

$\min_w \ell(w) = \sum_i \log(1 + e^{y_i x_i^T w})$

better formulation

SVM $\min_w \|w\|, \quad y_i^T w \geq 1$

Soft Margin SVM $\min_{w, \beta} \frac{1}{2} \|w\|^2 + \sum_i \xi_i$

$J(w, \beta, \alpha) = \frac{1}{2} \|w\|^2 + \beta^T \xi + \alpha^T (1 - \xi - \tilde{X}^T w)$

$\alpha \geq 0 \quad \begin{bmatrix} \xi \\ \tilde{X}^T w \end{bmatrix}$

$g(\alpha) = \min_{w, \beta} J(w, \beta, \alpha)$

$\beta(\alpha) = 1^T \alpha + \min_w \frac{1}{2} \|w\|^2 - \alpha^T \tilde{X}^T w + \min_{\beta} \beta^T (1 - \alpha)$

$\beta(\alpha) = -\frac{1}{2\lambda} \alpha^T \tilde{X} \tilde{X}^T \alpha - \delta_{(-\infty, 1]}(\alpha) + 1^T \alpha$

$\max_{\alpha \geq 0} g(\alpha) = \max_{0 \leq \alpha \leq 1} \frac{1^T \alpha - \alpha^T \tilde{X} \tilde{X}^T \alpha}{2\lambda}$

Projected GD $\max_{0 \leq \alpha \leq 1} -\frac{\alpha^T \tilde{X} \tilde{X}^T \alpha}{2\lambda} + 1^T \alpha$

$\alpha_{t+1} = \text{proj}_{[0, 1]^n}(\alpha_t + \eta_t \left(-\frac{\tilde{X} \tilde{X}^T}{\lambda} \alpha_t + 1 \right))$

$\hat{\theta}_{\text{TV}} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2} \|y - A\theta\|^2 + \frac{\lambda}{2} \|\theta\|^2$

$\hat{\theta}_{\text{TV}} = (A^T A + \lambda I_d)^{-1} A^T y, \quad y = A\theta + \xi$

$E\|\hat{\theta}_{\text{TV}} - \theta^*\|^2 = \|(A^T A + \lambda I_d)^{-1} A^T A \theta^* - \theta^*\|^2$

$+ E\|(A^T A + \lambda I_d)^{-1}\|^2$

use trace trick \rightarrow

+ SVD + assumption $E[\xi_i \xi_j] = \mathbb{1}(i=j)$

$\text{Var} = \sigma^2 \sum_{i=1}^d \frac{s_i^2}{s_i^2 + \lambda^2}$

$\text{bias}^2 = \sum_{i=1}^d \frac{s_i^2}{(s_i^2 + \lambda)^2}$

$\lim_{\lambda \rightarrow 0} \text{bias}^2 = \sum_{i=1}^d s_i^2 = \|V_{d-r}^T \theta^*\|^2$

$M = U \Sigma V^T$

$\sum_{i=1}^d s_i^2 = \|M\|_F^2 = \text{tr}(M M^T) = \text{tr}(\Sigma \Sigma^T)$

Column of squares of singular values

$\min_w \|y - Xw\|^2 + \|w\|_L^2$

$\|w\|_L^2 = \sum_{j=1}^d (w_j)^T L w_j$

Sylvester eq: $X^T X w^T + w^T L = X^T y$

$w = \begin{bmatrix} -w^{(1)} \\ -w^{(2)} \\ \vdots \\ -w^{(d)} \end{bmatrix}$

$f_y(y) = \frac{1}{\sqrt{2}} (y - H)^T \Sigma^{-1} (y - \mu)$

$y^T \Sigma^{-1} y = \langle yy^T, \Sigma^{-1} \rangle$

$= \langle \text{Vec}(yy^T), \text{Vec}(\Sigma^{-1}) \rangle$

$\min_{\theta} \frac{1}{2} \|y - A\theta\|^2 + \lambda f(\theta)$

$A = USV^T \quad \tilde{y} = U^T y$

$\min_{\theta, \beta} \frac{1}{2} \|\tilde{y} - S\beta\|^2 + \lambda f(\theta)$

s.t. $V^T \beta = \beta$

init r_t, s_t

$\theta_t = \text{prox}_{\frac{\lambda F}{\delta t}}(r_t) \quad \beta_t = \text{prox}_{\frac{\lambda f}{\delta t}}(s_t)$

$\eta_t = \frac{p \delta_t}{\nabla \cdot \text{prox}_{\frac{\lambda F}{\delta t}}(r_t)} \quad r_t = p \delta_t$

$\delta_t = \eta_t - s_t$

$\delta_{t+1} = \eta_t - \delta_t$

$\eta_t = V^T \left(\frac{\eta_t + \theta_t - \beta_t}{\delta_t} \right) \quad \eta_{t+1} = V \left(\frac{\eta_t + \theta_t - \beta_t - \eta_t}{\delta_{t+1}} \right)$

Mean Estimation in High Dimensions
 $y_i \stackrel{\text{iid}}{\sim} N(\theta^*, I_{d_n})$ $\bar{Y}_n = \frac{1}{n} \sum_{j=1}^n Y_j$
 $E[\|\bar{Y}_n - \theta^*\|^2] = E \sum_{i=1}^{d_n} (\bar{y}_{ni} - \theta^*)^2 = \sum_{i=1}^{d_n} \text{Var}(\bar{y}_{ni})$
 $\bar{Y}_n \xrightarrow{L^2} \theta^* \Rightarrow \bar{Y}_n \xrightarrow{P} \theta^*$
 $\bar{Y}_n \not\xrightarrow{P} \theta^* \Rightarrow \bar{Y}_n \xrightarrow{L^2} \theta^*$
 $A \Rightarrow B \Leftrightarrow \delta \geq \gamma^* \text{ and } P(A) \leq P(B)$
 $n \|\bar{Y}_n - \theta^*\|^2 \sim \chi^2(d_n)$
 $P(\|\bar{Y}_n - \theta^*\| > t) = P(n \|\bar{Y}_n - \theta^*\|^2 > nt^2) = 1 - F_{\chi^2(d_n)}(nt^2)$
 $\lim_{n \rightarrow \infty} Q\left(\frac{d_n}{2}, \frac{nt^2}{2}\right) = Q\left(\frac{d_n}{2}, \frac{nt^2}{2}\right)$

RIP A G R^{mxd} satisfies RIP of order S with const $\delta_S(A)$ if
 $\left\| \frac{A_S^T A_S - I_S}{n} \right\|_2 \leq \delta_S(A)$
 for all subsets S of size at most S.
 Let $\theta \in \mathbb{R}^d$ have support S.
 $-\delta_S \|\theta_S\|^2 \leq \theta_S^T \left(\frac{A_S^T A_S - I_S}{n} \right) \theta_S \leq \delta_S \|\theta_S\|^2$
 $(1 - \delta_S) \|\theta\|^2 \leq \frac{1}{n} \|A \theta\|^2 \leq (1 + \delta_S) \|\theta\|^2$
 $\frac{1}{n} \|\theta\|^2 \leq \|\theta\|^2 \leq \frac{1}{1 - \delta_S} \|\theta\|^2$

Convergence in Probability in High Dimensions
 $X_n \in \mathbb{R}^d$ $X_n \xrightarrow{P} X$ if
 $\lim_{n \rightarrow \infty} P(\sup_i |X_{ni} - X_i| > t) = 0$ ($\lim_{n \rightarrow \infty} P(\|X_n - X\|_\infty > t) = 0$) uniform restricted nullspace property holds for any subset S

Maximal Inequality Let x_1, \dots, x_N be $N \in \mathbb{R}^d$ s.t. cardinality $|S| \leq 8$
 $x_i \sim \text{SubG}(\sigma^2)$
 Then $E \left[\max_{1 \leq i \leq N} x_i \right] \leq \sigma \sqrt{2 \log N}$

Conc inequality for SubGaussian R.V.
 y with p, if $E[e^{t(y-\mu)}] \leq e^{\frac{3t^2}{2}}$
 $\theta \sim \text{SubG}(F)$ $P(|y - \mu| > \varepsilon) \leq 2 \exp\left(-\frac{\varepsilon^2}{2F^2}\right)$

Sparse Mean Estimation $y_i \sim N(\theta^*, I_d)$
 θ^* is S-approximately sparse $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$
 $\hat{\theta}: \hat{\theta}(i) = \bar{y}_n(i) \mathbf{1}(i \in \text{Top } S \text{ coordinates of } \bar{Y}_n)$
S1 Begin with $|\bar{y}_{ni} - \theta^*| \leq ct$ $\forall i$ & find upper bound on $|\hat{\theta} - \theta^*| \leq u$

S2 $\|\hat{\theta} - \theta^*\| \geq 3 \sqrt{\frac{28 \log 2d}{n} + \|\theta_S^*\|_1 \sqrt{2}}$
 $\Rightarrow \exists i: |\bar{y}_{ni} - \theta^*| \geq ct$
 $A \Rightarrow B$
 $P(A) \leq P(B)$
 $\leq \sum_i P(|\bar{y}_{ni} - \theta^*| \geq ct) = 2dc^2 \quad (= \delta)$

$y \stackrel{\text{iid}}{\sim} N(\lambda \theta^*, I_m)$
 Basis Pursuit (BP) $\min_{\theta \in \mathbb{R}^d} \|\theta\|_1$ such that $X\theta = y$
 $T(\theta^*) = \{ \Delta \in \mathbb{R}^d \mid \|\theta^* + \Delta\|_1 \leq \|\theta^*\|_1, t > 0 \}$
 $C(S) = \{ \Delta \in \mathbb{R}^d \mid \|\Delta_S\|_1 \leq \|\theta_S^*\|_1 \}$

Def: The matrix X satisfies the Restricted Nullspace Property wrt S if $C(S) \cap \text{null}(X) = \{0\}$
 The following are equivalent
 a) For any vector $\theta \in \mathbb{R}^d$ with support S, BP applied with $y = X\theta^*$ has unique soln $\hat{\theta} = \theta^*$
 b) The matrix X satisfies Restricted Nullspace property wrt X.
 $\hat{\theta} = \theta - \theta^*$
 $\Rightarrow a) \|\theta_S^*\|_1 = \|\theta^*\|_1 \geq \|\theta^* + \hat{\theta}\|_1$
 $\geq \|\theta_S^*\|_1 - \|\hat{\theta}_S\|_1 + \|\hat{\theta}_S^*\|_1$

Conclude $\hat{\theta} \in C(S)$ but also $\hat{\theta} \in \text{null}(X)$
 Thus $\hat{\theta} = 0$ & $\theta = \theta^*$

$a \Rightarrow b$ it suffices to show that if BP relaxation succeeds w.r.t sparse vectors then $\text{null}(X) \cap C(S) = \emptyset$. Choose $\beta^* \in \text{null}(X)$
 $\min_{\beta} \|\beta\|_1, \text{s.t. } X\beta = X \begin{bmatrix} \theta_S^* \\ 0 \end{bmatrix}$ by ass. unique optimal soln $\beta = [\theta_S^* \ 0]^T$
 $\|\beta\|_1 \leq \|\theta_S^*\|_1$, Thus $\theta^* \notin C(S)$