

High Dimensional Statistics

Yanbo Tang

Imperial College London

Feb. 2025

Content of this week

- Concentration for functions of random variables
- Performance bounds for linear regression
- LASSO

Sub-Gaussian random vectors

Definition

A random variable X is called sub-Gaussian with proxy variance σ^2 if there exists a $\sigma^2 > 0$:

$$\mathbb{E}[\exp(\lambda(X - \mu))] \leq \exp\left(\frac{\lambda^2 \sigma^2}{2}\right)$$

for all $\lambda \in \mathbb{R}$. A random vector $X \in \mathbb{R}^d$ is sub-Gaussian with proxy variance σ^2 if $u^\top X$ is sub-Gaussian with proxy variance σ^2 for all $u \in S^{d-1}$.


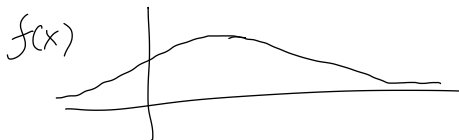
Concentration for functionals

$$\frac{\sum_{i=1}^n X_i - E(X)}{n}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f(X_1, X_2, X_3, \dots, X_n) - E[f(X_1, X_2, \dots, X_n)]$$

not good $f(x_i)$

Bounded differences

Suppose we have a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ for all

$x_1, x_2, \dots, x_d, x'_1, x'_2, \dots, x'_d \in \mathbb{R}$

$$|f(x'_1, x_2, \dots, x_j, \dots, x_d) - f(x_1, x_2, \dots, x_j, \dots, x_d)| \leq L_1 < \infty$$

$$\vdots$$

$$|f(x_1, x_2, \dots, x'_j, \dots, x_d) - f(x_1, x_2, \dots, x_j, \dots, x_d)| \leq L_j$$

$$\vdots$$

$$|f(x_1, x_2, \dots, x_j, \dots, x'_d) - f(x_1, x_2, \dots, x_j, \dots, x_d)| \leq L_d.$$

Bounded differences/McDiaramids

Proposition

(Bounded differences inequality/McDiaramids) Suppose that f satisfies the bounded difference property with parameters (L_1, \dots, L_n) and that the random vector

$$X = (X_1, X_2, \dots, X_n)$$

has independent components. Then

$$\mathbb{P}[|f(X) - \mathbb{E}[f(X)]| \geq t] \leq 2e^{-\frac{2t^2}{\sum_{k=1}^n L_k^2}} \quad \text{for all } t \geq 0.$$

U-Statistics

"Commaid" $\frac{\sum_{i=1}^n X_i}{n} \approx E(X)$

$$E[|X_1 - X_2|] \approx \frac{\sum_{i < j} |X_i - X_j|}{\binom{n}{2}}$$

$$U = \frac{1}{\binom{n}{2}} \sum_{i < j} g(X_i, X_j)$$

Imagine $\|g\|_\infty < b$ (or $X_i \in [-b, b]$ a.s.)

$$\begin{aligned} \forall j \in \{1, \dots, n\} |U(X_1, \dots, X_n) - U(X_1, \dots, X_j', \dots, X_n)| &\leq \frac{1}{\binom{n}{2}} \sum_{i \neq j} |g(X_i, X_j) - g(X_i, X_j')| \\ &\leq \frac{(n-1) \cdot 2b}{\binom{n}{2}} = \frac{4b}{n} \end{aligned}$$

$$\therefore L_j = \frac{4b}{n} \quad \forall j \in \{1, 2, \dots, n\}$$

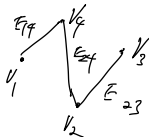
$$P(|U - E(U)| \geq t) \leq 2 \exp\left(-\frac{nt^2}{8b^2}\right)$$

Asymptotic statistics
Van der Vaart

Clique sizes for random graphs

$$G = (V, E)$$

A "clique" is a set of fully connected vertices.



Maximal clique size $\in \{1, \dots, d\}$

Erdős-Rényi: Random graph



$$X_{ij} = \text{Bernoulli}(p)$$

$$C(G) = C(X_{11}, X_{12}, \dots, X_{ij}, \dots, X_{dd})$$

$$|C(X_{11}, \dots, X_{ij}', \dots, X_{dd}) - C(X_{11}, \dots, X_{dd})| \leq 1$$

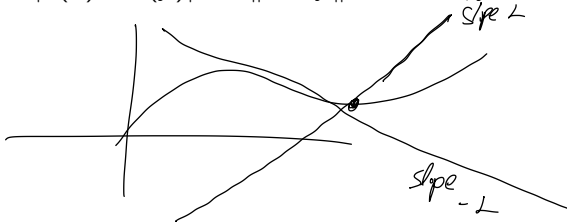
$$p\left[\frac{1}{n} |C(G) - E[C(G)]| \geq t\right] \leq 2e^{-nt^2}$$

$$n = \binom{d}{2}$$

Lipschitz property

We say that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -Lipschitz with respect to the Euclidean norm if:

$$|f(x) - f(y)| \leq L\|x - y\|_2 \text{ for all } x, y \in \mathbb{R}^n.$$



Gaussian concentration

\uparrow
 n

Theorem

Let (X_1, \dots, X_n) be a vector of IID standard Gaussian random variables and let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a L -Lipschitz function with respect to the Euclidean norm. Then $f(X) - E[f(X)]$ is sub-Gaussian with parameter at most L and

$$\mathbb{P}[|f(X) - E[f(X)]| \geq t] \leq 2 \exp\left(\frac{-t^2}{2L^2}\right) \text{ for all } t \geq 0.$$

Singular value decomposition

As a reminder for a real matrix $A \in \mathbb{R}^{n \times d}$, the singular value decomposition is:

$$A = \sum_{i=1}^r s_i(A) u_i v_i^T, \text{ where } r = \text{rank}(A).$$

The non negative numbers $s_i(A)$ are called the singular values of A , the vectors $u_i \in \mathbb{R}^n$ are the left singular vectors of A , and $v_i \in \mathbb{R}^d$ are the right singular vectors of A .

Singular values are eigenvalues $\sqrt{\lambda(A^T A)}$, $\sqrt{\lambda(A A^T)}$

Wely's theorem

$$s_1 \geq s_2 \geq \dots \geq s_d$$

Theorem

Given two matrices X and Y in $\mathbb{R}^{n \times d}$, we have

$$\max_{i=1, \dots, d} |s_i(X) - s_i(Y)| \leq s_1(X - Y) \leq \|X - Y\|_F,$$

where $\|\cdot\|_F$ is the Frobenius norm of a $\mathbb{R}^{n \times d}$ matrix:

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d a_{ij}^2} = \sqrt{\text{Trace}(A^\top A)} = \sqrt{\sum_{i=1}^{\min(n,d)} s_i(A)^2}.$$

Singular values of Gaussian random matrices

$X \in \mathbb{R}^{m \times d}$ where $x_j \sim N(0, 1)$

$$\max_{i=1, \dots, d} |S_i(x) - S_i(y)| \leq \|X - Y\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^d (x_{ij} - y_{ij})^2}$$

$$\forall_{k=1, \dots, d} \quad P(|S_k(X) - \mathbb{E}[S_k(X)]| \geq \delta) \leq \exp\left(-\frac{\delta^2}{2}\right)$$

Terence Tao notes on Random matrix theory

What about non-Gaussians?

We need some other assumptions:

Definition

A distribution supported in \mathbb{R}^n with density $p(x) = \exp(-\psi(x))$ is said to be γ strongly log concave if there exists a $\gamma > 0$ such that:

$$\lambda\psi(x) + (1 - \lambda)\psi(y) - \psi(\lambda x + (1 - \lambda)y) \geq \frac{\gamma}{2}\lambda(1 - \lambda)\|x - y\|_2^2,$$

for all $\lambda \in [0, 1]$ and $x, y \in \mathbb{R}^n$.

Theorem

Let \mathbb{P} be any strongly log-concave distribution with parameter $\gamma > 0$. Then for any L -Lipschitz function with respect to the Euclidean norm:

$$\mathbb{P}[|f(X) - \mathbb{E}[f(X)]| \geq t] \leq 2 \exp\left(\frac{-\gamma t^2}{4L^2}\right).$$

Concentration inequalities
by MASSART

Linear Regression

Prediction

Low MSE
for new obs.

↑

this week

Adaptation

sparsity pattern
recovery

Inference

Confidence interval
for θ

Introduction

Most regression models can be written in the form of:

$$Y_i = f(x_i) + \epsilon_i, i = 1, \dots, n,$$

where $f(\cdot)$ is some functional relationship and ϵ_i are some centred error terms.

Here, we assume the data generating model is:

$$Y_i = x_i^\top \theta^* + \epsilon_i, i = 1, \dots, n,$$

$$\epsilon \sim \text{i.i.d. } N(0, \sigma^2)$$

MSE

We first consider the performance of our estimated models in terms of the *Mean Square Error* (MSE), for a general regression problem this is:

$$MSE(\hat{f}_n) = \frac{1}{n} \sum_{i=1}^n \left(\hat{f}_n(x_i) - f(x_i) \right)^2,$$

for us this will simplify to

$$MSE(X\hat{\theta}) = \frac{1}{n} \|X(\hat{\theta}_n - \theta^*)\|_2^2,$$

where $\hat{\theta}_n$ is some estimated value for the regression parameter and θ^* is the true data-generating value of θ .

Unconstrained estimation

Proposition

The least squares estimator $\hat{\theta}^{LS} \in \mathbb{R}^d$ satisfies

$$X^T X \hat{\theta}^{LS} = X^T Y.$$

Moreover, $\hat{\theta}^{LS}$ can be chosen to be

$$\hat{\theta}^{LS} = (X^T X)^\dagger X^T Y,$$

where $(X^T X)^\dagger$ denotes the Moore-Penrose pseudoinverse of $X^T X$.

$$A \in \mathbb{R}^{m \times n}$$

$$Ax = b$$

A^\dagger is such that

$$\forall x \in \mathbb{R}^n \quad \|Ax - b\|_2$$

$$\geq \|A z - b\|_2$$

$$z = A^\dagger b$$

Proof: $\theta \mapsto \|y - X\theta\|_2^2$ is convex so all minima
satisfy $\nabla_\theta \|y - X\theta\|_2^2 = 0$

$$\nabla_\theta \|y - X\theta\|_2^2 = -2(X^T X - \theta^T X^T X)^T = 0$$

$$\Rightarrow X^T X \theta = X^T y$$

Cont.

\lesssim

$$5 \log(n) \cdot p \lesssim \log(n) \cdot p$$

Theorem

Assume that the linear model holds where $\varepsilon \sim \text{subG}_n(\sigma^2)$. Then the least squares estimator $\hat{\theta}^{LS}$ satisfies

$$\mathbb{E}[MSE(X\hat{\theta}^{LS})] = \frac{1}{n} \mathbb{E} \|X\hat{\theta}^{LS} - X\theta^*\|_2^2 \lesssim \sigma^2 \frac{r}{n},$$

where $r = \text{rank}(X^\top X)$. Moreover, for any $\delta > 0$, with probability at least $1 - \delta$,

$$MSE(X\hat{\theta}^{LS}) \lesssim \sigma^2 \frac{r + \log(1/\delta)}{n}.$$

Small fact that we'll need:

Moments of sub-Gaussian random variables. Let X be any random variable such that

$$\mathbb{P}[|X| > t] \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right),$$

then for any positive integers $k \geq 1$,

$$E[|X|^k] \leq (2\sigma^2)^{k/2} k\Gamma(k/2).$$

Proof

$$\|\gamma - X\hat{\theta}\|_2^2 \leq \|\gamma - X\theta^*\|_2^2 = \|\varepsilon\|_2^2$$

$$\begin{aligned} \|\varepsilon\|_2^2 &\geq \|\gamma - X\hat{\theta}\|_2^2 = \|X\hat{\theta} - X\theta^*\|_2^2 - 2\varepsilon^T X(\hat{\theta} - \theta^*) + \|\varepsilon\|_2^2 \\ &= \|X\hat{\theta} - X\theta^*\|_2^2 - 2\varepsilon^T X(\hat{\theta} - \theta^*) \end{aligned}$$

$$0 \geq \|X\hat{\theta} - X\theta^*\|_2^2 - 2\varepsilon^T X(\hat{\theta} - \theta^*)$$

$$\|X\hat{\theta} - X\theta^*\|_2^2 \leq 2\varepsilon^T X(\hat{\theta} - \theta^*) = 2\|X\hat{\theta} - X\theta^*\|_2 \cdot \left(\frac{\varepsilon^T X(\hat{\theta} - \theta^*)}{\|X\hat{\theta} - X\theta^*\|_2} \right)$$

Let $\Phi = [\phi_1, \dots, \phi_n] \in \mathbb{R}^{p \times n}$ be an orthonormal basis for X .

$$\exists V: X(\hat{\theta} - \theta^*) = \Phi V$$

$$\frac{\varepsilon^T X(\hat{\theta} - \theta^*)}{\|X\hat{\theta} - X\theta^*\|_2} = \frac{\varepsilon^T \Phi V}{\|\Phi V\|_2} = \frac{\varepsilon^T \Phi V}{\|V\|_2} = \frac{\tilde{\varepsilon}^T V}{\|V\|_2} \leq \sup_{u \in \mathbb{B}_2^n} \tilde{\varepsilon}^T u$$

$$\tilde{\varepsilon} = \Phi \varepsilon \quad \text{Sub-Gaussian (Exercise)} \quad (\mathbb{R}^2)$$

$$\mathbb{E}[\|X\hat{\theta} - X\theta^*\|_2^2] \leq \mathbb{E}\left[4 \sup_{u \in \mathbb{B}_2^n} (\tilde{\varepsilon}^T u)^2\right] = 4 \sum_{i=1}^n \mathbb{E}[\tilde{\varepsilon}_i^2] \leq 16\sigma^2 n$$

Proof Cont.

$$\sup_{\text{solve}} (\hat{\epsilon}' u)^2 \leq 8 \log(d) \sigma^2 \cdot r + 2 \sigma^2 \log\left(\frac{1}{\delta}\right)$$

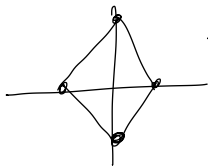
L^1 constraint

$$\hat{\theta}_k \in \argmin_{\theta \in K} \|y - x\theta\|_2^2$$

$$\|x\hat{\theta}_k - x\theta^*\|_2^2 \leq 2\varepsilon^T x(\hat{\theta}_k - \theta^*)$$

$$\leq 2 \sup_{\theta \in K-K} (\varepsilon^T x\theta) \quad K-K = \{x-y, x, y \in K\}$$

$$\text{if } K \text{ symmetric} \quad = 4 \sup_{v \in K} (\varepsilon^T v)$$



L^1 ball

has 2-d vertices

$$V = \{e_1, -e_1, e_2, -e_2, \dots\}$$

$$XV = \{x_1, -x_1, x_2, -x_2, \dots\}$$

Sub-Gaussian width of L^1 ball

Theorem

Let P be a polytope with N vertices $v^{(1)}, \dots, v^{(N)} \in \mathbb{R}^d$ and let $X \in \mathbb{R}^d$ be a random vector such that $[v^{(i)}]^\top X, i = 1, \dots, N$, are sub-Gaussian random variables with variance proxy σ^2 . Then

$$\mathbb{E} \left[\max_{\theta \in P} \theta^\top X \right] \leq \sigma \sqrt{2 \log(N)},$$

and

$$\mathbb{E} \left[\max_{\theta \in P} |\theta^\top X| \right] \leq \sigma \sqrt{2 \log(2N)}.$$

Moreover, for any $t > 0$,

$$\mathbb{P} \left(\max_{\theta \in P} \theta^\top X > t \right) \leq N e^{-\frac{t^2}{2\sigma^2}},$$

and

$$\mathbb{P} \left(\max_{\theta \in P} |\theta^\top X| > t \right) \leq 2N e^{-\frac{t^2}{2\sigma^2}}.$$

Performance

Theorem

Let \mathcal{B}_1 be the unit ℓ_1 ball of \mathbb{R}^d , $d \geq 2$ and assume that $\theta^* \in \mathcal{B}_1$. Moreover, assume the conditions of Theorem 6 and that the columns of \mathbb{X} are normalized such that $\max_j |\mathbb{X}_j|_2 \leq \sqrt{n}$. Then the constrained least squares estimator $\hat{\theta}_{\mathcal{B}_1}^{LS}$ satisfies

$$\mathbb{E}[MSE(\mathbb{X}\hat{\theta}_{\mathcal{B}_1}^{LS})] = \frac{1}{n} \mathbb{E}|\mathbb{X}\hat{\theta}_{\mathcal{B}_1}^{LS} - \mathbb{X}\theta^*|_2^2 \lesssim \sigma \sqrt{\frac{\log d}{n}}.$$

Moreover, for any $\delta > 0$, with probability $1 - \delta$, it holds

$$MSE(\mathbb{X}\hat{\theta}_{\mathcal{B}_1}^{LS}) \lesssim \sigma \sqrt{\frac{\log(d/\delta)}{n}}.$$

Proof

$$\|X\hat{\theta} - X\theta\|_2^2 \leq 4 \sup_{V \in XK} (C\varepsilon^T V)$$

suppose $\varepsilon \sim \text{sub}_2(\sigma^2)$ $\forall x_j$ $\|x_j\|_2^2 = n$ $\varepsilon^T x_j \sim \text{sub}_2(n\sigma^2)$

since $\varepsilon^T x_j = \|x_j\|_2 \frac{\varepsilon^T x_j}{\|x_j\|_2}$

apply our previous bound on polynomial

$$\begin{aligned} P(|\text{MSE}(\hat{\theta}_n)| > t) &\leq P\left(\sup_{V \in XK} (C\varepsilon^T V) > \frac{nt}{4}\right) \\ &\leq 2d \exp(-nt^2/32\sigma^2) \end{aligned}$$

$$2d e^{-nt^2/32\sigma^2} \leq \delta \iff$$

$$t^2 \geq 32 \frac{\log(2d)}{n} + 32 \frac{\sigma^2 \log(1/\delta)}{n}$$

L^0 performance

Denote by $B_0(k)$ the ℓ_0 "ball" of \mathbb{R}^d , i.e., the set of k -sparse vectors, defined by

$$B_0(k) = \{\theta \in \mathbb{R}^d : |\theta|_0 \leq k\}.$$

Our goal is to control the MSE of $\hat{\theta}_K^{IS}$ when $K = B_0(k)$. Note that computing $\hat{\theta}_{B_0(k)}^{IS}$ defined as:

$$\hat{\theta}_{B_0(k)}^{LS} \in \operatorname{argmin}_{\theta \in B_0(k)} \|Y - \mathbb{X}\theta\|_2^2,$$

but this would require computing $\binom{d}{k}$ least squares estimators since this loss is no longer smooth due to the constraint

Best subset selection

Theorem

Fix a positive integer $k \leq d/2$. Let $K = B_0(k)$ be set of k -sparse vectors of \mathbb{R}^d and assume that $\theta^* \in B_0(k)$. Moreover, assume the conditions of Theorem 6. Then, for any $\delta > 0$, with probability $1 - \delta$, it holds

$$MSE(\mathbb{X}\hat{\theta}_{B_0(k)}^{IS}) \lesssim \frac{k\sigma^2}{n} \log\left(\frac{ed}{2k}\right) + \log(6) \frac{\sigma^2 k}{n} + \frac{\sigma^2}{n} \log(1/\delta).$$

Matching rates for LASSO?

Assumption INC(k) We say that the design matrix \mathbb{X} has incoherence k for some integer $k > 0$ if

$$\left| \frac{\mathbb{X}^T \mathbb{X}}{n} - I_d \right|_{\infty} \leq \frac{1}{32k}$$

where the $|A|_{\infty}$ denotes the largest element of A in absolute value. Equivalently,

- ① For all $j = 1, \dots, d$,

$$\left| \frac{\|\mathbb{X}_j\|_2^2}{n} - 1 \right| \leq \frac{1}{32k}.$$

- ② For all $1 \leq i, j \leq d, i \neq j$, we have

$$\frac{|\mathbb{X}_i^T \mathbb{X}_j|}{n} \leq \frac{1}{32k}.$$

Technical Lemma

Lemma

Fix a positive integer $k \leq d$ and assume that \mathbb{X} satisfies assumption $\text{INC}(k)$. Then, for any $S \in \{1, \dots, d\}$ such that $|S| \leq k$ and any $\theta \in \mathbb{R}^d$ that satisfies the cone condition

$$|\theta_{S^c}|_1 \leq 3|\theta_S|_1,$$

it holds

$$|\theta|_2^2 \leq 2 \frac{|\mathbb{X}\theta|_2^2}{n}.$$

We will interpret the cone condition more carefully next week when we consider sparse recovery.

Theorem

Fix $n \geq 2$. Assume that the linear model (2.2) holds where $\varepsilon \sim \text{sub}G_n(\sigma^2)$. Moreover, assume that $\|\theta^*\|_0 \leq k$ and that X satisfies assumption $\text{INC}(k)$. Then the Lasso estimator $\hat{\theta}^{\mathcal{L}}$ with regularization parameter defined by

$$2\tau = 8\sigma\sqrt{\frac{\log(2d)}{n}} + 8\sigma\sqrt{\frac{\log(1/\delta)}{n}}$$

satisfies

$$\text{MSE}(X\hat{\theta}^{\mathcal{L}}) = \frac{1}{n}\|X\hat{\theta}^{\mathcal{L}} - X\theta^*\|_2^2 \lesssim k\sigma^2 \frac{\log(2d/\delta)}{n}.$$

and

$$\|\hat{\theta}^{\mathcal{L}} - \theta^*\|_2^2 \lesssim k\sigma^2 \frac{\log(2d/\delta)}{n}.$$

with probability at least $1 - \delta$.

Proof

Proof Cont.