

SDS 384 11: Theoretical Statistics Lecture 7: Talagrand's inequality

Purnamrita Sarkar Department of Statistics and Data Science The University of Texas at Austin

Theorem

Consider a convex function $f : \mathbb{R}^n \to \mathbb{R}$ with Lipschitz constant L. Also consider n iid random variables $X_1, \ldots, X_n \in [0, 1]$. We have for t > 0

$$\mathcal{P}(|f(X) - M_f| \ge t) \le 4 \exp\left(-\frac{t^2}{16L^2}\right)$$

where M_f is the median of f.

- $P(f(X) \ge M_f) \ge 1/2$ and $P(f(X) \le M_f) \ge 1/2$
- Often the median can be replaced by the mean with a little give in the *t*.

- Let *d* denote the Euclidean distance
- Define $A = \{x : f(x) \le M_f\}$
- Define $d(x, A) = \inf_{y \in A} d(x, y)$
- Define $A_t = \{x : d(x, A) \le t\}$
- Since f is 1 Lipschitz (WLOG), $x \in A_t \Rightarrow f(x) \le M_f + t$
- So $P(x \in A_t) \le P(f(x) \le M_f + t)$
- All we need is to upper bound $P(x \notin A_t)$
- Since f is convex, A is a convex set.

Theorem

Let $A \subset \mathbb{R}^n$ be a convex set and $X \sim Unif(\{0,1\}^n)$. Then,

$$P(X \in A)P(X \notin A_t) \leq e^{-t^2/16}$$

 This is basically saying that if A is convex and and P(x ∈ A) is large then At takes up most of the space in the unit hypercube for t ≫ 1.

Is convexity needed?

Example

Let $A := \{x \in \{0,1\}^n : \sum_{i=1}^n x_i \le n/2\}$. Consider a product measure such that $X_i \sim Bernoulli(1/2)$. Let $X = (X_1, \dots, X_n)$. Then $P(X \in A)$ is large. But is $P(X \notin A_t)$ large?

- Note that A is not convex.
- Also see that

$$|y^{T}1 - x^{T}1| \le ||y - x||_{1} = ||y - x||_{2}^{2}$$
$$\{y \in A_{t}\} \subseteq \{y^{T}1 \le n/2 + t^{2}\}$$
$$P(Y \notin A_{t}) \ge P(Y^{T}1 \ge n/2 + t^{2})$$

- Now P(X ∉ A_t), which is large for t ≈ (log n)^{1/4}, contrary to the result of Talagrand.
- What if we define A as a subset of Rⁿ?

- Now A is convex.
- Distance to A of a point with more than n/2 ones is simply its distance to the hyperplane $x^T 1 n/2 = 0$
- Consider a point y with n/2 + k ones.
- The distance to the previous nonconvex A is \sqrt{k}
- But distance to the convex A is $|y^T 1 n/2|/\sqrt{n} = k/\sqrt{n}$

$$\{y \in A_t^{(conv)}\} = \{y^T 1 - n/2 \le t\sqrt{n}\}$$
$$P(Y \notin A_t^{(conv)}) = P(Y^T 1 \ge n/2 + t\sqrt{n})$$

• Here, everything is fine since this is indeed large when $t \gg 1$

- First note that $E[(f(X) M_f)^2] \le CL^2$ by using Talagrand's inequality. (How?)
- Now note that $var(f(X)) \leq E[(f(X) M_f)^2] \leq CL^2$
- Finally $P(|f(X) E[f(X)]| \ge 2\sqrt{\operatorname{var}(f(X))}) \le 1/4$.
- So we must have $M_f \in [E[f(X)] \pm cL]$
- So, $P(|f(X) E[f(X)]| \ge cL + t) \le 4e^{-t^2/16L^2}$

Consider a random matrix $M = [X_{ij}] \in [a, b]^{n \times m}$ where X_{ij} are independent random variables.

$$P(||M||_{op} \ge E[||M||_{op}] + c\sqrt{\log n}) = o(1)$$

- For $E[X_{ij}] = 0$ and $var(X_{ij}) = \sigma^2$, it can be shown that $E[||M||_{op}] \le 2\sigma\sqrt{n}$.
- $||M||_{op}$ is 1 Lipschitz and convex. (how?)

Consider a random matrix $M = [X_{ij}] \in [a, b]^{n \times m}$ where X_{ij} are independent random variables.

$$P(||M||_{op} \ge E[||M||_{op}] + c\sqrt{\log n}) = o(1)$$

- For $E[X_{ij}] = 0$ and $var(X_{ij}) = \sigma^2$, it can be shown that $E[||M||_{op}] \le 2\sigma\sqrt{n}$.
- $||M||_{op}$ is 1 Lipschitz and convex. (how?)

Consider a iid sequence $X = \{X_i\}_{i=1}^n$. We will bound $f(X) := \sup_{a \in \mathcal{A}} a^T X$ where \mathcal{A} is a compact subset of \mathbb{R}^n such that $\mathcal{W} = \sup_{a \in \mathcal{A}} ||a||_2 < \infty$.

- Why cant we just use Chernoff?
- First let us check if f(X) is Lipschitz. Let a_* and a'_* be the maximizers of f(X) and f(X'). $f(X) - f(X') = a_*^T X - a'_*^T X' \le a_*^T (X - X')$ • $\le \sup_{a \in \mathcal{A}} a^T (X - X') \le \mathcal{W} ||X - X'||_2$
- How about convex? Consider the set $S_c = \{x : f(x) \le c\}$.
 - consider $x, y \in S_c$. Then

$$f(\lambda x + (1 - \lambda)y) \le f(\lambda x) + f((1 - \lambda)y) \le c$$

Consider a iid sequence $X = \{X_i\}_{i=1}^n$. We will bound $f(X) := \sup_{a \in \mathcal{A}} a^T X$ where \mathcal{A} is a compact subset of \mathbb{R}^n such that $\mathcal{W} = \sup_{a \in \mathcal{A}} ||a||_2 < \infty$.

- If $X_i \sim N(0,1)$ using Gaussian+Lipschtz $P(|f(X) - E[f(X)]| \ge t) \le 2e^{-\frac{t^2}{2W^2}}$
- If X_i are bounded, then Talagrand gives us the same thing (modulo constants).
- How about McDiarmid?

Complexity

Example

Consider a iid Rademacher sequence $X = \{X_i\}_{i=1}^n$. We will bound $f(X) := \sup_{a \in \mathcal{A}} a^T X$ where \mathcal{A} is a compact subset of \mathbb{R}^n such that $\mathcal{W} = \sup_{a \in \mathcal{A}} ||a||_2 < \infty$.

- Consider $X(k) \in [0,1]$
- Consider X and X' differing in the k-th coordinate, $f(X) - f(X') = a_*^T X - a_*'^T X' \le a_*^T (X - X')$ $\leq \sup_{a \in A} a_k(X(k) - X'(k)) \le \sup_{a \in A} |a_k|$
- So McDiarmid gives:

$$P(|f(X) - E[f(X)]| \ge t) \le 2 \exp(-\frac{t^2}{2\sum_i \sup_{a \in \mathcal{A}} |a_i|^2})$$