

Maximal inequalities, covering nets, norms of random matrices

Setting: Random process $\{X_t\}_{t \in T}$. A few examples:

- $\{g^T t\}_{t \in T}$, $g \sim N(0, I)$ a Gaussian vector
- $\{u^T X v\}_{\substack{u \in \mathbb{R}^n: \|u\|_2 \leq 1 \\ v \in \mathbb{R}^m: \|v\|_2 \leq 1}}, X \in \mathbb{R}^{n \times m}$ a random matrix
- $\left\{ \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X_i) \right\}_{f \in \mathcal{F}}, \mathcal{F} \text{ some function class}, X_1, \dots, X_n \text{ iid r.v.'s.}$

Goal: (1) Sharp upper bounds for $\mathbb{E} \sup_{t \in T} X_t$

(2) Sharp tail bounds for $\mathbb{P}[\sup_{t \in T} X_t \geq u]$

These questions are often equivalent if $\sup_{t \in T} X_t$ concentrates around its mean.

Principle: If $\{X_t\}_{t \in T}$ "sufficiently continuous," then size of $\sup_{t \in T} X_t$ is controlled by "complexity" of T .

Basic idea: If T has finite cardinality and $X_t \geq 0 \forall t \in T$,

$$\text{then } \mathbb{E} \sup_{t \in T} X_t \leq \sum_{t \in T} \mathbb{E} X_t.$$

Example: For any r.v.'s X_1, X_n :

$$\mathbb{E} \sup_{i=1}^n |X_i| \leq \sum_{i=1}^n \mathbb{E}|X_i| \leq n \cdot \sup_{i=1}^n \mathbb{E}|X_i|$$

$$\mathbb{E} \sup_{i=1}^n |X_i| \stackrel{\text{(Jensen's)}}{\leq} \left(\mathbb{E} \sup_{i=1}^n |X_i|^p \right)^{1/p} \leq n^{1/p} \sup_{i=1}^n (\mathbb{E}|X_i|^p)^{1/p}$$

Lemma (maximal inequality): Suppose $|T|$ is finite and

$\log \mathbb{E} e^{\lambda X_t} \leq \psi(\lambda)$ for all $\lambda \geq 0$, $t \in T$, where ψ is convex, $\psi(0) = \psi'(0) = 0$. Set $\psi^*(t) = \sup_{\lambda \geq 0} \lambda t - \psi(\lambda)$.

Then

$$\mathbb{E} \sup_{t \in T} X_t \leq \psi^{*-1}(\log |T|).$$

Cor: If $\mathbb{E} e^{\lambda X_t} \leq \frac{\lambda^2 \sigma^2}{2}$ for all $t \in T$ and $\lambda \geq 0$, then

$$\mathbb{E} \sup_{t \in T} X_t \leq \sqrt{2\sigma^2 \log |T|}$$

(Apply Lemma w/ $\psi(\lambda) = \frac{\lambda^2 \sigma^2}{2}$, $\psi^*(t) = \frac{t^2}{2\sigma^2}$.)

Proof: $\mathbb{E} \sup_{t \in T} X_t \stackrel{\text{(Jensen's)}}{\leq} \frac{1}{\lambda} \log \mathbb{E} e^{\lambda \sup_{t \in T} X_t} \quad \forall \lambda > 0$

$$= \frac{1}{\lambda} \log \mathbb{E} \sup_{t \in T} e^{\lambda X_t}$$

$$\leq \frac{1}{\lambda} \log \sum_{t \in T} \mathbb{E} e^{\lambda X_t}$$

$$\leq \frac{1}{\lambda} \log (|T| \cdot \sup_{t \in T} \mathbb{E} e^{\lambda X_t}) \leq \frac{\log |T| + \psi(\lambda)}{\lambda}.$$

$$\Rightarrow \log |T| \geq \sup_{\lambda \geq 0} (\lambda \mathbb{E} \sup_{t \in T} X_t - \psi(\lambda)) = \psi^*(\mathbb{E} \sup_{t \in T} X_t).$$

$$\Rightarrow \mathbb{E} \sup_{t \in T} X_t \leq \psi^{*-1}(\log |T|). \quad (\psi^* \text{ convex, increasing on } [0, \infty)).$$

Remark: This is closely related to the union bound

$$\mathbb{P}\left[\sup_{t \in T} X_t \geq u\right] \leq \sum_{t \in T} \mathbb{P}[X_t \geq u] \leq |T| \cdot e^{-\lambda u + \psi(\lambda)} \quad \forall \lambda \geq 0$$

$$\Rightarrow \mathbb{P}\left[\sup_{t \in T} X_t \geq u\right] \leq e^{\log |T| - \psi^*(u)}$$

$$\Rightarrow \mathbb{P}\left[\sup_{t \in T} X_t \geq \psi^{*-1}(\log |T| + s)\right] \leq e^{-s}$$

Lemma can be proven (up to a constant) by integrating this tail.

Intuition: If $\{X_t\}_{t \in T}$ are independent and $\sum_{t \in T} \mathbb{P}[X_t \geq u]$ is small,

$$\mathbb{P}\left[\sup_{t \in T} X_t \geq u\right] = 1 - \prod_{t \in T} (1 - \mathbb{P}[X_t \geq u])$$

$$\geq 1 - \prod_{t \in T} e^{-\mathbb{P}[X_t \geq u]} \approx \sum_{t \in T} \mathbb{P}[X_t \geq u]$$

so union bound is tight. Bound is loose when $\{X_t\}_{t \in T}$ are dependent.

Covering nets

Idea: Trade off dependence of $\{X_t\}_{t \in T}$ with the continuity

$$X_t \approx X_s \text{ for } t \approx s.$$

Def: Let (T, d) be a metric space. $N \subseteq T$ is an ε -net of T if, for all $t \in T$, there exists $\pi(t) \in N$ s.t. $d(t, \pi(t)) \leq \varepsilon$. The covering number is

$$N(T, d, \varepsilon) = \inf \{ |N| : N \text{ is } \varepsilon\text{-net of } T \}.$$

Prop: If $\{X_t\}_{t \in T}$ is L -Lipschitz, i.e. there exists (possibly random) $L \geq 0$ for which $|X_t - X_s| \leq L \cdot d(s, t)$ $\forall s, t \in T$, and by $\mathbb{E} e^{\lambda X_t} \leq \frac{\lambda^2}{2} + \mathbb{E} X_t^2$ $\forall \lambda \geq 0$ and $t \in T$, then

$$\mathbb{E} \sup_{t \in T} X_t \leq \inf_{\varepsilon > 0} \varepsilon \cdot \mathbb{E} L + \sqrt{\mathbb{E} X_t^2 \log N(T, d, \varepsilon)}$$

Proof: For any ε -net N ,

$$\begin{aligned} \sup_{t \in T} X_t &= \sup_{t \in T} (X_t - X_{\pi(t)}) + X_{\pi(t)} \\ &\leq L \cdot d(t, \pi(t)) + \sup_{s \in N} X_s. \end{aligned}$$

$$\text{Thus } \mathbb{E} \sup_{t \in T} X_t \leq \varepsilon \cdot \mathbb{E} L + \sqrt{2\varepsilon^2 \log |N|}$$

Take inf over $\varepsilon > 0$ and all ε -nets N . □

Remark: Sometimes $T \subseteq S$ a larger space, and it's simpler to construct an ε -net of T with points in S . Let

$$N^{ext}(T, d, \varepsilon) = \inf \{ |N| : N \subseteq S \text{ is } \varepsilon\text{-net of } T \}.$$

$$\text{Then } N^{ext}(T, d, \varepsilon) \leq N(T, d, \varepsilon) \leq N^{ext}(T, d, \varepsilon/2)$$

(See Homework 7).

Example: X_1, \dots, X_n iid r.v.'s on $[0, 1]$,

$$\tilde{\mathcal{F}} = \{ f : [0, 1] \rightarrow [0, 1], \text{ 1-Lipschitz} \}.$$

$$\text{Let } X_f = \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X_i), \quad W = \sup_{f \in \tilde{\mathcal{F}}} |X_f|$$

(1-Wasserstein distance between $\frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ and law of X_i)

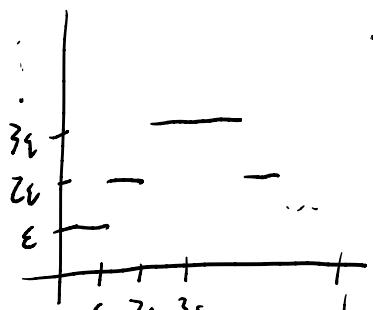
- $|X_f - X_g| \leq \frac{1}{n} \sum_{i=1}^n |f(X_i) - g(X_i)| \leq \frac{1}{n} \sum_{i=1}^n \|f - g\|_\infty$.

- Each $f(X_i) - \mathbb{E} f(X_i) \in [-\mathbb{E} f(X_i), \mathbb{E} f(X_i)]$, so

$$\log \mathbb{E} e^{\lambda X_f} \leq \frac{\lambda^2}{8n} \quad \forall \lambda \geq 0 \text{ by Hoeffding's lemma +}$$

Hoeffding's inequality

- Let $\mathcal{N} = \{\text{piecewise constant functions on } [0, 1] \text{ w/ jumps at } \varepsilon, 2\varepsilon, 3\varepsilon, \dots, \text{ values in } 0, \varepsilon, 2\varepsilon, \dots, \text{ where value changes by at most } \varepsilon \text{ at each jump}\}$



For any $f \in \mathcal{S}$: If $f(k\varepsilon) \in [j\varepsilon, (j+1)\varepsilon]$

then $\forall x \in [k\varepsilon, (k+1)\varepsilon]$,

$$|f(x) - j\varepsilon| \leq |f(x) - f(k\varepsilon)| + |f(k\varepsilon) - j\varepsilon| \leq 2\varepsilon.$$

Also, $f((k+1)\varepsilon) \in [(j-1)\varepsilon, j\varepsilon] \cup [j\varepsilon, (j+1)\varepsilon] \cup [(j+1)\varepsilon, (j+2)\varepsilon]$

$\Rightarrow \exists g \in \mathcal{N}$ (round value of f at left endpoint of each interval down to nearest $j\varepsilon$) s.t. $\|f-g\|_\infty \leq 2\varepsilon$.

$$|\mathcal{N}| \leq \left(\frac{1}{\varepsilon}\right) \cdot 3^{1/\varepsilon} \Rightarrow N(\mathcal{S}, \|\cdot\|_\infty, \varepsilon) \leq C^{1/\varepsilon} \quad \forall \varepsilon \in (0, 1).$$

$$\text{So, } \mathbb{E} W \leq \inf_{\varepsilon > 0} \left(2\varepsilon + \sqrt{\frac{1}{2n} \log C^{1/\varepsilon}} \right) \leq C n^{-1/3}.$$

In later lectures, we'll:

(i) Improve this bound to $C n^{-1/2}$, using the typical size

$$|X_f - X_g| \leq \frac{1}{\sqrt{n}} \|f-g\|_\infty \text{ w.h.p.}$$

(ii) Extend to non-Lipschitz classes \mathcal{F} for which $N(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)$ is infinite, by using other metrics for \mathcal{F} .

Norm of random matrices

Let $X \in \mathbb{R}^{n \times m}$. $\|X\|_{op} = \sup_{u \in B^n, v \in B^m} \underbrace{u^T X v}_{:= X_{u,v}}$,

$$B^n := \{u \in \mathbb{R}^n : \|u\|_2 \leq 1\},$$

Then for any $u, u' \in B^n, v, v' \in B^m$,

$$\begin{aligned} |X_{u,v} - X_{u',v'}| &= |u^T X v - u'^T X v'| \\ &\leq |u^T X(v-v')| + |(u-u')^T X v'| \\ &\leq \|X\|_{op} \cdot (\|v-v'\|_2 + \|u-u'\|_2) \end{aligned}$$

$\Rightarrow \{X_{u,v}\}_{u,v}$ is $\|X\|_{op}$ -Lipschitz w.r.t. the metric

$$d((u,v), (u',v')) = \|u-u'\|_2 + \|v-v'\|_2.$$

To bound $N(B^n \times B^m, d, \varepsilon)$:

Def: Let (T, d) be a metric space. $\mathcal{D} \subseteq T$ is an ε -packing of T if $d(s, t) > \varepsilon \quad \forall s \neq t \in \mathcal{D}$.

The packing number is

$$D(T, d, \varepsilon) = \sup \{|D| : D \text{ is } \varepsilon\text{-packing of } T\}.$$

Prop: For any $\varepsilon > 0$,

$$D(T, d, 2\varepsilon) \leq N(T, d, \varepsilon) \leq D(T, d, \varepsilon).$$

Proof: Let \mathcal{D} be any 2ε -packing, \mathcal{N} any ε -net.

For any $t \in \mathcal{D}$, $\exists \pi(t) \in \mathcal{N}$ s.t. $d(t, \pi(t)) \leq \varepsilon$.

Furthermore $\pi(t) \neq \pi(s)$ for $t \neq s \in \mathcal{D}$, because

$$d(\pi(t), \pi(s)) \geq d(t, s) - d(t, \pi(t)) - d(s, \pi(s)) > 0.$$

So $|\mathcal{D}| \leq |\mathcal{N}|$, implying $D(T, d, 2\varepsilon) \leq N(T, d, \varepsilon)$.

Let \mathcal{D} be any maximal ε -packing, i.e. $\mathcal{D} \cup \{t\}$ is not a ε -packing for any $t \in T \setminus \mathcal{D}$. Then for any $t \in T \setminus \mathcal{D}$, $\exists s \in \mathcal{D}$ s.t. $d(t, s) \leq \varepsilon$. So \mathcal{D} is also a ε -net, implying $N(T, d, \varepsilon) \leq |\mathcal{D}| \leq D(T, d, \varepsilon)$. □

Prop: Let $B^n = \{u \in \mathbb{R}^n : \|u\|_2 \leq 1\}$. Then for any $\varepsilon \in (0, 1)$,

$$\left(\frac{1}{\varepsilon}\right)^n \leq N(B^n, \|\cdot\|_2, \varepsilon) \leq \left(\frac{2}{\varepsilon} + 1\right)^n < \left(\frac{3}{\varepsilon}\right)^n$$

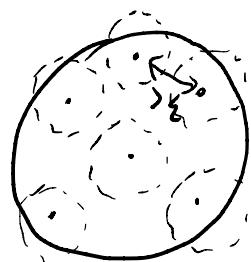
Proof: Let $\text{Vol}(\cdot)$ be the volume in \mathbb{R}^n . If \mathcal{N} is an ε -net,

$$\text{then } B^n = B(0, 1) \subseteq \bigcup_{t \in \mathcal{N}} B(t, \varepsilon)$$

$$\Rightarrow \text{Vol}(B(0,1)) \leq \sum_{t \in \mathcal{N}} \text{Vol}(B(t, \varepsilon)) \\ = |\mathcal{N}| \cdot \varepsilon^n \cdot \text{Vol}(B(0,1))$$

$$\Rightarrow |\mathcal{N}| \geq \left(\frac{1}{\varepsilon}\right)^n, \text{ implying } N(B, \|\cdot\|_2, \varepsilon) \geq \left(\frac{1}{\varepsilon}\right)^n.$$

If \mathcal{D} is an ε -packing, then $\{B(t, \frac{\varepsilon}{2}): t \in \mathcal{D}\}$ are disjoint balls contained in $B(0, 1 + \frac{\varepsilon}{2})$.



$$\Rightarrow \text{Vol}(B(0, 1 + \frac{\varepsilon}{2})) \geq \sum_{t \in \mathcal{D}} \text{Vol}(B(t, \frac{\varepsilon}{2})) \\ \Rightarrow (1 + \frac{\varepsilon}{2})^n \cdot \text{Vol}(B(0,1)) \geq |\mathcal{D}| \left(\frac{\varepsilon}{2}\right)^n \cdot \text{Vol}(B(0,1)) \\ \Rightarrow |\mathcal{D}| \leq \left(\frac{2}{\varepsilon} + 1\right)^n$$

$$\text{Thus } N(B, \|\cdot\|_2, \varepsilon) \leq D(B, \|\cdot\|_2, \varepsilon) \leq \left(\frac{2}{\varepsilon} + 1\right)^n. \quad \square$$

Thm: Suppose $X \in \mathbb{R}^{n \times m}$ has independent mean-zero σ^2 -subgaussian entries.

Then $\mathbb{E} \|X\|_{op} \leq C\sigma(\sqrt{n} + \sqrt{m})$ for a universal constant $C > 0$.

Proof: Recall $\|X\|_{op} = \sup_{u \in B^n, v \in B^m} |X_{u,v}| = \sup_{u \in B^n, v \in B^m} u^T X v$

where $|X_{u,v}|$ is $\|X\|_{op}$ -Lipschitz w.r.t. $d((u,v), (u',v')) = \|u-u'\|_2 + \|v-v'\|_2$.

If \mathcal{N}_u is an ε -net of B^n , \mathcal{N}_v an ε -net of B^m

(in $\|\cdot\|_2$), then $N_n \times N_r$ is 2ε -net of $B^n \times B^m$ in d .

$$\begin{aligned} \Rightarrow N(B^n \times B^m; d, \varepsilon) &\leq N(B; \|\cdot\|_2, \frac{\varepsilon}{2}) \cdot N(B^m; \|\cdot\|_2, \frac{\varepsilon}{2}) \\ &\leq \left(\frac{4}{\varepsilon} + 1\right)^{n+m} \end{aligned}$$

For each $(u, v) \in B^n \times B^m$, $u^T X v$ is σ^2 -subgaussian by Hoeffding's inequality. So

$$\mathbb{E} \|X\|_{op} \leq \underbrace{\varepsilon \cdot \mathbb{E} \|X\|_{op}}_{\text{Lipschitz constant of } \{X_{u,v}\}} + \sqrt{2\sigma^2 \log \left(\frac{4}{\varepsilon} + 1\right)^{n+m}}. \quad \forall \varepsilon > 0.$$

Pick $\varepsilon = 1/2$ and rearrange: $\mathbb{E} \|X\|_{op} \leq C\sigma(\sqrt{n} + \sqrt{m})$

Remark: $\|X\|_{op}$ is a convex, 1-Lipschitz function of its entries.

If the entries are bounded $X_{ij} \in [-C_0, C_0]$, or satisfy a LSI $\text{Ent } f(X_{ij})^2 \leq C\sigma^2 \mathbb{E} f'(X_{ij})^2$, then by concentration

$$P[\|X\|_{op} \geq \mathbb{E} \|X\|_{op} + \sigma \cdot t] \leq e^{-ct^2}.$$

Alternatively, assuming only X_{ij} are σ^2 -subgaussian, the above argument may be adapted into a union bound:

$$\begin{aligned}
& \mathbb{P}[\|X\|_{op} \geq \varepsilon \cdot \|X\|_{op} + C\sigma(\sqrt{n} + \sqrt{m}) + \sigma \cdot t] \\
& \leq \mathbb{P}[\exists (u,v) \in \mathcal{N}_u \times \mathcal{N}_v : X_{u,v} \geq C\sigma(\sqrt{n} + \sqrt{m}) + \sigma \cdot t] \\
& \leq \left(\frac{4}{\varepsilon} \cdot 1\right)^{n+m} e^{-\frac{1}{2}(C(\sqrt{n} + \sqrt{m}) + t)^2} \\
& \leq \left(\frac{4}{\varepsilon} \cdot 1\right)^{n+m} e^{-\frac{C^2}{2}(n+m) - \frac{t^2}{2}} \quad \forall t \geq 0
\end{aligned}$$

$$\text{Take } \varepsilon = 1/2, \quad C/2 = \log 9$$

$$\Rightarrow \mathbb{P}[\|X\|_{op} \geq 2C\sigma(\sqrt{n} + \sqrt{m}) + 2\sigma \cdot t] \leq e^{-t^2/2}.$$

Thm: Let $X \in \mathbb{R}^{n \times m}$ have independent rows $x_1, \dots, x_n \in \mathbb{R}^m$ s.t.
 $\mathbb{E}x_i = 0$, $\mathbb{E}x_i x_i^\top = I$, $u^\top x_i$ is σ^2 -subgaussian for every
(fixed) unit vector $u \in \mathbb{R}^m$. Let $s_{\min}(X)$ and
 $s_{\max}(X) = \|X\|_{op}$ be the smallest/largest singular values. Then
for a universal constant $C > 0$,

$$\begin{aligned}
& \mathbb{P}[\sqrt{n} - C\sigma^2(\sqrt{m} + t) \leq s_{\min}(X) \leq s_{\max}(X) \leq \sqrt{n} + C\sigma^2(\sqrt{m} + t)] \\
& \geq 1 - 2e^{-t^2}.
\end{aligned}$$

Interpretation: With probability $1 - e^{-cm}$, $\forall u \in \mathbb{R}^m$ w/ $\|u\|_2 = 1$:
 $\|Xu\|_2 \in \sqrt{n} \pm C\sqrt{m}$, i.e. $\frac{1}{\sqrt{n}}X$ is near isometry if $m \ll n$.

$$\text{Proof: Consider } \left\| \frac{1}{n} X^T X - I \right\|_{op} = \sup_{\substack{u \in \mathbb{R}^m \\ \|u\|_2=1}} |u^T (\frac{1}{n} X^T X - I) u|$$

$$= \sup_{\substack{u \in \mathbb{R}^m \\ \|u\|_2=1}} \underbrace{\left| \frac{1}{n} \|Xu\|_2^2 - 1 \right|}_{:= X_u}$$

X_u is $2 \left\| \frac{1}{n} X^T X - I \right\|_{op}$ -Lipschitz wrt $\|\cdot\|_2$. Let N be $\frac{1}{4}$ -net:

$$\left\| \frac{1}{n} X^T X - I \right\|_{op} \leq \frac{1}{2} \cdot \left\| \frac{1}{n} X^T X - I \right\|_{op} + \sup_{u \in N} X_u$$

$$\Rightarrow \left\| \frac{1}{n} X^T X - I \right\|_{op} \leq 2 \sup_{u \in N} X_u$$

$$\Rightarrow \mathbb{P}\left[\left\| \frac{1}{n} X^T X - I \right\|_{op} \geq \delta\right] \leq \mathbb{P}\left[2 \sup_{u \in N} X_u \geq \delta\right] \leq \underbrace{|N|}_{\leq q^m} \cdot \sup_{u \in N} \mathbb{P}[X_u \geq \frac{\delta}{2}]$$

Here: $\frac{1}{n} \|Xu\|_2^2 - 1 = \frac{1}{n} \sum_{i=1}^n (x_i^T u)^2 - \mathbb{E}(x_i^T u)^2$, where (by assumption)

$$\|x_i^T u\|_{\chi_2} \leq C\sigma \Rightarrow \|(x_i^T u)^2 - \mathbb{E}(x_i^T u)^2\|_{\chi_2} \leq C'\sigma^2$$

$$\Rightarrow \mathbb{P}\left[\left| \frac{1}{n} \|Xu\|_2^2 - 1 \right| \geq \frac{\delta}{2}\right] \leq 2e^{-c_1 n \cdot \min\left(\frac{\delta^2}{\sigma^4}, \frac{\delta}{\sigma^2}\right)}$$

(Bennet's inequality)

$$\text{Set } \delta = \sigma^2 \max(s, s^2)$$

$$\Rightarrow \mathbb{P}\left[\left\| \frac{1}{n} X^T X - I \right\|_{op} \geq \sigma^2 \max(s, s^2)\right] \leq q^m \cdot 2e^{-c_1 ns^2}$$

Note we must have $\sigma^2 \geq 1$, and for any unit vector $u \in \mathbb{R}^m$,

$$\left| \frac{1}{\sqrt{n}} \|Xu\|_2 - 1 \right| \geq \sigma^2 s \Rightarrow \left| \frac{1}{n} \|Xu\|_2^2 - 1 \right| \geq \max(\sigma^2 s^2, \sigma^4 s^2) \geq \sigma^2 \max(s, s^2)$$

(See Lecture 2.) Thus

$$\mathbb{P} \left[\sup_{u: \|u\|_2=1} \left| \frac{1}{\sqrt{n}} \|Xu\|_2 - 1 \right| \geq \sigma^2 s \right]$$

$$\leq \mathbb{P} \left[\left\| \frac{1}{n} X^T X - I \right\|_{op} \geq \sigma^2 \max(s, s^2) \right] \leq 9^m 2 e^{-c_1 n s^2}$$

$$\Rightarrow \mathbb{P} \left[\sqrt{n} - C \sigma^2 (\sqrt{n} + t) \leq s_{min}(X) \leq s_{max}(X) \leq \sqrt{n} + C \sigma^2 (\sqrt{n} + t) \right]$$

$$= \mathbb{P} \left[\sup_{u: \|u\|_2=1} \left| \frac{1}{\sqrt{n}} \|Xu\|_2 - 1 \right| \geq \frac{C \sigma^2}{\sqrt{n}} (\sqrt{n} + t) \right] \leq 2 e^{-t^2}$$

for $C > 0$ large enough.

□