## Q1 : Sparse principal component analysis

(a) Under $H_1$, $A = \lambda u v^T + B$,

$$P\left[ \text{val}(A) > \lambda - \frac{t}{\sqrt{d}} \right]$$

$$= P\left[ \sup_{\substack{x,y \in B_2^d \\ \|x\|_0 = \|y\|_0 = k}} x^T A y > \lambda - \frac{t}{\sqrt{d}} \right]$$

$$\leq P\left[ \sup u^T u v^T v \cdot \lambda + u^T B v > \lambda - \frac{t}{\sqrt{d}} \right]$$

$$\leq P\left[ \lambda + u^T B v > \lambda - \frac{t}{\sqrt{d}} \right]$$

$$= P\left[ u^T (B\sqrt{d}) v > -t \right]$$

$$= 1 - P\left[ u^T (B\sqrt{d}) v \leq -t \right]$$

But $u^T \sqrt{d} B v = \sum_{i,j=1}^{d} u_i v_j \sqrt{d} B_{ij}$ and

$u_i v_j \sqrt{d} B_{ij}$ is $|u_i v_j|$-subgaussian by HW 6, 5a)

So $u^T \sqrt{d} B v$ is the sum of iid $|u_i v_j|$-subgaussian RVs
and hence it is a $\sum_{i,j=1}^{d} u_i^2 v_j^2 = 1$-subgaussian RV, with mean zero.

Hence $\mathbb{P}\left[u^T B v \le -t\right] \le e^{-t^2/2}$ and hence

$$\mathbb{P}\left[u^T B v > -t\right] \ge 1 - e^{-t^2/2}.$$

(b) WTS $\exists\, c < 10$ s.t. if $A = B$,

$$\mathbb{P}\left[\text{val}(A) > c\sqrt{\frac{k}{d}\left(1 + \log\frac{d}{k}\right) + \frac{1}{d}\log\left(\frac{1}{r}\right)'}\right] \le d.$$

Let $\tilde{B}^2 = B_0^d(k) \wedge B_1^d(1)$. We have $\text{val}(A) = \sup\limits_{(u,y) \in \tilde{B}^2} u^T A y$.

Consider $N$ an $\varepsilon$-net of $\tilde{B}$ for the $2$-norm, s.t. $\forall\, x \in \tilde{B}$, $\exists\, y \in N$ s.t. $\|x - y\| \le \varepsilon$ and $\text{supp}(u) = \text{supp}(y)$, and s.t. $N$ is minimal. Then

$$|N| \le \binom{d}{k}\left(1 + \frac{2}{\varepsilon}\right)^k$$

since there are maximum $\binom{d}{k}$ combinations of $k$ points from $d$-dimensional vectors and using the upper bound from lectures on the size of an $\varepsilon$-Net (minimal).

$$|N| \le \left(\frac{ed}{k}\right)^k \left(1 + \frac{2}{\varepsilon}\right)^k \qquad (\text{stirling's})$$

Then if we take $(x, y) \in \tilde{B}_1^2$ such that $x^T A y = \text{val}(A)$, and $(\tilde{x}, \tilde{y}) \in N^2$ the corresponding vectors in $N$, then

$$x^T A y - \tilde{x}^T A \tilde{y} = (x - \tilde{x}) A y + \tilde{x}^T A (y - \tilde{y})$$

$$= \|x - \tilde{x}\|_2 \frac{(x - \tilde{x})^T}{\|x - \tilde{x}\|_2} A y + \tilde{x}^T A \frac{(y - \tilde{y})}{\|y - \tilde{y}\|_2} \|y - \tilde{y}\|_2$$

$$\leq \|x - \tilde{x}\|_2 \, \text{val}(A) + \text{val}(A) \|y - \tilde{y}\|_2$$

$$\leq 2 \varepsilon \, x^T A y$$

Then $\mathbb{P}\left[\text{val}(A) > t\right] \leq \mathbb{P}\left[\max_{\tilde{x}, \tilde{y} \in N} \tilde{x}^T A \tilde{y} \geq t(1 - 2\varepsilon)\right]$

$$\leq |N|^2 \sup_{\tilde{x}, \tilde{y} \in N} \mathbb{P}\left[\underbrace{\tilde{x}^T B \tilde{y}}_{} \geq t(1 - 2\varepsilon)\right] \quad (\text{union bound})$$

$$\text{(d - subgaussian by (a))}$$

$$\leq |N|^2 \sup_{\tilde{x}, \tilde{y} \in N} \exp\left\{-t^2 (1 - 2\varepsilon)^2 d / 2\right\}$$

$$\leq \left(\frac{ed}{k}\right)^{2k} \left(1 + \frac{2}{\varepsilon}\right)^{2k} \exp\left\{-t^2 (1 - 2\varepsilon)^2 \frac{d}{2}\right\}$$

Take $\varepsilon = \frac{1}{4}$ Then

$$\mathbb{P}_{H_0}\left[\text{val}(A) > t\right] \leq \left[\left(\frac{ed}{k}\right)^{2k} \cdot 9^{2k} \cdot \exp\left\{\frac{-t^2}{8} d\right\}\right]$$

$$\leq \left[\left(\frac{9ed}{k}\right) \exp\left\{\frac{-t^2}{16} \frac{d}{k}\right\}\right]^{2k}$$

To have $P_{H_0}\left[\text{val}(A) > t\right] \le \delta$, we need $\delta \ge \left(\left(\frac{9ed}{k}\right) \cdot \exp\left\{-\frac{t^2}{16}\frac{d}{k}\right\}\right)^{2k}$ i.e.

It is sufficient to choose

$$t^2 \ge 16\left(\frac{k}{d}\right)\left(\log 9e + \log\frac{d}{k} + \frac{1}{2k}\log\frac{1}{\delta}\right)$$

So taking $t(k,d) = 4\sqrt{\log 9e}\sqrt{\frac{k}{d}\left(1+\log\frac{d}{k}\right) + \frac{1}{d}\log\frac{1}{\delta}}$

is enough, hence we get the result with

$$c = 4\sqrt{\log 2e} \le 7.2 < 10$$

(c) For a fixed threshold $\tau = O(1)$, we want to have non-negligible power and level, power > level.

i.e. $\begin{cases} \tau = 1 - \frac{t}{\sqrt{d}} = c\sqrt{\frac{k}{d}\left(1+\log\frac{k}{d}\right) + \frac{1}{d}\log\left(\frac{1}{\delta}\right)} \\[2mm] 1 - e^{-t^2/2} > \delta \end{cases}$

$1 - e^{-t^2/2} > \delta$ gives us $\log(1-\delta) > -\frac{t^2}{2}$ but $\log(1-\delta) \simeq -\delta$ for small $\delta$ so $t \ge \sqrt{2\delta}$. Plugging into the first equation gives

$$1 - \sqrt{\frac{2\delta}{d}} > c\sqrt{\frac{k}{d}\left(1+\log\left(\frac{k}{d}\right)\right) + \frac{1}{d}\log\frac{1}{\delta}}$$

## Q2: Rademacher Complexity

(a) Let $\bar{\mathcal{F}} := \{ f - Pf \mid f \in \mathcal{F}\}$.

WTS $\frac{1}{2} R_n(\bar{\mathcal{F}}) \le \mathbb{E}\| P_n - P \|_{\mathcal{F}}$

$$\frac{1}{2} \mathbb{E}_{X,\varepsilon} \sup_{f \in \bar{\mathcal{F}}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(X_i) \right|$$

$$= \frac{1}{2} \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(X_i) - \mathbb{E}_{Y \sim P} f(Y) \right) \right|$$

$$\le \frac{1}{2} \mathbb{E}_{X,Y,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(X_i) - f(Y_i) \right) \right| \qquad \text{Jensen's inequality}$$

$$\le \frac{1}{2} \mathbb{E}_{X,Y} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} f(X_i) - f(Y_i) \right| \qquad \text{symmetrization}$$

$$= \frac{1}{2} \mathbb{E}_{X,Y} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} f(X_i) - \mathbb{E}f + \mathbb{E}f - f(Y_i) \right|$$

$$\le \frac{1}{2} \mathbb{E}_{X,Y} \sup_{f \in \mathcal{F}} \left( \left| \frac{1}{n} \sum_{i=1}^{n} f(X_i) - \mathbb{E}f \right| + \left| \frac{1}{n} \sum_{i=1}^{n} f(Y_i) - \mathbb{E}f \right| \right)$$

$$\le \mathbb{E}_{X} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} f(X_i) - \mathbb{E}f \right| \qquad \text{due to } X, Y \text{ iid}$$

$$= \mathbb{E}_{X} \| P_n - P \|_{\mathcal{F}}$$

(b) WTS: $R_n(\bar{\mathcal{F}}) \ge R_n(\mathcal{F}) - \frac{1}{\sqrt{n}} \| P \|_{\mathcal{F}}$

$$\mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(X_i) - \mathbb{E}f(X_i) \right) \right| \ge \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(X_i) \right| -$$

$$- \frac{1}{\sqrt{n}} \mathbb{E}_{X} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} f(X_i) - \mathbb{E}f \right|$$

By expanding $R_n(\mathcal{F})$,

$$R_n(\mathcal{F}) = \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(x_i) \right|$$

$$= \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(x_i) - \mathbb{E}f + \mathbb{E}f \right) \right|$$

$$\leq \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left( \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(x_i) - \mathbb{E}f \right) \right| + \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \mathbb{E}f \right| \right)$$

$$= \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(x_i) - \mathbb{E}(f) \right) \right| + \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \mathbb{E}f \right|$$

$$= R_n(\bar{\mathcal{F}}) + \mathbb{E}_{\varepsilon} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \right| \cdot \sup_{f \in \mathcal{F}} |\mathbb{E}f|$$

Since $\varepsilon_i \sim \text{unif}(\{\pm 1\})$,

$$\mathbb{E}_\varepsilon \left| \sum_{i=1}^{n} \varepsilon_i \right| = \mathbb{E}_\varepsilon \sqrt{\overline{\left( \sum_{i=1}^{n} \varepsilon_i \right)^2}} = \mathbb{E}_\varepsilon \sqrt{\overline{n + 2 \sum_{1 \leq i < j \leq n} \varepsilon_i \varepsilon_j}} \overset{(*)}{\leq} \sqrt{\overline{n + \mathbb{E}_\varepsilon \sum_{1 \leq i < j \leq n} \varepsilon_i \varepsilon_j}} = \sqrt{n}$$

$(*)$ follows from Jensen's inequality

Hence $\quad R_n(\mathcal{F}) \leq R_n(\bar{\mathcal{F}}) + \dfrac{\sup_{f \in \mathcal{F}} |\mathbb{E}f|}{\sqrt{n}} = R_n(\bar{\mathcal{F}}) + \dfrac{1}{\sqrt{n}} \|P\|_{\mathcal{F}}$

(c)   From (a) and (b)

$$\mathbb{E}_X \| P_n - P \|_{\widetilde{\mathcal{F}}} \geq \frac{1}{2} R_n(\widetilde{\mathcal{F}}) - \frac{1}{2\sqrt{n}} \| P \|_{\widetilde{\mathcal{F}}}$$

WTS (sufficient to show)

$$\| P_n - P \|_{\mathcal{F}} \geq \mathbb{E}_X \| P_n - P \|_{\mathcal{F}} - \delta \quad \text{w.p. at least} \quad 1 - e^{-n\delta^2/2b^2}$$

Now, if we write $y = \| P_n - P \|_{\mathcal{F}}$, we have

$$\mathbb{P} \left[ y - \mathbb{E} y > -\delta \right] = \mathbb{P} \left[ \mathbb{E} y - y \leq \delta \right] = 1 - \mathbb{P} \left[ \mathbb{E} y - y > \delta \right]$$

Note $y$ has bounded difference property with parameter $L = \frac{2b}{n}$ (shown in class)

Hence $\mathbb{P} \left[ \mathbb{E} y - y > \delta \right] \leq \exp \left( \frac{-n\delta^2}{2b^2} \right)$

This is what we wanted to show.

## Q3.

First we show $\quad \mathbb{E}_x \sup_{f \in F} \pm (P_n f - Pf) \leq 2 \tilde{R}_n(F)$

$$\mathbb{E}_x \sup_{f \in F} \pm (P_n f - Pf) = \mathbb{E}_x \sup_{f \in F} \left[ \frac{1}{n} \sum_{i=1}^{n} \pm \left( f(X_i) - \mathbb{E}_x f(X_i) \right) \right]$$

$$= \mathbb{E}_x \sup_{f \in F} \mathbb{E}_y \left[ \frac{1}{n} \sum_{i=1}^{n} \pm \left( f(X_i) - f(Y_i) \right) \right]$$

because $\mathbb{E}_\theta \sup_\theta \ell(\theta) \geq \sup_\theta \mathbb{E}_\theta \ell(\theta)$
$$\leq \mathbb{E}_{x,y} \left[ \sup_{f \in F} \frac{1}{n} \sum_{i=1}^{n} \pm \left( f(X_i) - f(Y_i) \right) \right]$$

symmetrization, $\varepsilon_i, -\varepsilon_i$ same distribution.
$$= \mathbb{E}_{x,y,\varepsilon} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \left( f(X_i) - f(Y_i) \right)$$

and $\varepsilon_i \left( f(X_i) - f(Y_i) \right) \overset{d}{=} \left( f(X_i) - f(Y_i) \right)$
$$\leq \mathbb{E}_{x,y,\varepsilon} \sup_{f \in F} \left( \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(X_i) \right) + \sup_{f \in F} \left( \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(Y_i) \right)$$

$$= \mathbb{E}_{x,\varepsilon} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(X_i) + \mathbb{E}_{y,\varepsilon} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(Y_i)$$

$$= 2 \tilde{R}_n(F).$$

So $\mathbb{E}_x \sup_{f \in F} P_n f - Pf \leq 2 \tilde{R}_n(F)$ and $\mathbb{E}_x \sup_{f \in F} - P_n f + Pf \leq 2 \tilde{R}_n(F)$

Now we show the second part.

Consider $\tilde{g}^z(X_1, \ldots X_n) = \sup\limits_{f \in \mathcal{F}} \frac{1}{n} \sum\limits_{i=1}^{n} \bar{f}(x_i) = \sup\limits_{f \in \mathcal{F}} \frac{1}{n} \sum\limits_{i=1}^{n} f(x_i) - E\left[f(V_i)\right]$

Define $X^{(i)}$ as the same vector as $X$, but index $i$ is changed.

Take $f \in \mathcal{F}$ s.t. $\frac{1}{n} \sum\limits_{i=1}^{n} f(X_i) \geq \tilde{g}(X) - \varepsilon$ , $\varepsilon > 0$

$\tilde{g}(X) - \tilde{g}(X^{(i)}) \leq \frac{1}{n} \sum\limits_{i=1}^{n} \bar{f}(X_i) + \varepsilon - \sup\limits_{f \in \mathcal{F}} \left( \frac{1}{n} \sum\limits_{i=2}^{n} \bar{f}(X_i^{(i)}) + \frac{1}{n} \bar{f}(X_1^{(i)}) \right)$

$\leq \frac{1}{n} \sum\limits_{i=1}^{n} \bar{f}(X_i) + \varepsilon - \left( \frac{1}{n} \sum\limits_{i=2}^{n} \bar{f}(X_i^{(i)}) + \frac{1}{n} \bar{f}(V_1^{(i)}) \right)$

$= \frac{1}{n} \left( \bar{f}(X_i) - \bar{f}(V_i^{(i)}) \right) + \varepsilon \leq \frac{2b}{n} + \varepsilon$

Similarly, taking $f \in \mathcal{F}$ s.t. $\hat{h}(X) - \varepsilon \leq \frac{1}{n} \sum\limits_{i=1}^{n} -\bar{f}(X_i)$, we have

$\hat{h}(X) - \hat{h}(X^{(i)}) \leq \frac{1}{n} \sum\limits_{i=1}^{n} -\bar{f}(X_i) + \varepsilon - \sup\limits_{f \in \mathcal{F}} \frac{1}{n} \sum\limits_{i=1}^{n} -\bar{f}(X_i^{(i)})$

$\leq \frac{1}{n} \sum\limits_{i=1}^{n} -\bar{f}(X_i) + \varepsilon - \frac{1}{n} \sum\limits_{i=1}^{n} -\bar{f}(X_i^{(i)})$

$= \frac{1}{n} \sum\limits_{i=1}^{n} -\bar{f}(X_i) + \varepsilon - \frac{1}{n} \left( \sum\limits_{i=2}^{n} -\bar{f}(X_i^{(i)}) + \bar{f}(X_1^{(i)}) \right)$

$= \frac{1}{n} \left( \bar{f}(X_i^{(i)}) - \bar{f}(X_i) \right) + \varepsilon \leq \frac{2b}{n} + \varepsilon$

Taking $\varepsilon \downarrow 0$, $\tilde{g}$ and $\hat{h}$ both satisfy the bounded differences property with parameter $\frac{2b}{n}$.

By the bounded differences inequality, with $\mathbb{E}[\tilde{g}(x)] = \mathbb{E}[\tilde{h}(x)] = 0$,

$$\mathbb{P}\left[\tilde{g}(x) > \mathbb{E}[\tilde{g}(x)] + \delta\right] = \mathbb{P}\left[\tilde{g}(x) > \delta\right] \leq \exp\left(-\frac{\delta^2 n}{2b^2}\right)$$

$$\mathbb{P}\left[\tilde{h}(x) > \mathbb{E}[\tilde{h}(x)] + \delta\right] = \mathbb{P}\left[\tilde{h}(x) \geq \delta\right] \leq \exp\left(-\frac{\delta^2 n}{2b^2}\right)$$

Hence $\sup\limits_{f \in F} (P_n f - Pf) \leq \delta$ with probability $\geq 1 - e^{-\delta^2 n/2b^2}$

and $\sup\limits_{f \in F} (Pf - P_n f) \leq \delta$ with probability $\geq 1 - e^{-\delta^2 n/2b^2}$

We work: let $m(f) := \sup\limits_{f \in F} P_n f - Pf$, $\quad m^-(f) = \sup\limits_{f \in F} Pf - P_n f$.

$$m(f) = \underbrace{m(f) - \mathbb{E}_x\, m(f)}_{\leq \delta \text{ wp at least } 1 - \exp\left(\frac{-\delta^2 n}{2b^2}\right)} + \underbrace{\mathbb{E}_x\, m(f)}_{\leq 2\hat{R}_n(F)}$$

Similarly, $$m^-(f) = \underbrace{m^-(f) - \mathbb{E}_x\, m^-(f)}_{\leq \delta \text{ wp at least } 1 - \exp\left(\frac{-\delta^2 n}{2b^2}\right)} + \underbrace{\mathbb{E}_x\, m^-(f)}_{\leq 2\hat{R}_n(F)}$$

Hence $m(f) \leq 2\hat{R}_n(F) + \delta$ w prob at least $1 - e^{-\delta^2 n/2b^2}$

$m^-(f) \leq 2\hat{R}_n(F) + \delta$ w' prob at least $1 - e^{-\delta^2 n/2b^2}$

Now
$$\sup_{f \in F} |Pf - P_n f| = \max\left( m(f), m^-(f) \right)$$

$$\Rightarrow \quad P\left( \sup_{f \in F} |P_n f - Pf| \geq 2\tilde{R}_n(F) + \delta \right)$$

$$= P\left( \{ m(f) > 2\tilde{R}_n(F) + \delta \} \cup \{ m^-(f) \quad 2\tilde{R}_n(F) + \delta \} \right)$$

$$\leq P\left( m(f) > 2\tilde{R}_n(F) + \delta \right) + P\left( m^-(f) > 2\tilde{R}_n(F) + \delta \right)$$

$$\leq 2\exp\left( \frac{-n\delta^2}{2b^2} \right).$$

Hence
$$\| P_n f - Pf \|_F \leq 2\tilde{R}_n(F) + \delta \quad w.p.$$

at least
$$1 - 2\exp\left( \frac{-n\delta^2}{2b^2} \right).$$

## Q4. Binary Regression

(a) $\mathcal{F}_{sgn} = \{ f_\theta(x) = sgn(\langle x, \theta \rangle) \mid \|\theta\| = 1, \theta \in \mathbb{R}^d \}$

Suppose $X = x_1, \cdots, x_n \in \mathbb{R}^d$ linly independent

WTS: $R(\mathcal{F}_{sgn}, X) = 1$

$$R(\mathcal{F}_{sgn}, X) = \mathbb{E}_\varepsilon \sup_{f \in \mathcal{F}_{sgn}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(x_i) \right|$$

Note $\frac{1}{n} \sum_{i=1}^{n} \varepsilon_i sgn(\langle x_i, \theta \rangle) \leq \frac{1}{n} \sum_{i=1}^{n} |\varepsilon_i| = \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i^2 = 1$

$\rightarrow \sup_{\substack{\theta \in \mathbb{R}^d \\ \|\theta\|=1}} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i sgn(\langle x_i, \theta \rangle) \leq 1$

Since $x_i \perp\!\!\!\perp x_j \; \forall i \neq j$, and $span\{x_i\}_{i=1}^{n} = \mathbb{R}^n$, $\exists \theta \in \mathbb{R}^d$ s.t.

$sgn(\langle x_i, \theta \rangle) = \varepsilon_i \; \forall i \in [n]$. (take $\theta = \frac{\sum_{i=1}^{n} x_i \varepsilon_i}{\|\sum_{i=1}^{n} x_i \varepsilon_i\|_2} \in \mathbb{R}^d$, nw. $\|\theta\|_2 = 1$ and

$\langle x_i, \theta \rangle = \langle x_i, \frac{\sum_{i=1}^{n} x_i \varepsilon_i}{\|\sum_{i=1}^{n} x_i \varepsilon_i\|_2} \rangle = \frac{\|x_i\|_2^2 \varepsilon_i}{\|\sum_{i=1}^{n} x_i \varepsilon_i\|_2} \Rightarrow sgn(\langle x_i, \theta \rangle) = sgn(\varepsilon_i)$

$\Rightarrow \sup_{\substack{\theta \in \mathbb{R}^d \\ \|\theta\|=1}} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i sgn(\langle x_i, \theta \rangle) = 1$, i.e.

$\sup_{f \in \mathcal{F}_{sgn}} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(x_i) \right| = \sup_{\theta \in \mathbb{R}^d} \left| \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(x_i) \right| = \frac{1}{n} \cdot n = 1$

and hence $R(F_{sgn}, x) = 1$.


This is a problem because it implies overfitting.
For $d \geq n$, and $x_i$'s linearly independent, we showed $R_n(F) = 1$.
So for the boolean loss function, we can choose the
joint distribution of $(X, Y)$ and $F$ such that $R_n(\ell \circ F) = O(1)$
which means $\exists$ a fixed proportion of the points and
$\|P_n - P\|_F$ won't decrease as $n \to \infty$ (still having $d \geq n$).


This is problematic in the context of binary regression, eg with
$X_i \sim N(0, \Sigma)$ and $\ell(z, y) = \frac{1}{2}(1 - zy) = \mathbb{1}\{y \neq z\}$
because we will likely encounter overfitting as the function will
perfectly fit random noise in a high-dimensional setting with
independent covariates.

(b) WTS it $h: \mathbb{R} \to \mathbb{R}$ is $\gamma$-Lipschitz, $\hat{F}$ any function class,

$h \circ \hat{F} = \{ h \circ f \mid f(-\hat{F}) \}$, then $R_n(h \circ \hat{F}) \leq \gamma R_n(\hat{F})$

Consider $\phi$ a 1-Lipschitz function, $T \subseteq \mathbb{R}^2$ bounded.

Take $t^1, t^2 \in T$.

$$t_1^1 + \phi(t_2^1) + t_1^2 + \phi(t_2^2) \leq t_1^1 + t_1^2 + |t_2^1 - t_2^2|$$
$$\leq \max\{ t_1^1 + t_1^2 + t_2^1 - t_2^2, \; t_1^1 - t_1^2 + t_2^1 - t_2^2 \}$$

wait, re-read: $\leq \max\{ t_1^1 + t_1^2 + t_2^1 - t_2^2, \; t_1^1 - t_1^2 + t_2^1 - t_2^2\}$

$$\leq \sup_{t \in T}(t_1 + t_2) + \sup_{t \in T}(t_1 + t_2)$$

Hence $\sup_{t^1, t^2 \in T} t_1^1 + \phi(t_2^1) + t_1^2 - \phi(t_2^2) \leq \sup_{t \in T}(t_1 + t_2) + \sup_{t \in T}(t_1 - t_2)$

i.e. $\sup_{t \in T}(t_1 + \phi(t_2)) + \sup_{t \in T}(t_1 - \phi(t_2)) \leq \sup_{t \in T}(t_1 + t_2) + \sup_{t \in T}(t_1 - t_2)$

Now let us consider $\gamma = 1$ and show by induction on $0 \leq j \leq n$,

$$\tilde{R}_n(h \circ f) \leq \mathbb{E}_{x, y, \varepsilon} \sup_{f \in \hat{F}} \frac{1}{n} \sum_{i=1}^{n-j} \varepsilon_i \, h(f(x_i), y_i) + \sum_{k=n-j+1}^{n} \varepsilon_k \cdot f(x_k)$$

But can $n = 0$ : Nothing to prove

$1 \leq j \leq n$: By induction hypothesis

$$\hat{R}_n(h \circ \hat{F}) \leq \mathbb{E}_{x, y, \varepsilon} \left[ \sup_{f \in \hat{P}} \frac{1}{n} \left( \sum_{i=1}^{n-j+1} \varepsilon_i \, h(f(x_i), y_i) + \sum_{k=n-j+2}^{n} \varepsilon_k \, f(x_k) \right) \right]$$

Let us define $t_{x,\varepsilon}^d := \sum_{i=1}^{n-j} \varepsilon_i \, h(f(x_i), y_i) + \sum_{k=n-j+1}^{n} \varepsilon_k \, f(x_k)$.

$t_{x,\varepsilon}^d \subseteq T$, $T$ bounded. Then

$$\tilde{R}_n(h \cdot \mathcal{F}) \leq \frac{1}{2n} \, \mathbb{E}_{x,y_i} \, \mathbb{E}_{\varepsilon / \{\varepsilon_{n-j+1}\}} \left[ \sup_{f \in \mathcal{F}} t_{x,\varepsilon}^d + h(f(x_{n-j+1}), y_{n-j+1}) \right.$$

all $\varepsilon$'s except $n-j+1$

$$\left. + \sup_{f \in \mathcal{F}} \left( t_{x,\varepsilon}^d - h(f(x_{n-j+1}), y_{n-j+1}) \right) \right)$$

By the hint, with $T = \{ t_{x,\varepsilon}^d \mid f \in \mathcal{F} \} \times \{ f(x_{n-j+1}) \mid f \in \mathcal{F} \} \subseteq \mathbb{R}^2$ bounded.

$$\leq \frac{1}{2n} \, \mathbb{E}_{x,y} \, \mathbb{E}_{\varepsilon \backslash \{\varepsilon_{n-j+1}\}} \left[ \sup_{f \in \mathcal{F}} \left( t_{x,\varepsilon}^d + f(x_{n-j+1}) \right) + \sup_{f \in \mathcal{F}} \left( t_{x,\varepsilon}^d - f(x_{n-j+1}) \right) \right]$$

$$= \frac{1}{n} \, \mathbb{E}_{x,y} \, \mathbb{E}_{\varepsilon} \sup_{f \in \mathcal{F}} \left[ t_{x,\varepsilon}^d + \varepsilon_j \, f(x_{n-j+1}) \right]$$

$$= \mathbb{E}_{x,y,\varepsilon} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^{n-j} \varepsilon_i \, h(f(x_i), y_i) + \sum_{k=n-j+1}^{n} \varepsilon_k \, f(x_k) \right]$$

which is what we wanted to prove. Now with $j = n$,

$$\tilde{R}_n(h \cdot \mathcal{F}) \leq \mathbb{E}_{x,y,\varepsilon} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{k=1}^{n} \varepsilon_k \, f(x_k) \right] = \tilde{R}_n(\mathcal{F}).$$

This shows for $h$ 1-Lipschitz in its first argument, $\hat{R}_n(h \circ F) \leq \tilde{R}_n(\mathcal{F})$.

Now let $h$ be $\gamma$-Lipschitz, $\gamma \geq 1$, in its first argument.

Take $\tilde{\phi}_h(f(x)) = \frac{1}{\gamma} \phi(f(x))$.

where $\phi_h(f(x)) = h(f(x), y)$ with $h$ $\gamma$-Lipschitz in first argument

Note $\tilde{\phi}_h(f(x))$ is 1-Lipschitz, and applying previous argument, $\forall n \in \mathbb{N}$

$$\mathop{\mathbb{E}}_{x,y,\varepsilon} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \phi(f(x_i)) = \gamma \mathop{\mathbb{E}}_{x,y,\varepsilon} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \tilde{\phi}(f(x_i))$$

$$\leq \gamma \mathop{\mathbb{E}}_{x,y,\varepsilon} \sup_{f \in \hat{F}} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(x_i)$$

Hence $\hat{R}_n(h \circ F) \leq \gamma \tilde{R}_n(\mathcal{F})$.

This is what we wanted to show

(c) Suppose $\ell: X \times Y \to Y$, $\ell((x,y)) \in [-b, b]$, $\ell$ is $\gamma$-Lipschitz in $z$
defined by $\ell_\theta((x,y)) = \ell(f_\theta(x), y)$.

Consider $\mathcal{A} := \{ \ell(f_\theta(x), y) \mid f_\theta \in \mathcal{F} \}$.

WTS: $\| L_n - L \|_{\mathcal{A}} \leq 2\gamma R_n(\mathcal{F})$, with

$$L_n = P_n \ell_\theta(X, Y), \quad L = \mathbb{E} \ell_\theta(X, Y)$$

From theorem proved in lecture, since $\mathcal{A}$ is class of $b$-bounded functions,
$\forall \delta \geq 0$, $n \geq 1$, we have $\| L_n - L \|_{\mathcal{F}} \leq 2 R_n(\mathcal{A}) + \delta$.

with probability at least $1 - \exp\left( \dfrac{-\delta^2 n}{2b^2} \right)$

Note by Question 3, we know $\| L_n - L \|_{\mathcal{F}} \leq 2 \tilde{R}_n(\mathcal{A}) + \delta$.

From above, with $\ell \circ \mathcal{F} := \{ \ell(f_\theta(x), Y) : f_\theta \in \mathcal{F} \} = \mathcal{A}$,

we have $\tilde{R}_n(\mathcal{A}) \leq \gamma \cdot \tilde{R}_n(\mathcal{F})$.

Hence $\| L_n - L \|_{\mathcal{A}} \leq 2 \gamma \cdot \tilde{R}_n(\mathcal{F})$ with probability $\geq 1 - \exp\left( \dfrac{-\delta^2 n}{2b^2} \right)$

(d) $\ell_c(z,y) = \max\left(0, \min\left(1, \frac{1}{2} - \frac{1}{2} cyz\right)\right)$. Let $|y| \leq 1$

- If $y = 0$, then $\ell_c(z,y) = \frac{1}{2}$, $\ell_c$ is trivially $\frac{c}{2}$-Lipschitz ($c \geq 0$ is assumed) since it is $0$-Lipschitz.

- If $y \neq 0$, we consider

Let $z_1, z_2 \in \mathbb{R}$.
$$|\ell_c(z_1,y) - \ell_c(z_2,y)| = \begin{cases} |\min\left(1, \frac{1}{2} - \frac{1}{2} cyz_2\right)| & \text{if} \quad (1) \\ |\min\left(1, \frac{1}{2} - \frac{1}{2} cyz_1\right) - \min\left(1, \frac{1}{2} - \frac{1}{2} cyz_2\right)| & \text{if} (2) \\ |\min\left(1, \frac{1}{2} - \frac{1}{2} cyz_1\right)| & \text{if} (3) \\ 0 & \text{if} (4) \end{cases}$$

(1) $\quad 0 \geq \min\left\{1, \frac{1}{2} - \frac{1}{2} cyz_1\right\}$ and $0 \leq \min\left(1, \frac{1}{2} - \frac{1}{2} cyz_2\right)$

(2) $\quad 0 \geq \min\left\{1, \frac{1}{2} - \frac{1}{2} cyz_1\right\}$ and $0 \geq \min\left(1, \frac{1}{2} - \frac{1}{2} cyz_2\right)$

(3) $\quad 0 \leq \min\left\{1, \frac{1}{2} - \frac{1}{2} cyz_1\right\}$ and $0 \geq \min\left\{1, \frac{1}{2} - \frac{1}{2} cyz_2\right\}$

(4) $\quad 0 \leq \min\left\{1, \frac{1}{2} - \frac{1}{2} cyz_1\right\}$ and $0 \leq \min\left\{1, \frac{1}{2} - \frac{1}{2} cyz_2\right\}$

- Case 1: $y > 0$ $\quad \frac{1}{2} - \frac{1}{2} cyz \leq 1 \iff z \geq \frac{-1}{cy}$

and $\frac{1}{2} - \frac{1}{2} cyz \geq 0 \iff z \leq \frac{1}{cy}$

then

$$\ell_c(z,y) = \begin{cases} 1 & \text{if} \quad z \leq -\frac{1}{cy} \\ \frac{1}{2} - \frac{1}{2} cyz & \text{if} \quad z \in \left]-\frac{1}{cy}, \frac{1}{cy}\right[ \\ 0 & \text{if} \quad z \geq \frac{1}{cy} \end{cases}$$

if $z \notin \left]-\frac{1}{cy}, \frac{1}{cy}\right]$, $\ell_c(z,y)$ is constant and hence trivially lipschitz $\forall c > 0$. If $z \in \left]-\frac{1}{cy}, \frac{1}{cy}\right[$, the

$\left|\left(\frac{1}{2} - \frac{1}{2} cyz_1\right) - \left(\frac{1}{2} - \frac{1}{2} cyz_2\right)\right| = \left|\frac{1}{2} cy (z_2 - z_1)\right| \leq \frac{c}{2} |z_2 - z_1|$.

Case 2: If $y < 0$, then $\frac{1}{2} - \frac{1}{2} czy \le 1 \Leftrightarrow z \le -\frac{1}{cy}$

$\frac{1}{2} - \frac{1}{2} cyz \ge 0 \Leftrightarrow z \ge \frac{1}{cy}$.

Here we have $\frac{1}{cy} < \frac{-1}{cy}$. and

$$\ell_c(z,y) = \begin{cases} 1 & \text{if } z > \frac{-1}{cy} \\ \frac{1}{2} - \frac{1}{2} cyz & \text{if } z \in [\frac{1}{cy}, \frac{-1}{cy}] \\ 0 & \text{if } z < \frac{1}{cy}. \end{cases}$$

Similarly as above, $\ell_c(z,y)$ is constant and hence $0$-Lipschitz if $z \notin [\frac{1}{cy}, \frac{-1}{cy}]$, and otherwise

$\left| (\frac{1}{2} - \frac{1}{2} cyz_1) - (\frac{1}{2} - \frac{1}{2} cyz_2) \right| \le \left| \frac{1}{2} cy(z_2 - z_1) \right| \le \frac{c}{2} |z_1 - z_2|$.

So for $|y| \le 1$, $\ell_c(z,y)$ is $\frac{c}{2}$-Lipschitz.

(e) $\mathcal{F}_{lin} \{ f_\theta(x) = \langle x, \theta \rangle \mid \theta \in \mathbb{R}^d, \|\theta\| = 1 \}$

WTS: $\hat{\mathcal{R}}_n(\mathcal{F}_{lin}) \underset{(1)}{=} \mathbb{E}_\theta \underset{x \in \mathcal{E}}{\sup} \left\| \frac{1}{n} \sum_{i=1}^{n} \xi_i x_i \right\| \underset{(2)}{\le} \sqrt{\frac{\mathbb{E}[\|x_i\|^2]}{n}}$

$$\hat{R}_n(\mathcal{F}_{lin}) = \mathbb{E}_{X,\varepsilon} \sup_{f \in \mathcal{F}_{lin}} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i f(X_i)$$

$$= \mathbb{E}_{X,\varepsilon} \sup_{\substack{\theta \in \mathbb{R}^d \\ \|\theta\|=1}} \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i \langle X_i, \theta \rangle$$

$$= \frac{1}{n} \mathbb{E}_{X,\varepsilon} \sup_{\|\theta\| \leq 1} \langle X^T \varepsilon, \theta \rangle$$

$$\leq \frac{1}{n} \left( \mathbb{E}_{X,\varepsilon} [\|X^T \varepsilon\|_2^2] \right)^{1/2} \qquad \text{Cauchy-Schwartz, Jensen's}$$

$$= \frac{1}{n} \sqrt{\sum_{i,j=1}^{n} \mathbb{E}\left[ \langle X_i, X_j \rangle \varepsilon_i \varepsilon_j \right]}$$

$$= \frac{1}{n} \sqrt{\sum_{i=1}^{n} \mathbb{E}\left[ \|X_i\|^2 \varepsilon_i^2 \right] + \sum_{\substack{i,j=1 \\ i \neq j}}^{n} \mathbb{E}\left[ \langle X_i, X_j \rangle \varepsilon_i \mathbb{E}[\varepsilon_j] \right]}$$

$$= \sqrt{\frac{\mathbb{E}[\|X_i\|^2]}{n}}$$

(f). Using part (e), $\forall n \geq 1$, $\delta \geq 0$, ($L_c$ is bounded by 1)
if $\|X_i\| \leq 1$, then we have for $|y_i| \leq 1 \; \forall Y$,

$$\|L_c - L\|_{\mathcal{F}_{lin}} \leq c \tilde{R}(\mathcal{F}_{lin}) + \delta \quad \text{w.p. at least} \quad 1 - \exp\left\{ -\frac{\delta^2 n}{2} \right\}$$

$$\leq \frac{c}{\sqrt{n}} + \delta. \quad \leq c + \delta.$$

This is due to $\|X_i\| \leq 1$ a.s. $\Rightarrow \sqrt{\frac{1}{n} \mathbb{E}(\|X_i\|^2} \leq \frac{1}{\sqrt{n}}$

We want $N \in \mathbb{N}$ s.t. $\forall n \geq N$, $\| L_n - L \|_{\hat{\mathcal{F}}_{lin}} \leq \varepsilon$ with probability at least $1 - e^{-\delta^2 n / 2} \geq 1 - \delta'$, for some $\delta' > 0$.

From the above, it suffices for $\delta^2 n \geq 2 \log(1/\delta')$ and $\frac{c}{\sigma_n} + \sigma \leq \varepsilon$.

i.e. we select $\delta(n)$ s.t. $\delta(n) \to 0$ and $\delta^2 n \to \infty$.

Take $\delta = n^{-1/4}$. This is enough to get uniform convergence in probability $\| L - L_n \|_{\hat{\mathcal{F}}_{lin}} \xrightarrow{P} 0$.

Now, we can bound excess risk for fixed $X, Y$ by

$$L(\hat{\theta}_n) - L(\theta^*) = |L(\hat{\theta}_n) - L_n(\hat{\theta}_n)| + |L_n(\hat{\theta}_n) - L_n(\theta^*)| + |L_n(\theta^*) - L(\theta^*)|$$

$$\leq 2 \| L - L_n \|_{\hat{\mathcal{F}}_{lin}}.$$

To bound this by $\varepsilon$, we pick $\| L - L_n \|_{\hat{\mathcal{F}}_{lin}} \leq \varepsilon/2$. To optimize the rate of convergence for the given probability $\delta'$, take

$$\delta = \sqrt{\frac{2 \log(1/\delta')}{\sqrt{n}}}, \quad \text{and observe}$$

$$\| L - L_n \|_{\hat{\mathcal{F}}_{lin}} \leq \frac{c + \sqrt{2 \log(1/\delta)}}{\sqrt{n}} \leq \varepsilon/2 \quad (*)$$

when $(*)$ is true for $n \geq 4 \left( \frac{c + \sqrt{2 \log(1/\delta')}}{\varepsilon} \right)^2$

When $\delta'$ is the $\delta$ used in the probability of deviation from mean of sample loss, and $\delta$ is in the bound on deviation from sample loss.

g) WTS $\quad$ minimize $\hat{P}_n \ell_c(z,y)$

$\qquad$ ($\iff$ minimize $\min_c \max(0, \min(1, \frac{1}{2} - \frac{1}{2}cyz))$)

We have $\quad s(z) = \text{sgn}(z) \mathbb{1}[|z| \geq 1] + z \mathbb{1}[|z| \leq 1]$ $\quad$ and let $\text{sgn}_c(z) = s(cy)$ be a "smooth classifier".

Now $\quad \ell(\text{sgn}_c(z), y) = \frac{1}{2}(1 - \text{sgn}_c(z)y)$

$\qquad\qquad = \frac{s}{2}(1 - s(cyz))$

$\qquad\qquad = \max(0, \min(\frac{1}{2}(1 - cyz), 1))$

$\qquad\qquad = \ell_c(z, y)$

This gives us the wanted mapping. Since $X_i \sim N(0, \frac{I_d}{d})$, and $\|\theta\| = 1$, we have (cf d)

$\langle X_i, \theta \rangle \sim N(0, \frac{1}{d})$. To determine the size of $c$, we

we would like that most $\langle X_i, \theta \rangle$ are s.t. $|\langle X_i, \theta \rangle| \geq \frac{1}{c}$. The mean of $|\langle X_i, \theta \rangle|$ is $\sqrt{\frac{2}{\pi d}}$

Hence to classify most points $\left(\text{proba} \geq \frac{1}{c}\right)$, we

need $c \geq \sqrt{\frac{\pi d}{2}}$, because, as $\left(\frac{u}{n}\mathbb{1}\{|\langle x, \theta \rangle| \geq 1/c\}\right)_{n \geq 1}$

a series of bounded variables, the concentration

inequalities will give us that at least $\frac{1}{2}$ of the

points will be classified w.h.p.

## Q5. Chaining Apetizer

(a) $\exp(s A_n(T)) = \exp\left(s \, \mathbb{E}_a \sup_{t \in T} \langle a, t \rangle\right)$

$$\leq \mathbb{E}_a \exp\left(s \cdot \sup_{t \in T} \langle a, t \rangle\right) \qquad \text{convexity}$$

$$\leq |T| \sup_{t \in T} \mathbb{E}_a \exp\left(s \langle a, t \rangle\right) \qquad |T| < \infty, \text{ union bound}$$

$$\leq |T| \sup_{t \in T} \prod_{i=1}^{n} \mathbb{E}_{a_i} \exp(s a_i t_i) \qquad a_i \text{ iid}$$

$$\leq |T| \sup_{t \in T} \prod_{i=1}^{n} \exp\left(\frac{\sigma^2 t_i^2 s^2}{2}\right) \qquad a_i \text{ subgaussian} \\ \text{\& mean-zero}$$

$$= |T| \sup_{t \in T} \exp\left(\frac{\sigma^2 s^2}{2} \|t\|_2^2\right)$$

$$= |T| \exp\left(\frac{\sigma^2 s^2}{2} \sup_{t \in T} \|t\|_2^2\right)$$

This holds $\forall s \in \mathbb{R}_+$, so $\forall s \in \mathbb{R}_+$

$$\exp(s A_n(T)) \leq |T| \exp\left(\frac{\sigma^2 s^2}{2} \sup_{t \in T} \|t\|_2^2\right)$$

$$\implies A_n(T) \leq \frac{1}{s}\left(\log|T| + \frac{s^2 \sigma^2}{2} \sup_{t \in T} \|t\|_2^2\right)$$

chose $s$ such that $\quad C\sigma \sup_{t \in T} \|t\| \sqrt{\log |T|} = \frac{1}{s}\left(\log|T| + \frac{\sigma^2 s^2}{2} \sup_{t \in T} \|t\|_2^2\right)$

$$\implies C\sigma \sup_{t \in T} \|t\|_2 \sqrt{\log|T|} = \frac{1}{s}\log|T| + \frac{\sigma^2 s}{2} \sup_{t \in T} \|t\|_2^2$$

Take $s = \dfrac{2}{\sigma \sup\limits_{t \in T} \|t\|_2} \sqrt{\log|T|}$, then for $c \geq \dfrac{3}{2}$, we obtain the result.

(b) $\forall t \in T$, $\exists t' \in N_\varepsilon$ s.t. $\|t - t'\| \leq \tilde{\varepsilon}$. (we use $\tilde{\varepsilon}$ for notation)

$$R_n(T) = \mathbb{E}_\varepsilon \sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \varepsilon_i t_i$$

$$\leq \mathbb{E}_\varepsilon \sup_{t' \in N} \sup_{\|\delta\| \leq \tilde{\varepsilon}} \frac{1}{n} \langle \varepsilon, t' + \delta \rangle,$$

$$\leq \frac{1}{n} \mathbb{E}_\varepsilon \left[ \sup_{t' \in N} \langle \varepsilon, t' \rangle + \sup_{\|\delta\| \leq \tilde{\varepsilon}} \langle \varepsilon, \delta \rangle \right]$$

$$\leq R_n(N) + \frac{1}{n} \mathbb{E}_\varepsilon \sup_{\|\delta\| \leq \tilde{\varepsilon}} \langle \varepsilon, \delta \rangle$$

$$\mathbb{E}_\varepsilon \sup_{\|\delta\| \leq \tilde{\varepsilon}} \frac{1}{n} \langle \varepsilon, \delta \rangle = \mathbb{E}_\varepsilon \, \tilde{\varepsilon} \sup_{\|\delta\|_2 \leq 1} \frac{1}{n} \langle \varepsilon, \delta \rangle$$

$$\leq \frac{1}{n} \mathbb{E}_\varepsilon \sup_{\|\delta\|_2 \leq 1} \tilde{\varepsilon} \|\varepsilon\|_2 \cdot \|\delta\|_2 \leq \frac{1}{n} \mathbb{E}_\varepsilon \, \tilde{\varepsilon} \|\varepsilon\|_2 = \tilde{\varepsilon} \sqrt{n}$$

using the fact that $\mathbb{E}_\varepsilon \, \tilde{\varepsilon} \|\varepsilon\|_2 = \mathbb{E}\, \tilde{\varepsilon} \sqrt{n}$

Hence $R_n(T) \leq R_n(N) + \tilde{\varepsilon} \sqrt{n}$

$$\leq \frac{3}{2} \sup_{t \in N} \|t\| \cdot \sqrt{\log|N|} + \varepsilon \sqrt{n}$$

$$\leq \frac{3}{2} \sup_{t \in T} \|t\| \sqrt{\log|N|} + \varepsilon \sqrt{n}$$

(c) We have $\sigma_{max}(B) = \|B\|_{op} = \sup\limits_{\substack{\|x\|=1 \\ \|y\|=1}} y^T B x = \sup\limits_{\substack{\|x\|\leq 1 \\ \|y\|\leq 1}} y^T B x$

We write $\|B\|_{op} = \sup\limits_{x,y \in S^{d-1}} \sum\limits_{i,j=1}^{d} B_{ij} y_i x_j$

Now consider $t_{ij} = x_i y_j$, and consider $\tilde{B}, \tilde{t} \in \mathbb{R}^{d^2}$ s.t.

$$\begin{cases} B_{ij} = \tilde{B}_{i+(d-1)j} \\ t_{ij} = \tilde{t}_{i+(d-1)j} \end{cases} \quad \forall i,j \in [d].$$

Then $\|x\| \leq 1, \|y\| \leq 1$

$$\Rightarrow \|\tilde{t}\|_2^2 = \sum_{i=1}^{d^2} |t_{ij}|^2$$

$$= \sum_{i,j=1}^{d^2} |x_i y_j|^2$$

$$= \sum_{i=1}^{d} |x_i|^2 \sum_{j=1}^{d} |y_j|^2$$

$$= \|x\|_2^2 \cdot \|y\|_2^2 \leq 1$$

Hence $\|B\|_{op} \leq \sup\limits_{\tilde{t}} \langle \tilde{B}, \tilde{t} \rangle$

with $T = \{ t \in B_2^{d^2}, \text{ s.t. } \exists x,y \in S^{d-1} \text{ s.t. } t_{i+(d-1)j} = y_i x_j \}$

each $B_{ij}$ is $\text{unif}(\{\pm 1\})$ hence a Rademacher RV. To build an $\varepsilon$-net for $T$, it suffices to build separate $\frac{\varepsilon}{2}$ nets

$N_1, N_2$ of $B_2^d$ and look at the vectors with entries mapped to the entries of $xy^T$ $\forall x \in N_1, y \in N_2$.

Indeed, if $t \in T$, $\exists x \in N_1, y \in N_2$, $\|\delta_1\|_2 \leq \varepsilon/4$, $\|\delta_2\| \leq \varepsilon/4$ with $\delta_1, \delta_2 \in \mathbb{R}^d$ s.t.

$$t_{i+(d-1)j} = (x_i + \delta_{1j})(y_j + \delta_{2j}),$$

and if we look at the vector $z \in B_2^{d^2}$, which entries are mapped to the ones in $xy^T$

$$\|t \cdot z\|_2^2 = \sum_{i,j=1}^{d} \left((x_i + \delta_{1i})(y_j + \delta_{2j}) - x_i x_j\right)^2$$

$$= \sum_{i,j=1}^{d} \left(x_i \delta_{2j} + y_j \delta_{1i} + \delta_{1i}\delta_{2j}\right)^2$$

$$= \|x\|_2^2 \|\delta_2\|_2^2 + \|y\|_2^2 \|\delta_1\|_2^2 + 2\langle x, \delta_1\rangle\langle y, \delta_2\rangle$$
$$+ 2\langle y, \delta_2\rangle\|\delta_1\|_2^2 + 2\langle x, \delta_1\rangle\|\delta_2\|_2^2 + \|\delta_1\|_2^2\|\delta_2\|^2$$

$$\leq \frac{\varepsilon^2}{16} + \frac{\varepsilon^2}{16} + \frac{2\varepsilon^2}{16} + \frac{2\varepsilon^3}{64} + \frac{2\varepsilon^3}{64} + \frac{\varepsilon^4}{128}$$

$$\leq \frac{\varepsilon^2}{4} + \frac{\varepsilon^3}{2} + \frac{\varepsilon^4}{16} \leq \varepsilon^2$$

We can build $\frac{\varepsilon}{4}$ nets of $B_2^d$ of size $\left(1 + \frac{3}{\varepsilon}\right)^d$ & $\exists$ an $\varepsilon$-net of $T$ of size $\left(1 + \frac{3}{\varepsilon}\right)^{2d}$. Now applying ($\delta$), we get

$$\mathbb{E}_{\tilde{B}} \|B\|_{op} = \mathbb{E}_{\tilde{B}} \sup_{\tilde{t} \in T} \langle \tilde{B}, \tilde{t} \rangle$$

$$\leq c \cdot \underbrace{\sup_{\tilde{t} \in T} \|\tilde{E}\|}_{\leq 1 \text{ shown above}} \sqrt{\log\left(1 + \frac{2}{\varepsilon}\right)^{2d}} + \varepsilon\sqrt{d^2}$$

$$= c \sqrt{d \log\left(1 + \frac{3}{\varepsilon}\right)^2} + \varepsilon d$$

By taking $\varepsilon = \frac{1}{\sqrt{d}}$, $\quad \mathbb{E}_{\tilde{B}} \|B\|_{op} \leq c \sqrt{2d \log\left(1 + 64d + 16\sqrt{d}\right)} + \sqrt{d}$

$$= O\left(\sqrt{d \log d}\right)$$

cd) The bound from (b) is not tight. The lossy term is $\varepsilon\sqrt{n}$. We obtained it by applying Cauchy-Schwartz to $\sum_{i=1}^{n} \varepsilon_i \delta_i$ to bound it by $\|\varepsilon\|_2 \|\delta\|_2$. But the bound is achieved only when $\delta$ is colinear to $\varepsilon$. However, taking the sup over all $\delta \in B_2^d [\varepsilon]$ was already lossy, as our covering might include include points out of our set $T$ ( in (b) ). Then we might not have $\delta$ colinear to $\varepsilon$ in the set which creates loo iness in the Cauchy-Schwartz inequality.