

Suppose X is uncountable. Can we still bound

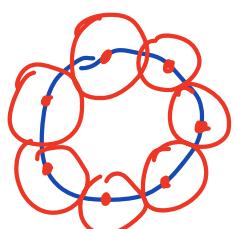
Choose x_1, \dots, x_T measure $y_\epsilon = \langle x_\epsilon, \theta_0 \rangle + \zeta_\epsilon$
(G-optimal) $\zeta_\epsilon \stackrel{\text{independent}}{\sim} \text{Normal}$

$$\begin{aligned}\hat{\theta} &= \left(\sum_t x_t x_t^\top \right)^{-1} \sum_t x_t y_\epsilon \\ &= \theta_0 + \left(\sum_t x_t x_t^\top \right)^{-1} \sum_t x_t \zeta_\epsilon \quad A = \sum_t x_t x_t^\top\end{aligned}$$

$$\begin{aligned}\max_{x \in X} |\langle \hat{\theta} - \theta_0, x \rangle| &= \max_x \left\| (\hat{\theta} - \theta_0)^\top A^{-1/2} \bar{A}^{1/2} x \right\|_2 \\ &\leq \max_x \left\| A^{1/2} (\hat{\theta} - \theta_0) \right\|_2 \cdot \left\| \bar{A}^{1/2} x \right\|_2 \\ &= \left\| \hat{\theta} - \theta_0 \right\|_A \cdot \underbrace{\max_x \|x\|_A}_{\leq \sqrt{d/T}} \cdot \sqrt{d + 2 \log(1/\delta)} \cdot \sqrt{d/T} \geq \sqrt{\frac{d}{T}}\end{aligned}$$

$$\left\| \hat{\theta} - \theta_0 \right\|_A = \left\| A^{1/2} (\hat{\theta} - \theta_0) \right\|_2 = \sup_{u: \|u\|_2 \leq 1} u^\top A^{1/2} (\hat{\theta} - \theta_0)$$

Lemma \exists cover $C \subset \mathbb{R}^d : |C| \leq (3/\varepsilon)^d$ and $x \in C, \|x\|_2 = 1$
 and $\forall x$ w/ $\|x\|_2 = 1$ we have $\|y - x\|_2 \leq \varepsilon$ for some $y \in C$.



In other words

$$\bigcup_{x \in C} B(x, \varepsilon) \supset \mathbb{S}^{d-1}$$

sphere
in d -dim

$$\begin{aligned}\left\| \hat{\theta} - \theta_0 \right\|_A &= \sup_{u: \|u\|_2 \leq 1} u^\top A^{1/2} (\hat{\theta} - \theta_0) \\ &= \sup_u \min_{y \in C} (u - y)^\top A^{1/2} (\hat{\theta} - \theta_0) + y^\top A^{1/2} (\hat{\theta} - \theta_0)\end{aligned}$$

$$\begin{aligned}&\leq \sup_u \min_{y \in C} \|u - y\|_2 \left\| \hat{\theta} - \theta_0 \right\|_A + \underbrace{y^\top A^{1/2} (\hat{\theta} - \theta_0)}_{\leq \varepsilon} \\ &\leq \sqrt{2 \log(\frac{|C|}{\delta})}\end{aligned}$$

$$\leq \varepsilon \|\hat{\theta} - \theta_*\|_A + \sqrt{2 \log(\frac{|C|}{\delta})}$$

$$\Rightarrow \|\hat{\theta} - \theta_*\|_A \leq \frac{1}{1-\varepsilon} \cdot \sqrt{2 \log(\frac{|C|}{\delta})} \leq \frac{1}{1-\varepsilon} \cdot \sqrt{2d \log(\frac{3/\varepsilon}{\delta}) + 2 \log(1/\delta)}.$$

$$\mathbb{E}[(y^T A'^{1/2} (\hat{\theta} - \theta_*))^2] = \mathbb{E}[y^T A'^{1/2} \bar{A}' (\sum_t x_t z_t) (\sum_s x_s z_s)^T \bar{A}' A'^{1/2} y]$$

$$\begin{aligned} &= y^T \bar{A}'^{1/2} \underbrace{\mathbb{E}\left[\sum_t x_t x_t^T z_t^2\right]}_{=A} \bar{A}'^{1/2} y \\ &= y^T y \\ &= 1 \end{aligned}$$

$$\mathbb{E}[z_t^2] = 1$$

Stochastic Processes

$$\|\mathbb{E}[\sum_t x_t x_t^T z_t^2] \bar{A}'^{-1/2} y\|_2$$

\leq

Flip two coins

Sample / outcome space $\mathcal{S} = \{HH, HT, TH, TT\}$

σ -algebra on a set \mathcal{S} called \mathcal{F} satisfies

1. $\mathcal{S} \in \mathcal{F}$,

2. $F \in \mathcal{F} \Rightarrow F^c \in \mathcal{F}$

3. $F_n \in \mathcal{F}$ $\forall n$ then $\cup_n F_n \in \mathcal{F}$.

Probability measure P assign numbers to \mathcal{F} ($P: \mathcal{F} \rightarrow \mathbb{R}$)

$P(F) \geq 0 \quad \forall F \in \mathcal{F}$

$P(\emptyset) = 0$

F_n are disjoint $\forall n$

$$P(\cup_n F_n) = \sum_n P(F_n)$$

Random variable X is a function mapping $\mathbb{R} \rightarrow \mathbb{R}$.

Conditional expectation

$$P(Y|X) = \frac{P(Y, X)}{P(X)}$$

(discrete)

If X is finite then

$$E[Y|X=x] = \sum_y y P(Y=y|X=x)$$

but otherwise

$$E[Y|X](\omega) = g(\omega) \text{ where } g \text{ is any function}$$

satisfying

$$\int_{\omega \in H} g(\omega) dP(\omega) = \int_{\omega \in H} Y(\omega) dP(\omega)$$

for all $H \in \sigma(X)$

\nexists measurable function f : $f(X(\omega)) = E[Y|X](\omega)$

X_1, X_2, X_3, \dots

$$\mathcal{F}_t = \sigma(\{X_s\}_{s \leq t}). \quad \mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \dots$$

Filtration: $\{\mathcal{F}_t\}_{t=1}^{\infty} = \mathcal{F}$

X_t is \mathcal{F} -adapted if X_t is \mathcal{F}_t measurable.

$$\mathbb{E}[X_t | \mathcal{F}_t] = X_t$$

$$\mathbb{E}[X_{t+1}] = \mathbb{E}[\mathbb{E}[X_{t+1} | \mathcal{F}_t]] = \mathbb{E}[X_t] = \mathbb{E}[X_0]$$

An \mathcal{F} -adapted sequence of R.V.s is a martingale if $\mathbb{E}[X_{t+1} | \mathcal{F}_t] = X_t \quad \forall t$, and $\mathbb{E}[|X_t|] < \infty$.

If $\mathbb{E}[X_{t+1} | \mathcal{F}_t] \leq X_t$ we say its supermartingale

If $\mathbb{E}[X_{t+1} | \mathcal{F}_t] \geq X_t$ " sub-martingale.

Ex. $Z_t \sim N(0, 1)$, then $X_t = \sum_{s=1}^t Z_s$ is a martingale.

$$\mathbb{E}[X_{t+1} | \mathcal{F}_t] = \mathbb{E}\left[\sum_{s=1}^{t+1} Z_s | \mathcal{F}_t\right] = \mathbb{E}[Z_{t+1} | \mathcal{F}_t] + \sum_{s=1}^t Z_s$$

$$\text{Ex. } Z_t = \begin{cases} 2 & \text{w.p. } 1/2 \\ 0 & \text{w.p. } 1/2 \end{cases}, \quad X_t = \prod_{s=1}^t Z_s, \quad X_0 = 1 = 0 + X_t$$

Ex. $Z_t \sim N(0, 1)$, $X_t = \prod_{s=1}^t \exp(\lambda Z_s)$, $\lambda > 0$

$$\mathbb{E}[X_{t+1} | \mathcal{F}_t] = \prod_{s=1}^t \exp(\lambda Z_s) \mathbb{E}\left[\exp(\lambda Z_{t+1}) | \mathcal{F}_t\right] e^{\lambda^2/2}$$

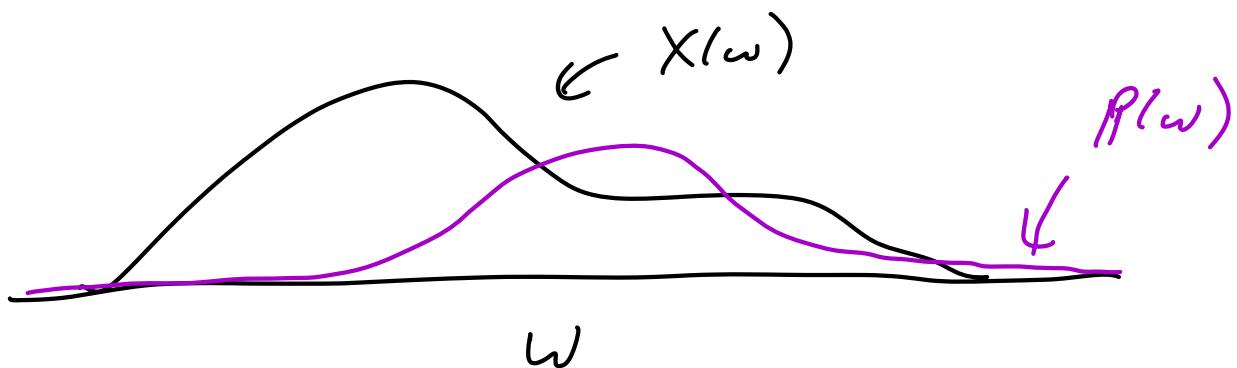
$$X_0 = 1 = X_t \cdot e^{\lambda^2/2}$$

$\geq X_t \Rightarrow X_t$ is sub-martingale

$X_t = \prod_{s=1}^t \exp(\lambda Z_s - \lambda^2/2) \Rightarrow X_t$ is super-martingale

$X(\omega) \quad \omega \in \Omega$

$$\mathbb{E}[X] = \sum_{\omega} X(\omega) P(\omega)$$



Def] A random variable $\tau \in \mathbb{N}$ is a stopping time if $\{\tau \leq t\}$ is \mathcal{F}_t measurable a.s.

Ex. $S_t = \sum_{s=1}^t Z_s$ $Z_s \sim N(0, 1)$ then

$\tau = \min\{t : S_t \geq \varepsilon\}$ is a stopping time

Ex. $\tau = \max\{t : S_t \geq \varepsilon\}$ is not a stopping time.

Lemma] Doob's Optional Stopping. Let τ be a stopping time and X_t be \mathcal{F} -adapted martingale. If

- $\exists N : P(\tau \leq N) = 1$, or
- $E[\tau] < \infty$ and $E[|X_{\tau \wedge n} - X_n|] \leq c$ for some c

then X_τ is well-defined and $E[X_\tau] = E[X_0]$.

Furthermore if X_τ is supermartingale $E[X_\tau] \leq E[X_0]$.
 " " sub-martingale $E[X_\tau] \geq E[X_0]$.

Lemma] Maximal inequality. Let X_t be \mathcal{F} -adapted sequence w/ $X_t \geq 0$ a.s. Then for any $\varepsilon > 0$

- $P(\max_{t \in [0, \tau]} X_t \geq \varepsilon) \leq \frac{E[X_\tau]}{\varepsilon}$ if X_τ is supermartingale
- $P(\max_{t \in [0, \tau]} X_t \geq \varepsilon) \leq \frac{E[X_\tau]}{\varepsilon}$ if X_τ is sub-martingale

Let Z_t be IID mean-zero 1-sub-Gaussian R.V. so that

$$\mathbb{E}[\exp(\lambda Z_t)] \leq e^{\lambda^2/2}. \quad S_t = \sum_{s=1}^t Z_s$$

$$M_t(\lambda) = \exp(\lambda S_t - t\lambda^2/2).$$

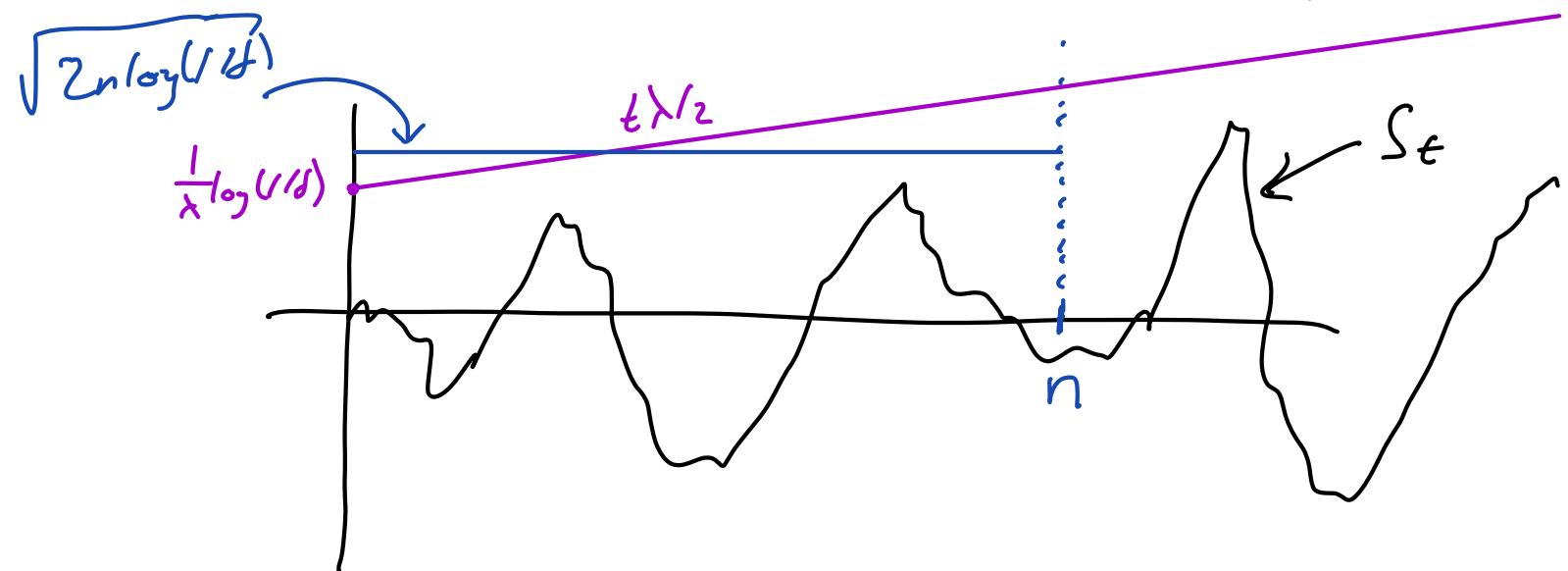
$$\begin{aligned} M_t(\lambda) &\text{ is supermartingale. } \mathbb{E}[M_{t+1}(\lambda) | \mathcal{F}_t] = \mathbb{E}[\exp(\lambda S_{t+1} - (t+1)\lambda^2/2) | \mathcal{F}_t] \\ &= \exp(\lambda S_t - t\lambda^2/2) \underbrace{\mathbb{E}[e^{\lambda Z_{t+1} - \lambda^2/2} | \mathcal{F}_t]}_{\leq 1} \\ &\leq M_t(\lambda) \end{aligned}$$

$$M_0(\lambda) := 1$$

By maximal inequality we have

$$\mathbb{P}\left(\max_{t \in \mathbb{N}} \exp(\lambda S_t - t\lambda^2/2) \geq \frac{1}{\delta}\right) \leq \frac{\mathbb{E}[M_0(\lambda)]}{1/\delta} = \delta$$

$$\Rightarrow \mathbb{P}\left(\exists t : S_t \geq t\lambda/2 + \frac{1}{\lambda} \log(1/\delta)\right) \leq \delta.$$



$$\text{Fix } n \in \mathbb{N} \text{ set } \lambda = \sqrt{\frac{2 \log(1/\delta)}{n}}$$

$$\Rightarrow \mathbb{P}(\exists t: S_t \geq (\epsilon/\sqrt{n} + \sqrt{n})\sqrt{\log(1/\delta)/2}) \leq \delta$$

$$\Rightarrow \mathbb{P}\left(\max_{t=1,\dots,n} : S_t \geq \sqrt{2n\log(1/\delta)}\right) \leq \delta$$

Chernoff Bound says $\mathbb{P}(S_n \geq \sqrt{2n\log(1/\delta)}) \leq \delta$

Method of Mixtures

Let $h(\lambda)$ be a probability distribution over λ

$$\text{Define } \bar{M}_t = \mathbb{E}_{\lambda} [M_t(\lambda)] = \int_{\lambda} M_t(\lambda) h(\lambda) d\lambda$$

\bar{M}_t is a supermartingale if $M_t(\lambda)$ is a superm.

$$\mathbb{E}[\bar{M}_{t+1} | \mathcal{F}_t] = \mathbb{E}\left[\int_{\lambda} M_{t+1}(\lambda) h(\lambda) d\lambda | \mathcal{F}_t\right]$$

$$= \int_{\lambda} \mathbb{E}[M_{t+1}(\lambda) | \mathcal{F}_t] h(\lambda) d\lambda$$

$$= \int_{\lambda} M_t(\lambda) h(\lambda) d\lambda$$

$$= \bar{M}_t$$

S_t is defined as above (sum 1-skell-Gaussian)

and define $h(\lambda) = \frac{1}{\sqrt{2\pi\nu^2}} e^{-\frac{\lambda^2}{2\nu^2}}$

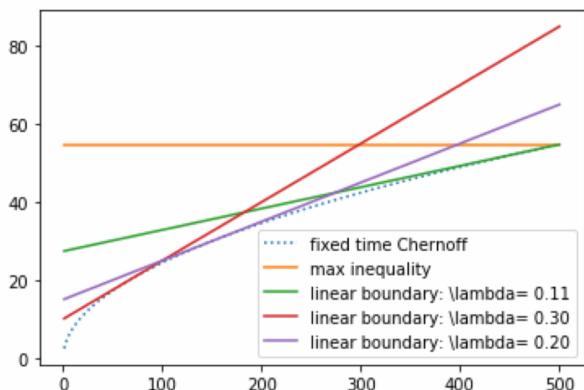
$$\bar{M}_t \doteq \int M_t(\lambda) h(\lambda) d\lambda$$

$$= \int \exp(\lambda S_t - \lambda^2/2) \frac{1}{\sqrt{2\pi\nu^2}} \exp(-\lambda^2/2\nu^2) d\lambda$$

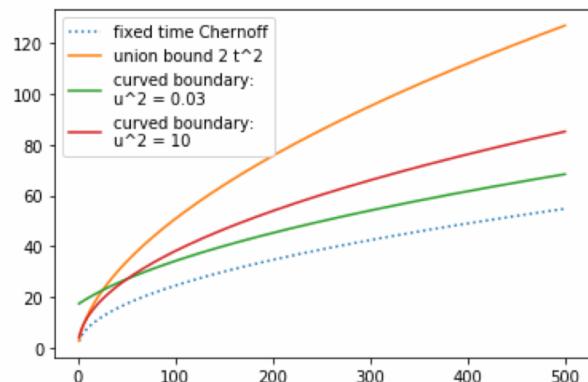
$$= \sqrt{\frac{\nu^{-2}}{t + \nu^{-2}}} \cdot \exp\left(S_t^2 (t + \nu^{-2})^{-1} / 2\right)$$

$$P(\exists t: |S_t| \geq \sqrt{2(t+\nu^{-2}) \left(\log(1/\delta) + \frac{1}{2} \log\left(\frac{t+\nu^{-2}}{\nu^{-2}}\right) \right)})$$

$$= P(\exists t: \bar{M}_t \geq \frac{1}{\delta}) \leq \delta.$$



(a) Fix $\delta = 0.05$. The ‘fixed time Chernoff’ represents $\sqrt{2t \log(1/\delta)}$ which holds at each t but not all $t \leq 500$ simultaneously (which is why it is dotted). The ‘max inequality’ holds for all $t \leq 500$, and the linear boundaries hold for all $t \in \mathbb{N}$ simultaneously.



(b) Fix $\delta = 0.05$. The ‘fixed time Chernoff’ represents $\sqrt{2t \log(1/\delta)}$ which holds at each t but not all $t \in \mathbb{N}$ simultaneously (which is why it is dotted). All other curves do hold for all $t \in \mathbb{N}$ simultaneously. “union bound $2t^2$ ” plots $\sqrt{2 \log(2t^2/\delta)}$.

Z_1, Z_2, \dots is an \mathcal{F} -adapted sequence. We say

α_t is predictable if α_t is \mathcal{F}_{t-1} -measurable.

Suppose for any λ we have $\mathbb{E}[Z_t | \mathcal{F}_{t-1}] = 0$

$$\mathbb{E}[\exp(\lambda Z_t) | \mathcal{F}_{t-1}] \leq \exp(\lambda^2 \alpha_t^2 / 2)$$

and define $S_t = \sum_{s=1}^t Z_s$. Then $V_t = \sum_{s=1}^t \alpha_s^2$

$$M_t(\lambda) = \exp(\lambda S_t - \lambda^2 V_t / 2)$$

is a super martingale.

$$\mathbb{P}\left(\exists t: S_t \geq \lambda V_t / 2 + \frac{1}{\lambda} \log(1/\delta)\right) \leq \delta$$

Define $\bar{M}_t = \mathbb{E}_{\lambda \sim h} [M_t(\lambda)]$ $h(\lambda) = \frac{1}{\sqrt{2\pi}} e^{-\lambda^2/2}$

$$\mathbb{P}\left(\exists t: |S_t| \geq \sqrt{(V_t+1) \log\left(\frac{V_t+1}{\delta}\right)}\right) \leq \delta.$$

Ex. Gradient Descent Analysis.

$$w_{t+1} = w_t - \gamma \nabla \ell_t(w_t) \quad \mathbb{E}[\ell_t(w)] \stackrel{?}{=} \ell(w).$$

$$\mathbb{E}[f(w_t) - f(w_*)] \leq \frac{RG}{\sqrt{t}} \quad \text{where} \quad R = \|w_* - w_0\|_2 \\ G = \|\nabla \ell(w)\|_2$$

$Z_t = \ell(w_t) - \mathbb{E}[\ell(w_t) | \mathcal{F}_{t-1}]$ is a martingale.

$$Z_t = \nabla \ell(w_t)(w_t - w_*) - \nabla \ell(w_t)(w_t - w_*)$$

$$\|\lambda\|_{\Sigma_t}^2 = \lambda^T \Sigma_t^{-1} \lambda$$

Suppose that $z_1, z_2, \dots \in \mathbb{R}^d$ is \mathcal{F} -adapted sequence.

$$\text{satisfying } \mathbb{E}[\exp(\langle \lambda, z_t \rangle) | \mathcal{F}_{t-1}] \leq \exp(\|\lambda\|_{\Sigma_t}^2 / 2)$$

for all $\lambda \in \mathbb{R}^d$, where $\sum_t \in \mathbb{R}^{d \times d}$ is a predictable sequence.

$$M_t(\lambda) = \exp(\langle \lambda, s_t \rangle - \|\lambda\|_{V_t}^2 / 2)$$

$$s_t = \sum_{s=1}^t z_s$$

$$V_t = \sum_{s=1}^t \Sigma_s$$

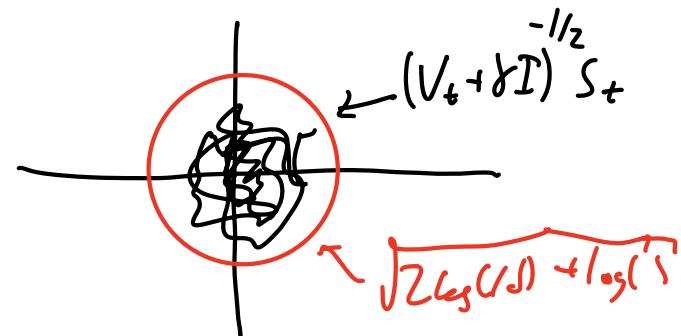
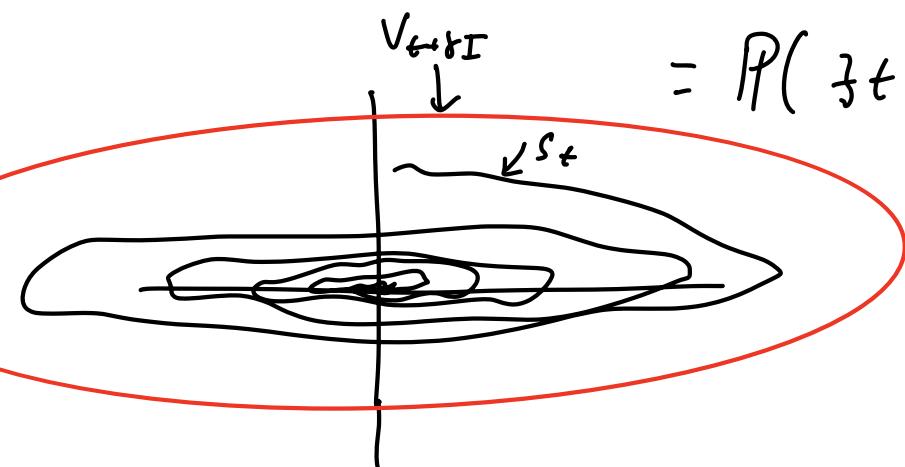
$$\mathbb{P}(\exists t: \langle \lambda, s_t \rangle - \|\lambda\|_{V_t}^2 / 2 \geq \log(1/\delta)) \leq \delta$$

$$\bar{M}_t = \mathbb{E}_{\lambda \sim h}[M_t(\lambda)] \quad h(\lambda) = \frac{1}{(2\pi\gamma)^{d/2}} \exp(-\|\lambda\|^2\gamma/2)$$

$$h(\lambda) \equiv \mathcal{N}(0, \frac{1}{\gamma} I_d)$$

$$\mathbb{P}(\exists t: \|s_t\|_{(V_t + \gamma I)^{-1}} \geq \sqrt{2\log(1/\delta) + \log\left(\frac{\|V_t + \gamma I\|}{\gamma}\right)})$$

$$= \mathbb{P}(\exists t: \bar{M}_t \geq \frac{1}{\delta}) \leq \delta$$



Online Linear Regression.

for $t=1, 2, \dots$

Observe $x_t \in \mathbb{R}^d$

Player predicts $\hat{y}_t = \langle x_t, \hat{\theta}_t \rangle$

Observes $y_t = \langle x_t, \theta_* \rangle + \gamma_t \leftarrow \gamma_t$ mean-zero sub-Gaussian.

$$\hat{\theta}_{t+1} = \underset{\theta}{\operatorname{argmin}} \sum_{s=1}^t (y_s - \langle \theta, x_s \rangle)^2 + \gamma \|\theta\|_2^2$$

$$= \left(\underbrace{\sum_{s=1}^t x_s x_s^\top}_{V_t} + \gamma I \right)^{-1} \sum_{s=1}^t x_s y_s$$

$$S_t = \sum_{s=1}^t x_s \gamma_s$$

$$= (V_t + \gamma I)^{-1} V_t \theta_* + (V_t + \gamma I)^{-1} S_t$$

x_t is \mathcal{F}_{t-1} measurable

from above

$$z_s \equiv x_s \gamma_s \quad \mathbb{E}[\exp(\langle \lambda, x_s \gamma_s \rangle) | \mathcal{F}_{s-1}] = \mathbb{E}[\exp(\langle \lambda, x_s z_s \rangle) | \mathcal{F}_{s-1}]$$

$$\sum_s \equiv x_s x_s^\top \leq \exp(\langle \lambda, x_s \rangle^2 / 2)$$

$$= \exp(\lambda^\top x_s x_s^\top \lambda / 2)$$

$$= \exp(\|\lambda\|_{x_s x_s^\top}^2 / 2)$$

$$\|\hat{\theta}_{t+1} - \theta_*\|_{(V_t + \gamma I)} = \|S_t - \gamma \theta_*\|_{(V_t + \gamma I)^{-1}}$$

$$\leq \|S_t\|_{(V_t + \gamma I)^{-1}} + \gamma \|\theta_*\|_{(V_t + \gamma I)^{-1}}$$

$$\leq \|S_t\|_{(V_t + \gamma I)^{-1}} + \sqrt{\gamma} \|\theta_a\|_2$$

$$\mathbb{P}\left(\exists t: \|\hat{\theta}_{t+1} - \theta_a\|_{(V_t + \gamma I)} \geq \sqrt{\gamma} \|\theta_a\|_2 + \sqrt{2 \log(1/\delta) + \log\left(\frac{|V_t + \gamma I|}{\gamma^d}\right)}\right) \leq \delta$$

$$VCB: \quad \mathcal{X}_t = \underset{x \in \mathcal{X}}{\text{argmax}} \quad VCB_t(x)$$

$$VCB_t(x) = \underset{\theta \in C_t}{\text{argmax}} \langle \theta, x \rangle$$

$$C_t = \left\{ \theta: \|\theta - \hat{\theta}_t\|_{(V_t + \gamma I)} \leq R_t \right\}$$