CS 59800 - RL
Linear Regression.

Agenda:
- Ridge regression
- Concentration bound

---

Consider a linear model of

$$X_t = \langle A_t, \theta_* \rangle + \eta_t \quad \longrightarrow \text{scalar}$$

scalar

$$A_t \in R^d \quad \theta_* \in R^d$$

and a sequence of $A_1, X_1, \cdots, A_n, X_n$, with a
filtration $\mathbb{F} = \{ \mathcal{F}_t \}$, such that, $\mathcal{F}_t = \sigma \left( A_1, X_1, \cdots A_{t+1} \right)$

Note that, $X_t$ is $\mathcal{F}_t$ measurable.
For noise $\eta_t$, we have:

$$E\left[ \exp(\alpha \eta_t) \mid \mathcal{F}_{t-1} \right] \leq \exp\left( \frac{\alpha^2}{2} \right)$$

for all $\alpha \in R$ and $t \in [n]$.
This is the 1-sub-Gaussian assumption on the noise.

---

The process is as follows:
- At each time step $t$, someone chooses $A_t$ from a
set $D_t$, and we observe $X_t$.
Given $A_1, X_1, \cdots A_t, X_t$,
Can we estimat $\theta_*$?

- How about solving a ridge regression for $\theta_*$?
→ at time $t$,

$$\min_{\theta \in R^d} \sum_{s=1}^{t} (X_s - \langle A_s, \theta \rangle)^2 + \|\theta\|_V^2$$

For positive Definite matrix $V$

$$\|\theta\|_V^2 = \theta^T V \theta$$

Let's define

- $S_t = \sum_{s=1}^{t} \eta_s A_s$

- $V_t = \sum_{s=1}^{t} A_s A_s^T$

- $V_t(V) = V + V_t = V_t(V)$

As you may have seen, we usually set $V = \lambda I$ for $\lambda > 0$

- $V_t(\lambda) = \lambda I + V_t$

Therefore, the minimizer of ridge regression problem is

$$\hat{\theta}_t = V_t(v)^{-1} \sum_{s=1}^{t} X_s A_s \qquad \text{(why?)}$$

we are interested is how good is this $\hat{\theta}_t$ estimate.

To simplify the notation, we use $V = \lambda I$

Let's define • $M_t(\alpha) = \exp\left(\langle \alpha, S_t \rangle - \frac{1}{2} \|\alpha\|_{V_t}^2\right)$

for $\alpha \in \mathbb{R}^d$

Lemma: For any $\alpha \in \mathbb{R}^d$, the process $M_t(\alpha)$ is an $\mathbb{F}$-adapted supermartingale.

Proof:

— It is clear that $M_t(\alpha)$ is $\mathcal{F}_t$-measurable for all $t$ by definition. We are left to show

$$E\left[M_t(\alpha) \mid \mathcal{F}_{t-1}\right] \leq M_{t-1}(\alpha) \quad \text{a.s.}$$

Let's expand $M_t(\alpha) \to$

$$E\left[M_t(\alpha) \mid \mathcal{F}_{t-1}\right] = E\left[\exp\left(\langle \alpha, S_t \rangle - \frac{1}{2} \|\alpha\|_{V_t}^2\right) \mid \mathcal{F}_{t-1}\right]$$

$$= E\left[\exp\left(\langle \alpha, S_{t-1}\rangle - \frac{1}{2}\|\alpha\|_{V_{t-1}}^2\right) \exp\left(\eta_t \langle \alpha, A_t\rangle - \frac{1}{2}\|\alpha\|_{A_t A_t^T}\right) \mid \mathcal{F}_{t-1}\right]$$

$$= M_{t-1}(\alpha) \underbrace{E\left[\exp \eta_t \langle \alpha, A_t\rangle - \frac{1}{2}\|\alpha\|_{A_t A_t^T} \mid \mathcal{F}_{t-1}\right]}_{\leq 1 \quad \text{a.s.}}$$

$$\leq M_{t-1}(\alpha) \quad \text{a.s.}$$

↳ we concluded that $M_t(\alpha)$ is a supermartingale sequence.

For a Gaussian measure $h$ with covariance $V$, let's
define:

$$\overline{M}_t = \int_{\mathbb{R}^d} M_t(\alpha)\, dh(\alpha)$$

now using Radon-Nikodym derivative, we have:

$$\overline{M}_t = \frac{1}{(2\pi)^{d/2} \det(V^{-1})^{\frac{1}{2}}} \int_{\mathbb{R}^d} \exp\left(\langle \alpha, S_t\rangle - \frac{1}{2}\|\alpha\|_{V_t}^2 - \frac{1}{2}\|\alpha\|_V^2\right) d\alpha$$

Note that:

$$\underbrace{\|S_t\|_{(V+V_t)^{-1}}^2}_{V_t(V)} - \|\alpha - (V+V_t)^{-1} S_t\|_{V+V_t}^2$$

$$= 2\langle \alpha, S_t\rangle - \|\alpha\|_{V_t}^2 - \|\alpha\|_V^2 \qquad (\text{why?})$$

$$\Rightarrow \overline{M}_t = \frac{1}{(2\pi)^{d/2}\det(V^{-1})^{\frac{1}{2}}} \int_{\mathbb{R}^d} \exp\left(\|S_t\|_{V_t(V)^{-1}}^2 \times \frac{1}{2} - \frac{1}{2}\|\alpha - V_t(V)^{-1} S_t\|_{V_t(V)}^2\right) d\alpha$$

$$= \frac{1}{(2\pi)^{d/2}\det(V^{-1})^{\frac{1}{2}}} \underbrace{\exp\left(\frac{1}{2}\|S_t\|_{V_t(V)^{-1}}^2\right)} \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\|\alpha - V_t(V)^{-1} S_t\|_{V_t(V)}\right) d\alpha$$

$$= \frac{\det(V_t(V)^{-1})^{\frac{1}{2}}}{\det(V^{-1})^{\frac{1}{2}}} \exp\left(\frac{1}{2}\|S_t\|_{V_t(V)^{-1}}^2\right) \times \overbrace{\left( \frac{1}{(2\pi)^{d/2}\det(V_t(V)^{-1})^{\frac{1}{2}}} \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\|\alpha - V_t(V)^{-1} S_t\|_{V_t(V)}^2\right) d\alpha \right)}^{\text{equal to } 1}$$

Note that $\bar{M}_t$ is also super martingale (why?)

Using the general form of maximal inequality, we have:

$$\mathbb{P}\left(\sup_t \bar{M}_t > \frac{1}{\delta}\right) \leq \delta$$

we know that $\{t \in [n]: \frac{\det(V_t(V'))^{\frac{1}{2}}}{\det(V^{-1})^{\frac{1}{2}}} \exp\left(\frac{1}{2}\|S_t\|_{V_t(V)}^{-1}\right) > \frac{1}{\delta}\}$

$$\subset \{\sup_t \bar{M}_t > \frac{1}{\delta}\}$$

$$\leadsto \mathbb{P}\left(t \in [n]: \frac{\det(V_t(V)^{-1})^{\frac{1}{2}}}{\det(V^{-1})^{\frac{1}{2}}} \exp\left(\frac{1}{2}\|S_t\|_{V_t(V)}^{-1}\right) > \delta\right) \leq \delta$$

$$\leadsto \mathbb{P}\left(t \in [n]: \|S_t\|_{V_t(V)^{-1}} > 2\log\frac{1}{\delta} + \log\left(\frac{\det(V_t(V))}{\det(V)}\right)\right) \leq \delta$$

**Theorem:** For $\delta \in (0,1)$, with probability at least $1-\delta$, for $t \in [n]$, we have; for $V = \lambda I$

$$\|\hat{\theta}_t - \theta_*\|_{V_t(V)} \leq \sqrt{\beta_t(\delta)} := \sqrt{\lambda}\|\theta_*\| + \sqrt{2\log\left(\frac{1}{\delta}\right) + \log\left(\frac{\det(V_t(V))}{\det(V)}\right)}$$

Furthermore if $\|\theta_*\| \leq S$, define confidene interval/set

$$C_t(\delta) = \{\theta \in \mathbb{R}^d: \|\hat{\theta}_{t-1} - \theta\|_{V_{t-1}(V)} \leq \sqrt{\lambda}S + (\overset{\text{for t-1}}{\frown})\}$$
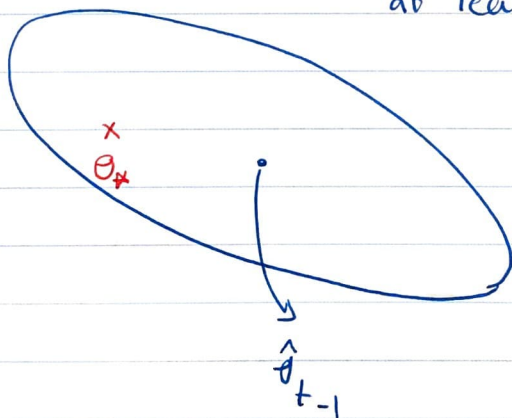
Then $\mathbb{P}\left(\text{exists } t \in [n]; \; \theta_* \notin C_t^{(\delta)}\right) \leqslant \delta$

where $C_t(\delta) = \Big\{ \theta \in \mathbb{R}^d : \|\theta_{t-1} - \theta\|_{V_t(v)}$

$$\leqslant \sqrt{\lambda} \, S + \sqrt{2 \log\left(\frac{1}{\delta}\right) + \log\left(\frac{\det(V_{t-1}(v))}{\det(V)}\right)}$$

what does it mean?

$$\Big\{ \theta \in \mathbb{R}^d; \; \|\hat{\theta}_{t-1} - \theta\|^2_{V_{t-1}(\lambda)} \, < \, \beta_{t-1}(\delta) \Big\}$$

is an elipse such that $\theta_*$ is in it, always, with propability at least $1 - \delta$



$\hat{\theta}_{t-1}$

Can we simplify $\sqrt{\beta_t(\delta)}$? So far, we did not use the fact that we want to set $V = \lambda I$. The results holds for any positive definit $V$. For $V = \lambda I \rightarrow \det(V) = \lambda^d$