

Guest lecture: Adaptive Online Learning and distribution shift
(Continuous)

Setting: Nature chooses P_1, P_2, \dots, P_n

1. $(x_t, y_t) \sim P_t$, x_t is revealed to learner
2. learner generates $\hat{y}_t = f_{\theta_t}(x_t) = \theta_t^T x_t$
chooses θ_t
3. incurs loss $\frac{(\hat{y}_t - y_t)^2}{2}$. $\mathcal{L}(x_t, y_t, \hat{y}_t)$

No assumption on $P_1 \dots P_n$

Example 0: $P_1 = P_2 = \dots = P_n$ iid setting

$$\underline{\sigma_1^*} = \sigma_{gm} E (x_{01}^T, y)^2$$

Example (Concurrence shift): $P_t(xy) = P_t(x) \cdot P(y|x)$

Example (concept shift): $P_+(x, y) = \overbrace{P(x)}^{\text{prior}} \cdot P(y|x)$

Example (Label shift) $P_+(x, y) = \underbrace{P_+(y)}_{\text{prior}} \cdot \underbrace{P(x|y)}_{\text{likelihood}}$

No regret Online learning:

1. (x_t, y_t) is chosen by Nature. (adversarial setting)
2. $\sim \theta_t \rightarrow \hat{y}_t$ _____
3. $\text{loss}(\hat{y}_t - y_t)^2$ _____

Static Regret: $\sum_{t=1}^n (x_t^\top \theta - y_t)^2 - \sum_{t=1}^n (x_t^\top u - y_t)^2 = \underbrace{\sum_{t=1}^n (x_t^\top \theta - y_t)^2}_{\text{your performance}} - \underbrace{\sum_{t=1}^n (x_t^\top u - y_t)^2}_{u = \arg \min_{\theta} \sum_{t=1}^n (x_t^\top \theta - y_t)^2} = \underbrace{O(n)}_{O(\log n)}$

Alge: Vork - ~~Azory~~ ~~Hermit~~ ~~Hermit~~ Forester
VdW

Alg 2: Online Newton Step

It doesn't work well with non-stationarity.

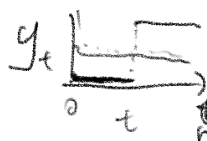
Example: $x_t = 1 \quad \forall t, \quad y_t = 0 \text{ for } t = 1, 2, \dots, \frac{n}{2}$
 $y_t = 1 \text{ for } t = \frac{n}{2} + 1, \dots, n$

$$l(\theta) = (x^T \theta - y)^2$$

$$\left[\nabla^2 l(\theta) \right] \geq \underline{xx^T} \geq 0$$

$$> \alpha \cdot \nabla l(\theta) \cdot \nabla l(\theta)$$

epiconvex



$$\min_{\theta} \frac{\sum_{t=1}^n (\theta - y_t)^2}{\frac{1}{4}n} \quad \theta^* = \frac{1}{2}$$

$$\sum_{t=1}^n (\theta - y_t)^2 \leq \frac{1}{4}n + o(\log n)$$

Dynamic Regret (Universal Dynamic Regret)

$$\text{DynReg}_{(U_{1:n})} = \sum_{t=1}^n (x_t^\top \theta_t - y_t)^2 - \sum_{t=1}^n (x_t^\top u_t - y_t)^2$$

$U_{1:n}$ is a seq of Comparators

(*) Sublinear Dyn Regret is impossible.

Example $x_t = 1 \forall t$, $y_t \sim \text{Ber}(0.5)$ i.i.d., $U_t = y_t$

$$\mathbb{E}[\text{Dynamic Regret}] \geq \Omega(n)$$

$$\mathbb{E}\left[\sum_{t=1}^n (y_t - y_t)^2 - \underbrace{(y_t - y_t)^2}_0\right] \geq \sum_{t=1}^n \mathbb{E}[(0.5 - y_t)^2] = \frac{1}{4}n$$

What can we still do?

Def: $TV(U_{1:n}) = \sum_{t=2}^n \|u_t - u_{t-1}\|_1$

we can parameterize Regret with $TV(U_{1:n})$

(Zinkevich, 2003) $\sqrt{n(1 + TV(U_{1:n}))}$

(Zhang, Zhou, 2016) $\sqrt{n(1 + TV(U_{1:n}))}$ is optimal for convex-losses

(Babu and W., 2021) $n^{\frac{1}{3}} (TV(U_{1:n}))^{\frac{2}{3}} + \log n$ optimal for exp-concave losses

$\text{poly}(d)$ dependence omitted

other description of $U_{1:n}$

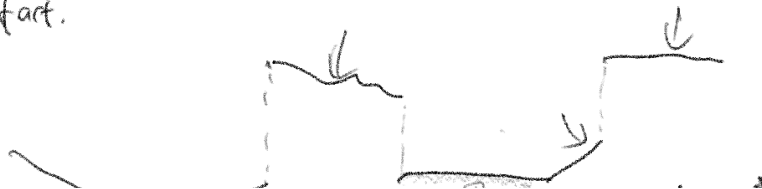
changes: $\ell_0\text{-TV}_{(U_{1:n})} = \sum_{t=2}^n \mathbb{1}(u_{t-1} \neq u_t)$

$\ell_2\text{-TV}_{(U_{1:n})} = \sum_{t=2}^n \|u_{t-1} - u_t\|_2$ "dimension-free"

In fact Babu and W.

D. Regret $\leq \log n \cdot (1 + \# \text{ of changes}) \wedge (n^{\frac{1}{3}} TV(U_{1:n})^{\frac{2}{3}} + \log n)$

in fact.



$3 \log n + n^{\frac{1}{3}} TV(\text{sequence, excluding change points})$

- nonparametric Regression

$$l_t(\theta) = (f_0(x_t) + \epsilon_t - \theta_t)^2$$

$$x_{t,n} = 0, \frac{1}{n}, \frac{2}{n}, \dots, \frac{n-1}{n}, 1$$

$$\underbrace{E\left[\sum_{t=1}^n (f_0(x_t) - \theta_t)^2\right]}_{MSE} = E\left[\sum_{t=1}^n (f_0(x_t) + \epsilon_t - \theta_t)^2\right] - \underbrace{E\left[\sum_{t=1}^n \epsilon_t^2\right]}_{n\sigma^2}$$

Optimal Rate for Hodder (dual) / Schabender, Boucheras, TV class $\theta_t = f_0(x_t)$

- Apply to non-stochastic control (LQR) Baby and W. (Maurits 22)
 \times Optimal Regret

Apply to nonstationary supervised learning

$$(x_t, y_t) \sim P_t, \quad \theta_t^* = E_{(x,y) \sim P_t} [(y - x^T \theta)^2]$$

"Covariate Shift" $y_t \sim N(\theta_0^T x_t, \sigma^2) \Rightarrow \theta_t^* \equiv \theta_0$
 also known give $O(\log n)$ static Regret

"Concept Shift" $x_t \sim P, y_t \sim N(\theta_t^T x_t, \sigma^2)$
 (quiz: concept + covariate shift) $\theta_t = \theta_t^*, \text{ Dyn. Regret } (\theta_t^* \cdot \theta_n^*) \leq O(n^{\frac{1}{3}} TV(\theta_{1:n}^*)^{\frac{2}{3}} + \log n)$

"label shift"

$$\theta_t^* = \underbrace{E_{y \sim P_t}}_{\arg \min_{\theta}} \underbrace{E_{x \sim P(x|y)}}_{\theta} [l((x,y), \theta)]$$

Assume L_t is λ -strongly convex $L_t(\theta)$

Lemma from Convex Opt: $\|\theta_t^* - \theta_{t-1}^*\| \leq \frac{\|\nabla_{\theta} (L_t - L_{t-1})(\theta_t^*)\|_*}{\lambda}$

$$L_t(\theta) = \sum_{y \in Y} P(y) \cdot E_x[d(x,y), \theta] = \langle P_t, E_x[d(x, \cdot), \theta] \rangle$$

$$\nabla_{L_t} L_t(\theta) = \langle P_t - P_{t-1}, \nabla_{\theta} E_x[d(x, \cdot), \theta] \rangle$$

$$\|\nabla(L_t - L_{t-1})(\theta)\| \leq \|P_t - P_{t-1}\| \cdot \max_y \|\nabla_{\theta} E[d(x,y), \theta_t^*]\|_*$$

$$TV(\theta_{t:n}^*) \leq \underbrace{\sum_t \|P_t - P_{t+1}\|_1}_{\text{cumulative label shift}} \cdot G \leq G$$

$$D. \text{Regret} \leq O\left(n^{\frac{1}{3}} \cdot TV(P_{t:n})^{\frac{2}{3}} \cdot G^{\frac{2}{3}}\right)$$

what if $\lambda = 0$?

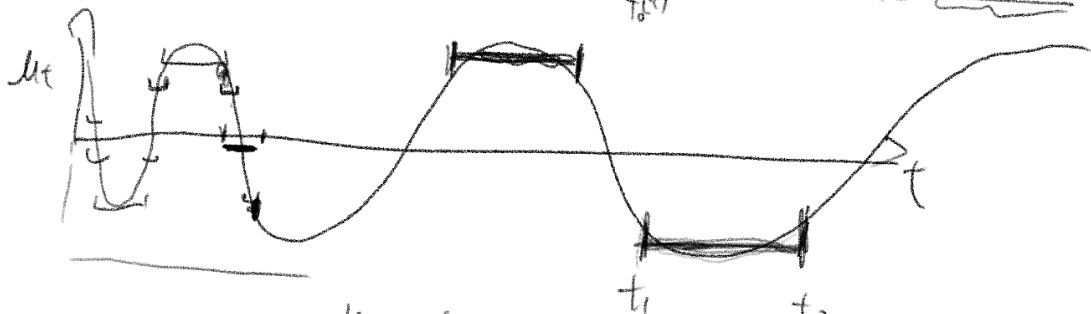
Compare with $\theta_{t+1}^* = \arg\min_{\theta} \mathbb{E}_{y \sim P_t} \mathbb{E}_X [\ell(x, y, \theta)] + \lambda \|\theta\|_2$

instantiate $u_t = \theta_t^*$ (θ_t^* is optimal for P_t)

$$\underbrace{\sum_{t=1}^n \ell(\theta_t, (x_t, y_t))}_{\text{your learner}} \leq \min_{u_{1:n}} \underbrace{\sum_t \ell(u_t, (x_t, y_t))}_{\downarrow} + \underbrace{\text{Regret}(u_{1:n})}_{\uparrow}$$

Algorithmic Idea (Strongly Adaptive Online Learner)

$$\ell(\theta_t) = (y_t - \theta_t)^2 \quad y_t \sim \mathcal{N}(\mu_t, \sigma^2) \quad f(t)$$



$$\text{Dynamic Regret} \leq \sum_{i=1}^k \sum_{t_i \leq t \leq t_{i+1}} \text{Dynamic Regret}(t_i : t_{i+1})$$

$$C_i \geq \frac{B^2 \epsilon}{\sqrt{n_i}}$$

Boudely

length of interval i

instantaneous regret bound

$$k \leq n^{\frac{1}{3}} TV(P_{t:n})^{\frac{2}{3}}$$

$$C_n \geq \sum_{i=1}^k C_i \geq \sum_{i=1}^k \frac{B^2}{\sqrt{n_i}} = k \sum_{i=1}^k \frac{1}{k} \frac{B^2}{\sqrt{n_i}} \geq k \frac{B^2}{\sqrt{n}}$$

1
"magical Rule"

$$k \leq \underline{\underline{C n^{\frac{2}{3}} n^{\frac{1}{3}} B^{\frac{4}{3}}}}$$

$$\begin{aligned} \frac{1}{\sqrt{x}} & \text{Expected} = \sqrt{\frac{1}{k}} \\ \text{Concave function} & = \frac{k}{\sqrt{k^2}} \sqrt{n} \\ \text{Convex} & = k^{\frac{2}{3}} B^{\frac{4}{3}} \\ & = \underline{\underline{\sqrt{n}}} \end{aligned}$$

① Adaptive Restarting: "change point detector" checking "magical Rule"

② "l-edge" over many leaves

$$\sqrt{n} \sqrt{\ln n}$$

len1 len2
↓
t=1

