

HW2. Problem:

$$\|x_i\|_{\infty} \leq X_{\infty}$$

* Thm: If A is OARO-stable w/ rate g , then

$$\mathbb{E}_{S \sim D^m} [L_D(A(S)) - L_S(A(S))] \leq g(m).$$

* l_2 -regularization gives stability:

Assume: ① $l(w, z)$ is ρ -Lipschitz wrt w , for any z

② $l(w, z)$ is convex wrt w , for any z .

③ $\hat{w} = A(S) = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left(\frac{\lambda}{2} \|w\|_2^2 + L_S(w) \right)$

① ② ③ $\Rightarrow A$ is $g(m) = \frac{2\rho^2}{\lambda m}$ - OARO stable.

want to find λ . such that it minimizes $\mathbb{E}_{S \sim D^m} L_D(w^*) + \frac{\lambda}{2} \|w^*\|_2^2$

$$\mathbb{E}_{S \sim D^m} L_D(A(S)) = \underbrace{\mathbb{E}_{S \sim D^m} L_S(A(S))}_{\text{can increase when } \lambda \uparrow} + \underbrace{\mathbb{E}_{S \sim D^m} [L_D(A(S)) - L_S(A(S))]}_{\leq \frac{2\rho^2}{\lambda m} \downarrow \text{ with } \lambda \uparrow}$$

output of \hat{w} (output of A).

optimizing w.r.t

$$L_S(A(S)) \leq F_S(A(S)) \leq F_S(w^*)$$

factory expectation:

$$\begin{aligned} \mathbb{E}_S L_S(A(S)) &\leq \mathbb{E}_S [F_S(w^*)] \\ &= \mathbb{E}_S \left[\frac{\lambda}{2} \|w^*\|^2 + \underbrace{L_S(w^*)}_{\frac{1}{m} \sum_{i=1}^m \ell(w^*; z_i)} \right] \\ &= \frac{\lambda}{2} \|w^*\|^2 + L_D(w^*) \end{aligned}$$

In summary:

$$\forall w^*, \quad \mathbb{E}_S L_D(A(S)) \leq L_D(w^*) + \underbrace{\frac{\lambda}{2} \|w^*\|^2}_{\uparrow} + \underbrace{\frac{2\rho^2}{\lambda m}}_{\downarrow} \Rightarrow 1$$

Interpretations:

① $\lambda = \frac{1}{\sqrt{m}}$ (or any other fn of m $\frac{1}{m} \ll \lambda \ll 1$)

$$\frac{1}{\sqrt{m}} \|w^*\|^2 + \frac{\rho^2}{\sqrt{m}}$$

$z \in \mathbb{N}$

$$\mathbb{E} L_D(A(S)) \leq \min_{w^*} \left(L_D(w^*) + \frac{\|w^*\|^2 + \rho^2}{\sqrt{m}} \right)$$



$$\frac{\beta^2 + \rho^2}{\sqrt{m}}$$

$$\geq \frac{\beta \cdot \rho}{\sqrt{m}}$$

This is a model selection result.

② Fix hypothesis class

$$\mathcal{H} = \{ w \in \mathbb{R}^d : \|w\|_2 \leq B \}.$$

$$\mathbb{E}_S L_D(A(S)) \leq \min_{w \in \mathcal{H}} L_D(w) + \underbrace{\frac{\lambda}{2} B^2 + \frac{2\rho^2}{\lambda m}}_{2 \cdot \rho \cdot B \cdot \sqrt{\frac{1}{m}}}.$$

$$\lambda^* = 2 \frac{\rho}{B} \sqrt{\frac{1}{m}}.$$

This is a PAC-like guarantee, except that the guarantee is in expectation, not in high probability.

$$m \geq \frac{4\rho^2 B^2}{\epsilon^2} \Rightarrow \mathbb{E}_S [L_D(A(S))] - \min_{w \in \mathcal{H}} L_D(w) \leq \epsilon.$$

Online learning:

Ex: 1. spam detection

At each time step $t=1, 2, \dots, T$

— receive an email $x_t \in \mathcal{X}$

— predict $\hat{y}_t \in \{ \underbrace{-1}_{\text{non-spam}}, \underbrace{+1}_{\text{spam}} \}$

— see $y_t \in \{-1, +1\}$. $f_t(w) = I(y_t \langle w, x_t \rangle \leq 0)$.

$$\text{minimize } \sum_{t=1}^T I(\hat{y}_t \neq y_t) = \ln(1 + e^{-y_t \langle w, x_t \rangle})$$

online classification. [regression.]

$$\Omega = \{ w : \|w\|_2 \leq B \}$$

$$f_w(x_t)$$

2. sequential investment (portfolio selection).

$W_1 =$ initial capital.

For $t = 1, 2 \dots T$

— decide $P_t \in \Delta^{N-1} = \{P \in \mathbb{R}_+^N : \sum_i P(i) = 1\}$
 for asset $i \in \{1, \dots, N\}$ investing $W_t \cdot P_t(i)$

— observe relative prices $r_t \in \mathbb{R}_+^N$.

$\prod_{t=1}^T \langle P_t, r_t \rangle$

maximizing

$f_t(P) = -\ln \langle P, r_t \rangle$

asset i : $W_t \cdot P_t(i) \cdot r_t(i)$

$W_{t+1} = \sum_i (W_t \cdot P_t(i) \cdot r_t(i)) = W_t \cdot \langle P_t, r_t \rangle$

Goal: maximize W_{T+1} .

3. aggregating weather prediction:

For each day $t = 1, 2 \dots T$

— obtain weather temperature predictions from N experts (model).

— make final prediction by following an expert drawn from $P_t \in \Delta^{N \times T} = \Omega$

— observe the losses of each model

$l_t \in [0, 1]^N$

$f_t(P) = \langle P, l_t \rangle$

Goal: minimize $\sum_t \langle p_t, l_t \rangle$

4. production recommendation (multi-armed bandits)

MAB.

$$f_t(a) = l_t(a)$$

②

$$p_t \in \Delta^{K-1}$$

For $t = 1, 2, \dots, T$

$$f_t(p) = \langle p, l_t \rangle$$

a new customer

— recommend $a_t \in \{1, \dots, K\}$, (product) to

— observe loss of a_t :

(e.g. clicked $\Rightarrow 0$
not clicked $\Rightarrow 1$)

$l_t(a_t)$
bandit feedback

$l_t \in [0, 1]^K$
full-info feedback.

Goal: minimize $\sum_t l_t(a_t)$

5. personalized product recommendation:

given N policies: $\pi^1 \dots \pi^N$ each $\pi^i: \mathcal{X} \rightarrow \{1, \dots, K\}$.

(take $N = K$. $\pi^i(x) \equiv i$)
 \Rightarrow MAB.

For $t = 1, 2, \dots, T$

— observe contextual info of customer x_t .

— random selection of N policies $\sim W_t \in \Delta^{N-1}$.

selected policy π_t . $f_t(w) = \mathbb{E}_{\pi \sim w} l_t(\pi(x_t))$

$$= \sum_{i=1}^N w_i l_t(\pi^i(x_t))$$

recommend $\pi_t(x_t)$

— observe loss $l_t(\pi_t(x_t))$ (we don't observe other entries of $l_t \in [0, 1]^K$)

Contextual bandits.

Online (convex) optimization:

- decision set Ω (action space) often convex

(Ω is a convex set: $\forall u, v \in \Omega, \forall \lambda \in (0,1)$
 $\lambda u + (1-\lambda)v \in \Omega$)



(convex)



(non convex)

For $t = 1, 2, \dots, T$.

- learner picks $w_t \in \Omega$
- environment picks loss fn $f_t: \Omega \rightarrow \mathbb{R}$.
- learner suffers loss $f_t(w_t)$,
& observes information on f_t .

performance measure: $\sum_{t=1}^T f_t(w_t)$.

setting of OL:

- stochastic: $\{f_t\}_{t=1}^T \stackrel{iid}{\sim} D$.
- oblivious adversary: $\{f_t\}_{t=1}^T$ chosen ahead of time.
- adaptive adversary: $\forall t, f_t$ can depend on $\{w_s\}_{s=1}^t$.

categorized by feedback model:

- full info: see f_t (or $\nabla f_t(w_t)$).

- bandit : $f_t(w_t) \in \mathbb{R}$
- other feedback settings.

key performance measure:

regret

$$R_T = \sum_{t=1}^T f_t(w_t) - \min_{w \in \Omega} \sum_{t=1}^T f_t(w)$$

(want = $o(T)$)

Next class:

Outline to batch conversion.

small $R_T \rightarrow$ get stat-learning alg. w/ small excess loss.

Next Thursday (Mar. 25): pre-recorded lecture.