

Online Optimization/Learning

• Learner vs. Environment

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

L selects $f_t \in \mathcal{F}$

E selects $l_t : \mathcal{F} \rightarrow \mathbb{R}$, reveals l_t^* to L

L incurs loss $l_t(f_t)$

*: full info feedback

Regret of the Learner after T rounds:

$$R_T := \sum_{t=1}^T l_t(f_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T l_t(f)$$

Want sublinear regret: $R_T/T \rightarrow 0$ as $T \rightarrow \infty$

• Online Conv Opt. (Zinkevich; Hazan, Abernethy, Rakhlin)

\mathcal{F} : closed bdd conv set in \mathbb{R}^d

$$D := \max \{ \|f - f'\| : f, f' \in \mathcal{F} \} < \infty$$

$(l_t)_{t \geq 0}$: conv, L -Lip, m -strongly convex

Projected GD:

$$f_{t+1} = \Pi(f_t - \alpha_t \nabla l_t(f_t))$$

$$R_T = \begin{cases} O(\sqrt{T}) & \text{for conv, Lip } (\alpha_t = \frac{1}{\sqrt{t}}) \\ O(\log T) & \text{for } m\text{-s.c., Lip. } (\alpha_t = \frac{c}{t}) \end{cases}$$

— can be shown that this is "optimal" ...

Perceptron Mistake Bound

(F. Rosenblatt, 1958)

data: $\tilde{z}_i = (\tilde{x}_i, \tilde{y}_i)$ $i \in [n]$
 $\tilde{x}_i \in \mathbb{R}^d$, $\tilde{y}_i \in \{\pm 1\}$

classifiers: $x \mapsto \text{sgn} \langle f, x \rangle$ ($f \in \mathbb{R}^d$)

Assume the data are separable: $\exists f^* \in \mathbb{R}^d$ s.t.
 $\tilde{y}_i \langle f^*, \tilde{x}_i \rangle \geq 1$ $\forall i \in [n]$

Goal: find $\hat{f} \in \mathbb{R}^d$ s.t. all examples are classified correctly

Perceptron: ($\alpha > 0$: parameter)

• $f_1 = 0$ (init)

• for $t=1, 2, \dots$:

$I_t \in [n] \leftarrow$ if $\exists i \in [n]$ s.t. $\tilde{y}_i \langle f_t, \tilde{x}_i \rangle < 0$,
then $I_t = i$
arbitrary o/w

$z_t = \tilde{z}_{I_t} = (\tilde{x}_{I_t}, \tilde{y}_{I_t})$

$f_{t+1} = \begin{cases} f_t & \text{if } y_t \langle f_t, x_t \rangle \geq 0 \\ f_t + \alpha y_t x_t & \text{if } y_t \langle f_t, x_t \rangle < 0 \end{cases}$

• if f_{t+1} makes no mistakes on data, STOP

Suppose L, B are chosen s.t.:

$$L \geq \max_{i \in [n]} \|\tilde{x}_i\|$$

$$B \geq \max \{ \|f^*\| : \min_{i \in [n]} \tilde{y}_i \langle f^*, \tilde{x}_i \rangle \geq 1 \}$$

[Then Perceptron will find a separating classifier in at most $L^2 B^2$ steps.]

Proof

1) map to Online Convex opt

time t :

- Learner generates f_t (acc. to algo)

- Environment responds with

$$l_t(f) := \begin{cases} 0 & \text{if } f_t \text{ is correct on } z_t \\ (1 - y_t \langle f, x_t \rangle)_+ & \text{o/w} \end{cases}$$

\downarrow
 $\in \mathbb{R}^d$

Note $l_t(\cdot)$ are convex, L -Lip.

$$\nabla l_t(f) = \begin{cases} 0, & \text{if } f_t \text{ is correct on } z_t \\ -y_t x_t & \text{o/w} \end{cases}$$

$$f_{t+1} = f_t - \alpha \nabla l_t(f_t) \quad : \text{perceptron } f_1 = 0$$

2) Let f^* be s.t. $\|f^*\| \leq B$

and $\tilde{y}_i \langle f^*, \tilde{x}_i \rangle \geq 1 \quad \forall i \in [n]$

Regret (in T rounds) w.r.t. f^*

$$V_t := \|f_t - f^*\|^2$$

$$V_1 = \|f^*\|^2 \leq B^2$$

$$V_{t+1} = \|f_{t+1} - f^*\|^2$$

$$= \|f_t - \alpha g_t - f^*\|^2$$

$$g_t := \nabla l_t(f_t)$$

$$= V_t - 2\alpha \langle g_t, f_t - f^* \rangle + \alpha^2 \|g_t\|^2$$

$$\leq V_t - 2\alpha \langle g_t, f_t - f^* \rangle + \alpha^2 L^2$$

$$\Rightarrow 2 \langle g_t, f_t - f^* \rangle \leq \frac{V_t - V_{t+1}}{\alpha} + \alpha L^2$$

$l_t(\cdot)$ convex \Rightarrow

$$l_t(f^*) - l_t(f_t) \geq \langle g_t, f^* - f_t \rangle$$

$$2(l_t(f_t) - l_t(f^*)) \leq \frac{V_t - V_{t+1}}{\alpha} + \alpha L^2$$

$$3) \quad l_t(f^*) = 0 \quad \forall t$$

suppose that f_1, \dots, f_T do not separate the data

$$\Rightarrow l_t(f_t) \geq 1 \quad \forall t$$

$$2 \sum_{t=1}^T [l_t(f_t) - l_t(f^*)] \geq 2T$$

$$2 \sum_{t=1}^T [l_t(f_t) - l_t(f^*)] \leq \frac{V_1}{\alpha} + \alpha L^2 T \leq \frac{B^2}{\alpha} + L^2 T$$

∴ if none of f_1, \dots, f_T separate the data, then

$$2T \leq \frac{B^2}{\alpha} + \alpha L^2 T \quad \forall \alpha > 0$$

$$2T \leq \inf_{\alpha > 0} \left(\frac{B^2}{\alpha} + \alpha L^2 T \right) = 2LB\sqrt{T}$$

$$\Rightarrow \sqrt{T} \leq LB \quad \square$$

Outline to Batch Generalization

$z_1, z_2, \dots, z_n \stackrel{iid}{\sim} P$

• $T = n$ rounds

• $f_1 \in \mathcal{F}$ (init)

• $t = 1, 2, \dots, n$:

L generates f_t

E "reveals" z_t

L incurs loss $l(f_t, z_t)$

$$f_{t+1} = A(f_1, \dots, f_t, z_1, \dots, z_t)$$

$$L(f) := \mathbb{E}_P[l(f, z)]$$

$$f_t = A(f_1, \dots, f_{t-1}, z_1, \dots, z_{t-1}) \perp\!\!\!\perp z_t$$

$$\Rightarrow \mathbb{E}[l(f_t, z_t) | z_1, \dots, z_{t-1}] = L(f_t) \quad \forall t$$

Now we can analyze $\frac{1}{T} \sum_{t=1}^T L(f_t)$

vs. $\frac{1}{T} \sum_{t=1}^T l(f_t, z_t)$

• Assume $l(f, z) \in [0, 1]$

• $(Y_t)_{t \geq 0}$: $Y_0 := 0$
 $Y_t := \sum_{s=1}^t [L(f_s) - l(f_s, z_s)]$

$$\mathbb{E}[L(f_s) - l(f_s, z_s) | z_1, \dots, z_{s-1}] = 0$$

$\Rightarrow (Y_t)_{t \geq 0}$ is a martingale

$$\mathbb{E}[Y_t | z_1, \dots, z_s] = Y_s \quad \text{if } 0 \leq s < t$$

By Azuma-Hoeffding inequality,

$$\mathbb{P}\{Y_T \geq \varepsilon T\} \leq \exp\left(-\frac{T\varepsilon^2}{2}\right) \quad \forall \varepsilon > 0$$

Set $\varepsilon = \sqrt{\frac{2 \log(1/\delta)}{T}}$, so

$$\frac{1}{T} Y_T \leq \sqrt{\frac{2 \log(1/\delta)}{T}} \quad \text{w.p. } \geq 1 - \delta$$

or

$$\frac{1}{T} \sum_{t=1}^T L(f_t) \leq \frac{1}{T} \sum_{t=1}^T l(f_t, z_t) + \sqrt{\frac{2 \log(1/\delta)}{T}}$$

Regret: $\frac{1}{T} \sum_{t=1}^T l(f_t, z_t) \leq \underbrace{\inf_{f \in \mathcal{F}} \frac{1}{T} \sum_{t=1}^T l(f, z_t)}_{\text{min-empirical risk}} + \frac{R_T}{T}$

• w.p. $\geq 1 - \delta$,

$$\frac{1}{T} \sum_{t=1}^T L(f_t) \leq \frac{1}{T} \inf_{f \in \mathcal{F}} L_T(f) + \frac{R_T}{T} + \sqrt{\frac{2 \log(1/\delta)}{T}}$$

• take f^* s.t. $L(f^*) = \min_{f \in \mathcal{F}} L(f)$

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T L(f_t, z_t) &\leq \frac{1}{T} \sum_{t=1}^T L(f^*, z_t) + \frac{R_T}{T} \\ &= L_T(f^*) + \frac{R_T}{T} \end{aligned}$$

$$\leq L(f^*) + \sqrt{\frac{2 \log(1/\delta)}{T}} + \frac{R_T}{T}$$

w.p. $\geq 1 - \delta$ by Hoeffding (z_1, \dots, z_T iid)

\Rightarrow w.p. $\geq 1 - 2\delta$,

$$\frac{1}{T} \sum_{t=1}^T L(f_t) \leq L(f^*) + \frac{R_T}{T} + \sqrt{\frac{8 \log(1/\delta)}{T}}$$

• if $f \mapsto L(f, z)$ is convex, then

$$\bar{f}_T := \frac{1}{T} \sum_{t=1}^T f_t \quad \text{satisfies}$$

$$L(\bar{f}_T) \leq \frac{1}{T} \sum_{t=1}^T L(f_t)$$

$$\leq L(f^*) + \frac{R_T}{T} + \sqrt{\frac{8 \log(1/\delta)}{T}} \quad \text{w.p.} \\ \geq 1 - 2\delta.$$