## 11.3 Generalization + online learning

$(\mathcal{F}, \mathcal{Z}, l)$          $l(f_t, z_t) =$ loss at time

Suppose $z_1, \dots z_n$ are independent, identically dist$^{d}$

Use projected gradient descent alg:

$f_0$ - fixed

$$f_{t+1} = \Pi\left(f_t + \alpha_t \nabla l(f_t, z_t)\right)$$

$\left(\begin{array}{c} \text{as in} \\ \text{Zinkevich} \end{array}\right)$

(Use each data sample once.)

**Thm 11.2** <u>Generalization</u> ability $\left[ l \in [0,1] \right]$
of on-line algorithm. Then for
any $T \geq 1$, with probability at least $1-\delta$:

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

$\underset{\substack{\uparrow \\ \text{generalization} \\ \text{risk}}}{}$          $\underset{\substack{\text{empirical} \\ \text{risk} \\ \text{(observed)}}}{}$

(b)
If $l(f,z)$ is convex in $f$ and $\bar{f}_T = \dfrac{1}{T}\displaystyle\sum_{t=1}^{T} f_t$

then          $L(\bar{f}_T) \leq \dfrac{1}{T}\displaystyle\sum_{t=1}^{T} l(f_t, z_t) + \sqrt{\dfrac{2 \log \frac{1}{\delta}}{T}}$

<u>Proof</u> Let $Y_0 = 0$, $Y_t = \sum_{s=1}^{t} \underbrace{L(f_s) - l(f_s, Z_s)}$

(conditional Mean zero)
given ft

so $E\left(Y_{t+1} - Y_t \,\middle|\, \begin{array}{c} f_1 \cdots f_t \\ Z_1 \cdots Z_t \\ Y_1 \cdots Y_t \end{array}\right) = 0$ i.e. $Y$ is a martingale

$\underbrace{Y_{t+1} - Y_t}$ martingale difference sequence with values in $\underline{[-1, 1]}$

$\left(\text{Hoeffding lemma: } E\left[e^{s(A - E[A])}\right] \le e^{-\frac{2s^2}{(a-b)^2}}\right)$

if $A \in [a, b]$

$E\left[e^{s(Y_{t+1} - Y_t)} \,\middle|\, past\right] \le e^{-\frac{2s^2}{4}} = e^{-s^2/2}$

$4 = (1 - (-1))^2$

So Azuma Hoeffding bound:

$E\left[Y_T \ge tT\right] \le e^{-\frac{2t^2T^2}{T \cdot 4}} = e^{-\frac{t^2T}{2}} = \underline{\delta}$

$P\left[\frac{Y_T}{T} \ge \sqrt{\frac{2 \log \frac{1}{\delta}}{T}}\right]$

$\boxed{\sqrt{\frac{2 \log \frac{1}{\delta}}{T}}} = t$

<u>Proof of (b)</u> $L(f) = E_Z \, l(f, Z)$ is

also convex in $f$.

So (b) follows from (a) by Jensen's inequality.

## Corollary 11.1

$$\bar{f}_n = \frac{1}{n} \sum_{t=1}^{n} f_t$$

Use projected gradient descent

$$\alpha_t = \frac{D}{L\sqrt{2t}}$$

$D =$ diameter of $\mathcal{F}$

$$= \max_{f, f' \in \mathcal{F}} \| f - f' \|$$

Assume $\ell(\cdot, z)$ is $L$-Lipschitz and convex for each $z$.

Then with probability at least $1 - 2\delta$

$$L(\bar{f}_n) \leq L^* + DL \sqrt{\frac{2}{n}} + \sqrt{\frac{8 \log 1/\delta}{n}}$$

---

**Proof**

(1)
$$L(\bar{f}_n) \leq \frac{1}{n} \sum_{t=1}^{n} \ell(f_t, z_t) + \sqrt{\frac{2 \log 1/\delta}{n}} \qquad \text{w.p.} \geq 1 - \delta$$

(generalization discussed above)

(2)
$$\frac{1}{n} \sum_{t=1}^{n} \ell(f_t, z_t) \leq \frac{1}{n} \sum_{t=1}^{n} \ell(f^*, z_t) + DL \sqrt{\frac{2}{n}}$$

(by Zinkevich)

(3)
$$\frac{1}{n} \underbrace{\sum_{t=1}^{n} \ell(f^*, z_t)}_{\substack{\text{independent} \\ \text{mean } L^*}} \leq L^* + \sqrt{\frac{2 \log 1/\delta}{n}} \qquad \text{w.p.} \geq 1 - \delta$$

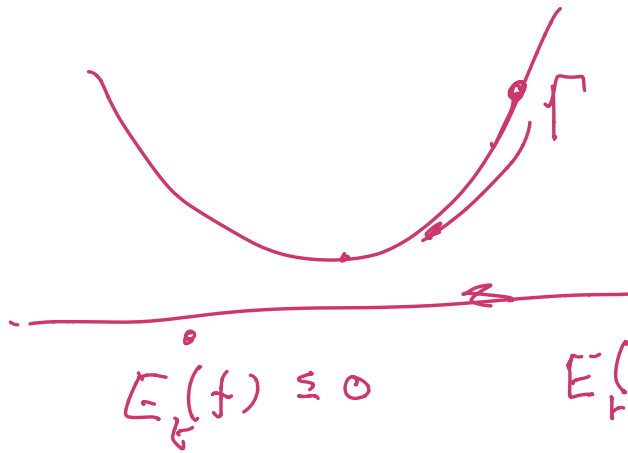$\longleftarrow$ Hoeffding inequality

Combine (1) − (3) to complete proof.

More fresh samples:

$Z_t, Z_{t+1}, \ldots Z_n$ are fresh samples

for $f_t$ so $L(f_t) \approx \frac{1}{n-t+1} \sum_{s=t}^{n} l(f_t, Z_s)$

# #3

(a) $\dot{f} = -\nabla \Gamma(f)$



$f(z_t) - f^* = O(\frac{1}{t})$

$\dot{E}_t(f) \leq 0$ $\qquad$ $E_t(f_t) \leq E_0(f_0)$

(b) — Accelerated gradient descent
(convex optimization) — rate $O(\frac{1}{t^2})$
convergence.

4. Zinkevich with $f^{*}$ — time verying

$$(f_t^*)$$

Constraint on
expert $\underline{\quad}$ $\sum_{t=1}^{T-1} \| f_{(t+1)}^* - f_t^* \| \leq W$

5. $\sqrt{T}$ is best possible for Zinkevich

6. Python SGD

# Chapter 12

## 12.1    Bound on average error probability for binary hypothesis testing

(Background)

**Observe** $Y$

$H_0:$ $Y$ has pmf $P_0$    $\leftarrow$ Bayesian prior $\pi_0 = \pi_1 = \frac{1}{2}$
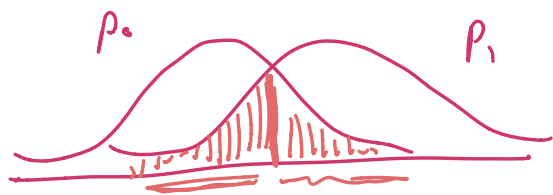
$H_1:$ $Y$ has pmf $P_1$    $\swarrow$

$f^*(y) \rightarrow \{0,1\}$    decision rule (classifier)

$$f^*(y) = \begin{cases} 1 & \text{if } P_1(y) > P_0(y) \\ 0 & \text{if } P_1(y) < P_0(y) \end{cases} \quad \left( \begin{array}{c} \text{Bayesian} \\ \text{rule} \end{array} \right)$$

$$\left( \overline{P_e} = \frac{1}{2} \sum_y P_0(y)\, f(y) + P_1(y)\,(1 - f(y)) \right.$$

$$\left( P\{ f(y) = H \} \qquad H = H_0 \text{ or } H_1 \right.$$

**Then**    $\overline{P_e^*} = \frac{1}{2} \sum_y P_0(y) \wedge P_1(y)$

$P_0$  $P_1$

$$\rho = \sum_y \sqrt{P_0(y) P_1(y)} \quad \leftarrow \text{Bhattacharyya coefficient}$$

$$\underline{\text{Lemma } 12.1} \qquad \frac{\rho^2}{4} \leq \boxed{P_e^*} \leq \frac{\rho}{2}$$
$$(c)$$

(b)

$$Y_1, \ldots, Y_n$$

$$H_0 \quad Y_i \sim P_{0,i} \quad 1 \leq i \leq n$$

$$H_1 \quad Y_i \sim P_{1,i} \quad 1 \leq i \leq n$$

$$\rho\left( P_{0,1}(y_1) \cdots P_{0,n}(y_n), \; P_{1,1}(y_1) \cdots P_{1,n}(y_n) \right)$$
$$= \prod_{j=1}^n \rho\left( P_{0,j}, P_{1,j} \right)$$