

$A_1, \dots, A_k$  events

$$P(\cup_i A_i) \leq P(A_1) + \dots + P(A_k)$$

---

Hoeffding inequality

$Z_1, \dots, Z_m$  independent, identically distributed  
(iid)

drawn from Bernoulli dist.

$$P(Z_i = 1) = \phi, \quad P(Z_i = 0) = 1 - \phi$$

$$\text{let } \hat{\phi} = \frac{1}{m} \sum_i Z_i$$

For  $\delta > 0$

$$P(|\phi - \hat{\phi}| > \delta) \leq 2e^{-2\delta^2 m}$$

Training set  $S = \{(x^i, y^i), i=1, \dots, m\}$

Samples drawn iid from prob dist  $D$

Training error for hypothesis  $h$

$$\hat{\epsilon}(h) = \frac{1}{m} \sum 1\{h(x^i) \neq y^i\}$$

Generalization error

$$\epsilon(h) = P_{(x,y) \sim D} (h(x) \neq y)$$

---

Ex linear classifier

$$h_{\theta}(x) = 1\{\theta^T x \geq 0\}$$

Pick  $\hat{\theta} = \operatorname{argmin}_{\theta} \hat{\epsilon}(h_{\theta})$

---

In general hypothesis class  $H$

# Finite $\mathcal{H}$

$$\mathcal{H} = \{h_1, \dots, h_k\}$$

Draw a sample, and let  $z$  denote whether  $h_i(x)$  misclassifies it, i.e.

$$z = 1\{h_i(x) \neq y\}$$

---

$$\hat{E}(h_i) = \frac{1}{m} \sum_j z_j$$

But this is just the est. mean of a Bernoulli distribution!

So

$$P(|E(h_i) - \hat{E}(h_i)| > \gamma) \leq 2e^{-\gamma^2 m}$$

Want this to be true for any  $h_i \in H$

$$P(\exists h \in H |\epsilon(h) - \hat{\epsilon}(h_i)| > \gamma) = P(A_1 \cup \dots \cup A_k)$$

~~###~~

$$\leq \sum_{i=1}^k P(A_k)$$

$$\leq \sum_{i=1}^k 2e^{-2\gamma^2 m}$$

$$= 2ke^{-2\gamma^2 m}$$

Subtract 1 from both sides...

$$\begin{aligned} & P(\neg \exists h \in H |\epsilon(h_i) - \hat{\epsilon}(h_i)| > \gamma) \\ &= P(\forall h \in H |\epsilon(h_i) - \hat{\epsilon}(h_i)| \leq \gamma) \geq 1 - 2ke^{-2\gamma^2 m} \end{aligned}$$

Q: Given  $\gamma$  and some  $\delta > 0$

how ~~much~~ large must  $m$  be  
to guarantee that w. Prob  $1 - \delta$   
the training error will be within  $\gamma$   
of generalization error?

$$A: m > \frac{1}{2\delta^2} \log \frac{2k}{\delta}$$

sample  
complexity

---

$$Q: |\hat{\epsilon}(h) - \epsilon(h)| \leq ???$$

---

$$h^* = \arg \min_{h \in \mathcal{H}} \epsilon(h)$$

$$\begin{aligned} \epsilon(\hat{h}) &\leq \hat{\epsilon}(\hat{h}) + \delta \\ &\leq \hat{\epsilon}(h^*) + \delta \\ &\leq \epsilon(h^*) + 2\delta \end{aligned}$$

---

Thm

$$\epsilon(\hat{h}) \leq \min_{h \in \mathcal{H}} \epsilon(h) + 2 \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$$

$k$  hypotheses result

## Case of $\infty$ h

Not the right argument but has the right intuitions

- $H$  parameterized by  $d$  real numbers
- Assume 64 bits per real #.

Then  $2^{64d}$  possibilities

So using sample complexity

$$m \geq O\left(\frac{1}{\gamma^2} \log_2 \frac{2^{64d}}{\delta}\right) = O\left(\frac{d}{\gamma^2} \log \frac{1}{\delta}\right) \\ = O_{\gamma, \delta}(d)$$

---

These notes follow

Andrew Ng's notes

## VC Dimension

$$|\epsilon(h) - \hat{\epsilon}(h)| \leq \sqrt{\left( \frac{d}{m} \log \frac{m}{d} + \frac{1}{m} \log \frac{1}{\delta} \right)}$$

$$() = \frac{1}{m} \log \left( \frac{m}{d} \right)^d + \frac{1}{m} \log \frac{1}{\delta}$$

$$= \frac{1}{m} \log \left( \frac{1}{\delta} \left( \frac{m}{d} \right)^d \right)$$