

§ 1

Motivation / Overview

1.

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\} \sim D^m, \quad y_i \in \{-1\}$$

H. Find good $h \in \mathcal{H}$ to reduce $R(h) = \mathbb{E}_{(x,y) \sim D} [1_{h(x) \neq y}]$.

2.

Reduce to ERM. $\underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{R}(h) = \frac{1}{m} \sum_{i=1}^m 1_{h(x_i) \neq y_i}$.

3.

Impose a structure on \mathcal{H} .

$\{\mathcal{H}_k\}_{k \in \mathbb{N}}$ be $\mathcal{H}_k \subseteq \mathcal{H}_{k+1}$, and $\mathcal{H} = \bigcup_{k \in \mathbb{N}} \mathcal{H}_k$.

Search for a good hypothesis set \mathcal{H}_k and also for a good hypothesis $h \in \mathcal{H}_k$.

- Some algorithm do this simultaneously:

- Structural Risk Minimization (SRM)
- Regularization

- Some does one after another

- Cross Validation (CV)
- n-fold CV

\mathcal{H}_k : model.

4.

Convex Surrogate Loss.

$h \in \mathcal{H}_k, f \in [X \rightarrow \mathbb{R}], h_f(x) = \operatorname{sgn}(f(x))$.

$$\hat{R}_s(h_f) \xrightarrow{?} \hat{R}_s^\Phi(f)$$

§2

Analysis of Error

1. $D \in \Pr(x|y), S \in D^m, H$.

2. $h_{\text{Bayes}} \in [g \rightarrow y]$, which has the minimum $R^* = R(h_{\text{Bayes}})$

$$h_{\text{Bayes}}(x) = \begin{cases} +1 & \text{if } P[y=+1|x] \geq P[y=-1|x] \\ -1 & \text{otherwise} \end{cases}$$

Bayes Error

3. $R(h_{\text{Bayes}}) = \underset{x \sim D_x}{\mathbb{E}} [\min \{P[y=+1|x], P[y=-1|x]\}]$.

4. $h \in H$... Result of algorithm

$$\text{Excess error : } R(h) - R^* = R(h) - \inf_{h \in H} R(h') + \inf_{h \in H} R(h') - R^*$$

$\underbrace{\quad}_{\text{estimation error}}$ $\underbrace{\quad}_{\text{approximation error}}$

§3

Empirical Risk Minimization Formally

1. $D \in \Pr(x|y), S \sim D^m, H$

2. $h_S^{\text{ERM}} = \underset{h \in H}{\operatorname{argmin}} \hat{R}_S(h)$.

3. What can we say about $R(h_S^{\text{ERM}})$?

Theorem

$$\Pr_{S \sim D^m} [R(h_S^{\text{ERM}}) \leq \inf_{h \in H} R(h) + 2R_m(H) + 2\sqrt{\frac{\log 2/\delta}{2m}}] \geq 1 - \delta.$$

$\delta = 2e^{-2m(\frac{\epsilon}{2} - R_m(H))}$

$$\Rightarrow R(h_S^{\text{ERM}}) - R^* \leq \inf_{h \in H} R(h) - R^* + (\quad) \text{ approx. error}$$

$$\Rightarrow \text{estimation} \leq 2R_m(H) + 2\sqrt{\frac{\log 2/\delta}{2m}} \text{ with prob } \geq 1 - \delta.$$

Equivalently, $\Pr_{S \sim D^m} [R(h_s^{ERM}) - \inf_{h \in H} R(h) > \varepsilon] \leq 2e^{-2m \left(\frac{C}{2} - R_m(H) \right)^2}$

Proposition 4.1

$$\Pr [R(h_s^{ERM}) - \inf_{h \in H} R(h) > \varepsilon] \leq \Pr [\sup_{h \in H} |R(h) - \hat{R}_s(h)| > \frac{\varepsilon}{2}]$$

$\exists h \in H \text{ s.t. } R(h) - \hat{R}_s(h) > \frac{\varepsilon}{2} \text{ or }$
 $R(h) - \hat{R}_s(h) < -\frac{\varepsilon}{2}.$

(Recall)

$$\Pr_{S \sim D^m} [\forall h \in H \quad R(h) \leq \hat{R}_s(h) + R_m(H) + \frac{\log(1/\delta)}{2m}] \geq 1 - \delta$$

$$\Pr_{S \sim D^m} [\forall h \in H \quad R(h) \geq \hat{R}_s(h) - R_m(H) - \frac{\log(1/\delta)}{2m}] \geq 1 - \delta$$

(pf of thm)

$$\Pr [R(h_s^{ERM}) - \inf_{h \in H} R(h) > \varepsilon] \leq \Pr [\exists h \in H \text{ s.t. } R(h) - \hat{R}_s(h) > \frac{\varepsilon}{2}]$$

$$+ \Pr [\exists h \in H \text{ s.t. } R(h) - \hat{R}_s(h) < -\frac{\varepsilon}{2}]$$

(By prop. 4.1)

$$= \Pr [\exists h \in H \text{ s.t. } R(h) > \hat{R}_s(h) + R_m(H) + \frac{\log(2/\delta)}{2m}]$$

$$+ \Pr [\quad " \quad R(h) < \hat{R}_s(h) - R_m(H) - \frac{\log(2/\delta)}{2m}]$$

$$\leq \frac{\delta}{2} + \frac{\delta}{2} = \delta \quad (\text{By thm 3.5}).$$

Thus, $\Pr [R(h_s^{ERM}) - \inf_{h \in H} R(h) \leq \varepsilon] \geq 1 - \delta$.

(pf of prop 4.1)

$$R(h_s^{ERM}) - \inf_{h \in H} R(h) \leq 2 \sup_{h \in H} |R(h) - \hat{R}_s(h)|$$

Take $K > 0$. $\exists h_k \in H$ s.t. $R(h_k) < \inf_{h \in H} R(h) + K$

$$R(h_s^{ERM}) - \inf_{h \in H} R(h) < R(h_s^{ERM}) - R(h_k) + K$$

$$= R(h_s^{ERM}) - \hat{R}_s(h_s^{ERM}) + \hat{R}_s(h_s^{ERM}) - R(h_k) + K$$

$$\leq R(h_s^{ERM}) - \hat{R}_s(h_s^{ERM}) + \hat{R}_s(h_k) - R(h_k) + K$$

$$\leq 2 \sup_{h \in H} |R(h) - \hat{R}_s(h)| + K$$

As $K > 0$ is arbitrary, done.

(another pf)

$$\begin{aligned}
 R(h_s^{\text{ERM}}) - R(h) &= R(h_s^{\text{ERM}}) - \hat{R}_s(h_s^{\text{ERM}}) + \hat{R}_s(h_s^{\text{ERM}}) - R(h) \\
 &\leq R(h_s^{\text{ERM}}) - \hat{R}_s(h_s^{\text{ERM}}) + \hat{R}_s(h) - R(h) \\
 &\leq 2 \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_s(h)|
 \end{aligned}$$

Take sup over $h \in \mathcal{H}$ on both side

$$\sup_{h \in \mathcal{H}} (R(h_s^{\text{ERM}}) - R(h)) = R(h_s^{\text{ERM}}) - \inf_{h \in \mathcal{H}} R(h) \leq 2 \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_s(h)|$$

§2.

Structural Risk Minimization (SRM)

①

$\mathcal{H} = \bigcup_{k=1}^K \mathcal{H}_k, \quad \mathcal{H}_k \subseteq \mathcal{H}_{k+1}$. Given $S \in (\mathcal{X} \times \mathcal{Y})^m$,

$$(k, h) \in \underset{k, h \in \mathcal{H}_k}{\operatorname{argmin}} (\hat{R}_s(h) + R_m(\mathcal{H}) + \sqrt{\frac{\log k}{m}}).$$

↳ Meaning of the term?

② Lemma

(Uniform Generalization Bound).

$$\underset{S \sim D^m}{\mathbb{P}} [\forall k, \forall h \in \mathcal{H}, \quad R(h) \leq \hat{R}_s(h) + R_m(\mathcal{H}) + \sqrt{\frac{\log k}{m}} + \sqrt{\frac{\log 2/\delta}{2m}}] \geq 1 - \delta.$$

③ Theorem 4.2

(Error Bound of SRM)

$$\underset{S \sim D^m}{\mathbb{P}} [R(h_s^{\text{SRM}}) \leq \inf_{h \in \mathcal{H}} (R(h) + 2R_m(\mathcal{H}_{k(h)}) + 2\sqrt{\frac{\log k(h)}{m}} + 2\sqrt{\frac{\log 4/\delta}{2m}})] \geq 1 - \delta,$$

where $k(h) = \min \{k \mid h \in \mathcal{H}_k\}$.

(pf of lemma)

$$\begin{aligned}
 &\underset{S \sim D^m}{\mathbb{P}} [\exists k, \exists h \in \mathcal{H}, \quad R(h) > \hat{R}_s(h) + R_m(\mathcal{H}) + \sqrt{\frac{\log k}{m}} + \sqrt{\frac{\log 2/\delta}{2m}}] \\
 &\leq \sum_k \underset{S \sim D^m}{\mathbb{P}} [\exists h \in \mathcal{H}, \quad R(h) > \hat{R}_s(h) + R_m(\mathcal{H}) + \varepsilon], \text{ where } \varepsilon = \sqrt{\frac{\log k}{m}} + \sqrt{\frac{\log 2/\delta}{2m}} \\
 &= \sum_k e^{-2m\varepsilon^2} = \sum_k e^{-2m(\sqrt{\frac{\log k}{m}} + \sqrt{\frac{\log 2/\delta}{2m}})^2} \\
 &\leq \sum_k e^{-2m(\frac{\log k}{m} + \frac{\log 2/\delta}{2m})} = \sum_k \frac{1}{k^2} \left(\frac{\delta}{2}\right) = \frac{\pi^2}{6} \cdot \frac{\delta}{2} < \delta.
 \end{aligned}$$

(pf of thm 4.2)

$$\text{Let } \varepsilon = 2\sqrt{\frac{\log 4/\delta}{2m}}$$

$$\mathbb{R}[R(h_s^{\text{SRM}}) - \inf_{h \in \mathcal{H}} (R(h) + 2R_m(\mathcal{H}_{k(h)}) + 2\sqrt{\frac{\log k(h)}{m}}) > \varepsilon]$$

Suppose $\varphi_0 \Rightarrow \varphi_1 \Rightarrow \varphi_2 \Rightarrow \dots, \varphi = \bigvee_{i=1}^{\infty} \varphi_i$. ↳ $F(h)$

Then, $\lim_{n \rightarrow \infty} P[\varphi_n] = P[\varphi]$.

$\exists h_1, h_2, h_3, \dots \in \mathcal{H}$ s.t. $F(h_1) \geq F(h_2) \geq \dots$ and $F(h_n) \rightarrow \inf_{\mathcal{H}_{\text{eff}}} F(h)$.
 If $R(h_s^{\text{SRM}}) - F(h_n) > \varepsilon$ $\forall n$, then $R(h_s^{\text{SRM}}) - \inf_{\mathcal{H}_{\text{eff}}} F(h) \geq \varepsilon$

Suffices to show $\mathbb{P}_{\mathcal{H}_{\text{eff}}}[R(h_s^{\text{SRM}}) - F(h) > \varepsilon] \leq \delta$.

$$\mathbb{P}[R(h_s^{\text{SRM}}) - (R(h) + 2R_m(H_{k(h)}) + 2\sqrt{\frac{\log k(h)}{m}}) > \varepsilon]$$

$$\leq \mathbb{P}\left[R(h_s^{\text{SRM}}) - (\hat{R}_s(h_s^{\text{SRM}}) + R_m(H_{k_s^{\text{SRM}}}) + \sqrt{\frac{\log k_s^{\text{SRM}}}{m}}) + (\hat{R}_s(h) + R_m(H_{k(h)}) + \sqrt{\frac{\log k(h)}{m}}) - (*) > \varepsilon\right]$$

$$= \mathbb{P}\left[R(h_s^{\text{SRM}}) - (\hat{R}_s(h_s^{\text{SRM}}) + R_m(H_{k_s^{\text{SRM}}}) + \sqrt{\frac{\log k_s^{\text{SRM}}}{m}}) - R(h) + \hat{R}_s(h) - R_m(H_{k(h)}) - \sqrt{\frac{\log k(h)}{m}} > \varepsilon\right]$$

$$\leq \mathbb{P}[R(h_s^{\text{SRM}}) - (\hat{R}_s(h_s^{\text{SRM}}) + R_m(H_{k_s^{\text{SRM}}}) + \sqrt{\frac{\log k_s^{\text{SRM}}}{m}}) > \frac{\varepsilon}{2}] + \mathbb{P}[-R(h) + \hat{R}_s(h) - R_m(H_{k(h)}) - \sqrt{\frac{\log k(h)}{m}} > \frac{\varepsilon}{2}]$$

$$\leq 2 \cdot e^{-2m(\frac{\varepsilon}{2})^2} \times 2 \quad (\text{By lemma}), \quad \text{Using } \sqrt{\frac{\log 2/\delta}{m}} = \frac{\varepsilon}{2}$$

$$= 4e^{-2m(\frac{\varepsilon}{2})^2} = \delta.$$

§

Cross Validation

①

$\{H_k\}_{k \in \mathbb{N}}$ with $H_k \subseteq \mathcal{H}_{\text{eff}}$.

$$\text{SRM} \cdots \arg \min_{k \in \mathbb{N}, H_k} \left(\hat{R}_s(h) + R_m(H_k) + \sqrt{\frac{\log k}{m}} \right).$$

②

(Cross Validation.)

Given $S \in (\mathbb{X} \times \mathbb{Y})^m$, $S = S_1 \cup S_2$, $|S_1| = (1-\alpha)m$, $|S_2| = \alpha m$

$$k \in \arg \min_k \hat{R}_s(h_{S_1, k}^{\text{ERM}}).$$

Return $h_{S_1, k}^{\text{ERM}}$.

Theorem 4.4

$$\mathbb{P}_{S \sim D^m} [R(h_s^{cv}) \leq \inf_k (R(h_{S_1, k}^{ERM}) + 2\sqrt{\frac{\log \max(k, k(h_s^{cv}))}{\alpha m}}) + 2\sqrt{\frac{\log 4/\delta}{\alpha m}}] \geq 1 - \delta.$$

$$\geq 1 - \delta.$$

$$\text{Proposition 4.3 } \mathbb{P} \left[\sup_k |R(h_{S_1, k}^{ERM}) - \hat{R}(h_{S_1, k}^{ERM})| - \sqrt{\frac{\log k}{\alpha m}} > \varepsilon \right] \leq 4e^{-2\alpha m \varepsilon^2}.$$

(pf of thm) $\mathbb{P} \left[\cup_k |R(h_{S_1, k}^{ERM}) - \hat{R}(h_{S_1, k}^{ERM})| \geq \varepsilon + \sqrt{\frac{\log k}{\alpha m}} \right] \geq 1 - 4e^{-2\alpha m \varepsilon^2}$.
 (By proposition 4.3). $\frac{\log 4/\delta}{\alpha m} = 1 - \delta$.

$$\begin{aligned} & \mathbb{P}[R(h_s^{cv}) \leq R(h_{S_1, k}^{ERM}) + 2\sqrt{\frac{\log \max(k, k(h_s^{cv}))}{\alpha m}} + 2\sqrt{\frac{\log 4/\delta}{\alpha m}}] \\ & \geq \mathbb{P}[R(h_s^{cv}) - \hat{R}_{S_1}(h_s^{cv}) + \hat{R}_{S_1}(h_{S_1}^{ERM}) - R(h_{S_1, k}^{ERM}) \leq 0] \\ & \leq 2\sqrt{\frac{\log(\max(k, k(h_s^{cv}))}{\alpha m}} + 2\sqrt{\frac{\log 4/\delta}{\alpha m}} \\ & \geq \mathbb{P}[R(h_s^{cv}) - \hat{R}_{S_1}(h_s^{cv}) \geq \sqrt{\log}] \end{aligned}$$

thm 4.4 $\Leftrightarrow \mathbb{P}[\cup_k R(h_s^{cv}) \leq R(h_{S_1, k}^{ERM}) + 2\sqrt{\frac{\log \max(k, k(h_s^{cv}))}{\alpha m}} + 2\sqrt{\frac{\log 4/\delta}{\alpha m}}] \geq 1 - \delta$

If ① \Rightarrow ②, done.

$$\begin{aligned} R(h_s^{cv}) - R(h_{S_1, k}^{ERM}) & \leq R(h_s^{cv}) - \hat{R}_{S_2}(h_s^{cv}) + \hat{R}_{S_2}(h_{S_1, k}^{ERM}) - R(h_{S_1, k}^{ERM}) \\ & \leq 2\sqrt{\frac{\log \max(k, k(h_s^{cv}))}{\alpha m}} + 2\sqrt{\frac{\log 4/\delta}{\alpha m}} \quad (\text{By ②}) \end{aligned}$$

(pf of prop. 4.3)

$$\mathbb{P} \left[\exists_k \mathbb{P}_{S \sim D^m} [R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})] > \sqrt{\frac{\log k}{\alpha m}} + \varepsilon \right]$$

$$\leq \sum_k \mathbb{P}_{S \sim D^m} [R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})] > \sqrt{\frac{\log k}{\alpha m}} + \varepsilon$$

$$= \sum_k \mathbb{E}_{S_1, S_2} [\mathbb{P}_{S_1} [R(h_{S_1, k}^{ERM}) - \hat{R}_{S_2}(h_{S_1, k}^{ERM})] > \sqrt{\frac{\log k}{\alpha m}} + \varepsilon | S_1]$$

$$= \sum_k \mathbb{E}_{S_i} 2e^{-2\alpha m \left(\frac{\log k}{\sqrt{\alpha m}} + \varepsilon \right)^2}$$

$$\leq \sum_k \mathbb{E}_{S_i} \frac{2}{k^2} e^{-2\alpha m \varepsilon^2} = \frac{\pi^2}{3} e^{-2\alpha m \varepsilon^2} \leq 4 e^{-2\alpha m \varepsilon^2}.$$

§ n-fold CV.

① Given $S \in (\mathcal{X} \times \mathcal{Y})^m$, $\{H_k\}_{k \in N}$.

$S = S_1 \cup S_2 \cup \dots \cup S_n$, $|S_i| = m_i$.

$$\hat{R}_{cv}(H_k) = \frac{1}{n} \sum_{i=1}^n \hat{R}_{S_i}(h_{S_i S_i, k}^{\text{ERM}})$$

$k \in \arg \min_k \hat{R}_{cv}(H_k)$. Return $h_{S_k}^{\text{ERM}}$.

② Consider case $m_1 = m_2 = \dots = m_n = \alpha$.

$$\mathbb{E}_S [\hat{R}_{cv}(H_k)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S_i} [\hat{R}_{S_i}(h_{S_i S_i, k}^{\text{ERM}})]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S_i S_i} \mathbb{E}_{S_i} [\hat{R}_{S_i}(h_{S_i S_i, k}^{\text{ERM}}) | S_i S_i]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S_i S_i} R(h_{S_i S_i, k}^{\text{ERM}})$$

$$= \mathbb{E}_{S' \sim \text{uniform}} R(h_{S', k}^{\text{ERM}}).$$