

HW1

Convexify the constraint

$$\rho = \frac{1}{\|\omega\|}$$

$$\Rightarrow \|\omega\|^2 = \frac{1}{\rho^2}$$

$g(\omega) \leq 0$: $\lambda g(\omega)$ is added to Lagrangian
where $\lambda \geq 0$.

$$g(\omega) = 0$$

$$\begin{aligned} \hookrightarrow \quad & g(\omega) \geq 0 \Rightarrow -g(\omega) \leq 0 \\ & \text{and } g(\omega) \leq 0 \quad (\lambda g(\omega) \text{ to Lagrangian where } \lambda \leq 0) \\ & \quad \quad \quad \downarrow \quad \quad \quad \hookrightarrow \end{aligned}$$

$\lambda g(\omega)$ to Lagrangian
where $\lambda \geq 0$

For an equality constraint, we add
 $\lambda g(\omega)$ to Lagrangian, $\lambda \in \mathbb{R}$.

Set up

Empirical risk minimization

Ingredients

Data: $(x_i, y_i) \quad i \in [m] \quad i = 1, 2, \dots, m.$

$(x_i, y_i) \sim \mathcal{D}$ iid.

$$x \in \mathcal{X} \quad y \in \mathcal{Y}$$

Loss: $\ell: (\mathcal{X} \times \mathcal{Y}) \times \mathcal{H} \rightarrow \mathbb{R}^+$

↓
hypothesis class

or
space of functions over
which we want to learn to
fit the data.

Mean-square loss:

$$\text{squared loss} \quad \ell((x, y), h) = \|h(x) - y\|^2$$

ERM:

$$\frac{1}{m} \sum_{(x, y) \in S} \ell((x, y), h) = \frac{1}{m} \sum_{(x, y) \in S} \|h(x) - y\|^2$$

MSE (pytorch.funct.)
error

ERM: solve optimization:

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{(x, y) \in S} \ell((x, y), h)$$

↑
sample set

Prob. dist of sample

$$S \sim \mathcal{D}^m = \mathcal{D} \times \dots \times \mathcal{D}$$

↓

joint prob. dist of $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

Goal: $\min_{h \in \mathcal{H}} \mathbb{E}_{(x, y) \sim \mathcal{D}} \ell((x, y), h) = \text{generalization error over } \mathcal{H}$

ML: h that "generalizes".

\mathcal{D} : unknown!

R.V.

Replacing / approximating $\mathbb{E}_{(x, y) \sim \mathcal{D}} \ell((x, y), h)$ with sample average.

convex in SLLN.

Emp risk of $h \xrightarrow{m \rightarrow \infty} \text{gen. error of } h.$

$$\text{Emp. risk: } \hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \ell((x_i, y_i), h)$$

$$S = \{(x_i, y_i)\}_{i \in [m]} = \frac{1}{m} \sum_{(x, y) \in S} \ell((x, y), h)$$

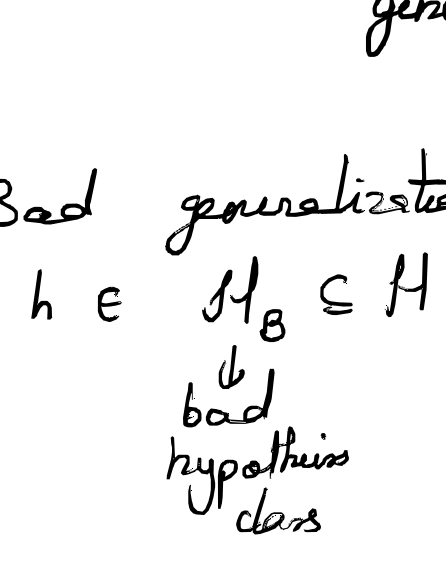
$$\text{Generalization error: } R(h) = \mathbb{E}_{(x, y) \sim \mathcal{D}} \ell((x, y), h)$$

Generalization gap:

$$R(h) - \hat{R}_S(h)$$

Overfitting: $\hat{R}_S(h) \approx 0$ but $R(h)$ is \uparrow .

Memorization: open questions! (generative models)



$$y = \begin{cases} +1 & \text{when } x \in \text{shaded square} \\ -1 & \text{o.w.} \end{cases}$$

$$\text{Ar}(\text{shaded square}) = \frac{1}{2} \text{Ar}(\text{square})$$

$$(x, y) \sim \text{Unif}(\text{square} \times \{-1, 1\})$$

My predictor:

$$h(x) = \begin{cases} y_i & x = x_i \\ 1 & \text{o.w.} \end{cases}$$

$$\text{Given } S = \{(x_i, y_i)\}_{i=1}^m, \hat{R}_S(h) = ? = 0.$$

$$R(h) = \frac{1}{2}.$$

Realizability assumption: $\exists h^* \in \mathcal{H}$ s.t.

$$R(h^*) = 0.$$

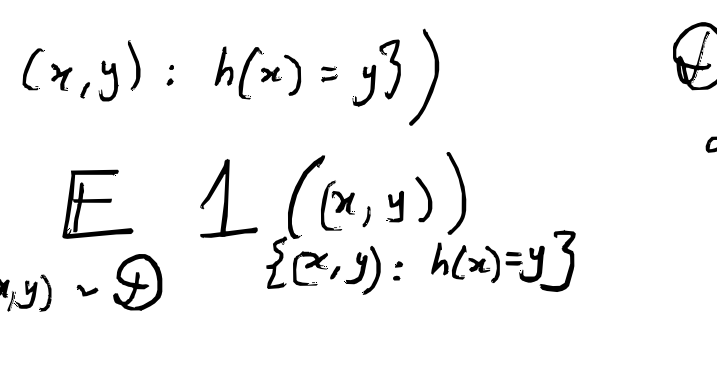
\Rightarrow For almost every S ,

$$\hat{R}_S(h) = 0 \text{ if } h_S \text{ is ERM solution for } S$$

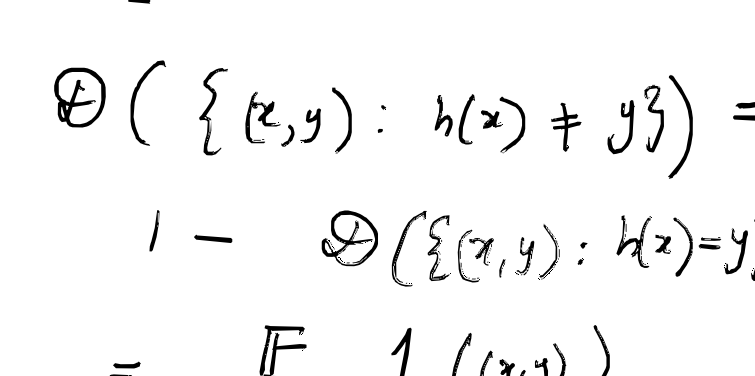
$$\text{i.e. } h_S = \arg \min_{h \in \mathcal{H}} \frac{1}{m} \sum_{(x, y) \in S} \ell((x, y), h)$$

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{(x, y) \in S} \ell((x, y), h) = 0.$$

$$S \supseteq \text{supp}(\mathcal{D}^m) \subseteq (\mathcal{X} \times \mathcal{Y})^m$$



For $h \in \mathcal{H}$,



When does solving an ERM lead to generalization?

Bad generalization: Fix $\epsilon > 0$.

$$h \in \mathcal{H}_B \subseteq \mathcal{H} \text{ if } R(h) > \epsilon.$$

↓
bad hypothesis class

$$\Pr(R(h) > \epsilon) = ?$$

$$\mathcal{D}^m(\text{bad samples}) =$$

$$\mathcal{D}^m(\{S \in (\mathcal{X} \times \mathcal{Y})^m : R(h_S) > \epsilon\}) = ?$$

$$\{S : R(h_S) > \epsilon\} \subseteq \{S : \hat{R}_S(h) = 0 \text{ and } h \in \mathcal{H}_B\}$$

$$\mathcal{D}^m(\text{bad samples}) \leq \mathcal{D}^m(\{S : \hat{R}_S(h) = 0 \text{ \& } h \in \mathcal{H}_B\})$$

Caveat: under realizability assumption, we have $\hat{R}_S(h) = 0$ a.s.

But in practice, "optimization error" needs to be considered.

$$\mathcal{D}^m(\text{bad samples}) \leq \bigcup_{h \in \mathcal{H}_B} \mathcal{D}^m(\{S : \hat{R}_S(h) = 0\})$$

$$\text{supp}(\mathcal{D}^m)$$

$$\text{Fix } h \in \mathcal{H}_B. \quad (\mathcal{H}_B = \{h \in \mathcal{H} : R(h) > \epsilon\})$$

$$\mathcal{D}^m(\{S : \hat{R}_S(h) = 0\})$$

$$= \mathcal{D}^m(\{S = \{(x_i, y_i)\}_{i \in [m]} : h(x_i) = y_i \forall i \in [m]\})$$

$$= (\mathcal{D}(\{(x, y) : h(x) = y\}))^m$$

(Read as Prob. of finding $(x, y) \sim \mathcal{D}$ s.t. $y = h(x)$)

\mathcal{D} : prob. measure

$$\mathcal{D}(\{ \}) = \text{Pr of "event" represented by set}$$

$$= \mathbb{E}_{(x, y) \sim \mathcal{D}} \mathbb{1}_{\{ \}}$$

$$h \in \mathcal{H}_B$$

$$\mathcal{X} \times \mathcal{Y}$$

$$\mathcal{D}(\{(x, y) : h(x) = y\}) \quad \mathcal{D}: \text{data dist.}$$

$$= \mathbb{E}_{(x, y) \sim \mathcal{D}} \mathbb{1}_{\{(x, y) : h(x) = y\}}$$

$$(\text{Indicator function } \mathbb{1}_A((x, y)) = \begin{cases} 1 & (x, y) \in A \\ 0 & \text{o.w.} \end{cases})$$

$$= \mathcal{D}(\{(x, y) : h(x) \neq y\}) =$$

$$1 - \mathcal{D}(\{(x, y) : h(x) = y\})$$

$$= \mathbb{E}_{(x, y) \sim \mathcal{D}} \frac{\mathbb{1}_{\{(x, y) : h(x) \neq y\}}}{\mathbb{1}_{\{(x, y) : h(x) \neq y\}}} \quad (\text{zero-one loss})$$

$$= \mathbb{E}_{(x, y) \sim \mathcal{D}} \ell((x, y), h)$$

$$= R(h) > \epsilon$$

$$\Rightarrow (\mathcal{D}(\{(x, y) : h(x) = y\}))^m \leq (1 - \epsilon)^m$$

$$\Rightarrow \bigcup_{h \in \mathcal{H}_B} (\mathcal{D}(\{(x, y) : h(x) = y\}))^m \leq |\mathcal{H}_B| (1 - \epsilon)^m$$

(Recall Union bound) cardinality of set \mathcal{H}_B

We have shown,

$$\mathcal{D}^m(\{\text{bad samples}\}) \leq |\mathcal{H}_B| (1 - \epsilon)^m$$

$$\leq |\mathcal{H}| e^{-\epsilon m}$$

$$(1 - \epsilon)^m \leq e^{-\epsilon m}.$$

In other words,

$$1 - \mathcal{D}^m(\{\text{bad samples}\}) \geq 1 - |\mathcal{H}| e^{-\epsilon m}$$

For every $\epsilon > 0$,

with Prob $\geq 1 - |\mathcal{H}| e^{-\epsilon m}$, the

generalization error $R(h_S) < \epsilon$ for the ERM h_S of S

$$\text{Set } \delta = |\mathcal{H}| e^{-\epsilon m}$$

$$\log \delta = \log |\mathcal{H}| - \epsilon m$$

$$\epsilon m = \frac{\log |\mathcal{H}|}{\delta}$$

$$\epsilon = \frac{1}{m} \frac{\log |\mathcal{H}|}{\delta}$$

Alt: For any $\delta \in (0, 1)$, with Prob $\geq 1 - \delta$ over training samples $S \sim \mathcal{D}^m$,

$$R(h_S) < \frac{1}{m} \frac{\log |\mathcal{H}|}{\delta} \text{ for an}$$

ERM h_S over the training samples S .

(PAC)

Probably approximately correct bounds