

Empirical risk minimization

loss function / risk : on data space \mathcal{X}

$l(x, y, h) \in \mathbb{R}^+$

↓
hypothesis
(function to be learned from data, to mimic/approx target)

target
↓
 $f(x) = y$

$x_i, y_i = \{-1, 1\}$
↓ ↑
inputs labels

Want: find h s.t.
 $h(x) \approx y$

Given: $(x_i, y_i = f(x_i))_{i \in [m]}$

$(x_i, y_i) \sim \mathcal{D}$ (data distribution)
probability measure

goal: Find an h
s.t. $\min_{h(x,y) \sim \mathcal{D}} \mathbb{E} l(x, y, h)$

loss fn: $l(x, y, h) = 0$ if $f(x) = h(x) = y$.

Hypothesis class : function class where h lives.

Goal: $\min_{h \in \mathcal{H}(x,y)} \mathbb{E} l(x, y, h)$
↓
hypothesis class

Examples of loss functions

(for binary classification)

- $l(x, y, h) = \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{o.w.} \end{cases}$

zero-one loss.

$\min_{h \in \mathcal{H}(x,y) \sim \mathcal{D}} \mathbb{E} \mathbb{1}_{\{h(x) \neq y\}}$

0-1 loss hinge loss

Recall

$h(x) = (\omega^T x + b)$

Soft SVM:

Constraints:

$$y_i (\omega^T x_i + b) \geq 1 - \xi_i$$

$$y_i h(x_i) \leq 1 - \xi_i$$

obj fn: $\min \|\xi\|$

$l(x, y, h) = \max \left\{ 0, 1 - \frac{y h(x)}{\rho} \right\}$

Goal: $\min_{h \in \mathcal{H}(x,y) \sim \mathcal{D}} \mathbb{E} l(x, y, h)$

Hypothesis classes

Bin. class: $h(x) = \text{sgn}(\omega^T x + b)$
SVM

parameterized by (ω, b)
 \mathbb{R}^d, \mathbb{R}

Goal: $\min_{h \in \mathcal{H}(x,y) \sim \mathcal{D}} \mathbb{E} l(x, y, h)$: generalization error

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

$\mathbb{E}_{(x,y) \sim \mathcal{D}} l(x, y, h)$: gene. err of h .

\mathcal{D} : unknown Empirical risk minimization (ERM)

Actual: $\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m l(x_i, y_i, h)$

$\frac{1}{m} \sum_{i=1}^m l(x_i, y_i, h) \xrightarrow{m \rightarrow \infty} \mathbb{E}_{(x,y) \sim \mathcal{D}} l(x, y, h)$ (LLN)

ML: Solving ERM's.

Supervised learning: data are of the form (x_i, y_i)
↑
 y_i : evaluations of the target function

Unsupervised learning: data are of the form $\{x_i\}_{i=1}^m$

how to solve ERM

Optimization:

- Parameterize \mathcal{H}
e.g. NNs, kernels
- w : set of parameters $\in \mathbb{R}^p$

ERM

$$\min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m l(x_i, y_i, h) = \min_{w \in \mathbb{R}^p} \frac{1}{m} \sum_{i=1}^m l(x_i, y_i, h(x_i, w))$$

$x \rightarrow h(x, w)$: on input space
 $w \rightarrow h(x, w)$: on parameter space

Overparameterize : $p \gg d \times m$

- GD : $R_S(h)$
 $R_S(w) = \frac{1}{m} \sum_{i=1}^m l(x_i, y_i, h(x_i, w))$
↑
empirical risk $S = \{(x_i, y_i)\}_{i \in [m]}$

GD Update : $w_{t+1} = w_t - \eta_t \nabla_w R_S(w_t)$
↓
learning rate

Stochastic gradient descent

SGD update : $w_{t+1} = w_t - \eta_t \hat{\nabla}_w R_S(w_t)$
↑
(noisy estimate of true gradient)

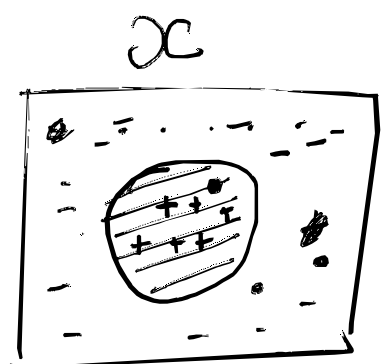
$\hat{\nabla}_w R_S(w_t) = \frac{1}{|I|} \sum_{i \in I} \nabla_w l(x_i, y_i, h(x_i, w))$

Workhorse
AdamW (momentum)
accelerated version of GD/SGD.

If you solve ERM, do you have guarantees on the generalization error?

Ben-David, Shalev-Schwartz

Example:



$\mathcal{D} \equiv \text{Unif.}$

$$A_r(\text{shaded circle}) = \frac{A_r(X)}{2}$$

$$\text{Target function: } f(x) = \begin{cases} 1 & x \in \text{shaded circle} \\ -1 & \text{o.w.} \end{cases}$$

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

$$\begin{matrix} \nearrow \\ \text{training} \\ \text{sample} \end{matrix} S = \{(x_i, y_i)\}_{i=1}^m \sim \mathcal{D}^m$$

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

Empirical risk

$$R_S(h) = 0.$$

Generalization error

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(x,y, h) = ? = 1/2.$$

Overfitting / memorization:

$R_S(h) = 0$ but gen. error is high.

$$\begin{aligned} R_S(h) &= \frac{1}{m} \sum_{i=1}^m \ell((x_i, y_i), h) \\ &= \frac{1}{m} \sum_{i=1}^m 0 = 0. \end{aligned}$$