

16.1 Empirical risk

First, we introduce some notations,

- Data point $(x, y) \sim P \in \mathcal{P}$,
where P is the joint distribution and \mathcal{P} is the space of distributions.
- Loss function $\ell : \mathcal{Y} \times \mathcal{U} \rightarrow [0, 1]$,
 \mathcal{Y} is the space of true labels and \mathcal{U} is the space of predicted labels
- Hypothesis space $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{U}$
- Data $Z^n = (Z_1, \dots, Z_N) \stackrel{iid}{\sim} P$

Definition 16.1. *some function $f \in \mathcal{F}$, the empirical risk is*

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(y_i, f(x_i))] \quad (16.1)$$

Now we have,

$$\begin{aligned} \mathbb{E}[\hat{R}_n(f)] &= \frac{1}{n} \sum_{i=1}^n n \mathbb{E}[\ell(y_i, f(x_i))] \\ &= \mathbb{E}[\ell(y_i, f(x_i))] \\ &= R(f) \end{aligned}$$

So based on the Hoeffding inequality,

$$P(|\hat{R}_n(f) - R(f)| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

$$\implies \text{w.p. } 1 - 2e^{-2n\epsilon^2}$$

$$|\hat{R}_n(f) - R(f)| < \epsilon$$

We could find $\hat{R}_n(f)$ for any $f \in \mathcal{F}$, then why not $\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}_n(f)$

Here is how the logic goes: $f^* := \operatorname{argmin}_{f \in \mathcal{F}} R(f)$, and

$$R(\hat{f}_n) \approx \hat{R}_n(\hat{f}_n) \approx \hat{R}_n(f^*) \approx R(f^*)$$

However, this doesn't work. We need conditions on $\mathcal{P}, \mathcal{F}, \ell$. First we introduce the "induced loss class",

$$\mathcal{L}(\mathcal{F}) = \{\ell(\cdot, f(\cdot)) : f \in \mathcal{F}\}$$

and

$$q(n, \epsilon) := \sup_{P \in \mathcal{P}} P^n(Z^n \in \mathbb{Z}^n : \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)| \geq \epsilon)$$

This is corresponding to the worst case probability distribution over all \mathcal{P} of "bad" samples (those with worst/large risk deviation over \mathcal{F}).

We say $\mathcal{L}(\mathcal{F})$ has "uniform convergence of empirical means" (UCEM) property w.r.t. \mathcal{P} , if

$$\lim_{n \rightarrow \infty} q(n, \epsilon) = 0 \quad \forall \epsilon > 0$$

Theorem 16.2. *if $\mathcal{L}(\mathcal{F})$ has UCEM property the ERM is PAC.*

Proof.

$$\begin{aligned} R(\hat{f}_n) - R(f^*) &= R(\hat{f}_n) - \hat{R}_n(\hat{f}_n) + \hat{R}_n(\hat{f}_n) - \hat{R}_n(f^*) + \hat{R}_n(f^*) - R(f^*) \\ \hat{R}_n(\hat{f}_n) - \hat{R}_n(f^*) &\leq 0 \quad \text{by ERM} \\ R(\hat{f}_n) - \hat{R}_n(\hat{f}_n) &\leq \sup_f (R(f) - R_n(f)) \leq \sup_f |R(f) - \hat{R}_n(f)| \\ \hat{R}_n(f^*) - R(f^*) &\leq \sup_f (R(f) - R_n(f)) \leq \sup_f |R(f) - \hat{R}_n(f)| \end{aligned}$$

So,

$$R(\hat{f}) - R(f^*) \leq 2 \sup_f |R(f) - \hat{R}_n(f)|$$

then with the UCEM property we have, $\exists n_0(\epsilon, \delta)$ s.t.

$$q(n, \epsilon/2) \leq \delta \quad \forall n > n_0$$

■

16.2 ERM algorithm

Now, let's consider an algorithm,

$$\begin{aligned} r_{\mathcal{A}}(n, \epsilon) &= \sup_{P \in \mathcal{P}} P^n(Z^n : R(\hat{f}) \geq R(f^*) + \epsilon) \\ &\leq \sup_{P \in \mathcal{P}} P^n(Z^n : \sup_{f \in \mathcal{F}} |\hat{R}_n(\hat{f}) - R(f)| \geq \epsilon/2) \\ &= q(n, \epsilon/2) \leq \delta \end{aligned}$$

UCEM is sufficient for the ERM algorithm to "work". Then our question is, when do we have UCEM? It turns out $\mathbb{E} \left[\sup_{f \in \mathcal{F}} |R(f) - \hat{R}_n(f)| \right] = o(1)$ is sufficient for UCEM when $\ell : \mathcal{Y} \times \mathcal{U} \rightarrow [0, M]$.

Further more, we define

$$\Delta_n = \sup_f |R(f) - \hat{R}_n(f)|$$

Then we have the following with McDiarmid's inequality,

$$P^n(\Delta_n - E\Delta_n \geq t) \leq e^{-2nt^2} \quad \forall \ell \rightarrow [0, 1]$$

What we want is $P^n(\Delta_n \geq \epsilon) \rightarrow 0$

$$\begin{aligned}
P^n(\Delta_n \geq \epsilon) &= P^n(\Delta_n - E\Delta_n \geq \epsilon - E\Delta_n) \\
&\quad \exists \text{ n.s.t. } E\Delta_n < \epsilon/2 \\
&\leq P^n(\Delta_n - E\Delta_n \geq \epsilon/2) \leq e^{-n\epsilon^2/2} \\
&\leq f \quad \text{with n large enough}
\end{aligned}$$

When is $E\Delta_n$ is small?

16.3 Rademacher complexity

Theorem 16.3. *ERM satisfies*

- $R(\hat{f}_n) \leq R(f^*) + 2\Delta_n(\mathcal{F})$,
- $R(\hat{f}_n) \leq \hat{R}_n(\hat{f}_n) + \Delta_n(\mathcal{F})$.
- $\Delta_n = \sup_{f \in \mathcal{F}} |R(f) - \hat{R}_n(f)|$

We need to control Δ_n ,

Definition 16.4.

$$\hat{\mathfrak{R}}_n(\mathcal{F}) = \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right| \middle| Z^n \right] \quad (16.2)$$

σ_i 's are Rademacher variables, $P(\sigma_i = 1) = P(\sigma_i = -1) = \frac{1}{2}$

- Think of this as “projection” of \mathcal{F} on to Z^n
- “Correlation” of \mathcal{F} with noise
- $\sup_f |f| \leq M$ or this doesn't work
- “empirical Rademacher complexity”

Definition 16.5. *Expected Rademacher complexity*

$$\mathfrak{R}_n(\mathcal{F}) = \mathbb{E}_{P^n} \left[\hat{\mathfrak{R}}_n(\mathcal{F}) \right]$$

Theorem 16.6.

$$\mathbb{E}[\Delta_n(\mathcal{F})] \leq \mathfrak{R}_n(\mathcal{F})$$

Proof. (by symmetrization) First we need Z^n and “ghost” samples \tilde{Z}^n

$$\begin{aligned}
\Delta_n(\mathcal{F}) &= \sup_{f \in \mathcal{F}} |R(f) - \hat{R}_n(f)| \\
&= \sup_{f \in \mathcal{F}} \tilde{\mathbb{E}} \left[|\tilde{R}_n(f) - \hat{R}_n(f)| \right] \\
&\leq \tilde{E} \left[\sup_{f \in \mathcal{F}} |\tilde{R}_n(f) - \hat{R}_n(f)| \right]
\end{aligned}$$

So,

$$\mathbb{E}_{Z^n} \Delta_n(\mathcal{F}) = \mathbb{E}_{Z^n} \mathbb{E}_{\tilde{Z}^n} \left[\sup_{f \in \mathcal{F}} |\tilde{R}_n(f) - \hat{R}_n(f)| \right]$$

Because Z^n and \tilde{Z}^n are independent, with the symmetrization

$$\begin{aligned} \tilde{R}_n(f) - \hat{R}_n(f) &= \frac{1}{n} \sum_{i=1}^n f(\tilde{Z}_i) - f(Z_i) \\ &\stackrel{dist}{=} \frac{1}{n} \sum_{i=1}^n f(Z_i) - f(\tilde{Z}_i) \\ &\stackrel{dist}{=} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(Z_i) - f(\tilde{Z}_i)) \end{aligned}$$

So,

$$\begin{aligned} \mathbb{E}_{Z^n} \Delta_n(\mathcal{F}) &= \mathbb{E}_{Z^n} \mathbb{E}_{\tilde{Z}^n} \left[\left| \frac{1}{n} \sum_{i=1}^n \sigma_i(f(Z_i) - f(\tilde{Z}_i)) \right| \right] \\ &\leq \mathbb{E}_{Z^n \sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i(f(Z_i)) \right| \right] + \mathbb{E}_{\tilde{Z}^n \sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i(f(\tilde{Z}_i)) \right| \right] \\ &= 2 \mathbb{E}_{Z^n \sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i(f(Z_i)) \right| \right] \\ &= \mathbb{E}_{Z^n} \left[\mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i(f(Z_i)) \right| \middle| Z^n \right] \right] \\ &= \mathfrak{R}_n(\mathcal{F}) \end{aligned}$$

■

References