

Statistical Learning Theory

Lijun Zhang

Nanjing University, China

August 13, 2017

Outline

1 Introduction

2 Formalization

3 Analysis

4 Discussions

Outline

1 Introduction

2 Formalization

3 Analysis

4 Discussions

Introduction

Nothing is more practical than a good theory!

[Vapnik, 1998]

Introduction

Nothing is more practical than a good theory!

[Vapnik, 1998]

■ The Goal

- Provide a framework for studying the problem of inference

■ The Methodology

- Assumptions of statistical nature about the underlying phenomena

■ No Free Lunch

- If there is no assumption on how the **past** (i.e. training data) is related to the **future** (i.e. test data), prediction is impossible.
- If there is no **a priori restriction on the possible phenomena** that are expected, it is impossible to generalize and therefore is thus no better algorithm

Outline

1 Introduction

2 Formalization

3 Analysis

4 Discussions

Formalization [Bousquet et al., 2004]

■ Notations

- Input space: \mathcal{X} , Output Space $\mathcal{Y} = \{-1, 1\}$
- Function: $g : \mathcal{X} \mapsto \mathcal{Y}$
- Risk of g :

$$R(g) = \Pr(g(X) \neq Y) = \mathbb{E} [\mathbb{1}_{g(X) \neq Y}]$$

- Regression function:

$$\eta(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}] = 2 \Pr[Y = 1|X = \mathbf{x}] - 1$$

- Target function (or Bayes classifier):

$$t(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}))$$

Formalization [Bousquet et al., 2004]

■ Notations

- Input space: \mathcal{X} , Output Space $\mathcal{Y} = \{-1, 1\}$
- Function: $g : \mathcal{X} \mapsto \mathcal{Y}$
- Risk of g :

$$R(g) = \Pr(g(X) \neq Y) = \mathbb{E} [\mathbb{1}_{g(X) \neq Y}]$$

- Regression function:

$$\eta(\mathbf{x}) = \mathbb{E}[Y|X = \mathbf{x}] = 2\Pr[Y = 1|X = \mathbf{x}] - 1$$

- Target function (or Bayes classifier):

$$t(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}))$$

Theorem 1

$t(\cdot)$ is the optimal function, i.e.,

$$R(t) = \inf_g R(g), \quad t \in \operatorname{argmin}_g R(g)$$

Proof of Theorem 1

$$\begin{aligned} R(g) &= \Pr(g(X) \neq Y) \\ &= E[\mathbb{1}_{g(X) \neq Y}] \\ &= E_X [E_Y [\mathbb{1}_{g(X) \neq Y} | X]] \\ &= E_X [\Pr(g(X) \neq Y | X)] \end{aligned}$$

Proof of Theorem 1

$$\begin{aligned}
 R(g) &= \Pr(g(X) \neq Y) \\
 &= \mathbb{E} [\mathbb{1}_{g(X) \neq Y}] \\
 &= \mathbb{E}_X [\mathbb{E}_Y [\mathbb{1}_{g(X) \neq Y} | X]] \\
 &= \mathbb{E}_X [\Pr(g(X) \neq Y | X)]
 \end{aligned}$$

Then, $t(\cdot)$ minimizes $R(\cdot)$ if

$$\begin{aligned}
 t(\mathbf{x}) &= \operatorname{argmin}_{y \in \{-1, 1\}} \Pr(y \neq Y | X = \mathbf{x}) \\
 &= \begin{cases} 1, & \Pr[Y = 1 | X = \mathbf{x}] \geq \Pr[Y = -1 | X = \mathbf{x}] \\ -1, & \Pr[Y = 1 | X = \mathbf{x}] < \Pr[Y = -1 | X = \mathbf{x}] \end{cases} \\
 &= \operatorname{sign}(2 \Pr[Y = 1 | X = \mathbf{x}] - 1)
 \end{aligned}$$

Empirical Risk

■ The Dilemma

Although

$$t(\mathbf{x}) = \text{sign}(\eta(\mathbf{x})) = \text{sign}(2\Pr[Y=1|X=\mathbf{x}] - 1)$$

is optimal, the distribution is **unknown** in general.

■ Empirical Risk

$$R_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{g(\mathbf{x}_i) \neq y_i}$$

where $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ are i.i.d. in $\mathcal{X} \times \mathcal{Y}$

Practical Solutions

■ Empirical Risk Minimization

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}} R_n(g)$$

■ Structural Risk Minimization

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}_d, d \in \mathbb{N}} R_n(g) + \text{pen}(d, n)$$

■ Regularization

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}} R_n(g) + \lambda \|g\|^2$$

Bounds

- g_n is a function learnt from $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

- $g_n \in \mathcal{G}$

- Error Decomposition

$$\begin{aligned}
 & R(g_n) - \min_g R(g) \\
 &= \underbrace{R(g_n) - \min_{g \in \mathcal{G}} R(g)}_{\text{Estimation Error}} + \underbrace{\min_{g \in \mathcal{G}} R(g) - \min_g R(g)}_{\text{Approximation Error}}
 \end{aligned}$$

- Risk Bound

$$R(g_n) - \min_{g \in \mathcal{G}} R(g)$$

- Generalization Error Bound

$$R(g_n) - R_n(g_n)$$

Empirical Process

■ Loss class

$$\mathcal{F} = \{f : (\mathbf{x}, y) \mapsto \mathbb{1}_{g(\mathbf{x}) \neq y} | g \in \mathcal{G}\}$$

■ Define

$$Pf = \mathbb{E}[f(X, y)] = R(g), \quad P_n f = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i, y_i) = R_n(g)$$

■ Empirical Process

$$\{Pf - P_n f\}_{f \in \mathcal{F}}$$

■ Supremum of Empirical Process

$$\sup_{f \in \mathcal{F}} Pf - P_n f, \quad \sup_{f \in \mathcal{F}} |Pf - P_n f|$$

Empirical Process

■ Lipschitz Loss class

$$\mathcal{F} = \{f : (\mathbf{x}, y) \mapsto \ell(y, g(\mathbf{x})) | g \in \mathcal{G}\}$$

- $|\ell(y, g(\mathbf{x})) - \ell(y, g'(\mathbf{x}))| \leq L|g(\mathbf{x}) - g'(\mathbf{x})|$

■ Define

$$Pf = \mathbb{E}[f(X, y)] = R(g), \quad P_n f = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i, y_i) = R_n(g)$$

■ Empirical Process

$$\{Pf - P_n f\}_{f \in \mathcal{F}}$$

■ Supremum of Empirical Process

$$\sup_{f \in \mathcal{F}} Pf - P_n f, \quad \sup_{f \in \mathcal{F}} |Pf - P_n f|$$

Empirical Process

■ Lipschitz Loss class

$$\mathcal{F} = \{f : (\mathbf{x}, y) \mapsto \ell(y, g(\mathbf{x})) | g \in \mathcal{G}\}$$

- $\ell(y, g(\mathbf{x})) = \max(0, 1 - yg(\mathbf{x}))$, $\ell(y, g(\mathbf{x})) = \log(1 + e^{-yg(\mathbf{x})})$

■ Define

$$Pf = \mathbb{E}[f(X, y)] = R(g), \quad P_n f = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i, y_i) = R_n(g)$$

■ Empirical Process

$$\{Pf - P_n f\}_{f \in \mathcal{F}}$$

■ Supremum of Empirical Process

$$\sup_{f \in \mathcal{F}} Pf - P_n f, \quad \sup_{f \in \mathcal{F}} |Pf - P_n f|$$

Outline

1 Introduction

2 Formalization

3 Analysis

4 Discussions

Problem

- The key problem is to bound

$$Pf - P_n f$$

for some $f \in \mathcal{F}$.

Theorem 2 (Hoeffding's Inequality)

Let X_1, \dots, X_n be *independent* random variables. Assume that the X_i are almost surely *bounded*, that is,
 $\Pr(X_i \in [a_i, b_i]) = 1$, $1 \leq i \leq n$. Denote

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Then, we have

$$\Pr\left(\left|\bar{X} - \mathbb{E}[\bar{X}]\right| \geq t\right) \leq 2 \exp\left(-\frac{2n^2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

An Elementary Analysis

Theorem 3

Assume $0 \leq f(\mathbf{x}, y) = \ell(y, g(\mathbf{x})) \leq 1$. With probability at least $1 - \delta$, for any $f \in \mathcal{F}$ that is independent from $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$,

$$|P_n f - P f| \leq \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$$

In other words, with probability at least $1 - \delta$, for any $g \in \mathcal{G}$ that is independent from $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$,

$$|R(g) - R_n(g)| \leq \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$$

An Elementary Analysis

Theorem 3

Assume $0 \leq f(\mathbf{x}, y) = \ell(y, g(\mathbf{x})) \leq 1$. With probability at least $1 - \delta$, for any $f \in \mathcal{F}$ that is independent from $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$,

$$|P_n f - P f| \leq \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$$

In other words, with probability at least $1 - \delta$, for any $g \in \mathcal{G}$ that is independent from $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$,

$$|R(g) - R_n(g)| \leq \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$$

- The above result cannot be used for empirical risk minimizer

$$g_n = \operatorname{argmin}_{g \in \mathcal{G}} R_n(g)$$

Remedy

Theorem 4

Assume $0 \leq f(\mathbf{x}, y) = \ell(y, g(\mathbf{x})) \leq 1$. With probability at least $1 - \delta$, for any/all $f \in \mathcal{F}$

$$|P_n f - Pf| \leq \sqrt{\frac{1}{2n} \log \frac{2|\mathcal{F}|}{\delta}}$$

With probability at least $1 - \delta$,

$$|R(g_n) - R_n(g_n)| \leq \sqrt{\frac{1}{2n} \log \frac{2|\mathcal{F}|}{\delta}}$$

- If $|\mathcal{F}|$ is small, analysis is easy.

Boole's inequality or Union Bound

Formally, for a countable set of events A_1, A_2, \dots , we have

$$\Pr\left(\bigcup_i A_i\right) \leq \sum_i \Pr(A_i)$$

An Example

With a probability at least $1 - \delta$,

$$X \leq f(\delta)$$

With a probability at least $1 - \delta$,

$$X \geq -f(\delta)$$



Boole's inequality or Union Bound

Formally, for a countable set of events A_1, A_2, \dots , we have

$$\Pr\left(\bigcup_i A_i\right) \leq \sum_i \Pr(A_i)$$

An Example

With a probability at least $1 - \delta$,

$$X \leq f(\delta)$$

With a probability at least $1 - \delta$,

$$X \geq -f(\delta)$$

By the union bound, with a probability at least $1 - 2\delta$, we have

$$|X| \leq f(\delta)$$

By the union bound, with a probability at least $1 - \delta$, we have

$$|X| \leq f(\delta/2)$$



A More Advanced Analysis (I)

■ Generalization Error Bound

$$R(g_n) - R_n(g_n) = Pf_n - P_n f_n \leq \sup_{f \in \mathcal{F}} (Pf - P_n f)$$

where $f_n : (\mathbf{x}, y) \mapsto \ell(y, g_n(\mathbf{x}))$

Theorem 5 (McDiarmid's inequality)

Let X_1, \dots, X_n be independent random variables taking values in a set A , and assume that $f : A^n \mapsto \mathbb{R}$ satisfies

$$\sup_{x_1, \dots, x_n, x'_i \in A} |f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i$$

for every $1 \leq i \leq n$. Then, for every $t > 0$,

$$P\{f(X_1, \dots, X_n) - E[f(X_1, \dots, X_n)] \geq t\} \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

A More Advanced Analysis (II)

- We need to upper bound

$$\sup_{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}'_i, y'_i)} \left| \sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \right|$$

A More Advanced Analysis (II)

- We need to upper bound

$$\sup_{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}'_i, y'_i)} \left| \sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \right|$$

For any $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}'_i, y'_i)$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \leq \frac{1}{n} \sup_{f \in \mathcal{F}} (f(\mathbf{x}_i, y_i) - f(\mathbf{x}'_i, y'_i)) \leq \frac{1}{n}$$

$$\sup_{f \in \mathcal{F}} (Pf - P'_n f) - \sup_{f \in \mathcal{F}} (Pf - P_n f) \leq \frac{1}{n} \sup_{f \in \mathcal{F}} (f(\mathbf{x}'_i, y'_i) - f(\mathbf{x}_i, y_i)) \leq \frac{1}{n}$$

A More Advanced Analysis (II)

- We need to upper bound

$$\sup_{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}'_i, y'_i)} \left| \sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \right|$$

For any $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}'_i, y'_i)$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \leq \frac{1}{n} \sup_{f \in \mathcal{F}} (f(\mathbf{x}_i, y_i) - f(\mathbf{x}'_i, y'_i)) \leq \frac{1}{n}$$

$$\sup_{f \in \mathcal{F}} (Pf - P'_n f) - \sup_{f \in \mathcal{F}} (Pf - P_n f) \leq \frac{1}{n} \sup_{f \in \mathcal{F}} (f(\mathbf{x}'_i, y'_i) - f(\mathbf{x}_i, y_i)) \leq \frac{1}{n}$$

Thus

$$\left| \sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \right| \leq \frac{1}{n}$$

which implies

$$\sup_{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), (\mathbf{x}'_i, y'_i)} \left| \sup_{f \in \mathcal{F}} (Pf - P_n f) - \sup_{f \in \mathcal{F}} (Pf - P'_n f) \right| \leq \frac{1}{n}$$

Symmetrization Inequality

Theorem 6

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq E \left[\sup_{f \in \mathcal{F}} (Pf - P_n f) \right] + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

- We proceed to upper bound $E \left[\sup_{f \in \mathcal{F}} (Pf - P_n f) \right]$

$$\begin{aligned}
 & E \left[\sup_{f \in \mathcal{F}} (Pf - P_n f) \right] = E \left[\sup_{f \in \mathcal{F}} (E [P'_n f] - P_n f) \right] \leq E \left[\sup_{f \in \mathcal{F}} (P'_n f - P_n f) \right] \\
 & = E \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}'_i, y'_i) - f(\mathbf{x}_i, y_i)) \right] = E \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(\mathbf{x}'_i, y'_i) - f(\mathbf{x}_i, y_i)) \right] \\
 & \leq E \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}'_i, y'_i) \right] + E \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n -\sigma_i f(\mathbf{x}_i, y_i) \right] = \underbrace{2 E \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i, y_i) \right]}_{:= \mathcal{R}_n(\mathcal{F})} \text{ LAMDA}
 \end{aligned}$$

Comparison Inequality [Ledoux and Talagrand, 1991]

Theorem 7

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq 2\mathcal{R}_n(\mathcal{F}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

- How to upper bound $\mathcal{R}_n(\mathcal{F})$

Comparison Inequality [Ledoux and Talagrand, 1991]

Theorem 7

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq 2\mathcal{R}_n(\mathcal{F}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

- How to upper bound $\mathcal{R}_n(\mathcal{F})$

Theorem 8 (Theorem 7 of [Meir and Zhang, 2003])

Let $\{\phi_i\}_{i=1}^n$ be functions with Lipschitz constant L_i , then

$$\mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^n \sigma_i \phi_i(g(\mathbf{x}_i)) \right] \leq \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^n \sigma_i L_i g(\mathbf{x}_i) \right]$$

Comparison Inequality [Ledoux and Talagrand, 1991]

Theorem 7

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq 2\mathcal{R}_n(\mathcal{F}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

- How to upper bound $\mathcal{R}_n(\mathcal{F})$

Theorem 8 (Theorem 7 of [Meir and Zhang, 2003])

Let $\{\phi_i\}_{i=1}^n$ be functions with Lipschitz constant L_i , then

$$\mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^n \sigma_i \phi_i(g(\mathbf{x}_i)) \right] \leq \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^n \sigma_i L_i g(\mathbf{x}_i) \right]$$

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i, y_i) \right] \leq L \mathbb{E} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right] = L \mathcal{R}_n(\mathcal{G})$$

Bounding Rademacher Complexities

Theorem 9

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq 2L\mathcal{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

Bounding Rademacher Complexities

Theorem 9

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq 2L\mathcal{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

■ Suppose

$$\mathcal{G} = \left\{ \mathbf{w}^\top \mathbf{x} : \mathbb{R}^d \mapsto \mathbb{R} \mid \|\mathbf{w}\| \leq \gamma \right\}$$

$$E \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \gamma} \sum_{i=1}^n \sigma_i \langle \mathbf{x}_i, \mathbf{w} \rangle \right]$$

$$\leq E \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \gamma} \left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\| \|\mathbf{w}\| \right]$$

$$\leq \gamma E \left[\left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\| \right] \leq \gamma \sqrt{E \left[\|\mathbf{x}_i\|^2 + \sum_{u \neq v} \sigma_u \sigma_v \mathbf{x}_u^\top \mathbf{x}_v \right]} \leq \gamma D \sqrt{n}.$$

Bounding Rademacher Complexities

Theorem 10

With a probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} (Pf - P_n f) \leq \frac{2L\gamma D}{\sqrt{n}} + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

■ Suppose

$$\mathcal{G} = \left\{ \mathbf{w}^\top \mathbf{x} : \mathbb{R}^d \mapsto \mathbb{R} \mid \|\mathbf{w}\| \leq \gamma \right\}$$

$$\begin{aligned} & \mathbb{E} \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \gamma} \sum_{i=1}^n \sigma_i \langle \mathbf{x}_i, \mathbf{w} \rangle \right] \\ & \leq \mathbb{E} \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \gamma} \left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\| \|\mathbf{w}\| \right] \\ & \leq \gamma \mathbb{E} \left[\left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\| \right] \leq \gamma \sqrt{\mathbb{E} \left[\left\| \mathbf{x}_i \right\|^2 + \sum_{u \neq v} \sigma_u \sigma_v \mathbf{x}_u^\top \mathbf{x}_v \right]} \leq \gamma D \sqrt{n}. \end{aligned}$$

Risk Bound

■ Define

$$g_* \in \operatorname{argmin}_{g \in \mathcal{G}} R(g)$$

■ Risk Bound

$$\begin{aligned} & R(g_n) - \min_{g \in \mathcal{G}} R(g) \\ &= R(g_n) - R(g_*) \\ &\leq R(g_n) - R(g_*) + R_n(g_*) - R_n(g_n) \\ &= R(g_n) - R_n(g_n) + R_n(g_*) - R(g_*) \\ &= \underbrace{\sup_{f \in \mathcal{F}} (Pf - P_n f)}_{\text{Theorem 10}} + \underbrace{R_n(g_*) - R(g_*)}_{\text{Theorem 3}} \end{aligned}$$

Risk Bound

■ Define

$$g_* \in \operatorname{argmin}_{g \in \mathcal{G}} R(g)$$

■ Risk Bound

$$\begin{aligned} & R(g_n) - \min_{g \in \mathcal{G}} R(g) \\ &= R(g_n) - R(g_*) \\ &\leq R(g_n) - R(g_*) + R_n(g_*) - R_n(g_n) \\ &= R(g_n) - R_n(g_n) + R_n(g_*) - R(g_*) \\ &= \underbrace{\sup_{f \in \mathcal{F}} (Pf - P_n f)}_{\text{Theorem 10}} + \underbrace{R_n(g_*) - R(g_*)}_{\text{Theorem 3}} \\ &\leq 2 \sup_{g \in \mathcal{G}} |R(g) - R_n(g)| = 2 \sup_{f \in \mathcal{F}} |Pf - P_n f| \end{aligned}$$

Theorem 10

Outline

1 Introduction

2 Formalization

3 Analysis

4 Discussions

■ Faster Rates

- Local Rademacher Complexities [Bartlett et al., 2005]
- Smoothness [Srebro et al., 2010]
- Strong Convexity [Sridharan et al., 2009]

■ Recover Error Bound

- Suppose the observation (\mathbf{x}, y) is generated according to some model

$$\Pr[Y = y | X = \mathbf{x}] = \frac{1}{1 + e^{-y\mathbf{w}_*^\top \mathbf{x}}}$$

- Upper bound the difference between \mathbf{w}_n and \mathbf{w}_*
- Compressive sensing, matrix completion

■ Active Learning

- The learner can select the training data actively
- Establishing theory is extremely difficult

Reference I

-  Bartlett, P. L., Bousquet, O., and Mendelson, S. (2005).
Local rademacher complexities.
The Annals of Statistics, 33(4):1497–1537.
-  Bousquet, O., Boucheron, S., and Lugosi, G. (2004).
Introduction to statistical learning theory.
In *Advanced Lectures on Machine Learning*, pages 169–207.
-  Ledoux, M. and Talagrand, M. (1991).
Probability in Banach Spaces: Isoperimetry and Processes.
Springer.
-  Meir, R. and Zhang, T. (2003).
Generalization error bounds for bayesian mixture algorithms.
Journal of Machine Learning Research, 4:839–860.
-  Srebro, N., Sridharan, K., and Tewari, A. (2010).
Optimistic rates for learning with a smooth loss.
ArXiv e-prints, arXiv:1009.3896.
-  Sridharan, K., Shalev-shwartz, S., and Srebro, N. (2009).
Fast rates for regularized objectives.
In *Advances in Neural Information Processing Systems 21*, pages 1545–1552

Reference II



Vapnik, V. N. (1998).
Statistical Learning Theory.
Wiley-Interscience.