

9 Statistical Learning Theory

얼마 많은 data가 있어야 어떤 수준의 예측 정확도를 얻을 수 있는지 수학적으로 분석해보자.

Empirical Risk를 잘 줄이는 모델이 true risk도 잘 줄이나? (항상 참은 아님) → 그럼 언제? or 어떤 조건 있을때?

9.1 Introduction

9.1.1 Statistical Learning Theory & No-Free-Lunch Theorems.

• iid sample $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{0, 1\}$ 을 사용해 binary classification을 학습시켜보자.

→ goal: Risk $R(h) = E[\mathbb{1}\{h(x) \neq y\}]$ 를 최소화하는 classifier h 찾기.

• suppose y is Bernoulli distributed with mean $\mu(x)$ ~~✗, ✗~~ often ≈ 0.5

→ if noise = non-zero, error $\neq 0$.

⇒ class. error $R(h)$ 를 최소화하는 classifier는 Bayes. Class. 그런데 (x, y) 결합 분포 몰라 → Bayes 계산 불가.

$$h^*(x) = \mathbb{1}\{\mu(x) > 1/2\} \quad \mu(x) \text{ 몰라}$$

solution) $\hat{R}(h)$ Empirical Risk Minimization을 사용해보자.

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\}$$

근데 이걸 최소화하는 함수는 단순히 training data 만 이용하는 overfitting model 일 수도 있다.

= No-Free-Lunch; 반면적으로 좋은 학습 방법도 ✗. 그럼 어떻게? 1) generative

→ discriminative

9.1.2 Generative vs. Discriminative

• Generative: $P(x, y)$ or $P(x|y), P(y)$ 를 직접 모델링 e.g) Gaussian Naives Bayes

• Discriminative: $P(y|x)$ or 결정 함수만 모델링 e.g) Logistic Regression, SVM.

• Risk Analysis: $R(\hat{h}) - R(h^*) = R(\hat{h}) - \inf_h R(h) + \inf_h R(h) - R(h^*)$

estimation error
data가 완벽히 생겼는지

approximation error
함수 class H 가 충분히 넓지 않아 생기는 error

* Statistical Learning Theory에서 Empirical Risk Minimization (ERM)의 일반적인 오류를 분석함. 우리는 $h \in H$ 중에서 training data에 대해

Empirical Risk를 최소화하는 \hat{h} 를 선택함. 궁극적 목표는 True Risk가 작은 모델을 찾는 것

• $R(h)$: true risk (population risk)

: $R(h) = E[\text{loss}(h(x), y)]$ 전체 data 분포에 대한 avg. loss

• $\hat{R}(h)$: empirical risk

$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \text{loss}(h(x_i), y_i)$ training data에 대한 avg. loss

• \hat{h} : ERM에서 선택된 가설

$\hat{h} = \argmin \hat{R}(h)$

• h_H : H 에서 true risk를 최소화하는 opt. 가설

$h_H = \argmin R(h)$

• $R(\hat{h})$: 선택된 모델 \hat{h} 의 true risk

• $\hat{R}(\hat{h})$: 선택된 모델 \hat{h} 의 empirical risk

• $R(h_H)$: H 내의 최적 가설 h_H 의 true risk

• $\hat{R}(h_H)$: H 내의 최적 가설 h_H 의 empirical risk

$$R(\hat{h}) - R(h_H) = \underbrace{R(\hat{h}) - \hat{R}(\hat{h})}_{\text{일반화 오류}} + \underbrace{\hat{R}(\hat{h}) - \hat{R}(h_H)}_{\text{emp. risk 최소화}} + \underbrace{\hat{R}(h_H) - R(h_H)}_{\text{일반화 오류}} = \underbrace{R(\hat{h}) - \hat{R}(\hat{h})}_{\text{generation gap}} + \underbrace{\hat{R}(h_H) - R(h_H)}_{\text{generation gap}}$$

generation gap 항상 ≤ 0
($\because \hat{h}$ 는 \hat{R} 를 최소화하는 함수)

$$\rightarrow R(\hat{h}) - R(h_H) \leq |R(\hat{h}) - \hat{R}(\hat{h})| + |\hat{R}(h_H) - R(h_H)|$$

$$\rightarrow |R(\hat{h}) - R(h_H)| \leq \sup_{h \in H} |\hat{R}(h) - R(h)| \quad \because \hat{h} \in H \text{ 이므로 어떤 } h \in H \text{ 에 대해서도 generation gap은 supremum 보다 작다.}$$

9.2 Risk Minimization & Generation

$$R(f) = \underbrace{R(f) - R(f_H)}_{\text{estimation error}} + \underbrace{R(f_H) - R(f^*)}_{\text{approximation error}} + R(f^*)$$

- complex model \Rightarrow approx error \downarrow , estima. error \uparrow
- simple model \Rightarrow approx error \uparrow , estima. error \uparrow

* trade-off 설명

e.g.) Linear Regression $H = \{h_\theta(x) = \theta^T x\}$.

Neural network
더 큰 H .

- Class H : 우리가 모델을 고르는 후보군 집합.

5차 다항식 $H = \{h(x) = a_0 + a_1x + \dots + a_5x^5\}$

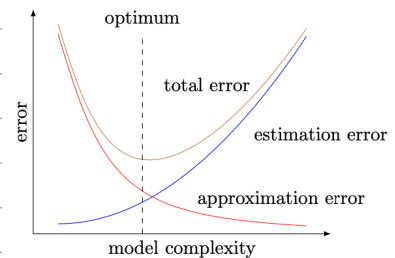
\rightarrow 100차 다항식은 1차 함수보다 더 많은 곡선 형태 내포.

= H 만이 진짜 최적함수 f^* 와 가까운 함수가 들어왔을 가능성 높음. \approx approximation 가까워짐 \Rightarrow approx. error \downarrow

\hookrightarrow But, data를 위변조 위험이 커짐 (Overfitting) data가 유한하면 일반화가 안될수도. e.g) 5차식 모델인데 data가 3개뿐이면...

(복잡한 모델은 학습해야 할 parameter 수 \uparrow , 즉 data 1개당 설명해야 할 정보량이 적어져서 noise 민감함)

| complexity | class H | Approx. Error | Estimation Error |
|------------|--------------|---------------------|--------------------------------------|
| complex | \downarrow | \uparrow (표현력 부족) | \downarrow |
| simple | \uparrow | \downarrow | \uparrow (추론이 잘 안될 확률 \uparrow) |



9.3 Uniform Convergence

Empirical Risk를 잘 줄이는 모델이 true risk도 잘 줄이나요? (항상 참은 아님) \rightarrow 증명해줄래.

- Definition: < Uniform Convergence >

$$\sup_{h \in H} |\hat{R}(h) - R(h)| \leq \epsilon \text{ (with high probability)}$$

Class H 의 모든 가설에 대해 $\hat{R}(h)$ 와 $R(h)$ 차이가 0이함을 보장.
(empirical risk와 true risk의 차이가 작다는 것을 균일하게 보장)

즉 ERM (Empirical Risk Minimization)으로 찾은 \hat{h} 도 일반화를 잘 하는거란 이론적 보장

- \hookrightarrow Prop 1: finite H , generalization 보장
- Theorem 1: Infinite H , covering number를 통해 보장
- Prop 2: Lipschitz loss + covering number 상한값에 bound 제공

Prop 1 + Theorem 1 + Prop 2 } uniform convergence $\checkmark \Rightarrow$ ERM 일반화 보장 \Rightarrow empirical risk 작으면 true risk 작다는 말이
Prop 1, Theorem 1, Prop 2를 만족할수록 성립.

• Prop 1: Finite Hypothesis Class

Suppose H finite, $0 \leq \text{loss}(f(x), y) \leq B$ is bounded,
with probability (at least) $1 - \delta$, for all $h \in H$:

$$R(f) \leq \hat{R}(f) + B \sqrt{\frac{\log(|H|/\delta)}{2n}}$$

$$P\left[\sup_h |R(h) - \hat{R}(h)| \geq \varepsilon\right] \leq 2|H|e^{-2n\varepsilon^2}$$

* finite H에 대해 Hoeffding's Inequality + union Bound \Rightarrow bound 가능

$$P\left[\max_i (R(f_i) - \hat{R}(f_i)) \geq t\right] = P\left[\bigcup_{i=1}^{|H|} (R(f_i) - \hat{R}(f_i) \geq t)\right]$$

$$\leq \sum_i P[R(f_i) - \hat{R}(f_i) \geq t] \leq \sum_i e^{-\frac{2nt^2}{B^2}} = |H|e^{-\frac{2nt^2}{B^2}}$$

* Hoeffding's Inequality (Theo 1)

• Theorem 1: Hoeffding's Inequality

Let z_1, \dots, z_n independent random variables taking values in $[a, b]$. Then for $\forall \beta > 0$:

$$P\left[\frac{1}{n} \sum_{i=1}^n (z_i - E[z_i]) \geq +\beta\right] \leq e^{-\frac{2n\beta^2}{(b-a)^2}}$$

infinite H에서는 union bound X 대신 복잡도 측정 도구 필요.

ε -covering number 개념 도입. 즉 H를 몇개의 함수로 근사적으로

덮을 수 있나? $N(H, \varepsilon, n)$

\rightarrow H가 작게 덮이면 (covering number ↓) uniform conv. 성립.

• Prop 2: Lipschitz loss + covering number

suppose loss is Lipschitz, for $\forall z, f, f' \in H$:

$$|\text{loss}(f, z) - \text{loss}(f', z)| \leq L \|f - f'\|$$

Moreover assume loss is bounded $\Leftrightarrow 0 \leq \text{loss}(f, z) \leq B$. Then,

$$\sup_f R(f) - \hat{R}(f) \leq \varepsilon \text{ with probability } 1 - N(H, \frac{\varepsilon}{4L}) \cdot e^{-\frac{2n\varepsilon^2}{B^2}}$$

Also,

$$P\left[\sup |R(f) - \hat{R}(f)| \geq \varepsilon\right] \leq N(H, \frac{\varepsilon}{4L}) \cdot e^{-\frac{2n\varepsilon^2}{B^2}}$$

Proof. Let $S = \frac{\varepsilon}{4L}$ for H . $f \in H: \exists f' \in S$ s.t.

$$|\text{loss}(f, z) - \text{loss}(f', z)| \leq L \|f - f'\| \leq \frac{\varepsilon}{4}.$$

$$\text{Then, } P\left[\exists f \in H: \frac{1}{n} \sum_{i=1}^n (E[\text{loss}(f, z_i)] - \text{loss}(f, z_i)) \geq \varepsilon\right]$$

$$\leq P\left[\exists f \in S: \frac{1}{n} \sum_{i=1}^n (E[\text{loss}(f, z_i)] - \text{loss}(f, z_i)) \geq \frac{\varepsilon}{2}\right]$$

$$\leq N(H, \frac{\varepsilon}{4L}) \cdot e^{-\frac{2n\varepsilon^2}{B^2}} \quad \square$$

* intuitive)

If $N(H, \varepsilon) = 100 \rightarrow \log N = \log(100) \approx 6.9 \Rightarrow$ 필요한 data $n \geq \frac{6.9}{\varepsilon^2}$

If $N(H, \varepsilon) = 10^6 \rightarrow \log N \approx 13.8 \Rightarrow$ 동일한 generation bound를 얻으려면 2배 data 필요