# Basic Concept of Ch.9

① 왜 어떤 ML model은 잘 generalization 되고, 어떤 모델은 overfitting 될까?

② 어떤 조건에서 empirical risk가 작다면, true risk도 작다고 보장할 수 있나?

#1. 정말 알고 싶은 것: true Risk $R(h)$ But 데이터분포 $D$모름 → empirical risk $\hat{R}(h)$ 를 최소화해보자!

근데.. $\hat{R}(h)$이 작다고 $R(h)$이 작은 증거 있어?

그럼 $R(h)$. $\hat{R}(h)$ 관계를 알아보자.

#2. Error $\quad R(\hat{h}) - R(f^*) = \underbrace{R(\hat{h}) - R(h_H)}_{\text{estimation error}} + \underbrace{R(h_H) - R(f^*)}_{\text{approximation error}}$ $\quad$ $f^*$: ideal. $h_H$: $H$중 가장 좋은 $h$, $\hat{h}$: ERM 선택된 $h$

분석해보니

#3 Complexity
- trade off

$H$가 단순하면 목값이 잘 못따라가 (approx. error ⬆) but 일반화 잘해 (estm. error ⬇)

$H$가 복잡하면 목값이 잘 따라가 (approx. error ⬇) but overfitting 될지도 (estm. error ⬆)

그럼... bias - variance tradeoff가 필요하게?

근데 ERM 잘 하려면 조건이 필요하대!

#4 Uniform Convergence $\quad \sup |R(h) - \hat{R}(h)| \le \varepsilon$ $\quad$ empirical risk 와 true risk의 차이중 가장 큰게

어떤 $\varepsilon$ 안에 있으면. empirical risk가 작으면 true risk도 작아.

③ 근데 Uniform convergence가 안되는 경우는?

#5 SGD / SGM 같은 stochastic optimization 알고리즘! 근데 왜?

#6 SGD는 왜 잘 작동? ∵ Algorithmic Stability 때문.

여기서는 어떤 algo가 입력 데이터를 약간 바꿔도 학습결과가 거의 안바뀌면 "stable"

stable algorithm은 generalization 잘해.

그럼 Stability가 수치적으로 어떻게 되는데?

↳ $E[R(\hat{h}) - \hat{R}(\hat{h})] \le \varepsilon$ 에서 strong convexity, Lipschitz 조건하에

$\varepsilon = \dfrac{2l^2}{\lambda n}$ 으로 stability 가져.

# 9 Statistical Learning Theory

얼마나 많은 Data가 있어야 어떤 수준의 예측 정확도를 얻을수 있는지 수학적으로 분석해보자.
Empirical Risk를 잘 줄이는 모델이 true risk도 잘 줄이나? (항상 참은 아님) → 그점 언제? or 어떤 조건 있을때?

## 9.1 Introduction

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

- iid sample $(x_1, y_1), \cdots (x_n, y_n) \in \mathbb{R}^d \times \{0,1\}$ 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 $\mu(x)$ $\varkappa, \emptyset$ often $\approx 0.5$
  
  ↳ if noise = non-zero, error $\neq 0$.

  ⇒ class. error $R(h)$를 최소화하는 Classifier는 Bayes. Classf. 그런데 $(x, y)$ 결합 분포 몰라 → Bayes 계산 불가.
  
  $$h^*(x) = \mathbb{1}\{\mu(x) > 1/2\} \qquad \mu(x) \text{ 몰라}$$

  solution) $\hat{R}(h)$ Empirical Risk Minimization 을 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 안 외우는 Overfitting model 일 수도 있다.
  = No-Free-Lunch ; 보편적으로 좋은 학습 방법은 X. 그럼 어떻게? 1) generative
  2) 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 : $R(\hat{h}) - R(h^*) = \underbrace{R(\hat{h}) - \inf R(h)}_{\substack{\text{estimation error} \\ \text{data가 유한하니 생기는 오차}}} + \underbrace{\inf R(h) - R(h^*)}_{\substack{\text{approximation error} \\ \text{함수 class H가 충분히 넓지 않아 생기는 오차}}}$

---

* Statistical Learning Theory 에서 Empirical Risk Minimization (ERM)의 일반화의 오류를 분석함. 우리는 $h \in H$ 중에서 training data에 대한 Empirical Risk를 최소화하는 $\hat{h}$를 선택함. 궁극적 목표는 True Risk가 작은 모델을 찾는 것

- $R(h)$ : true risk (population risk) : $R(h) = E[loss(h(x), y)]$ 전체 data 분포에 대한 avg. loss
- $\hat{R}(h)$ : empirical risk $\qquad \hat{R}(h) = \frac{1}{n}\sum_{i=1}^{n} loss(h(x), y)$ traing data에 대한 avg. loss
- $\hat{h}$ : ERM 에서 선택된 가설 $\qquad \hat{h} = \arg\min \hat{R}(h)$
- $h_H$ : H 에서 true risk를 최소화하는 opt. 가설 $\qquad h_H = \arg\min R(h)$
- $R(\hat{h})$ : 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

$$R(\hat{h}) - R(h_H) = \underbrace{R(\hat{h}) - \hat{R}(\hat{h})}_{\substack{\text{일반화 오류} \\ \text{generation gap}}} + \underbrace{\hat{R}(\hat{h}) - \hat{R}(h_H)}_{\substack{\text{emp. risk 최소화} \\ \text{항상} \leq 0 \\ (\because \hat{h}는 리스크를 \\ 최소화하는 함수)}} + \underbrace{\hat{R}(h_H) - R(h_H)}_{\substack{\text{일반화 오류 generation gap}}} \qquad \leq |R(\hat{h}) - \hat{R}(\hat{h})| + |\hat{R}(h_H) - R(h_H)|$$

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

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

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

  $f_H = \underset{f \in H}{\arg\min} R(f)$ ,

  $f^* = \underset{f}{\arg\min} R(f)$

- complex model $\Rightarrow$ approx error $\Downarrow$, estima. error $\Uparrow$

- simple model $\Rightarrow$ approx error $\Uparrow$, estima. error $\Uparrow$

\* trade-off 설명

- Class H : 우리가 모델을 고르는 후보균 집합.

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

  5차 다항식 $H = \{h(x) = a_0 + a_1 x + \cdots + a_5 x^5\}$

  Neural network 더 큰 H.

  → 100차 다항식은 1차 함수보다 더 많은 곡선 형태 내포.

  = H안에 진짜 최적함수 $f^*$와 가까운 함수가 들어있을 가능성 높음. ≈ approximation 가까워짐 ⇒ approx. error $\Downarrow$

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

  (복잡한 모델은 학습해야하는 parameter 수 $\theta_d$, 즉 data 1개당 설명해야하는 정보들이 적어져서 noise 민감해짐)

| Complexity | Class H | Approx. Error | Estimation Error |
|---|---|---|---|
| Complex | $\downarrow$ | $\uparrow$ (표현력 부족) | $\downarrow$ |
| Simple | $\uparrow$ | $\downarrow$ | $\uparrow$ (추정이 잘 안됨 학습 $\uparrow$) |



## 9.3 Uniform Convergence

Empirical Risk를 잘 줄이는 모델이 true risk도 잘 줄이나? (항상 참은 아님) → 증명해줄게.

- Definition : < Uniform Convergence >

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

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

즉 ERM (Empirical risk Minimization)으로 찾은 $\hat{h}$도 일반화를 잘 한거라고 이론적 보장

| | |
|---|---|
| ↳ Prop 1 | finite H, generalization 보장 |
| Theorem 1 | Infinite H, covering number를 통해 보장 |
| Prop 2 | Lipschitz loss + covering number 상황일때 bound 제공 |

Prop 1
+
Theorem 1   uniform   $\checkmark$ ⇒ ERM 일반화 보장 ⇒   empirical risk 작으면 true risk 작다는 말이
+          convergence                                  Prop1, Theorem1, Prop2 를 만족할때나 성립.
Prop 2

- **Prop 1 : Finite Hypothese Class**

  Suppose H finite, $0 \le \text{loss}(f(x), y) \le B$ is bounded,

  with probability (at least) $1 - \delta$, for all $h \in H$ :

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

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

  \* finite H에서는 Hoeffding's Inequality + union Bound 로 bound 가능

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

  \* $\le \sum_i^{|H|} P\left[R(f_i) - \hat{R}(f_i) \ge t\right] \le \sum_i^{|H|} e^{-\frac{2nt^2}{B^2}} = |H| e^{-\frac{2nt^2}{B^2}}$ \*

  \* Hoeffding's Inequality (Theo 1)

- **Theorem 1 : Höffding's Inequality**

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

  $$P\left[\frac{1}{n}\sum_i^n (z_i - E[z_i]) \ge +\beta\right] \le e^{-\frac{2n\beta^2}{(b-a)^2}}$$

  infinite H에서는 union bound X 대신 복잡도 측정 도구 필요.

  $\varepsilon$ - covering number 개념 도입. 즉 H를 몇개의 항수로 근사적으로

  덮을 수 있나? $N(H, \varepsilon, n)$

  → H가 작게 덮이면 (covering number 小) uniform conv. 성립.

- **Prop 2 : Lipschitz loss + covering number**

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

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

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

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

  ALSO:

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

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

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

  Then). $P\left[\exists f \in H : \frac{1}{n}\sum_i^n (E[\text{loss}(f, z)] - \text{loss}(f, z_i)) \ge \varepsilon\right]$

  $\le P\left[\exists f \in S : \frac{1}{n}\sum_i^n (E[\text{loss}(f, z)] - \text{loss}(f, z_i)) \ge \frac{\varepsilon}{2}\right]$

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

  \* intuitive )

  If $N(H, \varepsilon) = 100$, → $\log N = \log(1000) \approx 6.9$ ⇒ 필요한 data $n \ge \frac{6.9}{\varepsilon^2}$

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

왜 Stochastic optimization이 중요한가? 많은 ML 문제가 학습 = 최적화 문제로 표현됨.

$$\min_{w \in \Theta} R(w) := E_{z \sim D}\left[ \text{loss}(w,z) \right]$$

$R(w)$ : True risk (Population Risk) / $D$ : data 분포 / $z=(x,y)$ : data point

하지만 우리는 $D$를 모르고, 대신 data $z_1, \cdots, z_n$ 만 가짐. 그래서 ERM 사용.

$$\hat{R}(w) = \frac{1}{n}\sum_{i}^{n} \text{loss}(w, z_i) \text{ 를 최소화!}$$

• Stochastic Gradient Method (SGM)    현실에서 큰 data를 다루기위해 stochastic gradient method 사용

$$w_{k+1} = \prod_{\theta}\left( w_k - \alpha_k \nabla \text{loss}(w_k, z_k) \right)$$

- $f_w$ : model parameterized by $w \in \Theta$,
- $R(f_w)$ : risk $= E\left[ \text{loss}(f_w(x), y) \right] = R(w) = E\left[ \text{loss}(w,y) \right]$
- $z_k$ : mini batch.
- $\alpha_k$ : 학습률 stepsize
- $\prod_{\theta}$ : 제약조건 영역 (projection). optional.

✱ 성능 분석 : SGD는 실제로 얼마나 잘 일반화 할까?

Suppose norm of SG : $E\left[ \| \nabla \text{loss}(w_z, z_k) \|_2^2 \right] \leq B^2$ , which $B = \sup \| \nabla \text{loss}(w, z) \|_2$

SGD의 iterate avg $\bar{w}_n = \frac{1}{n}\sum w_i$ 에 대하, stepsize를 $\alpha_k = D/B\sqrt{n}$ 로 설정하면:

$$E\left[ R(\bar{w}_n) \right] - R(w^*) \leq \frac{BD}{\sqrt{n}}$$

→ 의미 : 평균 iterate risk는 $O\left(\frac{1}{\sqrt{n}}\right)$ 만큼 $w^*$ 보다 나쁘지 않음.

↑여기에 확률을 $\geq 1-\delta$ 로 설정하면:

$$R(\bar{w}_n) - R(w^*) \leq \frac{BD}{\sqrt{n}}\left( 1 + \sqrt{2\log(1/\delta)} \right)$$

✱ 결론)

SGD는 널리 쓰이는데 이 이론이 그 일반화 성능을 수학적으로 보장.
→ loss function의 convexity, Lipschitz, boundedness 가정함.

data 하나 바꾼다고 결과가 크게 바뀌지 않는 것. 왜? 일반적으로 stable algorithm은 generalization (일반화) 잘해.

✱ Flow   9.4) SGM이 왜 generalization 잘하는지 수학적으로 expectation, prob. bound 로 설명.
        9.5) 어떤 알고리즘이 generalization 잘하는가를 "얼마나 안정적인가?" 라는 관점으로 설명.

• Define : Uniform Stability

Consider learning algorithm $A$. given $Z = \{z_1, \cdots, z_n\} \in Z^n$ is $\varepsilon$-uniformly stable.

즉, training dataset $Z$와 그것과 단 한 sample만 가진 $Z^{(i)}$에 대하:

$$\forall_z, \quad E\left[ | \text{loss}(A(Z), z) - \text{loss}(A(Z^{(i)}), z) | \right] \leq \varepsilon$$

→ data 하나 바뀌어도 output function은 거의 바뀌지 않음 & 기댓값은 알고리즘의 randomness 에 대하 취함.

- Proposition 3 : Stability → Generalization

let A $\varepsilon$-uniformly stable, then generation error of A is bounded by $E\left[R(A(z)) - \hat{R}(A(z))\right] \leq \varepsilon$

⇒ 즉, 일반화 오차의 기댓값이 작다 ≈ Overfitting 위험이 작다.


- Proposition 4 : Regularized ERM is stable

let A(z) minimizing regularized empirical loss of model. 이 알고리즘이 다음 최적화 문제를 푼다면, :

$$\min_{w} \frac{1}{n} \sum^{n} loss(w, z_i) + \frac{\lambda}{2} \|w\|^2 \quad \leadsto \quad \text{이 algorithm은 } \varepsilon = \frac{2L^2}{\lambda n}.$$

- loss가 L-Lipschitz 이고,
- $\lambda$-strongly convex 라면

* interpretation
- $n$ ⇑ → good stability
- $\lambda$ ⇑ → good stability

* Conclusion

1) Stability는 data 수 $n$ 과 regularization $\lambda$ 에 의해 제어됨.

2) SGD나 Regularized ERM 처럼 실제 사용하는 알고리즘의 일반화 성능을 Stability 관점에서 해석할수 있음.

3) uniform convergence 없이 generalization을 설명