

모델은 단순히 training dataset에 잘 맞는 것만 부족, 새로운 데이터에 generalization이

4. Model Selection & Validation

잘 모르는 모델을 선택해야 한다 → Model을 선택하는 법, 검증 (Validation)하는 법.

4.1 The statistical learning setup

Problem Def)

- 우리는 어떤 입력 $x \in \mathcal{X}$ 와 출력 $y \in \mathcal{Y}$ 의 관계를 모델링해보고 싶어
- Data는 확률 분포 D 에 따라 생성된 $(x, y) \sim D$ 샘플들.
- But 실제로 우리는 D 를 모르고, Sample dataset $S = \{(x_i, y_i)\}_{i=1}^n$ 만 가지고 있다.

• Loss function: 예측이 얼마나 틀렸는지 확인하기 위해 $\text{loss}: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ 사용. $\text{ex}) \text{loss}(z, y) = (z - y)^2 \dots$

• Risk

1) Population Risk $R(h) = E[\text{loss}(h(x), y)]$ ← 우리가 진짜 알고 싶은 것. But D 를 모르면 어떻게 계산? $h^* = \arg\min R(h)$

2) Empirical Risk $\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \text{loss}(h(x_i), y_i)$ ← 우리가 실제로 계산할 수 있는 것. D 기점의 loss avg. (손실 평균)

기본 Idea: 모델 학습할 때 $\hat{R}(h)$ 를 최소화하는 방향으로 학습하되, $R(h)$ 가 작아야 함 ($\hat{R}(h)$ 으로 overfitting 하면 안 좋음)

4.2 Training, Validation, Test errors

ML에서 데이터를 3부분으로 나눠서 모델 성능 평가

1) Training Set - 모델 학습

2) Validation Set - 모델 선택 / Hyperparameter tuning 사용

3) Test Set - 최종 성능 평가 (unseen data)

• Risk 비교관점에서 본 Errors

(Model h 에 대해)

1) Training Error $\hat{R}(h, D_{\text{train}}) = \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \text{loss}(h(x_i), y_i)$

Empirical Risk on training set

모델이 training data에 얼마나 잘 맞는지 보여줌.

2) Validation Error $\hat{R}(h, D_{\text{val}}) = \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \text{loss}(h(x_i), y_i)$

여러 모델 중 어떤 모델이 generalization이 더 좋는지 비교할 때 사용

3) Test Error $\hat{R}(h, D_{\text{test}}) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \text{loss}(h(x_i), y_i)$

학습이 끝난 후 모델의 일반화 성능을 최종 평가하는 용도.

★ 왜 이렇게 나누는가?

- 모델이 training set에만 너무 잘 맞으면 (overfitting) test set 성능 떨어질 수 있음.
- Validation set이 이를 방지하고, 모델 선택 기준 제공
- test set은 한번도 보지 않은 data set, 즉 실제 상황처럼 test.

• Steps 1) 여러 모델(h_1, h_2, \dots)을 training set에 학습

2) Validation error $\hat{R}_{\text{val}}(h)$ 가 가장 낮은 모델 h^* 선택

3) 선택한 모델 h^* 에 대해 test error $\hat{R}_{\text{test}}(h)$ 평가.

True Risk $R(h) = E[\text{loss}(h(x), y)]$

Empirical Risk $\hat{R}(h) = \frac{1}{n} \sum \text{loss}(h(x), y)$

우리는 최종적으로 $R(h)$ 를 원하지만 바로 계산 불가. \rightarrow test set 위에서 $\hat{R}(h)$ 를 사용.

\rightarrow 그럼 질문: Empirical Risk $\hat{R}(h)$ 가 True Risk $R(h)$ 를 얼마나 잘 근사할까?

Hoeffding's Inequality

for sample (x_i, y_i) :

$$P[|\hat{R}(h, D) - R(h)| \leq O\left(\sqrt{\frac{\log(1/\delta)}{n}}\right)] \geq 1 - \delta$$

$\underbrace{\hspace{10em}}_{\substack{\text{test set 크기} \\ O(\dots) \\ \text{근사치}}}$ \uparrow 허용 가능한
실패 확률
e.g) $\delta = 0.05 \approx 95\%$ 신뢰도

해석) 좌항: test error, true risk 차이

우항: 그 차이가 $O\left(\sqrt{\frac{\log(1/\delta)}{n}}\right)$ 일 확률.

의미: test set이 충분히 크면 test error는 true risk에 가까워진다.

결론: Hoeffding's Inequality는 test error $\hat{R}(h)$ 가 true risk $R(h)$ 를 얼마나 잘 근사하는지 확률적으로 보장.

4.3 K-fold Cross Validation

• Goal: 여러 모델 중 generation 성능이 제일 좋은 모델 h 선택하자 (true risk $R(h)$ 가 가장 작은 모델 선택)

\rightarrow But, $R(h)$ 직접 계산 불가 \rightarrow validation set에서 측정된 $\hat{R}(h)$ 로 평가해보자.

• K-fold - Cross-Validation

1) Dataset D 를 크기가 같은 K 개의 부분집합으로 나눈다. D_1, D_2, \dots, D_K

2) 각 $i \in \{1, \dots, K\}$ 에 대해

- Train $D_{\text{train}} = D \setminus D_i$
- Validation $D_{\text{val}} = D_i$

3) 이때 모델 $h \in H$ 을 train set으로 학습하고, validation set으로 측정

4) 모든 fold에 대해 평균 Validation error 계산.

$$\hat{R}_i(h) = \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \text{loss}(h(x_i), y_i)$$

$$\hat{R}_{\text{cv}}(h) = \frac{1}{K} \sum_{i=1}^K \hat{R}_i(h)$$

\rightarrow 적용 목적. (왜 사용?)

모델 선택 시 단일 validation set가 아닌 여러 split을 사용 \rightarrow 보다 안정적인 선택 가능. 특히 Dataset 작을 때.

Validation error를 줄여줌. \star test set은 모델 성능 평가에만 사용. 모델 선택 (Model Selection)에서는 사용 X

5. Convex optimization → Basic Intro. • Convexity → Aspect) Optimization algorithm.

5.1 Optimization Problems

• Basic Structure

$\min_{x \in \mathbb{R}^d} f(x)$ • $f: \mathbb{R}^d \rightarrow \mathbb{R}$ Objective function
 • 목적: $f(x)$ 를 최소화하는 x 를 찾자.

• 최적해 x^* 는 $f(x^*) \leq f(x) \quad \forall x \in \mathbb{R}^d$ 를 만족해야 한다
 (= global minimum)

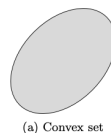
5.2 Convexity

• Definition 1: Convex Set

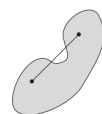
A convex set \mathcal{C} is any set s.t. for $x, y \in \mathcal{C}$ and $\forall \theta \in (0, 1)$:

$$\theta x + (1-\theta)y \in \mathcal{C}$$

≡ 두 점 사이를 잇는 선분이 모두 집합 안에 있다면 Convex set.



(a) Convex set



(b) Non-convex set

• Definition 2: Convex Function (Strictly convex)

그래프가 직선 연결보다 아래 있으면 convex

A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex

⇔ for all $x, y \in \mathbb{R}^d, \theta \in (0, 1)$ 에 대해

$$\theta f(x) + (1-\theta)f(y) \geq f(\theta x + (1-\theta)y)$$

$$\theta f(x) + (1-\theta)f(y) > f(\theta x + (1-\theta)y)$$

• Linear-function $f(x) = \langle a, x \rangle + b$

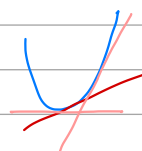
• Quadratic $f(x) = \frac{1}{2}x^T Q x + b^T x + c$

• Any Norm $f(x) = \|x\|$

• Prop 1: 함수 f 가 convex 하고, diff-able 하면 다음 성립 →

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$

→ 접선이 함수 그래프보다 아래에 있다.



• Prop 2: convex function f 에 대해, 어떤 점 $x^* \in \mathbb{R}^d$ 가 $\nabla f(x^*) = 0$ 이면 → $x^* = \text{global minimum}$ 이다.

convex function의 stationary point는 global minimum이다.

→ 기울기가 0인 지점이 global 최소값 (≈ 최적해를 가지는 x^*)

• 만약 f strictly convex 일 경우 x^* 는 유일한! global minimizer

5.3 Computational aspect of optimization algorithm

우리는 $\min f(x)$ 를 풀고 싶어. e.g) $f(x) = \|Ax - y\|$. 여기서 $A \in \mathbb{R}^{m \times n}$ 은 full column rank 행렬.

만약 A 가 full rank 라면 (해를 구할 수 있다면) 최적해 $\Rightarrow x^* = (A^T A)^{-1} A^T y$. 그런데 현실에선? 1) f 가 closed form

계산이 불가능한 경우에 어떻게 하지?

2) 최적해 구하기 어려움

→ 반복해서 최대한 비슷하게 만들어보라 (근사해 Approximate Solution 구해보라)

• Oracle Access - 실제 컴퓨터가 풀 수 가능한 것 (할 수 있는 것): 알고리즘이 접근 가능한 정보.

1) 0th-order oracle : 주어진 x 에 대해 $f(x)$ 만 반환

2) 1st-order oracle : $f(x)$ 와 $\nabla f(x)$ (Gradient)를 반환

3) 2nd-order oracle : $f(x)$, $\nabla f(x)$, $\nabla^2 f(x)$ (Hessian)까지 반환

+ 몇 번을 호출해야 (반복해야) 할까? = oracle complexity

• 우리가 어떤 정보를 가지고 최적화를 진행하나?

• 이 정보를 몇 번 요청해야 충분한 정확한 해를 얻을 수 있는가?

다음 chapter부터는 algorithm을 볼 텐데 이 이전에

우리는 $f(x)$, $\nabla f(x)$ 같은 Oracle Access 정보를 사용한다 가정!

ε-accurate solution을 얻기 위한 최소

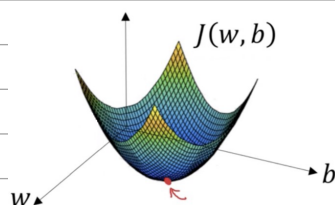
$$\|x^k - x^*\| \leq \varepsilon \text{ or } f(x^k) - f(x^*) \leq \varepsilon$$

이 중 하나를 만족할 때의 oracle 호출 횟수 = 복잡도.

6. Gradient Descent (안과수준 중급)

$$x^{k+1} = x^k - \alpha_k \nabla f(x^k) \quad \text{Gradient (도함수)를 } k \text{ 번 반복하면서 } \alpha_k \text{ 만큼씩 반복해서 경사를}$$

내려갈거야



- Goal: Gradient Descent 이 Quadratic function (2차 함수)에 대해 얼마나 빨리 최소값에 수렴하는지 수학적으로 분석해보자.

6.1 Convergence for quadratic function

• Quad. function :

$$f(x) = \frac{1}{2} x^T Q x - b^T x + c \rightarrow \nabla f(x) = Qx - b$$

$$\text{식 정리: } x^{k+1} - x^* = x^k - \alpha \nabla f(x^k) - x^*$$

• Q is symmetric, strictly positive definite.

• 이 함수 f는 convex

$$\text{Gradient Descent : } x^{k+1} = x^k - \alpha \nabla f(x^k)$$

$$\bullet \nabla f(x) = Qx - b. \rightarrow \text{최적화 식: } Qx^* = b.$$

$$= x^k - \alpha(Qx^k - b) - x^* = x^k - \alpha(Qx^k - Qx^*) - x^*$$

$$= x^k - \alpha Q(x^k - x^*) - x^* = (I - \alpha Q)(x^k - x^*)$$

$$\Rightarrow \text{It follows: } \|x^{k+1} - x^*\|_2 \leq \|I - \alpha Q\| \|x^k - x^*\|_2$$

Eigen Value: $\lambda_i = 1 - \alpha \lambda_i$; Symmetric (\approx invertible, singular value \checkmark , sym 행렬의 norm = 최대 Singular Value.)

$$\checkmark \|I - \alpha Q\| = \max(\lambda_{\max}(I - \alpha Q), \lambda_{\min}(I - \alpha Q)) = \max(\alpha M - 1, 1 - \alpha m) \quad M: \text{최대 SV}, m: \text{최소 SV}$$

이 값을 최소화하려면, $|\alpha M - 1| = |1 - \alpha m| \Rightarrow \text{stepsize } \alpha = \frac{2}{M+m}$

$$\text{값을 } \alpha = \frac{2}{M+m}, \kappa = \frac{M}{m} \text{로 설정하면,}$$

$$\bullet \lambda_{\min} = m: 1 - \frac{2}{M+m} \cdot m = \frac{M-m}{M+m} \quad \left. \begin{array}{l} \bullet \lambda_{\max} = M: 1 - \frac{2}{M+m} \cdot M = -\frac{M-m}{M+m} \end{array} \right\} \rightarrow \|I - \alpha Q\| = \max |1 - \alpha \lambda_i| = \frac{M-m}{M+m} = \frac{\kappa-1}{\kappa+1} = \frac{1-\frac{1}{\kappa}}{1+\frac{1}{\kappa}}$$

$$\bullet \lambda_{\max} = M: 1 - \frac{2}{M+m} \cdot M = -\frac{M-m}{M+m}$$

• Prop3: GD with $\alpha = \frac{2}{M+m}$ applied to a quadratic function obeys

$$\kappa = \frac{M}{m}$$

$$\|x^k - x^*\|_2 \leq \left(\frac{\kappa-1}{\kappa+1}\right)^k \|x^0 - x^*\|_2 \leq \varepsilon$$

$$\text{우리가 원하는 것: } \varepsilon\text{-accurate solution} \Leftrightarrow \|x^k - x^*\|_2 \leq \varepsilon, \quad \rho = \frac{\kappa-1}{\kappa+1} \Rightarrow \rho < 0.$$

$$\Rightarrow \rho^k \|x^0 - x^*\|_2 \leq \varepsilon \Rightarrow \rho^k \leq \frac{\varepsilon}{\|x^0 - x^*\|_2}$$

$$\log \Rightarrow k \log \rho \leq \log \left(\frac{\varepsilon}{\|x^0 - x^*\|_2} \right) \Rightarrow k \geq \frac{\log(\|x^0 - x^*\|_2 / \varepsilon)}{\log(1/\rho)}$$

여기서 최적 k를 N이라 해보자

$$N = \left(\log \frac{1}{\rho} \right)^{-1} \log \left(\frac{\|x^0 - x^*\|_2}{\varepsilon} \right)$$

9.1)

1) 어떤 행렬 Q의 condition number (κ)는 고유값의 비율

* Q is well-conditioned ($\kappa \approx 1$): 고유값이 서로 비슷 \Rightarrow 모든 방향에서 함수의 곡률이 비슷 \Rightarrow 수렴 fast, stable

* Q is ill-conditioned ($\kappa \gg 1$): 고유값 차이 \gg \Rightarrow 어떤 방향은 급경사, 평평 \Rightarrow 수렴 slow, zigzag 발생

$$2) \text{rate of convergence} = \left(1 - \frac{1}{\kappa}\right)^k$$

$\kappa \uparrow = \text{수렴률} \downarrow, \kappa \rightarrow \infty \Rightarrow \text{수렴률 매우 느림}$

3) Geometrical ...

Quadratic function (2차 함수)의 등고선이 Q에 따라 타원형의 모양. Q의 condition number가 클수록 길쭉한 타원
GD는 등고선을 따라 수직하강 방식으로 움직임. 곡률 차이가 크면 zigzag $ax^2 + by^2 = 1$