

14. Surrogate Models

- 전통적인 최적화 : 함수 f 를 알고 있음
- Surrogate model : 목적 함수 f 모름
- m 개의 계획점과 이에 대응하는 함수의 값 y
- 목적 함수 :
$$\min_{\theta} \|y - \hat{y}\|_p$$

L_p norm

- Surrogate Models

• m개의 계리점

$$X = \{x^{(0)}, x^{(1)}, \dots, x^{(m)}\}$$

• m개의 함수값

$$y = \{y^{(0)}, y^{(1)}, \dots, y^{(m)}\}$$

• m개의 예측값 : \hat{f}, θ (모수)

$$\hat{y} = \{\hat{f}_\theta(x^{(0)}), \hat{f}_\theta(x^{(1)}), \dots, \hat{f}_\theta(x^{(m)})\}$$

⇒ 최적화 문제는 계리점 찾기 아니고

함수라고 생각되는 surrogate 함수 구성하는
 θ 모수 찾기

• 목적 함수

$$: \underset{\theta}{\text{minimize}} \|y - \hat{y}\|_p \quad \begin{matrix} p=2 : \text{유클리디안} \\ \uparrow \\ L_p \text{ norm} \end{matrix}$$

• 선형 모형 (linear models)

$$\hat{f} = w_0 + w^T x, \quad \theta = f_{w_0, w}$$

$$\hat{f} = \theta^T x$$

• 선형회귀 (Linear regression)

$$: \underset{\theta}{\text{minimize}} \|y - \hat{y}\|_2^2 \rightarrow \text{최소 제곱 합}$$

• 행렬 표현

$$: \underset{\theta}{\text{minimize}} \|y - X\theta\|_2^2$$

$$\rightarrow \text{계리점 행렬} \quad X = \begin{bmatrix} (x^{(0)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad \text{계리점 행렬}$$

Newton 방법의 특수한 경우.

• Analytic Solution : 최소 제곱 추정량

$$\underset{\theta}{\text{minimize}} \|y - X\theta\|_2^2$$

$$\rightarrow \theta = (X^T X)^{-1} X^T y \quad \rightarrow \text{최소 제곱 추정량 구하기 어렵} \rightarrow \text{역행렬!}$$

→ 대규모 자료 : m이 매우 크고 X 차원 매우 높음

→ Newton 방법의 특수한 경우

↓

→ $(X^T X)^{-1}$: 이차 미분 정보 사용 X, gradient 사용하!

→ 경사도

$$2X^T(y - X\theta) \rightarrow (X\hat{\theta} - y)X^T$$

$$2X^T(y - X\theta) \rightarrow (\hat{y} - y)X^T$$

→ Full-batch Gradient Descent : X와 y 이용

$$\theta^{(t+1)} = \theta^{(t)} - \alpha (\hat{y} - y) X^T$$

계리점 행렬의 전체 data

- Stochastic Gradient Descent

• 한 개의 계리점 X와 함수값 y

• 경사도

$$(\hat{y}_{(i)} - y_{(i)}) \cdot X_{(i)}^T \Rightarrow (\hat{y}_{(i)} - y_{(i)}) \cdot X_{(i)}$$

• SGD

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \alpha (\hat{y}_{(i)} - y_{(i)}) \cdot X_{(i)}$$

- Minibatch Stochastic Gradient Descent

• 전체 자료를 m보다 작은 임의의 부분 집합들로 나눔

ex) 100개 자료 → 10개씩 10번의 미니배치

$$\theta_j^{(k+1)} = \theta_j^{(k)} - \alpha \sum_{k=1}^{(m/10) \times B} (\hat{y}_{(k)} - y_{(k)}) \cdot X_{(k)}$$

→ adam 방법도 사용