

10. Algorithmic Regularization [Idea!!!] = EARLY STOPPING

10.1 Early Stopping least-squares

• 우리가 관찰할 data  $y_i = \langle x_i, \theta^* \rangle + z_i, \quad z_i \sim \mathcal{N}(0, \sigma^2)$

- Goal:  $\theta^*$ 를 잘 추정해보자.  $\approx$  prediction risk를 줄이자.
- 이때 사용하는 estimator는 GD의 iterate  $\theta_t$ 를 충분히 학습시키지 않고 early stopping 하면 regularization 효과 有
- 즉: GD를 완전히 수렴시키지 않고 적절한 t에서 멈추는 걸로 overfitting 방지 가능 = Algorithmic Regularization (Implicit Regularization)

10.1.2 Iterates of gradient descent

GD의 iteration 식을 수학적으로 전개해서 다음 식 유도:

$$\theta_t = -\eta \sum_{k=0}^{t-1} V \Sigma^T (I - \eta \Sigma \Sigma^T)^k U^T r_0$$

전개한 식은 이후 risk decomposition에 사용

10.1.3 Risk of gradient descent iterates

$\theta_t$ 의 risk를 bias + variance로 나누자!

$$E[R(\theta_t)] = \underbrace{\sum_{i=1}^d (1 - \eta \sigma_i^2)^{2t}}_{\text{bias}^2} + \underbrace{\sigma^2 \sum_{i=1}^d \frac{(1 - (1 - \eta \sigma_i^2)^t)^2}{\sigma_i^2}}_{\text{variance}}$$

Insights: bias는 감소하고 variance는 증가함.  $\Rightarrow$  전체 risk는 U자형 (U-shaped) 곡선을 따라감.  
= 너무 작은 t는 underfitting, 너무 큰 t는 overfitting  
 $\rightarrow$  적절한 시점에 early stopping을 하는게 best risk.

10.1.4 Comparison to ridge regression estimator

일반 Ridge Regression과 비교할때?

- Ridge Regression Estimator

$$\hat{\theta}_{\text{ridge}} = V \cdot \text{diag}\left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right) V^T \theta^* + V \cdot \text{diag}\left(\frac{\sigma_i}{\sigma_i^2 + \lambda}\right) U^T z$$

- Risk

$$E[R(\theta_t)] = \sum_{i=1}^d \left(1 - \frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right)^2 + \sigma^2 \sum_{i=1}^d \left(\frac{\sigma_i}{\sigma_i^2 + \lambda}\right)^2$$

Ridge  $\equiv \lambda$ 로 Regularization  
GD  $\equiv t$ 로 Regularization  
ONLY) 모든  $\sigma_i$ 가 같을때만 동일 Risk가 됨

## 10.1 Early Stopping Least Squares

### • Problem State

$y_i = \langle x_i, \theta^* \rangle + z_i$ ,  $z_i \sim \mathcal{N}(0, \sigma^2)$ , linear model with noise 를 가정. 우리 goal: to estimate well  $\theta^*$   
then model risk  $\hat{h}(x) = \langle \hat{\theta}, x \rangle$  small

• Suppose feature vector  $x$  is Gaussian distributed ( $x \sim \mathcal{N}(0, I)$ ) then prediction risk is:

$$R(\hat{h}) = E[(\hat{h}(x) - y)^2] = E[(\langle \hat{\theta} - \theta^*, x \rangle - z)^2] = \|\hat{\theta} - \theta^*\|_2^2 + \sigma^2$$

(Interpret) 만약 estimate  $\hat{\theta}$  가 true model  $\theta^*$  와 유사하면 risk가 낮아지고 "좋은 모델" 이 됨.

### AS AN ESTIMATOR

we knew the iterates of gradient descent applied to the least-squares

$$L(\theta) = \frac{1}{2} \|X\theta - y\|_2^2$$

## 10.1.2 Iterates of gradient descent

GD의 Loss 부터 전개해보자.

$$\nabla L(\theta) = X^T(X\theta - y) = X^T r \quad \text{①} \quad r^{t+1} = (I - \eta XX^T) r^t \quad \& \quad r^t = (I - \eta XX^T)^t r^0 \quad \text{②}$$

$r := X\theta - y$  residual

① gradient descent Update  $\theta^{t+1} = \theta^t - \eta \nabla L(\theta^t) = \theta^t - \eta X^T r^t$

$$X\theta^{t+1} = X\theta^t - \eta XX^T r^t$$

$$r = X\theta - y$$

$$\Rightarrow r^{t+1} = X\theta^{t+1} - y = (X\theta^t - \eta XX^T r^t) - y = r^t - \eta XX^T r^t$$

$$\Rightarrow X\theta^t = r^t + y$$

$$r^{t+1} = r^t - \eta XX^T r^t$$

②

$$t=0 : r^1 = (I - \eta XX^T) r^0$$

$$t=1 : r^2 = (I - \eta XX^T) r^1 = (I - \eta XX^T)^2 r^0$$

$$t=2 : r^3 = (I - \eta XX^T) r^2 = (I - \eta XX^T)^3 r^0$$

$\vdots$

$\sum_{k=0}^{\infty} \eta^k$

$$r^t = (I - \eta XX^T)^t r^0$$

• GD iterate 정의:  $\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle x_i, \theta \rangle)^2 = \frac{1}{n} \|y - X\theta\|_2^2$

• GD update 식 정의:

$$\theta_{t+1} = \theta_t - \eta \cdot \nabla \hat{R}(\theta_t) = \theta_t - \frac{2\eta}{t} X^T (y - X\theta_t)$$

여기서 초기값  $\theta_0 = 0$  이라 가정하고 반복식 전개하면 closed-Form 사용.

$$\begin{aligned} \theta_t - \theta_0 &= (\theta_t - \theta_{t-1}) + (\theta_{t-1} - \theta_{t-2}) + \dots + (\theta_t - \theta_0) \\ &= -\eta \sum_{k=0}^{t-1} (X^T X)^k X^T r_0 \quad (\text{where } r_0 = y) \end{aligned}$$

여기서 SVD 이용하면  $\theta_t = -\eta \sum_{k=0}^{t-1} (X^T X)^k X^T r_0$ ,

$$X = U \Sigma V^T, \quad y = X\theta^* + z$$

$$X^T X = V \Sigma^2 V^T, \quad (X^T X)^k = V \Sigma^{2k} V^T \quad \theta_t = -\eta \sum_{k=0}^{t-1} (X^T X)^k X^T r_0 = -\eta \sum_{k=0}^{t-1} V \Sigma^{2k} V^T \cdot V \Sigma U^T \cdot y$$

$$= -\eta \sum_{k=0}^{t-1} V \Sigma^{2k+1} U^T y = -\eta \cdot V \left( \sum_{k=0}^{t-1} \Sigma^{2k+1} \right) \cdot U^T y$$

$\Sigma$  is diagonal matrix  $\rightarrow \Sigma^{2k+1} = \text{diag}(\dots) \Rightarrow$  각 singular vector 방향  $v_i$  에 대해 독립적으로 계산 가능

$$\theta_t = V \cdot \sum_{k=0}^{t-1} \left[ \left( \sum_{k=0}^{t-1} \eta \cdot \sigma_i^{2k+1} \right) \langle u_i, y \rangle \right] v_i$$

$$\sum_{k=0}^{t-1} \eta \sigma_i^{2k+1} = \sigma_i \cdot \eta \sum_{k=0}^{t-1} (\sigma_i^2)^k = \sigma_i \cdot \eta \cdot \frac{1 - (\sigma_i^2)^t}{1 - \sigma_i^2} \approx 1 - (1 - \eta \sigma_i^2)^t$$

attenuation factor

최종)  $\theta_t = V \cdot \sum_{i=1}^d (1 - (1 - \eta \sigma_i^2)^t) \cdot \langle u_i, y \rangle \cdot \frac{1}{\sigma_i} \cdot v_i$

반복 횟수에 따라 줄어드는 계수 (early stopping 핵심 역할)

label vector  $y$  가 data 방향  $u_i$  에 갖는 projection

→ 고유값 크기 (singular value) 에 대한 분포

### 10.1.3 Risk of gradient descent iterates.

우리가 진짜 관심 있는 것 = expected risk.

그리고 평균적으로 parameter  $\theta^* \sim N(0, I)$  라고 가정할 때 GD iterates  $\theta_t$  에 대한 risk = bias + Variance 가 됨.

$$E[R(\theta_t)] = \underbrace{\sum_{i=1}^d (1 - \eta \sigma_i^2)^{2t}}_{\text{bias}^2} + \underbrace{\sigma^2 \sum_{i=1}^d \frac{(1 - (1 - \eta \sigma_i^2)^t)^2}{\sigma_i^2}}_{\text{Variance}}$$

$\sigma_i$  = singular values of  $X$   
 $\eta$  = learning rate (step size)  
 $t$  = number of iterations

Interpretations:

- 1) 반복이 적으면 ( $t \downarrow$ ):  $(1 - \eta \sigma_i^2)^t \approx 1 \Rightarrow$  bias 가 크고, variance 작다 (underfitting)
  - 2) 반복이 많으면 ( $t \uparrow$ ):  $(1 - \eta \sigma_i^2)^t \rightarrow 0 \Rightarrow$  bias 가 작고, variance 커진다 (overfitting)
- $\therefore$  전체 risk는  $t$ 에 대해 (반복 횟수에 대해) U-shaped curve 를 가지며,  
 적절한  $t$ 에서 멈추는 것이 중요 !!! = **early stopping**

### 10.1.4 Comparison to ridge regression estimator

• Ridge Regression  $\hat{\theta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$

→ SVD

$$\hat{\theta}_{\text{ridge}} = V \cdot \text{diag}\left(\frac{\sigma_i}{\sigma_i^2 + \lambda}\right) V^T \theta^* + V \cdot \text{diag}\left(\frac{\sigma_i}{\sigma_i^2 + \lambda}\right) U^T z$$

$\Rightarrow$  Risk:

$$E[R(\hat{\theta}_{\text{ridge}})] = \underbrace{\sum_{i=1}^d \left(1 - \frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right)^2}_{\text{bias}^2} + \underbrace{\sigma^2 \sum_{i=1}^d \left(\frac{\sigma_i}{\sigma_i^2 + \lambda}\right)^2}_{\text{Variance}}$$

생각을 바꾸어야!

	< Ridge Regression >	< GD + Early Stopping >
• Bias-Variance Trade-off	둘다 $\text{bias}^2 + \text{Variance}$ 형태	
• SVD 표현	모든 $V_i$ -basis 에서 각 방향으로 shrinkage factor 작용	
• Shrinkage 효과	각 방향으로 적용된 weight 를 줄여서 overfitting 방지	
• High-frequency 방향 억제	작은 $\sigma_i$ 에서는 (noise가 심한 방향) weight 작아짐	
• Control Parameter	$\lambda$ (explicit)	$t$ (iteration 횟수, implicit)
• Shrinkage 방식	$\frac{\sigma_i^2}{\sigma_i^2 + \lambda}$ , $\lambda \rightarrow 0$ 일때 OLS	$1 - (1 - \eta \sigma_i^2)^t$ , $t \rightarrow \infty$ 일때 OLS
• Noise Sensitivity (Variance)	$\left(\frac{\sigma_i}{\sigma_i^2 + \lambda}\right)^2$	$\frac{(1 - (1 - \eta \sigma_i^2)^t)^2}{\sigma_i^2}$
• 학습 방식	단일 closed-form	반복 최적화 알고리즘 필요