

MACHINE 기계 학습 LEARNING

오일석 지음

5장. 딥러닝 최적화

PREVIEW

■ 과학과 공학에서 최적화

- 예) 웨이퍼에 최대 집적하는 반도체공학, 목적지까지 최단 경로 찾기 등
- 목적함수를 정의하고 최적해를 구하면 됨

가 , 가 .

■ 기계 학습의 최적화는 훨씬 복잡함

- 훈련집합으로 학습을 마친 후, 현장에서 발생하는 새로운 샘플을 잘 예측해야 함. 즉 일반화 능력이 뛰어나야 함
- 훈련집합이 테스트집합의 대리자 역할을 한다고 말할 수 있음
- 또한 MSE, 로그우도와 같은 목적함수가 궁극 목표인 정확률의 대리자 역할을 함



■ 기계 학습의 최적화가 어려운 이유

- 대리자 관계
- 목적함수의 비볼록 성질, 고차원 특징 공간, 데이터의 희소성 등
- 긴 시간 소요 (예, GPU 4대 장착한 컴퓨터에서 VGGNet을 훈련하는데 2~3주 소요)

가
가 . ().
MSE , .

■ 5장은 이런 어려움을 극복하는 여러 가지 효과적인 방안을 제시함

가
가 .

가 .

각 절에서 다루는 내용

- 5.1절 목적함수로서 평균제곱 오차의 단점을 지적하고, 딥러닝이 많이 활용하는 교차 엔트로피와 로그우도를 소개한다.
- 5.2절 스토캐스틱 경사 하강법의 성능 향상에 효과적인 전처리, 가중치 초기화, 모멘텀, 적응적 학습률, 활성화함수, 배치 정규화를 설명한다.
- 5.3절 딥러닝에서 규제의 필요성과 원리를 간략히 기술한다.
- 5.4절 많이 쓰는 규제 기법으로서 가중치 벌칙, 조기 멈춤, 데이터 확대, 드롭아웃, 앙상블을 설명한다.
- 5.5절 하이퍼 매개변수의 중요성과 최적값을 찾는 방법을 소개한다.
- 5.6절 2차 미분 정보를 활용하여 스토캐스틱 경사 하강법을 보완하는 기법을 기술한다.

5.1 목적함수: 교차 엔트로피와 로그우도

- 5.1.1 평균제곱 오차를 다시 생각하기
- 5.1.2 교차 엔트로피 목적함수
- 5.1.3 softmax 활성화함수와 로그우도 목적함수

시험에서는 틀린 만큼 합당한 벌점을 받는 것이 중요하다. 그래야 다음 시험에서 심기일전으로 공부하여 틀리는 개수를 줄일 가능성이 크기 때문이다. 틀린 개수에 상관없이 비슷한 벌점을 받는다면 나태해져 성적을 올리는 데 지연이 발생할 것이다. 이러한 원리가 기계 학습에도 적용될까?

5.1.1 평균제곱 오차 다시 생각하기

■ 평균제곱 오차(MSE) 목적함수

$$e = \frac{1}{2} \|\mathbf{y} - \mathbf{o}\|_2^2 \quad (5.1)$$

- 오차가 클수록 e 값이 크므로 벌점으로 훌륭함

■ 하지만 큰 허점

가 saturated
가 ? , MSE

- 왼쪽 상황은 $e = 0.2815$, 오른쪽 상황은 $e = 0.4971$ 이므로 오른쪽이 더 큰 벌점을 받아야 마땅함

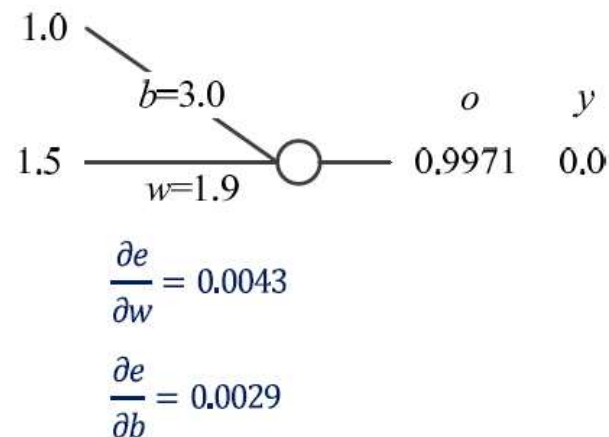
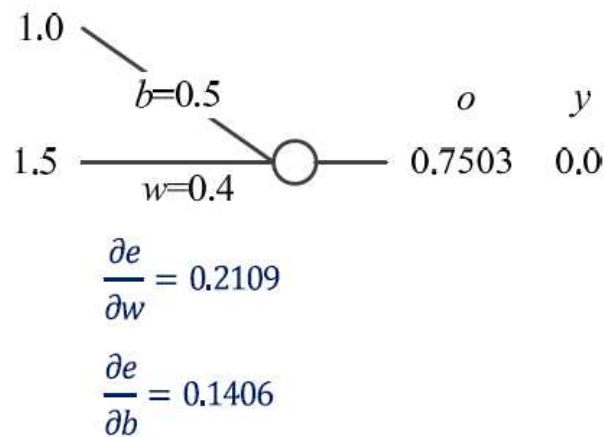


그림 5-1 MSE가 목적함수로서 부적절한 상황

5.1.1 평균제곱 오차 다시 생각하기

■ 큰 허점

- 식 (5.3)의 그레이디언트가 별점에 해당

$$e = \frac{1}{2}(y - o)^2 = \frac{1}{2}(y - \sigma(wx + b))^2$$

$$\left. \begin{aligned} \frac{\partial e}{\partial w} &= -(y - o)x\sigma'(wx + b) \\ \frac{\partial e}{\partial b} &= -(y - o)\sigma'(wx + b) \end{aligned} \right\} \quad (5.3)$$

- 그레이디언트를 계산해보면 왼쪽 상황의 그레이디언트가 더 큼 → 더 많은 오류를 범한 상황이 더 낮은 별점을 받은 꼴 → 학습이 더딘 부정적 효과

■ 이유

- $w x + b$ (아래 그래프의 가로축에 해당)가 커지면 그레이디언트가 작아짐

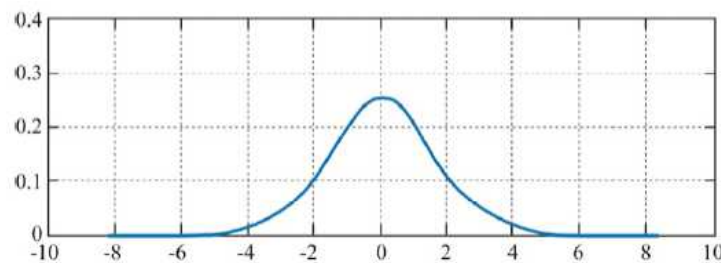
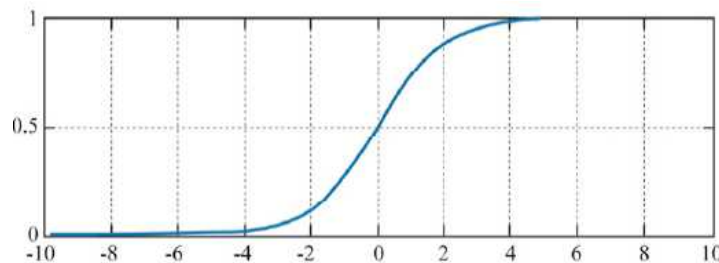


그림 5-2 로지스틱 시그모이드함수와 도함수

'o'		0	가	1	가
	saturated	.			
,					MSE
	.	(...)

(5.2)

5.1.2 교차 엔트로피 목적함수

■ 교차 엔트로피 MSE

- 레이블에 해당하는 y 가 확률변수 (부류가 2개라고 가정하면 $y \in \{0,1\}$)
- 확률 분포: P 는 정답 레이블, Q 는 신경망 출력

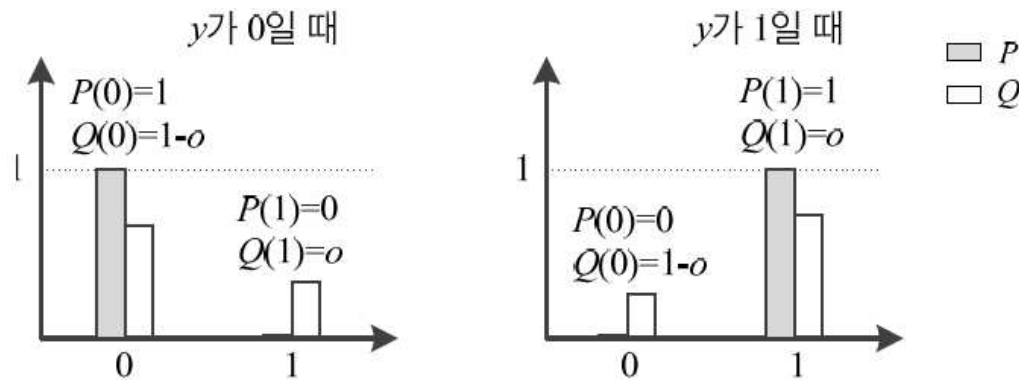


그림 5-3 레이블 y 가 0일 때와 1일 때의 P 와 Q 의 확률분포

- 확률분포를 통일된 수식으로 쓰면,

$$\begin{aligned} P(0) &= 1 - y & Q(0) &= 1 - o \\ P(1) &= y & Q(1) &= o \end{aligned}$$

- 교차 엔트로피 식은 $H(P, Q) = -\sum_{y \in \{0,1\}} P(y) \log_2 Q(y)$

가
가
saturated

5.1.2 교차 엔트로피 목적함수

■ 교차 엔트로피 목적함수

$$e = -(y \log_2 o + (1 - y) \log_2(1 - o)), \quad \text{이때, } o = \sigma(z) \text{이고 } z = wx + b \quad (5.4)$$

■ 제구실 하는지 확인

- y 가 1, o 가 0.98일 때 (예측이 잘된 경우)
 - 오류 $e = -(1 \log_2 0.98 + (1 - 1) \log_2(1 - 0.98)) = 0.0291$ 로서 낮은 값
- y 가 1, o 가 0.0001일 때 (예측이 엉터리인 경우)
 - 오류 $e = -(1 \log_2 0.0001 + (1 - 1) \log_2(1 - 0.0001)) = 13.2877$ 로서 높은 값

5.1.2 교차 엔트로피 목적함수

- 공정한 벌점을 부여하는지 확인 (MSE의 느린 학습 문제를 해결하나?)

- 도함수를 구하면,

$$\begin{aligned}
 \frac{\partial e}{\partial w} &= -\left(\frac{y}{o} - \frac{1-y}{1-o}\right) \frac{\partial o}{\partial w} \\
 &= -\left(\frac{y}{o} - \frac{1-y}{1-o}\right) x \sigma'(z) \\
 &= -x \left(\frac{y}{o} - \frac{1-y}{1-o}\right) o(1-o) \\
 &= x(o-y)
 \end{aligned}
 \quad \longrightarrow \quad
 \left. \begin{aligned}
 \frac{\partial e}{\partial w} &= x(o-y) \\
 \frac{\partial e}{\partial b} &= (o-y)
 \end{aligned} \right\} (5.5)$$

Saturated

- 그레이디언트를 계산해 보면, 오류가 더 큰 오른쪽에 더 큰 벌점(그레이디언트) 부과

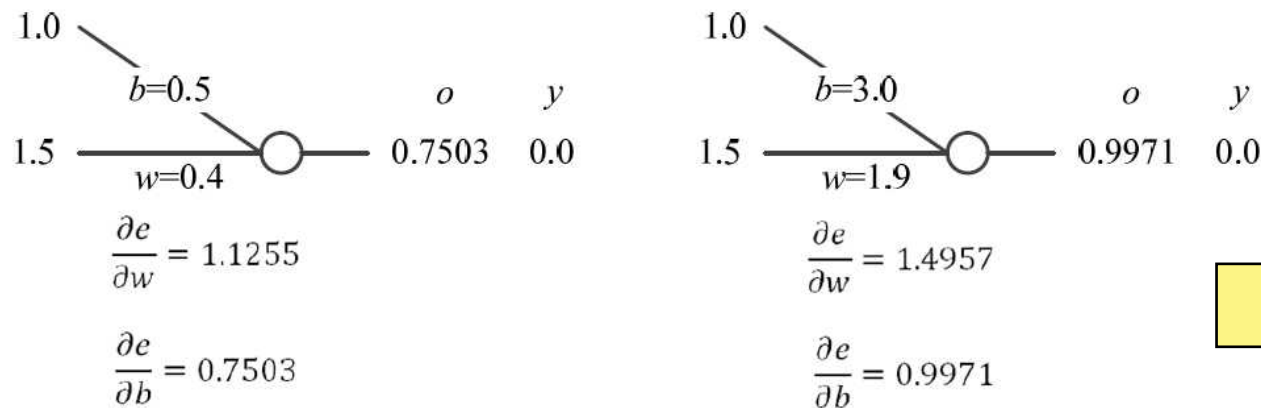


그림 5-4 교차 엔트로피를 목적함수로 사용하여 느린 학습 문제를 해결

5.1.2 교차 엔트로피 목적함수

- 식 (5.4)를 c 개의 출력 노드를 가진 경우로 확장
 - 출력 벡터 $\mathbf{o} = (o_1, o_2, \dots, o_c)^T$ 인 상황으로 확장 ([그림 4-3]의 DMLP)

$$e = - \sum_{i=1,c} (y_i \log_2 o_i + (1 - y_i) \log_2 (1 - o_i)) \quad (5.6)$$

5.1.3 softmax 활성화함수와 로그우도 목적함수

■ softmax 활성화함수

$$o_j = \frac{e^{s_j}}{\sum_{i=1,c} e^{s_i}} \quad (5.7)$$

■ 동작 예시

- softmax는 max를 모방(출력 노드의 중간 계산 결과 s_i^L 에서 최댓값은 더욱 활성화하고 작은 값은 억제)
- 모두 더하면 1이 되어 확률 모방

softmax
max

output
5 - 5
softmax output

max 0 1 , 가
softmax가

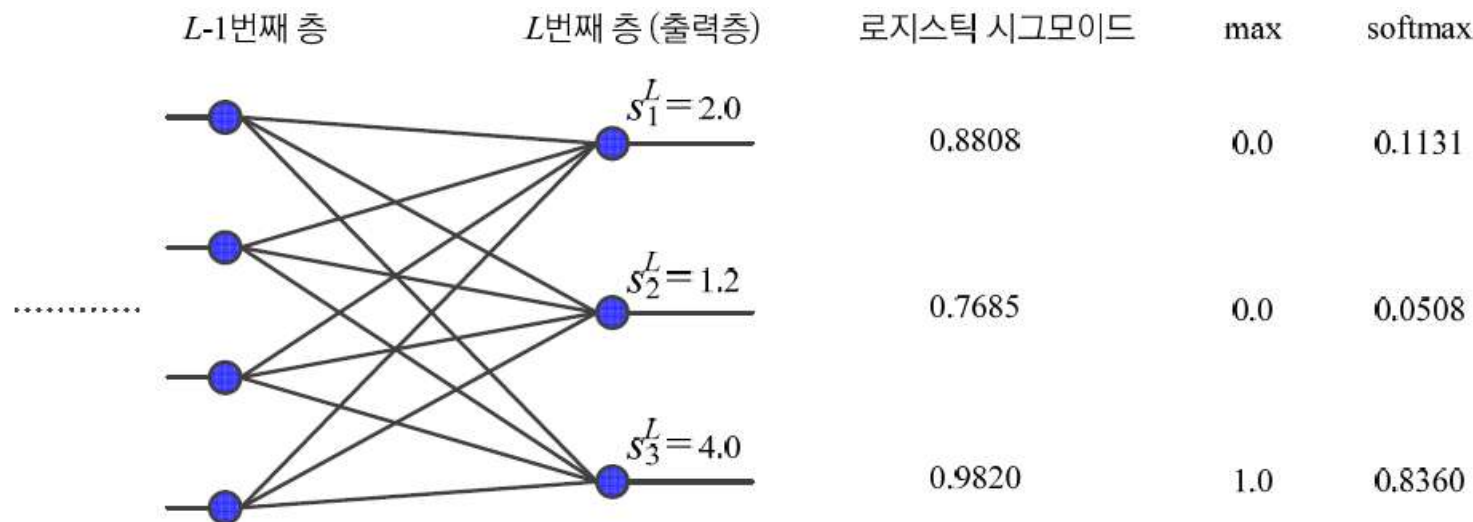
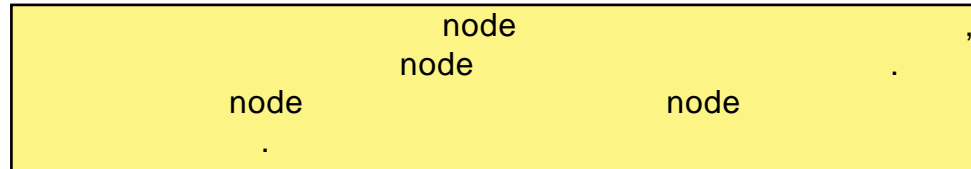


그림 5-5 출력층의 활성화함수로 로지스틱 시그모이드와 softmax 비교

5.1.3 softmax 활성화 함수와 로그우도 목적 함수

■ 로그우도 목적 함수

$$e = -\log_2 o_y$$



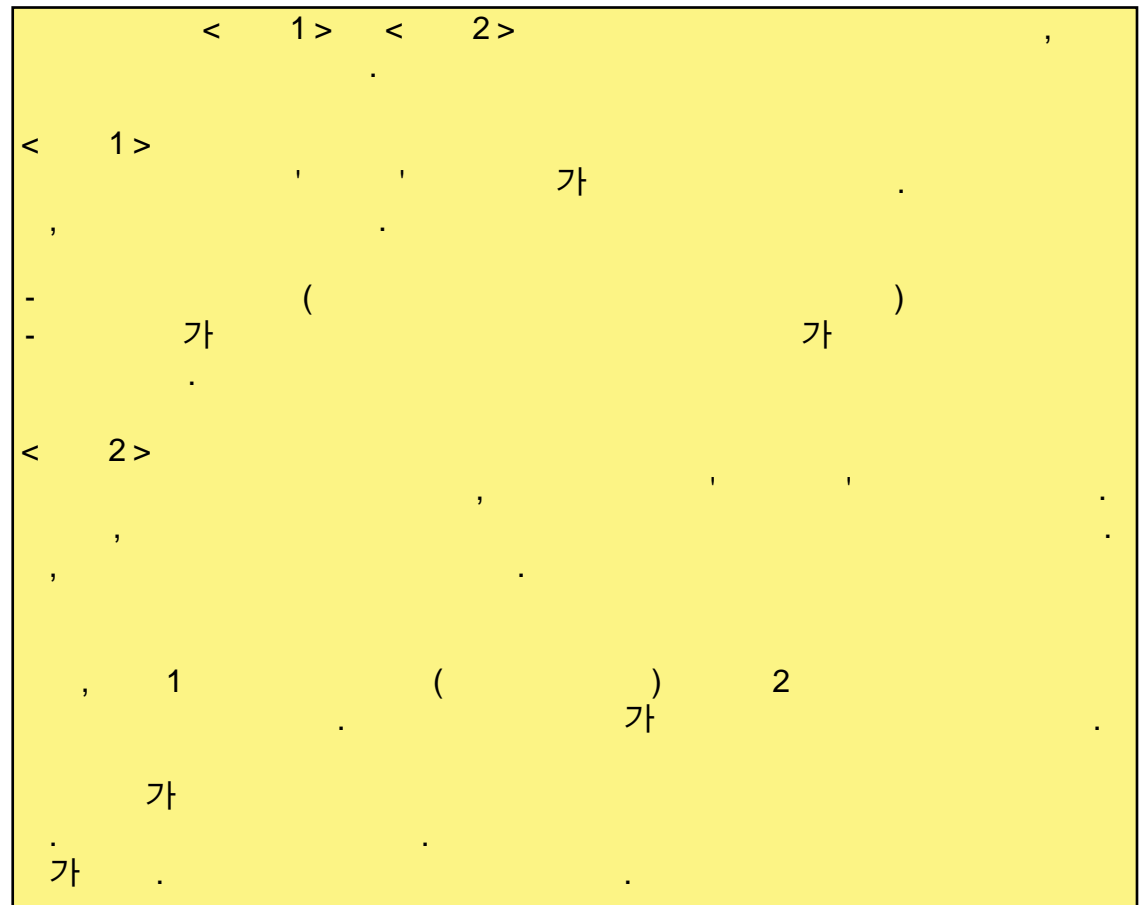
- 모든 출력 노드값을 사용하는 MSE나 교차 엔트로피와 달리 o_y 라는 하나의 노드만 사용
- o_y 는 샘플의 레이블에 해당하는 노드의 출력값
- 동작 예시1) [그림 5-5]에서 현재 샘플이 두 번째 부류라면 o_y 는 o_2 $e = -\log_2 0.0508 = 4.2990$. 잘못 분류한 셈이므로 목적함수값이 큼
- 동작 예시2) [그림 5-5]에서 현재 샘플이 세 번째 부류라면 o_y 는 o_3 $e = -\log_2 0.8360 = 0.2584$. 제대로 분류한 셈이므로 목적함수값이 작음

■ Softmax와 로그우도

- Softmax는 최댓값이 아닌 값을 억제하여 0에 가깝게 만든다는 의도 내포
- 학습 샘플이 알려주는 부류에 해당하는 노드만 보겠다는 로그우도와 잘 어울림
- 따라서 둘을 결합하여 사용하는 경우가 많음

5.2 성능 향상을 위한 요령

- 5.2.1 데이터 전처리
- 5.2.2 가중치 초기화
- 5.2.3 모멘텀
- 5.2.4 적응적 학습률
- 5.2.5 활성화함수
- 5.2.6 배치 정규화



5.2 성능 향상을 위한 요령

■ 여러 기관이 연구결과를 공유

- 『Neural Networks: Tricks of the Trade』는 연구결과를 한데 묶은 대표적인 책(1998년에 1판, 2012년에 2판)

■ 5.2절의 내용은 경험규칙

- 주어진 데이터에 잘 들어맞을지는 실험을 통해 신중히 확인해야 함
- Bengio의 권고 [Bengio2012]

“... the wisdom distilled here should be taken as a guideline, to be tried and challenged, not as a practice set in stone. ... 이 논문이 제시한 정제된 기법들은 자신의 문제에 적용한 다음 변형하여 새로운 기법을 만드는 길잡이 역할 정도로 받아들여야지 만고불변의 법칙으로 여겨서는 안 된다.”

5.2.1 데이터 전처리

■ 규모 문제

- 예) 건강에 관련된 데이터 (키(m), 몸무게(kg), 혈압)^T
 - 1.885m와 1.525m는 33cm나 차이가 나지만 특징값 차이는 불과 0.33
 - 65.5kg과 45.0kg은 20.5라는 차이
 - 첫 번째와 두 번째 특징은 대략 100배의 규모 차이
- $-\delta_j z_i$ 가 그레이디언트이기 때문에 첫 번째 특징에 연결된 가중치는 두 번째 특징에 연결된 가중치에 비해 100여 배 느리게 학습됨 → 느린 학습의 요인 .

gradient	가 100	, 10	,
- >	1.0	.	.

5.2.1 데이터 전처리

■ 모든 특징이 양수인 경우의 문제

- [그림 5-6]의 경우 ↑ 표시된 가중치는 모두 증가, ↓ 표시된 가중치는 모두 감소
- 이처럼 뭉치로 증가 또는 감소하면 최저점을 찾아가는 경로가 갈팡질팡하여 느린 수렴

, ?

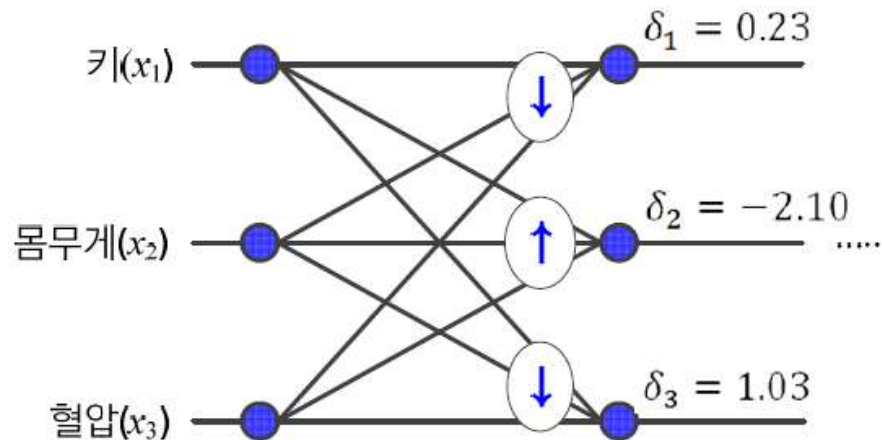


그림 5-6 특징이 모두 양수일 때 가중치가 뭉치로 갱신되는 효과

gradient가
가
gradient가
가
gradient가
가
- > 0 가

■ 식 (5.9)의 정규화는 규모 문제와 양수 문제를 해결해줌

$$x_i^{new} = \frac{x_i^{old} - \mu_i}{\sigma_i} \quad (5.9)$$

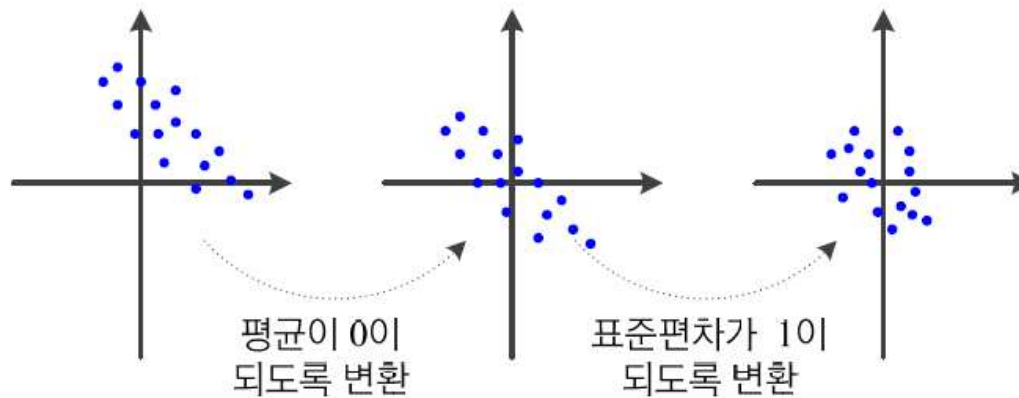


그림 5-7 표준점수로 변환

가

5.2.1 데이터 전처리

■ 명칭값을 nominal value 원핫 one-hot 코드로 변환

- 예) 성별의 남(1)과 여(2), 체질의 태양인(1), 태음인(2), 소양인(3), 소음인(4)
 - 명칭값은 거리 개념이 없음
- 원핫 코드는 값의 개수만큼 비트를 부여
 - 성별은 2비트, 체질은 4비트 부여
 - 예) 키 1.755m, 몸무게 65.5kg, 혈압 122, 남자, 소양인 샘플

$(1.755, 65.5, 122, 1, 3) \rightarrow (1.755, 65.5, 122, \underbrace{1, 0}_{\text{성별}}, \underbrace{0, 0, 1, 0}_{\text{체질}})$

5.2.2 가중치 초기화

■ 대칭적 가중치 문제

- [그림 5-8]의 대칭적 가중치에서는 z_1^{l-1} 과 z_2^{l-1} 가 같은 값이 됨. $-\delta_j z_i$ 가 그레이디언트기 때문에 u_{11}^l 과 u_{12}^l 이 같은 값으로 갱신됨 → 두 노드가 같은 일을 하는 **중복성 발생**
- **난수로 초기화**함으로써 대칭 파괴

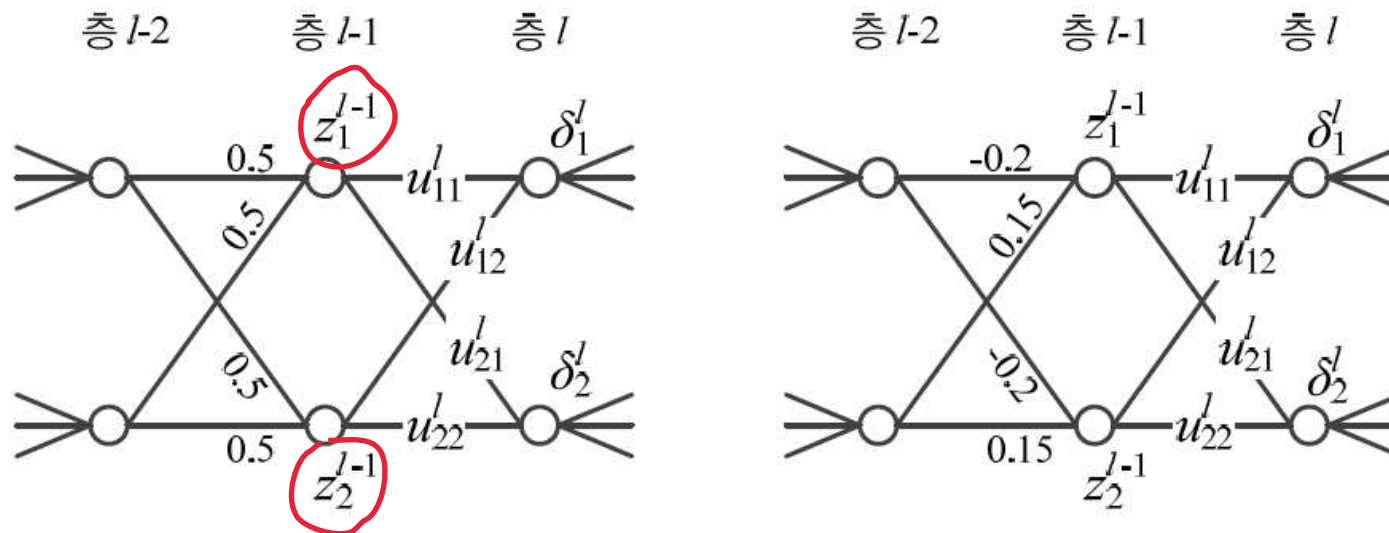


그림 5-8 대칭적 가중치로 초기화된 경우의 중복성 문제

가 가 ,
 l-1 가 .

5.2.2 가중치 초기화

■ 난수로 가중치 초기화

- 가우시언 분포 또는 균일 분포에서 난수 추출. 두 분포는 성능 차이 거의 없음
- 난수 범위는 무척 중요함 → 식 (5.10) 또는 식 (5.11)로 r 을 결정한 후 $[-r, r]$ 사이에서 난수 발생

$$r = \frac{1}{\sqrt{n_{in}}} \quad (5.10)$$

$$r = \frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}} \quad (5.11)$$

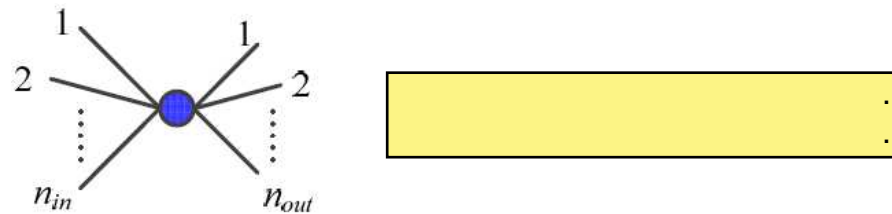


그림 5-9 노드로 들어 오는 에지 개수 n_{in} 과 노드에서 나가는 에지 개수 n_{out}

- 바이어스는 보통 0으로 초기화

가 가 0 가 , gradient가
가 .

■ 사례

- AlexNet [Krizhevsky2012]: 평균 0, 표준편차 0.001인 가우시언에서 난수 생성
- ResNet [He2016a]: 평균 0, 표준편차 $\sqrt{\frac{2}{n_{in}}}$ 인 가우시언에서 난수 생성

5.2.2 가중치 초기화

■ 또 다른 방법들

- [Saxe2014]: 가중치 벡터가 수직이 되도록 설정
- [Sussillo2014]: 임의 행로 random walk 활용하여 설정
- [Sutskever2013]: 가중치 초기화와 모멘텀을 동시에 최적화
- [Mishkin2016]: 가중치 분포가 아니라 노드의 출력값 분포가 일정하도록 강제화

[Saxe2014]
가 SVD , 가
가

[Mishkin2016]
1 : (Saxe2014) 가 .
2 : 1 () 가
가 (1
가 .
< >
- 가
- 가

5.2.3 모멘텀

■ 그레이디언트의 잡음 현상

- 기계 학습은 훈련집합을 이용하여 그레이디언트를 추정하므로 잡음 가능성 높음
- 모멘텀은 그레이디언트에 스무딩을 가하여 잡음 효과 줄임 → 수렴 속도 향상

■ 모멘텀을 적용한 가중치 갱신 수식

$$\left. \begin{aligned} \mathbf{v} &= \alpha \mathbf{v} - \rho \frac{\partial J}{\partial \Theta} \\ \Theta &= \Theta + \mathbf{v} \end{aligned} \right\} \begin{array}{l} \text{1 가} \\ \mathbf{v} \end{array} \quad (5.12)$$

- 속도 벡터 \mathbf{v} 는 이전 그레이디언트를 누적한 것에 해당함(처음에는 $\mathbf{v} = 0$ 로 출발)
- α 의 효과
 - $\alpha = 0$ 이면 모멘텀이 적용 안 된 이전 공식과 같음
 - α 가 1에 가까울수록 이전 그레이디언트 정보에 큰 가중치를 주는 셈 → Θ 가 그리는 궤적이 매끄러움
 - 보통 0.5, 0.9, 또는 0.99 사용 (또는 0.5로 시작하여 세대가 지남에 따라 점점 키워 0.99에 도달하는 방법)

5.2.3 모멘텀

vs.

- > gradient ? ? ?

■ 모멘텀의 효과

- 오버슈팅 현상 누그러뜨림

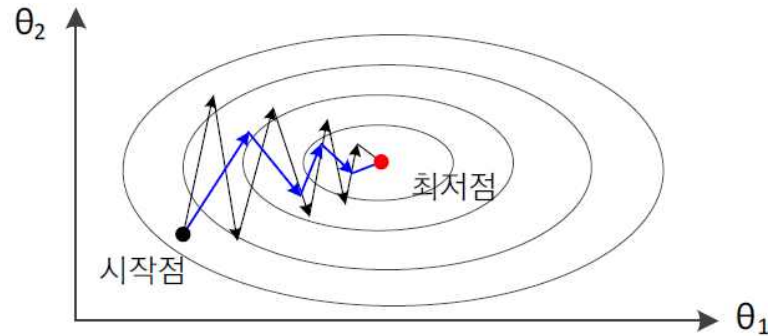


그림 5-10 모멘텀 효과

■ 네스테로프 모멘텀

- 현재 \mathbf{v} 값으로 다음 이동할 곳 $\tilde{\Theta}$ 를 예견한 후, 예견한 곳의 그레디언트 $\left. \frac{\partial J}{\partial \Theta} \right|_{\tilde{\Theta}}$ 를 사용

$$\begin{aligned}\tilde{\Theta} &= \Theta + \alpha \mathbf{v} \\ \mathbf{v} &= \alpha \mathbf{v} - \rho \left. \frac{\partial J}{\partial \Theta} \right|_{\tilde{\Theta}} \\ \Theta &= \tilde{\Theta} + \mathbf{v}\end{aligned} \quad (5.13)$$

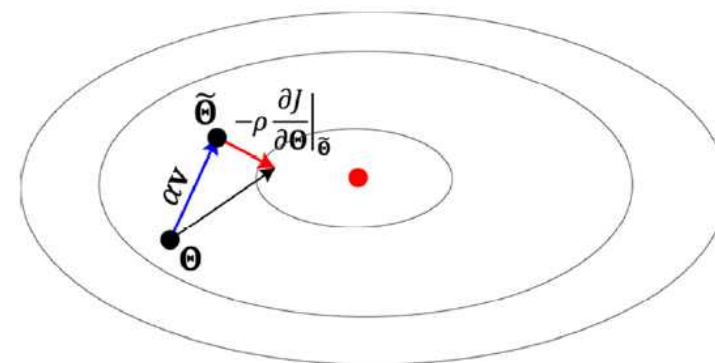


그림 5-11 네스테로프 모멘텀

가 gradient gradient 가
gradient .

5.2.3 모멘텀

■ 일반적인 경사 하강 알고리즘

알고리즘 5-1 경사 하강법

입력: 훈련집합 \mathbb{X} , \mathbb{Y} , 학습률 ρ

출력: 최적의 매개변수 $\hat{\theta}$

```
1  난수를 생성하여 초기해  $\theta$ 를 설정한다.
2  repeat
3       $\mathbf{g} = \frac{\partial J}{\partial \theta} \Big|_{\theta}$  //  $\theta$ 에서 그래디언트  $\frac{\partial J}{\partial \theta} \Big|_{\theta}$ 를 구해  $\mathbf{g}$ 라 한다.
4       $\theta = \theta - \rho \mathbf{g}$  // 식 (2.58)
5  until (멈춤 조건)
6   $\hat{\theta} = \theta$ 
```

네스테로프 모멘텀을 적용한 경사 하강 알고리즘

알고리즘 5-2 네스테로프 모멘텀을 적용한 경사 하강법

입력: 훈련집합 \mathbb{X} , \mathbb{Y} , 학습률 ρ , 속도 조절 계수 α

출력: 최적의 매개변수 $\hat{\theta}$

```
1  난수를 생성하여 초기해  $\theta$ 를 설정한다.
2   $\mathbf{v} = \mathbf{0}$  // 속도 항을  $\mathbf{0}$  벡터로 초기화
3  repeat
4       $\tilde{\theta} = \theta + \alpha \mathbf{v}$  // 예견
5       $\mathbf{g} = \frac{\partial J}{\partial \theta} \Big|_{\tilde{\theta}}$  // 예견한 곳의 그래디언트
6       $\mathbf{v} = \alpha \mathbf{v} - \rho \mathbf{g}$ 
7       $\theta = \theta + \mathbf{v}$ 
8  until (멈춤 조건)
9   $\hat{\theta} = \theta$ 
```


5.2.4 적응적 학습률

■ 학습률 ρ 의 중요성

- 너무 크면 오버슈팅에 따른 진자 현상, 너무 작으면 수렴이 느림

(learning rate)

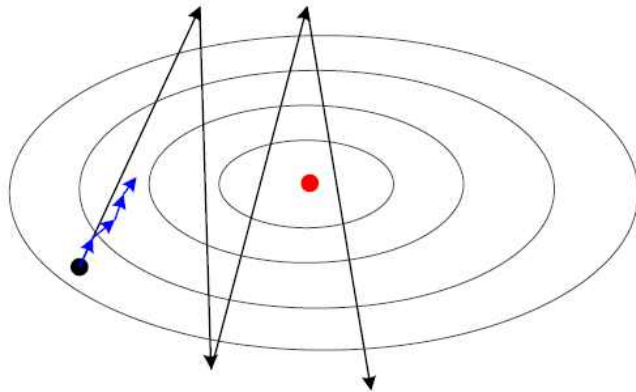


그림 5-12 학습률의 크기에 따른 최적화 알고리즘의 이동 궤적

■ 적응적 학습률

- 그레이디언트에 학습률 ρ 를 곱하면, $\rho \frac{\partial J}{\partial \theta} = \left(\rho \frac{\partial J}{\partial \theta_1}, \rho \frac{\partial J}{\partial \theta_2}, \dots, \rho \frac{\partial J}{\partial \theta_k} \right)^T$. 즉 모든 매개변수가 같은 크기의 학습률을 사용하는 셈
- 적응적 학습률은 매개변수마다 자신의 상황에 따라 학습률을 조절해 사용

가 (Jacobs1988) 가
gradient gradient

5.2.4 적응적 학습률

■ AdaGrad

- 라인 5~7을 자세히 쓰면

$$5. (r_1, r_2, \dots, r_k)^T = (r_1 + g_1^2, r_2 + g_2^2, \dots, r_k + g_k^2)^T$$

$$6. (\Delta\theta_1, \Delta\theta_2, \dots, \Delta\theta_k)^T = \left(-\frac{\rho g_1}{\epsilon + \sqrt{r_1}}, -\frac{\rho g_2}{\epsilon + \sqrt{r_2}}, \dots, -\frac{\rho g_k}{\epsilon + \sqrt{r_k}} \right)^T$$

$$7. (\theta_1, \theta_2, \dots, \theta_k)^T = (\theta_1 + \Delta\theta_1, \theta_2 + \Delta\theta_2, \dots, \theta_k + \Delta\theta_k)^T$$

- r**은 이전 그래디언트를 누적한 벡터
 - r_i 가 크면 $|\Delta\theta_i|$ 는 작아서 조금만 이동
 - r_i 가 작으면 $|\Delta\theta_i|$ 는 커서 많이 이동
 - 예) [그림 5-10]에서 θ_1 은 θ_2 보다 보폭이 큼
- 라인 6의 $\frac{\rho}{\epsilon + \sqrt{r_i}}$ 는 상황에 따라 보폭을 정해주는 적응적 학습률로 볼 수 있음

알고리즘 5-3 AdaGrad

입력: 훈련집합 \mathbb{X}, \mathbb{Y} , 학습률 ρ

출력: 최적의 매개변수 $\hat{\Theta}$

```
1  난수를 생성하여 초기해  $\Theta$ 를 설정한다.
2   $\mathbf{r} = \mathbf{0}$  // 그래디언트 누적 벡터 초기화
3  repeat
4      그래디언트  $\mathbf{g} = \frac{\partial J}{\partial \Theta} \Big|_{\Theta}$ 를 구한다.
5       $\mathbf{r} = \mathbf{r} + \mathbf{g} \odot \mathbf{g}$  //  $\odot$ 는 요소별 곱
6       $\Delta\Theta = -\frac{\rho}{\epsilon + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ 
7       $\Theta = \Theta + \Delta\Theta$ 
8  until (멈춤 조건)
9   $\hat{\Theta} = \Theta$ 
```

5.2.4 적응적 학습률

■ RMSProp

■ AdaGrad의 단점

- [알고리즘 5-3]의 라인 5는 단순히 제곱을 더함
- 따라서 오래된 그레디언트와 최근 그레디언트는 같은 비중의 역할 → **r이 점점 커져** 수렴 방해할 가능성

■ RMSProp은 가중 이동 평균 기법 적용

$$\mathbf{r} = \alpha \mathbf{r} + (1 - \alpha) \mathbf{g} \odot \mathbf{g} \quad (5.14)$$

- α 가 작을수록 최근 것에 비중을 둠
- 보통 α 로 0.9, 0.99, 0.999를 사용

AdaGrad	gradient
gradient	. RMSProp
gradient	.

알고리즘 5-4 RMSProp

입력: 훈련집합 \mathbb{X}, \mathbb{Y} , 학습률 ρ , 가중 이동 평균 계수 α

출력: 최적의 매개변수 $\hat{\Theta}$

```
1 난수를 생성하여 초기해  $\Theta$ 를 설정한다.
2  $\mathbf{r} = \mathbf{0}$  // 그레디언트 누적 벡터 초기화
3 repeat
4   그레디언트  $\mathbf{g} = \frac{\partial J}{\partial \Theta} \Big|_{\Theta}$ 를 구한다.
5    $\mathbf{r} = \alpha \mathbf{r} + (1 - \alpha) \mathbf{g} \odot \mathbf{g}$  //  $\odot$ 는 요소별 곱
6    $\Delta \Theta = -\frac{\rho}{\epsilon + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ 
7    $\Theta = \Theta + \Delta \Theta$ 
8 until (멈춤 조건)
9  $\hat{\Theta} = \Theta$ 
```

decay .

5.2.4 적응적 학습률

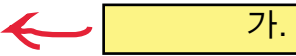
■ Adam

- RMSProp에 식 (5-12)의 모멘텀을 추가로 적용한 알고리즘

알고리즘 5-5 Adam

입력: 훈련집합 \mathbb{X}, \mathbb{Y} , 학습률 ρ , 모멘텀 계수 α_1 , 가중 이동 평균 계수 α_2

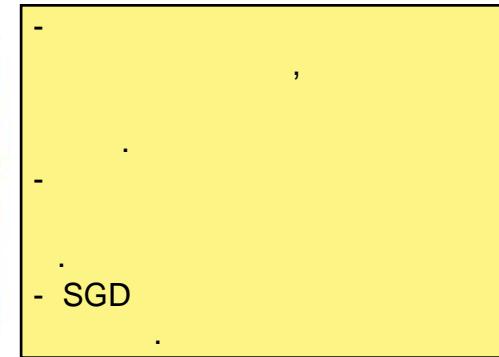
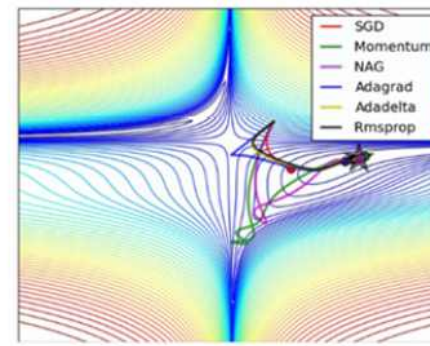
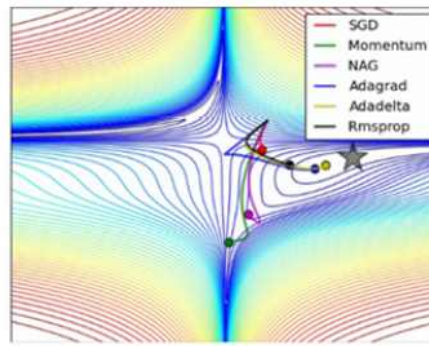
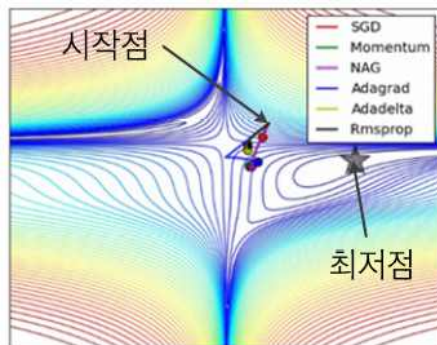
출력: 최적의 매개변수 $\hat{\theta}$

```
1  난수를 생성하여 초기해  $\theta$ 를 설정한다.
2   $\mathbf{v} = \mathbf{0}, \mathbf{r} = \mathbf{0}$ 
3   $t = 1$ 
4  repeat
5      그래디언트  $\mathbf{g} = \frac{\partial J}{\partial \theta} \Big|_{\theta}$ 를 구한다.
6       $\mathbf{v} = \alpha_1 \mathbf{v} + (1 - \alpha_1) \mathbf{g}$  // 속도 벡터 
7       $\mathbf{v} = \frac{1}{1 - (\alpha_1)^t} \mathbf{v}$ 
8       $\mathbf{r} = \alpha_2 \mathbf{r} + (1 - \alpha_2) \mathbf{g} \odot \mathbf{g}$  // 그래디언트 누적 벡터
9       $\mathbf{r} = \frac{1}{1 - (\alpha_2)^t} \mathbf{r}$ 
10      $\Delta \theta = -\frac{\rho}{\epsilon + \sqrt{\mathbf{r}}} \mathbf{v}$ 
11      $\theta = \theta + \Delta \theta$ 
12      $t++$ 
13 until (멈춤 조건)
14  $\hat{\theta} = \theta$ 
```

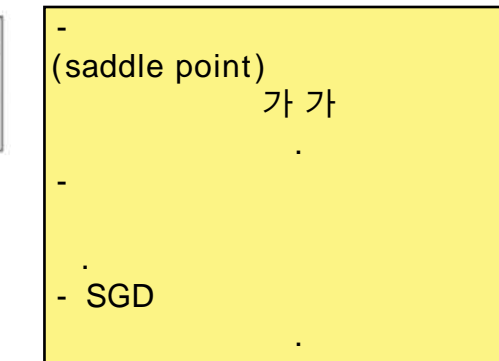
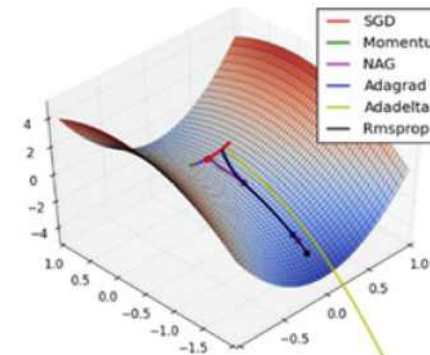
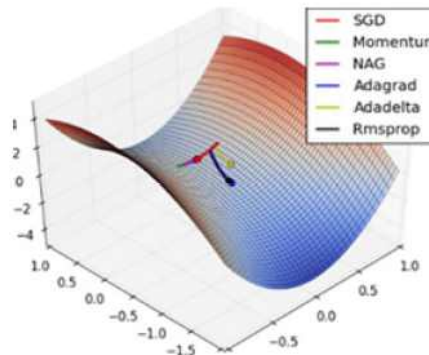
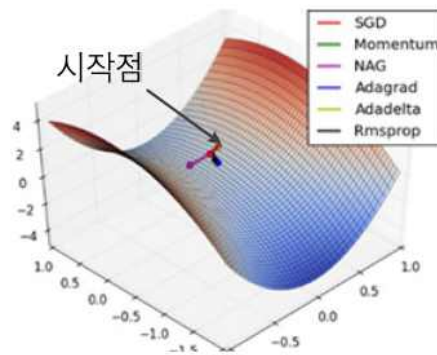
5.2.4 적응적 학습률

■ 동작 예시 실시간 애니메이션 <http://cs231n.github.io/neural-networks-3>

- [그림 5-13(a)]는 중앙으로 급강하하는 절벽 지형
- [그림 5-13(b)]는 중앙 부근에 안장점이 saddle point 있는 지형



1. (a) 협곡 지형



2. (b) 안장점 지형

그림 5-13 모멘텀과 적응적 학습률 기법의 수렴 특성을 보여주는 예제 4

5.2.5 활성화 함수

- 활성화값 z 를 계산하고 활성화함수 τ 를 적용하는 과정

$$\begin{aligned} z &= \mathbf{w}^T \tilde{\mathbf{x}} + b \\ y &= \tau(z) \end{aligned} \quad (5.15)$$

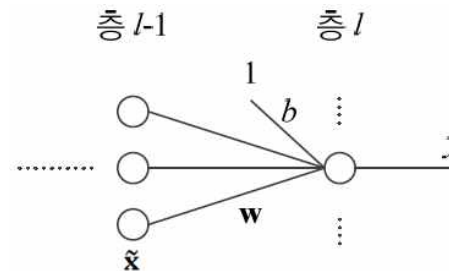
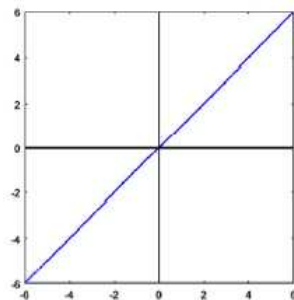
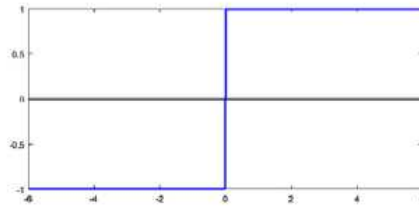


그림 5-14 신경망 노드의 연산

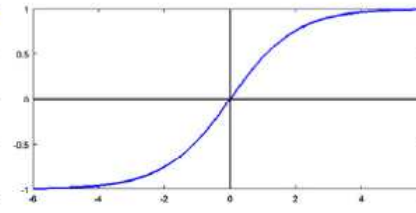
- 시대별 활성화함수



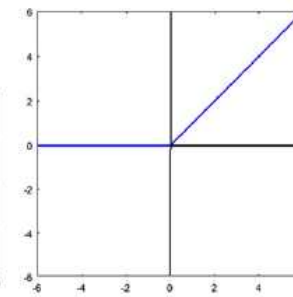
(a) 선형



(b) 계단(1950년대)



(c) tanh(1980년대)



(d) ReLU(2000년경~현재)

그림 5-15 활성화함수 τ

MLP

- X

- tanh는 활성화값이 커지면 포화 상태가 되고 그래디언트는 0에 가까워짐 → 매개변수 갱신(학습)이 매우 느린 요인

5.2.5 활성 함수

■ ReLU(Rectified Linear Unit) 활성화 함수

- 포화 문제 해소

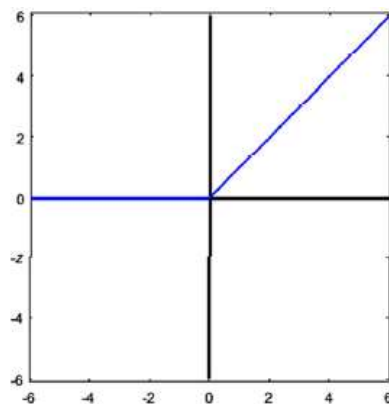
$$z = \mathbf{w}^T \tilde{\mathbf{x}} + b$$

$$y = \text{ReLU}(z) = \max(0, z)$$

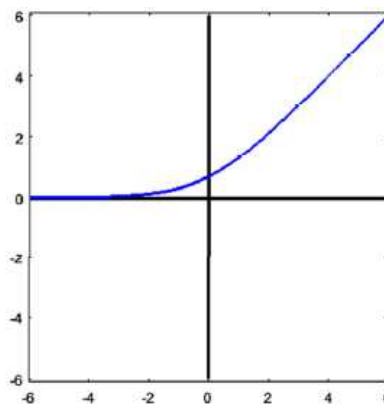
ReLU 0 가 가 .
 , , gradient가 가
 , 가 가
 . [Glorot2011]

■ ReLU의 변형

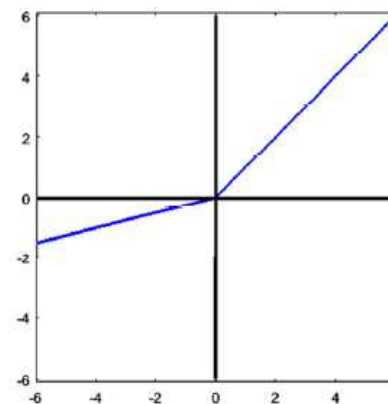
- Leaky ReLU (보통 $\alpha = 0.01$ 을 사용) $\text{leakyReLU}(z) = \begin{cases} z, & z \geq 0 \\ \alpha z, & z < 0 \end{cases} \quad (5.17)$
- PReLU (α 를 학습으로 알아냄)



(a) ReLU



(b) softplus



(c) leakyReLU와 PReLU

그림 5-16 ReLU의 변형

ReLU softplus
 , ReLU가 .

5.2.6 배치 정규화

■ 공변량 시프트 covariate shift 현상

- 학습이 진행되면서 층1의 매개변수가 바뀔에 따라 $\tilde{\mathbf{x}}^{(1)}$ 이 따라 바뀜 → 층2 입장에서 보면 자신에게 입력되는 데이터의 분포가 수시로 바뀌는 셈
- 층2, 층3, ...으로 깊어짐에 따라 더욱 심각
- 학습을 방해하는 요인으로 작용

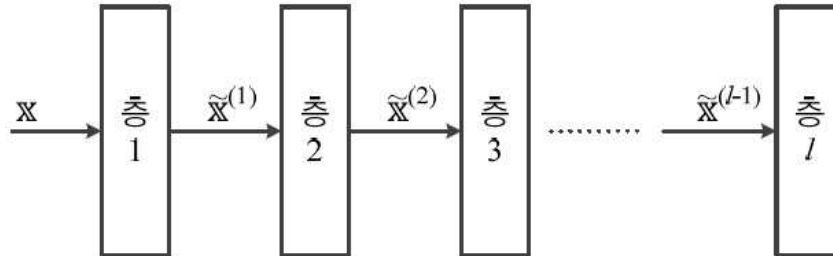


그림 5-17 공변량 시프트 현상

		?	
1	()	2

		?	
-		.	
-			.
-			.

5.2.6 배치 정규화

■ 배치 정규화

- 공변량 시프트 현상을 누그러뜨리기 위해 식 (5.9)의 정규화를 모든 층에 적용하는 기법

$$x_i^{new} = \frac{x_i^{old} - \mu_i}{\sigma_i} \quad , \quad \dots ? \quad (5.9)$$

- 정규화를 적용하는 곳이 중요
 - 식 (5.15)의 연산 과정 중 식 (5.9)를 어디에 적용하나?

$$\begin{aligned} z &= \mathbf{w}^T \tilde{\mathbf{x}} + b \\ y &= \tau(z) \end{aligned} \quad (5.15)$$

x ? z ? ?
[Ioffe2015] z

- 입력 $\tilde{\mathbf{x}}$ 또는 중간 결과 z 중 어느 것에 적용? → z 에 적용하는 것이 유리
- 훈련집합 전체 또는 미니배치 중 어느 것에 적용?
 - 미니배치에 적용하는 것이 유리 , .

5.2.6 배치 정규화

정규화 변환을 수행하는 코드

- 미니배치 $\mathbb{X}_B = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ 에 식 (5.15)를 적용하여 $\tilde{\mathbb{X}}_B = \{z_1, z_2, \dots, z_m\}$ 를 얻은 후, $\tilde{\mathbb{X}}_B$ 를 가지고 코드 1을 수행
- 노드마다 독립적으로 코드 1을 수행
- α 와 β 는 노드마다 고유한 매개변수로서 학습으로 알아냄

코드 1:

$$\mu_B = \frac{1}{m} \sum_{i=1}^m z_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (z_i - \mu_B)^2$$

$$\tilde{z}_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}, \quad i = 1, 2, \dots, m$$

$$z'_i = \gamma \tilde{z}_i + \beta, \quad i = 1, 2, \dots, m$$

가

$$\rightarrow z = \mathbf{w}^T \tilde{\mathbf{x}} + \beta$$

5.2.6 배치 정규화

■ 최적화를 마친 후 추가적인 후처리 작업 필요

- 각 노드는 전체 훈련집합을 가지고 독립적으로 코드2를 수행

코드 2:

$$\mu = \frac{1}{n} \sum_{i=1}^n z_i$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (z_i - \mu)^2$$

노드에 $\mu, \sigma^2, \gamma, \beta$ 를 저장한다. // 예측 단계에서 식 (5.18)로 변환을 수행하기 위함

가

■ 예측 단계

- 각 노드는 독립적으로 식 (5.18)을 적용(코드 1의 마지막 두 라인을 수행하는 셈)

$$z' = \frac{\gamma}{\sqrt{\sigma^2 + \epsilon}} z + \left(\beta - \frac{\gamma \mu}{\sqrt{\sigma^2 + \epsilon}} \right)$$

가

, z
z'

$$z = \mathbf{w}^T \tilde{\mathbf{x}} + b$$

$$y = \tau(z)$$

5.2.6 배치 정규화

■ CNN에서는

- 노드 단위가 아니라 특징 맵 단위로 코드 1과 코드 2를 적용

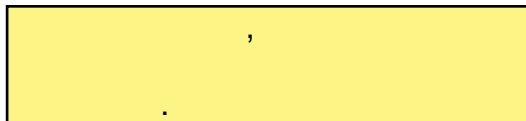


■ 배치 정규화의 긍정적 효과를 측정한 실험사례[Ioffe2015]

- 가중치 초기화에 덜 민감함
- 학습률을 크게 하여 수렴 속도 향상 가능
- 시그모이드를 활성화함수로 사용하는 깊은 신경망도 학습이 이루어짐

■ 배치 정규화는 규제 효과를 제공

- 드롭아웃이라는 규제 기법을 적용하지 않아도 높은 성능



5.3 규제의 필요성과 원리

- 5.3.1 과잉적합에 빠지는 이유와 과잉적합을 피하는 전략
- 5.3.2 규제의 정의

- 규제가 중요하기 때문에 1장에서 미리 소개한 내용
 - 1.5절의 과소적합과 과잉적합, 바이어스와 분산([그림 1-13], [그림 1-14])
 - 1.6절의 데이터 확대와 가중치 감소

5.3.1 과잉적합에 빠지는 이유와 과잉적합을 피하는 전략

■ 학습 모델의 용량에 따른 일반화 능력

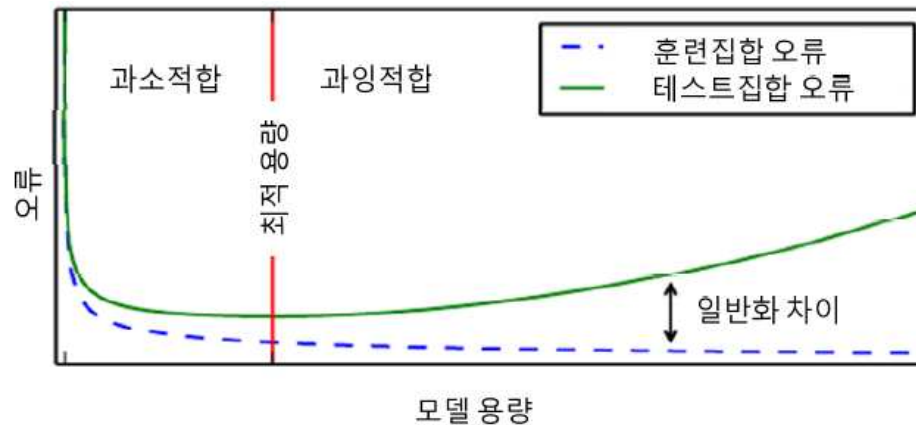


그림 5-18 학습 모델의 용량과 일반화 능력의 관계

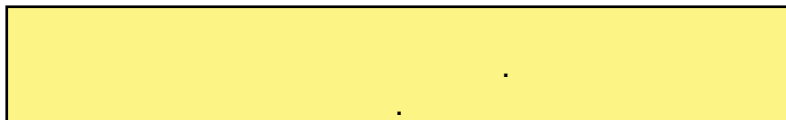
■ 대부분 가지고 있는 데이터에 비해 훨씬 **큰 용량의 모델**을 사용

- 예) VGGNet은 분류층에 1억 2천 1백만 개의 매개변수
- 훈련집합을 단순히 '암기'하는 과잉적합에 주의를 기울여야 함

■ **현대 기계 학습의 전략**



- **충분히 큰 용량의 모델을 설계한 다음, 학습 과정에서 여러 규제 기법을 적용**



5.3.2 규제 정의

- 규제는 오래 전부터 수학과 통계학에서 연구해온 주제
 - 모델 용량에 비해 데이터가 부족한 경우의 불량 문제를 ill-posed problem 푸는 데 사용
 - 적절한 가정을 투입하여 문제를 풀 → 입력과 출력 사이의 매핑은 **매끄럽다**는 사전 지식

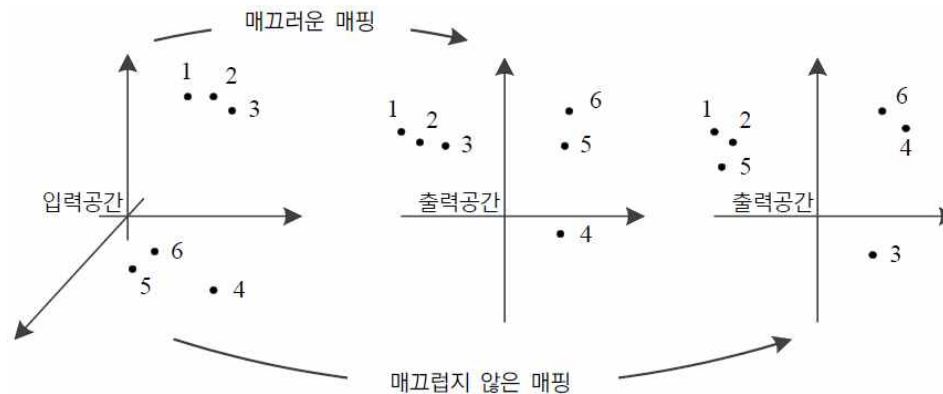


그림 5-20 사전 지식으로서 매끄러움의 특성

* (smooth function)
 가 가
 가 . (ex. 4
 , 12)
 * 가
 ()

- 티호노프의 규제 기법은 매끄러움 가정에 기반을 둔 식 (5.19)를 사용

$$\underbrace{J_{\text{regularized}}(\Theta)}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta)}_{\text{목적함수}} + \lambda \underbrace{R(\Theta)}_{\text{규제 항}} \rightarrow \boxed{} \quad (5.19)$$

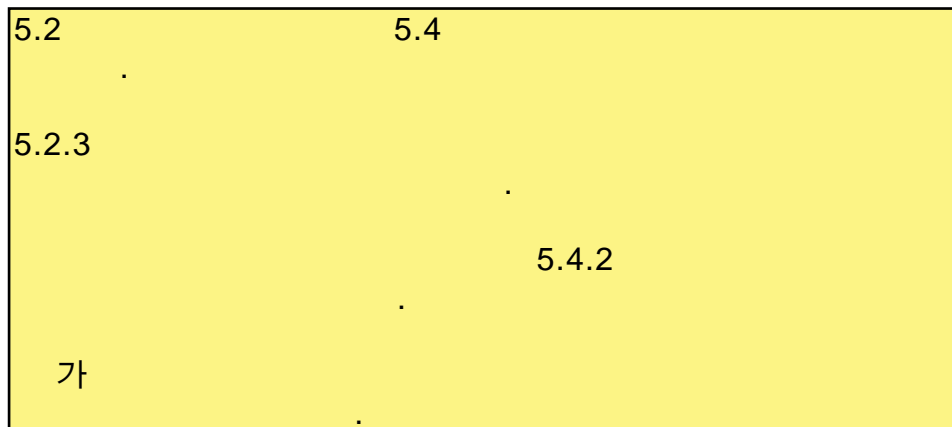
5.3.2 규제 정의

■ 현대 기계 학습도 매끄러움 가정을 널리 사용함

- 5.4.1절의 가중치 감쇠 기법
 - 모델의 구조적 용량을 충분히 크게 하고, '수치적 용량'을 제한하는 규제 기법
- 6장의 비지도 학습 등

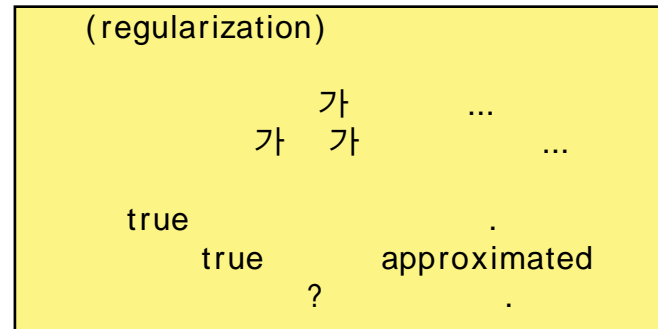
■ 『Deep Learning』 책의 정의

"...any modification we make to a learning algorithm that is intended to reduce its generalization error ... 일반화 오류를 줄이려는 의도를 가지고 학습 알고리즘을 수정하는 방법 모두"



5.4 규제 기법

- 5.4.1 가중치 벌칙
- 5.4.2 조기 멈춤
- 5.4.3 데이터 확대
- 5.4.4 드롭아웃
- 5.4.5 앙상블 기법



■ 명시적 규제와 암시적 규제

- **명시적 규제:** 가중치 감쇠나 드롭아웃처럼 목적함수나 신경망 구조를 직접 수정하는 방식
- **암시적 규제:** 조기 멈춤, 데이터 증대, 잡음 추가, 앙상블처럼 간접적으로 영향을 미치는 방식

5.4.1 가중치 벌칙

- 식 (5.19)를 관련 변수가 드러나도록 다시 쓰면,

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \underbrace{\lambda R(\Theta)}_{\text{규제 항}} \quad (5.20)$$

- 규제항은 훈련집합과 무관하며, 데이터 생성 과정에 내재한 사전 지식에 해당
- 규제항은 매개변수를 작은 값으로 유지하므로 모델의 용량을 제한하는 역할(수치적 용량을 제한함)

$$\lambda \left(\sum_{i=1}^n \sum_{j=1}^m \Theta_{ij}^2 \right) \quad \text{가}$$

- 규제항 $R(\Theta)$ 로 무엇을 사용할 것인가?

- 큰 가중치에 벌칙을 가해 작은 가중치를 유지하려고 주로 $L2$ 놈이나 $L1$ 놈을 사용

5.4.1 가중치 벌칙

■ $L2$ 놈 , ridge regression .

- 규제 항 R 로 $L2$ 놈을 사용하는 규제 기법을 '가중치 감쇠'라 weight decay 부름 → 식 (5.21)

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \underbrace{\lambda \|\Theta\|_2^2}_{\text{규제 항}} \quad (5.21)$$

- 식 (5.21)의 그레이디언트 계산

$$\nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) = \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda\Theta \quad (5.22)$$

5.4.1 가중치 벌칙

- 식 (5.22)를 이용하여 매개변수를 갱신하는 수식

$$\begin{aligned}\Theta &= \Theta - \rho \nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) \\ &= \Theta - \rho(\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda\Theta) \quad \longrightarrow \quad \underline{\Theta = (1 - 2\rho\lambda)\Theta - \rho\nabla J} \quad (5.23) \\ &= (1 - 2\rho\lambda)\Theta - \rho\nabla J(\Theta; \mathbb{X}, \mathbb{Y})\end{aligned}$$

- $\lambda = 0$ 으로 두면 규제를 적용하지 않은 원래 식 $\Theta = \Theta - \rho\nabla J$ 가 됨
- 가중치 감쇠는 단지 Θ 에 $(1 - 2\rho\lambda)$ 를 곱해주는 셈
 - 예를 들어, $\rho=0.01$, $\lambda = 2.0$ 이라면 $(1 - 2\rho\lambda)=0.96$
- 최종해를 원점 가까이 당기는 효과 (즉 가중치를 작게 유지함)

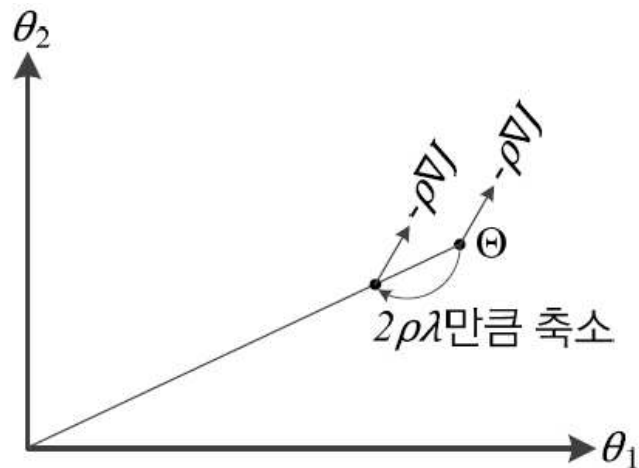


그림 5-21 L2 놈을 사용한 가중치 감쇠 기법의 효과

5.4.1 가중치 벌칙

■ 선형 회귀에 적용

- 선형 회귀는 훈련집합 $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, $\mathbb{Y} = \{y_1, y_2, \dots, y_n\}$ 이 주어지면, 식 (5.24)를 풀어 $\mathbf{w} = (w_1, w_2, \dots, w_d)^T$ 를 구하는 문제. 이때 $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$

$$\boxed{n} \left\{ \begin{aligned} w_1 x_{i1} + w_2 x_{i2} \cdots + w_d x_{id} = \mathbf{x}_i^T \mathbf{w} = y_i, \quad i = 1, 2, \dots, n \end{aligned} \right. \quad (5.24)$$

- 식 (5.24)를 행렬식으로 바꿔 쓰면,

$$\begin{pmatrix} 0 & \text{가} & . \end{pmatrix}$$

$$\boxed{1} \left\{ \mathbf{Xw} = \mathbf{y} \right. \quad (5.25)$$

- 가중치 감소를 적용한 목적함수 $\boxed{\text{MSE}}$.

$$J_{\text{regularized}}(\mathbf{w}) = \|\mathbf{Xw} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 = (\mathbf{Xw} - \mathbf{y})^T (\mathbf{Xw} - \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2 \quad (5.27)$$

ridge .

5.4.1 가중치 벌칙

- 식 (5.27)을 미분하여 0으로 놓으면,

$$\frac{\partial J_{regularized}}{\partial \mathbf{w}} = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} = \mathbf{0} \Rightarrow (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (5.28)$$

- 식 (5.28)을 정리하면,

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (5.29)$$

- 공분산 행렬 $\mathbf{X}^T \mathbf{X}$ 의 대각 요소가 2λ 만큼씩 증가 \rightarrow 역행렬을 곱하므로 가중치를 축소하여 원점으로 당기는 효과 ([그림 5-21])

weight decay

- 예측 단계에서는.

$$y = \mathbf{x}^T \hat{\mathbf{w}} \quad (5.30)$$

5.4.1 가중치 벌칙

예제 5-1 리지 회귀

훈련집합 $\mathbb{X} = \{\mathbf{x}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 3 \\ 3 \end{pmatrix}\}$, $\mathbb{Y} = \{y_1 = 3.0, y_2 = 7.0, y_3 = 8.8\}$ 이 주어졌다고 가정하자. 특징 벡터가 2차원이므로 $d=2$ 이고 샘플이 3개이므로 $n=3$ 이다. 훈련집합으로 설계행렬 \mathbf{X} 와 레이블 행렬 \mathbf{y} 를 다음과 같이 쓸 수 있다.

$$\mathbf{X} = \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 3.0 \\ 7.0 \\ 8.8 \end{pmatrix}$$

이 값들을 식 (5.29)에 대입하여 다음과 같이 $\hat{\mathbf{w}}$ 을 구할 수 있다. 이때 $\lambda = 0.25$ 라 가정하자.

$$\hat{\mathbf{w}} = \left(\begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{pmatrix} + \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{pmatrix} \begin{pmatrix} 3.0 \\ 7.0 \\ 8.8 \end{pmatrix} = \begin{pmatrix} 1.4916 \\ 1.3607 \end{pmatrix}$$

따라서 하이퍼 평면은 $y = 1.4916x_1 + 1.3607x_2$ 이다. 새로운 샘플로 $\mathbf{x} = (5 \ 4)^T$ 가 입력되면 식 (5.30)을 이용하여 12.9009를 예측한다.

5.4.1 가중치 벌칙

■ MLP와 DMLP에 적용

- 식 (3.21)에 식 (5.23)의 가중치 감쇠라는 규제 기법을 적용하면,

$$\left. \begin{array}{l} \mathbf{U}^1 = \mathbf{U}^1 - \rho \frac{\partial J}{\partial \mathbf{U}^1} \\ \mathbf{U}^2 = \mathbf{U}^2 - \rho \frac{\partial J}{\partial \mathbf{U}^2} \end{array} \right\} (3.21) \longrightarrow \left. \begin{array}{l} \mathbf{U}^1 = (1 - 2\rho\lambda)\mathbf{U}^1 - \rho \frac{\partial J}{\partial \mathbf{U}^1} \\ \mathbf{U}^2 = (1 - 2\rho\lambda)\mathbf{U}^2 - \rho \frac{\partial J}{\partial \mathbf{U}^2} \end{array} \right\} (5.31)$$

- [알고리즘 3-4]에 적용하면,

13. for ($k=1$ to c) for ($j=0$ to ρ) $u_{kj}^2 = u_{kj}^2 - \rho \Delta u_{kj}^2$ // 가중치 감쇠 적용하지 않은 원래 알고리즘

14. for ($j=1$ to ρ) for ($i=0$ to d) $u_{ji}^1 = u_{ji}^1 - \rho \Delta u_{ji}^1$

↓

13. for ($k=1$ to c) for ($j=0$ to ρ) $u_{kj}^2 = (1 - 2\rho\lambda)u_{kj}^2 - \rho \Delta u_{kj}^2$ // 가중치 감쇠 적용한 알고리즘

14. for ($j=1$ to ρ) for ($i=0$ to d) $u_{ji}^1 = (1 - 2\rho\lambda)u_{ji}^1 - \rho \Delta u_{ji}^1$

5.4.1 가중치 벌칙

- [알고리즘 3-6](미니배치 버전)에 적용하면,

$$14. \quad \mathbf{U}^2 = \mathbf{U}^2 - \rho \frac{\Delta \mathbf{U}^2}{t} \quad // \text{가중치 감쇠 적용하지 않은 원래 알고리즘}$$

$$15. \quad \mathbf{U}^1 = \mathbf{U}^1 - \rho \frac{\Delta \mathbf{U}^1}{t}$$

⇓

$$14. \quad \mathbf{U}^2 = (1 - 2\rho\lambda)\mathbf{U}^2 - \rho \frac{\Delta \mathbf{U}^2}{t} \quad // \text{가중치 감쇠 적용한 알고리즘}$$

$$15. \quad \mathbf{U}^1 = (1 - 2\rho\lambda)\mathbf{U}^1 - \rho \frac{\Delta \mathbf{U}^1}{t}$$

- DMLP를 위한 [알고리즘 4-1]에 적용하면,

$$16. \quad \text{for } (l=L \text{ to } 1) \quad // \text{가중치 감쇠 적용하지 않은 원래 알고리즘}$$

$$17. \quad \text{for } (j=1 \text{ to } n_l) \text{ for } (i=0 \text{ to } n_{l-1}) \quad u_{ji}^l = u_{ji}^l - \rho \left(\frac{1}{t}\right) \Delta u_{ji}^l$$

⇓

$$16. \quad \text{for } (l=L \text{ to } 1) \quad // \text{가중치 감쇠 적용한 알고리즘}$$

$$17. \quad \text{for } (j=1 \text{ to } n_l) \text{ for } (i=0 \text{ to } n_{l-1}) \quad u_{ji}^l = (1 - 2\rho\lambda)u_{ji}^l - \rho \left(\frac{1}{t}\right) \Delta u_{ji}^l$$

5.4.1 가중치 벌칙

■ $L1$ 놈

- 규제 항으로 $L1$ 놈을 적용하면, ($L1$ 놈은 $\|\Theta\|_1 = |\theta_1| + |\theta_2| + \dots$)

$$\underbrace{J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{\|\Theta\|_1}_{\text{규제 항}} \quad (5.32)$$

- 식 (5.32)를 미분하면,

$$\nabla J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y}) = \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + \lambda \mathbf{sign}(\Theta) \quad (5.33)$$

sign	1,
- 1	가 .

- 매개변수를 갱신하는 식에 대입하면,

$$\begin{aligned} \Theta &= \Theta - \rho \nabla J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y}) \\ &= \Theta - \rho (\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + \lambda \mathbf{sign}(\Theta)) \\ &= \Theta - \rho \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) - \rho \lambda \mathbf{sign}(\Theta) \end{aligned}$$

5.4.1 가중치 벌칙

- 매개변수를 갱신하는 식

$$\Theta = \Theta - \rho \nabla J - \rho \lambda \text{sign}(\Theta) \quad (5.34)$$

- 식 (5.34)의 가중치 감쇠 효과

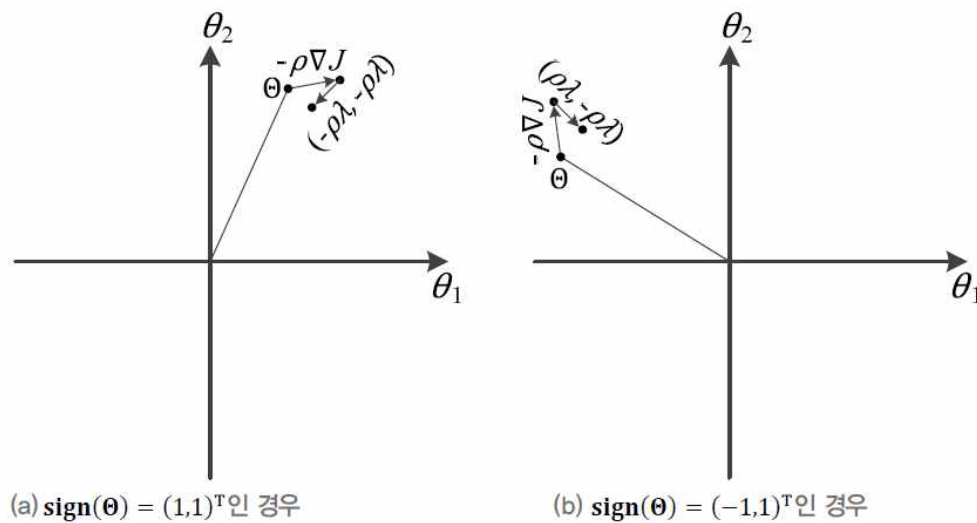


그림 5-22 L1 놈을 사용한 가중치 감쇠 기법의 효과

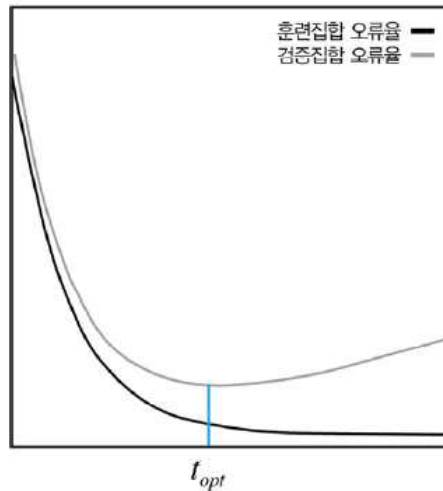
- L1 놈의 희소성 효과(0이 되는 매개변수가 많음)
 - 선형 회귀에 적용하면 특징 선택 효과

sign

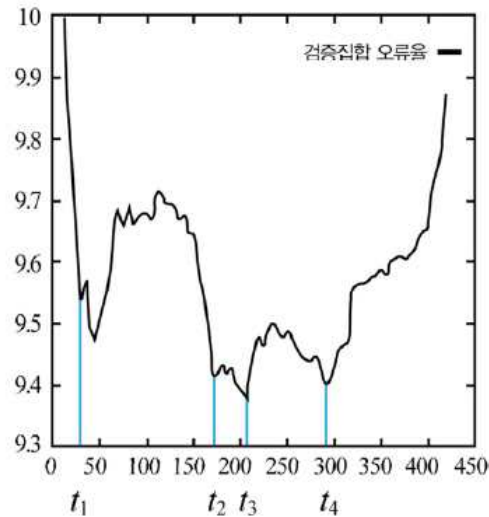
0 가

5.4.2 조기 멈춤

- 학습 시간에 따른 일반화 능력 [그림 5-23(a)]
 - 일정 시간(t_{opt})이 지나면 과잉적합 현상이 나타남 → 일반화 능력 저하
 - 즉 훈련 데이터를 단순히 암기하기 시작



(a) 개념적인 도표



(b) 실제 데이터에 나타나는 지그재그 현상

그림 5-23 학습 시간에 따른 성능 추이

t	p
(p:)

- 조기 멈춤이라는 규제 기법
 - 검증집합의 오류가 최저인 점 t_{opt} 에서 학습을 멈춤

5.4.2 조기 멈춤

- [알고리즘 5-6]은 현실을 제대로 반영하지 않은 순진한 버전
 - [그림 5-23(a)] 상황에서 동작

알고리즘 5-6 조기 멈춤을 채택한 기계 학습 알고리즘(지그재그 현상을 고려하지 않은 순진한 버전)

입력: 훈련집합 \mathbb{X} 와 \mathbb{Y} , 검증집합 \mathbb{X}' 와 \mathbb{Y}'

출력: 최적의 매개변수 $\hat{\theta}$, 최적해가 발생한 세대 \hat{t}

```
1  난수를 생성하여 초기해  $\theta_0$ 을 설정하고 오류율  $e_0 = 1.0$ 으로 설정한다. // 1.0은 오류율 최대치
2   $t=0$ 
3  while (true)
4      학습 알고리즘으로  $\theta_t$ 를 갱신하여  $\theta_{t+1}$ 을 얻는다.
5       $\theta_{t+1}$ 로 검증집합에 대한 오류율  $e_{t+1}$ 을 측정한다.
6      if( $e_{t+1} > e_t$ ) break
7       $t++$ 
8   $\hat{\theta} = \theta_t, \hat{t} = t$ 
```

5.4.2 조기 멈춤

■ 실제 세계는 [그림 5-23(b)]와 같은 상황

- 순진한 버전을 적용하면 t_1 에서 멈추므로 설익은 수렴
- 이에 대처하는 여러 가지 방안 중에서 [알고리즘 5-7]은 **참을성**을 반영한 버전

알고리즘 5-7 조기 멈춤을 채택한 기계 학습 알고리즘(참을성을 반영한 버전)

입력: 훈련집합 \mathbb{X} 와 \mathbb{Y} , 검증집합 \mathbb{X}' 와 \mathbb{Y}' , 참을성 인자 p , 세대 반복 인자 q

출력: 최적의 매개변수 $\hat{\theta}$, 최적해가 발생한 세대 \hat{t}

```
1  난수를 생성하여 초기해  $\theta_0$ 을 설정한다.
2   $\hat{\theta} = \theta_0, \hat{t} = 0$ 
3   $t = 0, \hat{e} = 1.0, j = 0$ 
4  while ( $j < p$ )
5      학습 알고리즘의 세대를  $q$ 번 반복하여  $\theta_{t+q}$ 를 얻는다.
6       $\theta_{t+q}$ 로 검증집합에 대한 오류율  $e_{t+q}$ 를 측정한다.
7      if ( $e_{t+q} < \hat{e}$ ) // 새로운 최적을 발견한 상황
8           $j = 0$  // 참는 과정을 처음부터 새로 시작
9           $\hat{\theta} = \theta_{t+q}, \hat{e} = e_{t+q}, \hat{t} = t + q$ 
10     else
11          $j = j + 1$ 
12      $t = t + q$ 
```

p 가
 q
()

5.4.3 데이터 확대

- 과잉적합 방지하는 가장 확실한 방법은 큰 훈련집합 사용

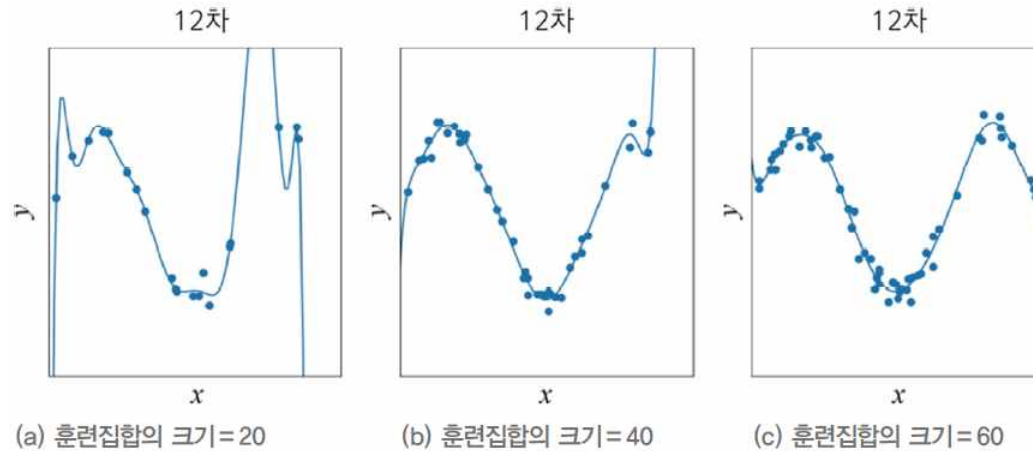


그림 1-17 데이터를 확대하여 일반화 능력을 향상함

- 하지만 데이터 수집은 비용이 많이 드는 작업
- 데이터 확대라는 규제 기법
 - 데이터를 인위적으로 변형하여 확대함
 - 자연계에서 벌어지는 잠재적인 변형을 프로그램으로 흉내 내는 셈

5.4.3 데이터 확대

- 예) MNIST에 어파인 변환(이동, 회전, 크기)을 적용

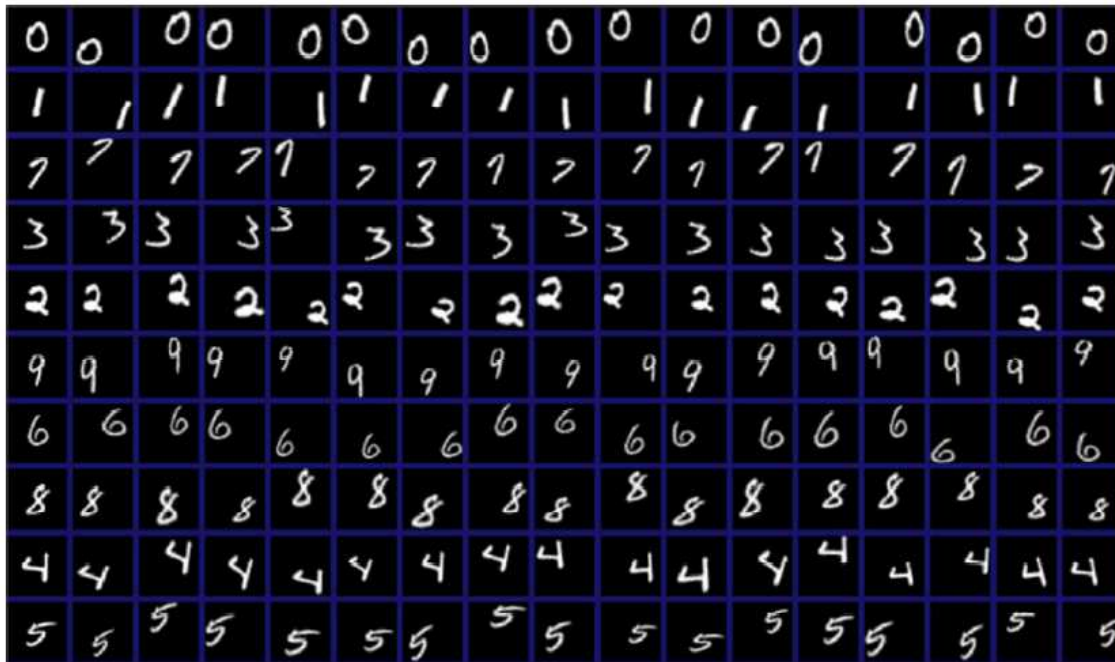
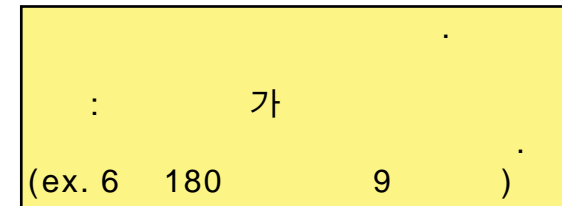
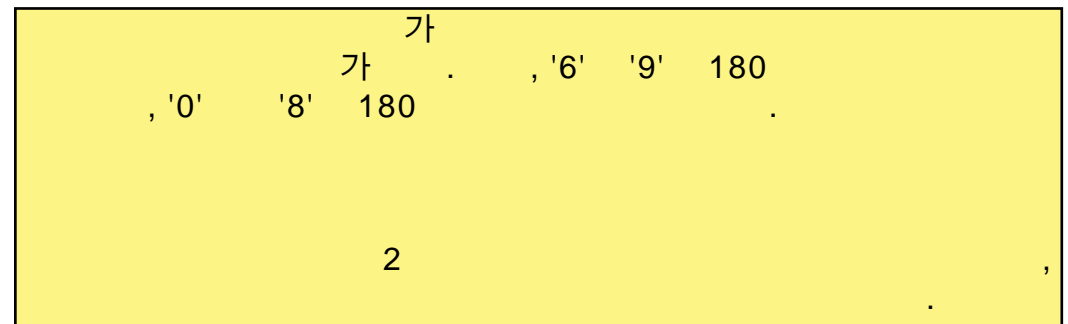


그림 5-24 필기 숫자 데이터의 다양한 변형⁸



- 한계
 - 수작업 변형
 - 모든 부류가 같은 변형 사용



5.4.3 데이터 확대

- 예) 모핑을 이용한 변형 [Hauberg2016]

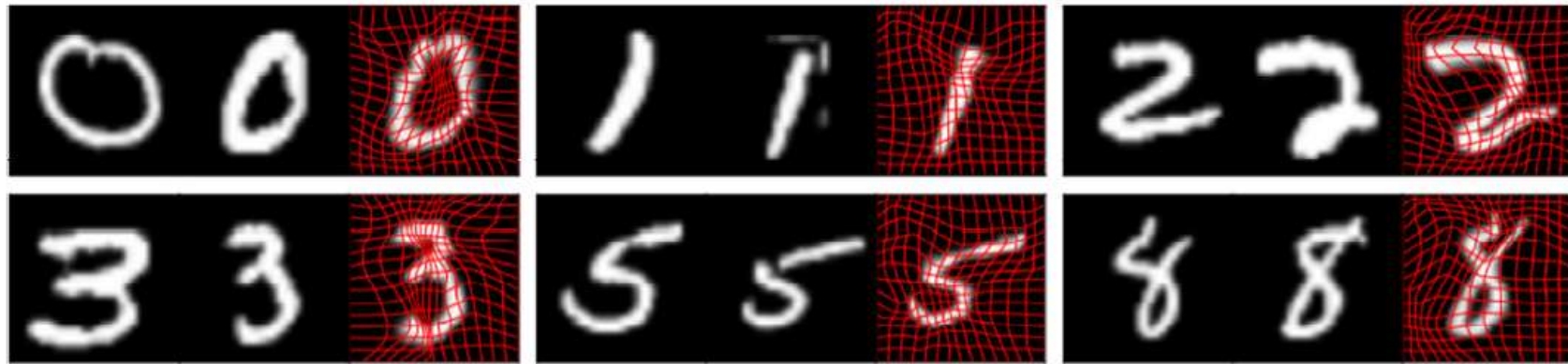


그림 5-25 비선형 변환 학습

- 비선형 변환으로서 어파인 변환에 비해 훨씬 다양한 형태의 확대
- 학습 기반: 데이터에 맞는 '비선형 변환 규칙을 학습'하는 셈

5.4.3 데이터 확대

■ 예) 자연영상 확대 [Krizhevsky2012]

- 256*256 영상에서 224*224 영상을 1024장 잘라내어 이동 효과. 좌우 반전까지 시도하여 2048배로 확대
- PCA를 이용한 색상 변환으로 추가 확대 (1%)
- 예측 단계에서는 [그림 5-26]과 같이 5장 잘라내고 좌우 반전하여 10장을 만든 다음 앙상블 적용

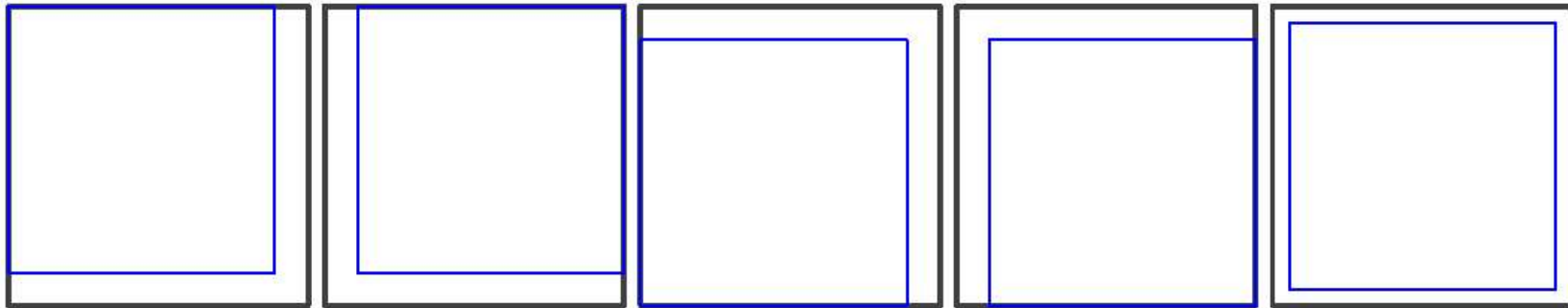


그림 5-26 예측 단계에서 영상 잘라내기

■ 예) 잡음을 섞어 확대하는 기법

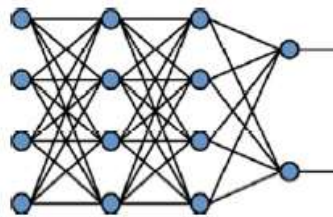
- 입력 데이터에 잡음을 섞는 기법 가
- 은닉 노드에 잡음을 섞는 기법 (고급 특징 수준에서 데이터를 확대하는 셈)

5.4.4 드롭아웃

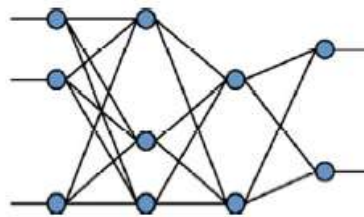
■ 드롭아웃 규제 기법

- 입력층과 은닉층의 노드 중 일정 비율을 임의로 선택하여 제거
- 남은 부분 신경망을 학습
- 많은 부분 신경망을 만들고, 예측 단계에서 앙상블 결합하는 기법으로 볼 수 있음

$$\frac{2^q}{n^k}$$



(a) 원래 신경망(4-4-4-2 구조)



(b) 드롭아웃된 3개의 신경망 예시

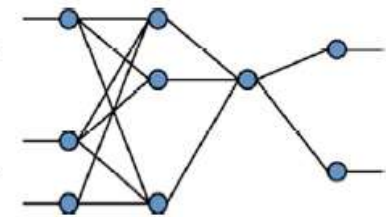
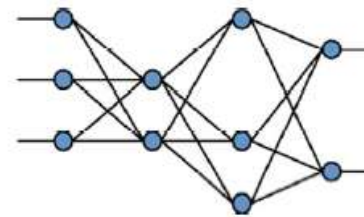


그림 5-27 드롭아웃된 신경망

- 많은 부분 신경망을 학습하고, 저장하고, 앙상블 결합하는 데 따른 계산 시간과 메모리 공간 측면의 부담

$$\left(\frac{2^q}{n^k} \right)$$

5.4.4 드롭아웃

■ 실제로는 가중치 공유 사용

- 하나의 신경망(하나의 가중치 집합)에 드롭아웃을 적용함([알고리즘 5-8])

알고리즘 5-8 드롭아웃을 채택한 기계 학습 알고리즘

입력: 드롭아웃 비율 p_{input} , p_{hidden}

출력: 최적해 $\hat{\Theta}$

```
1  난수를 생성하여 초기해  $\Theta$ 를 설정한다.
2  while (! 멈춤 조건) // 수렴 조건
3      미니배치  $\mathbb{B}$ 를 샘플링한다.
4      for ( $i=1$  to  $|\mathbb{B}|$ ) //  $\mathbb{B}$ 의 샘플 각각에 대해
5          입력층은  $p_{input}$ , 은닉층은  $p_{hidden}$  비율로 드롭아웃을 수행한다.
6          드롭아웃된 부분 신경망  $\Theta_i^{dropout}$ 로 전방 계산을 한다.
7          오류 역전파를 이용하여  $\Theta_i^{dropout}$ 를 위한 그레디언트  $\nabla_i^{dropout}$ 를 구한다.
8           $\nabla_1^{dropout}, \nabla_2^{dropout}, \dots, \nabla_{|\mathbb{B}|}^{dropout}$ 의 평균  $\nabla_{ave}^{dropout}$ 를 계산한다.
9           $\Theta = \Theta - \rho \nabla_{ave}^{dropout}$  // 가중치 갱신
10  $\hat{\Theta} = \Theta$ 
```

(gradient) .

5.4.4 드롭아웃

■ 라인 6의 전방 계산

$$\begin{array}{ll} l\text{번째 은닉층의 } j\text{번째 노드의 연산:} & \text{드롭아웃 적용:} \\ z_j^l = \tau_l(s_j^l) & z_j^l = \tau_l(s_j^l) \\ \text{이때 } s_j^l = \mathbf{u}_j^l \mathbf{z}^{l-1} & \Rightarrow \text{이때 } \begin{cases} \tilde{\mathbf{z}}^{l-1} = \mathbf{z}^{l-1} \odot \boldsymbol{\pi}^{l-1} \\ s_j^l = \mathbf{u}_j^l \tilde{\mathbf{z}}^{l-1} \end{cases} \end{array} \quad (5.35)$$

- 불린 배열 $\boldsymbol{\pi}$ 에 노드 제거 여부를 표시
- $\boldsymbol{\pi}$ 는 샘플마다 독립적으로 정하는데, 난수로 설정함
- 보통 입력층 제거 비율 $P_{input} = 0.2$, 은닉층 제거 비율 $P_{hidden} = 0.5$ 로 설정

5.4.4 드롭아웃

■ 예측 단계

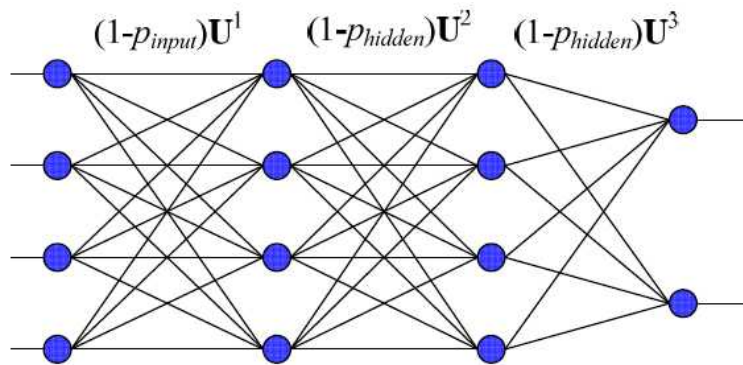


그림 5-28 드롭아웃의 예측 단계

- 앙상블 효과 모방
 - 가중치에 생존 비율 $(1-\text{드롭아웃 비율})$ 을 곱하여 전방 계산
 - 학습 과정에서 가중치가 $(1-\text{드롭아웃 비율})$ 만큼만 참여했기 때문

■ 메모리와 계산 효율

- 추가 메모리는 불린 배열 π , 추가 계산은 작음
- 실제 부담은 신경망의 크기에서 옴: 보통 은닉 노드 수를 $\frac{1}{p_{hidden}}$ 만큼 늘림

5.4.5 앙상블 기법

■ 앙상블

- 서로 다른 여러 개의 모델을 결합하여 일반화 오류를 줄이는 기법
- 현대 기계 학습은 앙상블도 규제로 여김

■ 두 가지 일

- 서로 다른 예측기를 학습하는 일
 - 예, 서로 다른 구조의 신경망 여러 개를 학습 또는 같은 구조를 사용하되 서로 다른 초깃값과 하이퍼 매개변수를 설정하고 학습
 - 예, 배깅(훈련집합을 여러 번 샘플링하여 서로 다른 훈련집합을 구성)
 - 예, 부스팅(i 번째 예측기가 틀린 샘플을 $i+1$ 번째 예측기가 잘 인식하도록 연계성을 고려)
- 학습된 예측기를 결합하는 일
 - 주로 투표 방식을 사용

■ 자세한 내용은 12장

5.5 하이퍼 매개변수 최적화

- 학습 모델에는 두 종류의 매개변수가 있음
 - 내부 매개변수
 - 신경망의 경우 에지 가중치로서 이 책은 θ 로 표기 (예, 식 (4.1)의 가중치 행렬 \mathbf{U}^l)
 - 학습 알고리즘이 최적화함
 - 하이퍼 매개변수
 - 모델의 외부에서 모델의 동작을 조정함
 - 예, 은닉층의 개수, CNN 마스크 크기와 보폭, 학습률, 모멘텀과 관련된 매개변수 등

5.5 하이퍼 매개변수 최적화

■ 하이퍼 매개변수 선택

- 표준 문헌이 제시하는 기본값을 사용하면 됨
 - 보통 여러 후보 값 또는 범위를 제시
- 후보 값 중에서 주어진 데이터에 최적인 값 선택 ← 하이퍼 매개변수 최적화

알고리즘 5-9 하이퍼 매개변수 최적화

입력: 훈련집합 \mathbb{X} , 검증집합 $\mathbb{X}_{\text{validate}}$, 하이퍼 매개변수집합 \mathcal{H}

출력: 최적 하이퍼 매개변수값 \hat{H}

```

1   $q_{\max} = 0$ 
2  repeat
3      하이퍼 매개변수 조합  $\tilde{H}$ 을 생성한다.
4       $\tilde{H}$ 으로 설정된 모델을  $\mathbb{X}$ 로 학습한다.
5      학습된 모델을  $\mathbb{X}_{\text{validate}}$ 로 성능  $q$ 를 측정한다.
6      if ( $q > q_{\max}$ )  $q_{\max} = q$ ,  $\hat{H} = \tilde{H}$ 
7  until(멈춤 조건)
```

- 라인 3을 구현하는 방법에 따라 수동 탐색, 격자 탐색, 임의 탐색

5.5 하이퍼 매개변수 최적화

■ 격자 탐색과 임의 탐색

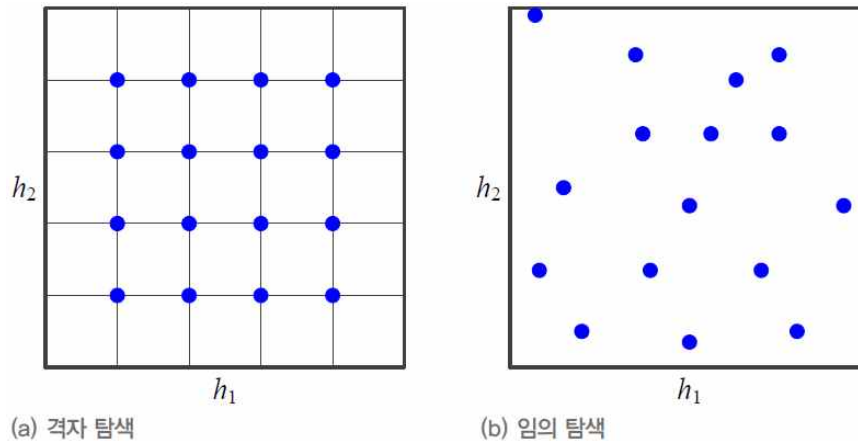


그림 5-29 하이퍼 매개변수 탐색 방법

- 임의 탐색은 난수로 매개변수 조합을 생성함

■ 로그 규모 간격

(ex. 2)

- 어떤 매개변수는 로그 규모를 사용해야 함
- 예, 학습률 범위가 [0.0001~0.1]일 때
 - 등간격은 0.0001, 0.0002, 0.0003, ..., 0.0998, 0.0999로 총 1000개 값을 조사
 - 로그 규모는 2배씩 증가시킴. 즉 0.0001, 0.0002, 0.0004, 0.0008, ..., 0.0256, 0.0512, ...를 조사함

, 가 가

5.5 하이퍼 매개변수 최적화

■ 차원의 저주 문제 발생

가 1~3

- 매개변수가 m 개이고 각각이 q 개 구간이라면 q^m 개의 점을 조사해야 함

■ 임의 탐색이 우월함

- [Bergstra2012]의 실험 (32개의 매개변수) → 임의 탐색이 유리함
- [그림 5-30]은 임의 탐색이 유리한 이유를 설명

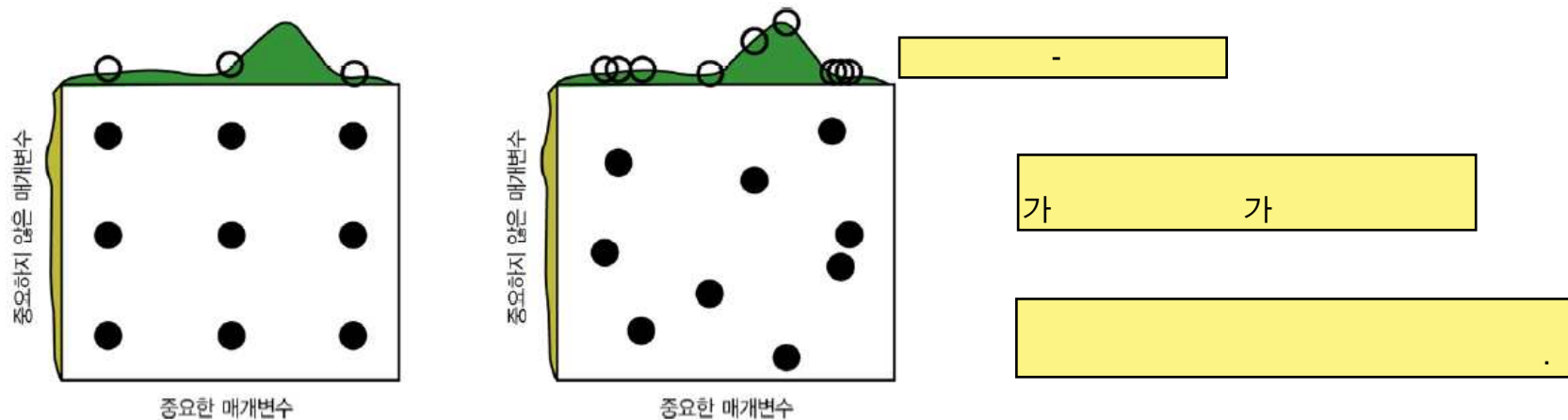


그림 5-30 격자 탐색과 임의 탐색의 성능 비교

5.6 2차 미분을 이용한 방법

- 5.6.1 뉴턴 방법
- 5.6.2 켈레 그래디언트 방법
- 5.6.3 유사 뉴턴 방법

```
-      ( - >      )  
-      ( - >      )
```

- 그래디언트(1차 미분)를 사용하는 경사 하강법
 - 현재 기계 학습의 주류 알고리즘
 - 두 가지 방향의 개선책 [Bottou2017]
 - 그래디언트의 잡음을 줄임 (예, 미니배치 사용 등)
 - 2차 미분 정보를 활용 → 5.6절의 주제

```
<      > (      )  
  
      (well - conditioned)  
      가      (ill -  
conditioned)      ...  
- 2      .  
-      ,      .
```

5.6 2차 미분을 이용한 최적화

■ 경사 하강법 다시 보기

예제 5-2 경사 하강법의 동작

변수가 2개인 아래 함수 J 의 최저점을 구하는 문제를 생각하자. 먼저 1차 도함수인 그래디언트 ∇J 를 구한다.

$$J(\mathbf{w}) = J(w_1, w_2) = (w_1 - 2)^2 + 4(w_2 - 1)^2$$
$$\nabla J(\mathbf{w}) = \left(\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2} \right)^T = (2(w_1 - 2), 8(w_2 - 1))^T$$

시작점이 $\mathbf{w}_1 = (4, 2)^T$ 라고 가정하면 \mathbf{w}_1 에서의 그래디언트는 $\nabla J(\mathbf{w}_1) = (4, 8)^T$ 이다. 식 (2.58)에 따라 새로운 점 \mathbf{w}_2 를 계산하면, $\mathbf{w}_2 = (3.2, 0.4)^T$ 가 된다. 이때 학습률 ρ 는 0.2라고 가정하자.

$$\mathbf{w}_2 = \mathbf{w}_1 - \rho \nabla J(\mathbf{w}_1) = \begin{pmatrix} 4 \\ 2 \end{pmatrix} - 0.2 \begin{pmatrix} 4 \\ 8 \end{pmatrix} = \begin{pmatrix} 3.2 \\ 0.4 \end{pmatrix}$$

\mathbf{w}_2 에서 그래디언트는 $\nabla J(\mathbf{w}_2) = (2.4, -4.8)^T$ 이고, 새로운 점은 $\mathbf{w}_3 = (3.2, 0.4)^T - 0.2(2.4, -4.8)^T = (2.72, 1.36)^T$ 가 된다. 비슷한 계산을 반복하면 다음과 같은 궤적을 그리며 최저점 $(2, 1)^T$ 로 접근하는 것을 알 수 있다.

$$\mathbf{w}_1 = \begin{pmatrix} 4 \\ 2 \end{pmatrix} \rightarrow \mathbf{w}_2 = \begin{pmatrix} 3.2 \\ 0.4 \end{pmatrix} \rightarrow \mathbf{w}_3 = \begin{pmatrix} 2.72 \\ 1.36 \end{pmatrix} \rightarrow \mathbf{w}_4 = \begin{pmatrix} 2.432 \\ 0.784 \end{pmatrix} \rightarrow \mathbf{w}_5 = \begin{pmatrix} 2.2592 \\ 1.1296 \end{pmatrix} \rightarrow \dots$$

5.6 2차 미분을 이용한 최적화

- 경사 하강법을 더 빠르게 할 수 있나
 - [그림 5-31]에서 파란 경로는 경사 하강법이 해를 찾아가는 과정
 - 1차 미분 정보로는 빨간 경로를 알 수 없음
 - 1차 미분은 현재 위치에서 지역적인 기울기 정보만 알려주기 때문

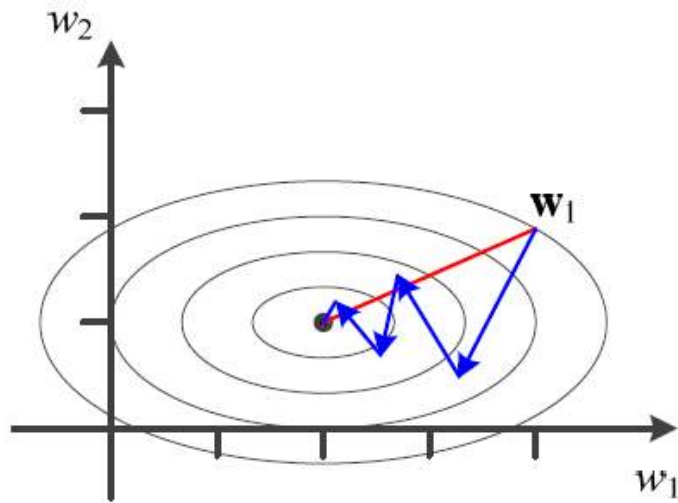


그림 5-31 1차 미분을 사용하는 경사 하강법을 더 빠르게 할 수 있는가

- 뉴턴 방법은 2차 미분 정보를 활용하여 빨간 경로를 알아냄

5.6.1 뉴턴 방법

- 테일러 급수를 적용하면,

$$J(w + \delta) \approx J(w) + J'(w)\delta + \frac{J''(w)}{2}\delta^2 \quad (5.36)$$

- 식 (5.37)은 변수가 여러 개일 때 (**H**는 헤시언 행렬)

$$J(\mathbf{w} + \boldsymbol{\delta}) \approx J(\mathbf{w}) + \nabla J^T \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T \mathbf{H} \boldsymbol{\delta} \quad (5.37)$$

- δ 로 미분하면,

w가 m
delta gradient - J' m
(gradient - J 1 , H 2 H m*m 가 .)

$$\frac{\partial J(w + \delta)}{\partial \delta} \approx J'(w) + \delta J''(w)$$

- [그림 5-32]처럼 $w + \delta$ 를 최소점이라 가정하면,

$$\frac{\partial J(w + \delta)}{\partial \delta} \approx J'(w) + \delta J''(w) = 0$$

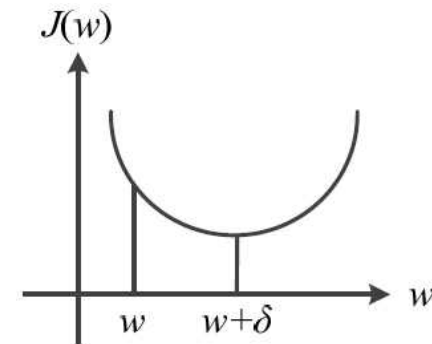


그림 5-32 변수가 하나인 상황의 뉴턴 방법

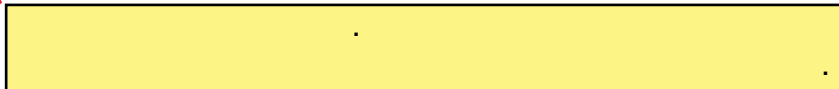
5.6.1 뉴턴 방법

- 식을 조금 정리하면,

$$\delta = -\frac{J'(w)}{J''(w)} = -(J''(w))^{-1}J'(w) \quad (5.38)$$

- 변수가 여러 개인 경우로 확장하면,

$$\delta = -\mathbf{H}^{-1}\nabla J \quad (5.39)$$



2 (3 0), (5.39) .
 , .
 MSE, , 2 ? 2 가 .
 , , .

5.6.1 뉴턴 방법

예제 5-3 뉴턴 방법

[예제 5-2]의 함수 $J(\mathbf{w}) = J(w_1, w_2) = (w_1 - 2)^2 + 4(w_2 - 1)^2$ 을 다시 활용한다. 그레이디언트와 헤시언 행렬을 구하면 아래와 같다.

$$\nabla J(\mathbf{w}) = \begin{pmatrix} 2w_1 - 4 \\ 8w_2 - 8 \end{pmatrix}$$

$$\mathbf{H} = \begin{pmatrix} 2 & 0 \\ 0 & 8 \end{pmatrix}$$

시작점 $(4, 2)^T$ 를 현재 점 \mathbf{w}_1 이라고 하면, \mathbf{w}_1 에서 그레이디언트는 $\nabla J(\mathbf{w}_1) = \begin{pmatrix} 4 \\ 8 \end{pmatrix}$ 이다. 이것을 식 (5.39)에 대입하면 아래와 같다.

$$\boldsymbol{\delta} = -\begin{pmatrix} 2 & 0 \\ 0 & 8 \end{pmatrix}^{-1} \begin{pmatrix} 4 \\ 8 \end{pmatrix} = -\begin{pmatrix} 1/2 & 0 \\ 0 & 1/8 \end{pmatrix} \begin{pmatrix} 4 \\ 8 \end{pmatrix} = -\begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

다음 점 \mathbf{w}_2 는 아래 식으로 계산할 수 있고, 이 점은 최저점임을 알 수 있다.

$$\mathbf{w}_2 = \mathbf{w}_1 + \boldsymbol{\delta} = \begin{pmatrix} 4 \\ 2 \end{pmatrix} - \begin{pmatrix} 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

5.6.1 뉴턴 방법

■ 뉴턴 방법의 적용

- [예제 5.3]은 2차 함수에 뉴턴 방법을 적용했으므로, 3차 항 이상을 무시한 식 (5.36)을 사용했음에도 최적 경우 제시
- 기계 학습이 사용하는 목적함수는 2차 함수보다 복잡한 함수이므로 한 번에 최적해에 도달 불가능 → [알고리즘 5-10]과 같이 반복하는 뉴턴 방법을 사용해야 함

알고리즘 5-10 뉴턴 방법

입력: 훈련집합 \mathbb{X} , \mathbb{Y}

출력: 최적해 $\hat{\theta}$

```
1  난수를 생성하여 초기해  $\theta$ 를 설정한다.
2  repeat
3     $\delta = -\mathbf{H}^{-1}\nabla J$ 를 계산한다. // 식 (5.39)
4     $\theta = \theta + \delta$ 
5  until (멈춤 조건)
6   $\hat{\theta} = \theta$ 
```

$$O(m^3) \quad (5.39)가$$

- >

- 라인 3에서 헤시언 \mathbf{H} 를 구해야 함 → 매개변수의 개수를 m 이라 할 때 $O(m^3)$ 이라는 과다한 계산량 → 켈레 그래디언트 방법이 대안 제시

5.6.2 켈레 그래디언트 방법

■ 직선 탐색 line search

- [그림 5-33(a)]는 경사 하강법이 학습률을 직선 탐색으로 찾는 상황을 예시

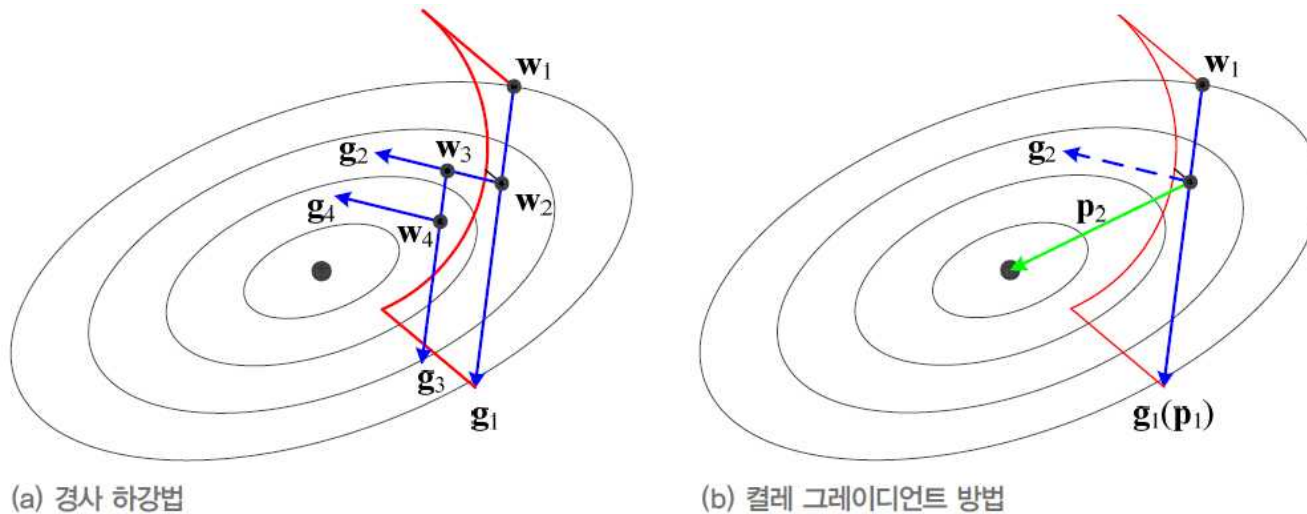


그림 5-33 경사 하강법과 켈레 그래디언트 방법의 비교

- 켈레 그래디언트 방법
 - [그림 5-33(a)]의 경사 하강법은 근본적으로 이전 경사 하강법과 같음
 - g_2 로 직선 탐색할 때 직전에 사용한 g_1 정보를 전혀 고려하지 않음
 - [그림 5-33(b)]의 켈레 그래디언트는 직전 정보를 사용하여 해에 빨리 접근

5.6.2 켈레 그레디언트 방법

■ 켈레 그레디언트

- 직전에 사용한 직선이 \mathbf{p}_{t-1} 이라면 다음 순간에는 식 (5.40)을 만족하는 \mathbf{p}_t 를 사용(\mathbf{p}_{t-1} 과 \mathbf{p}_t 를 켈레라 부름). 그런데 식 (5.40)은 여전히 헤시언 행렬을 포함

$$\mathbf{p}_t^T \mathbf{H} \mathbf{p}_{t-1} = 0 \quad (5.40)$$

- 켈레 그레디언트 방법에서는 식 (5.41)로 대체하여 계산
 - \mathbf{g}_{t-1} 과 \mathbf{g}_t 는 그레디언트로서, 그레디언트를 이용하여 2차 미분 정보를 근사하는 셈

$$\mathbf{p}_t = -\mathbf{g}_t + \beta_t \mathbf{p}_{t-1} \quad (5.41)$$

$$\text{Polak - Ribiere: } \beta_t = \frac{(\mathbf{g}_t - \mathbf{g}_{t-1})^T \mathbf{g}_t}{\mathbf{g}_{t-1}^T \mathbf{g}_{t-1}} \quad (5.42)$$

2 gradient gradient

$$\text{Fletcher - Reeves: } \beta_t = \frac{\mathbf{g}_t^T \mathbf{g}_t}{\mathbf{g}_{t-1}^T \mathbf{g}_{t-1}} \quad (5.43)$$

(5.42) (5.43) 2 가
 . 가 2 가 .
 - > 5 - 11 6

5.6.2 켈레 그레이디언트 방법

알고리즘 5-11 켈레 그레이디언트

입력: 훈련집합 \mathbb{X} , \mathbb{Y} , β 를 리셋하는 주기 q

출력: 최적의 매개변수 $\hat{\Theta}$

```
1   $\mathbf{p}_0 = \mathbf{0}$  // 라인 탐색에 사용할 벡터를  $\mathbf{0}$  벡터로 초기화
2   $t=1$ 
3  난수를 생성하여 초기해  $\Theta_1$ 을 설정한다.
4  repeat
5       $\Theta_t$ 에서 그레이디언트  $\mathbf{g}_t = \frac{\partial J}{\partial \Theta} \Big|_{\Theta_t}$ 를 구한다.
6      if  $(t \bmod q=1) \beta_t = 0$  //  $q$ 번에 한 번씩  $\beta$ 를 리셋
7      else  $\beta_t = \frac{(\mathbf{g}_t - \mathbf{g}_{t-1})^T \mathbf{g}_t}{\mathbf{g}_{t-1}^T \mathbf{g}_{t-1}}$  // 식 (5.42)
8       $\mathbf{p}_t = -\mathbf{g}_t + \beta_t \mathbf{p}_{t-1}$  // 식 (5.41)
9       $\Theta_t$ 와  $\Theta_t + \varepsilon \mathbf{p}_t$ 를 잇는 직선에서 직선 탐색을 수행하여 최소에 해당하는  $\hat{\varepsilon}$ 을 구한다.
10      $\Theta_{t+1} = \Theta_t + \hat{\varepsilon} \mathbf{p}_t$ 
11      $t++$ 
12 until (멈춤 조건)
13  $\hat{\Theta} = \Theta_t$ 
```

8

gradient

9

가

식 (5.42)와 식 (5.43)은 2차 함수를 가정하고 유도한 식임.
라인 6은 이 때문에 발생하는 오차를 보정하는 역할

5.6.3 유사 뉴턴 방법

■ 유사 뉴턴 방법의 기본 아이디어

- 헤시언 \mathbf{H} 의 역행렬을 근사하는 행렬 \mathbf{M} 을 사용 gradient .
- 처음에 단위 행렬 \mathbf{I} 로 시작하여 그레디언트 정보를 이용하여 점점 개선
- LFGS가 많이 사용됨
- 기계 학습에서는 \mathbf{M} 을 저장하는 메모리를 적게 쓰는 L-BFGS를 주로 사용함

■ 기계 학습에서 2차 미분 정보의 활용

- 현재 널리 활용되지는 않지만 연구 계속되고 있음

“... These methods, such as those built on noise reduction and second-order techniques, offer the ability to attain improved convergence rates, overcome the adverse effects of high nonlinearity and ill-conditioning, and exploit parallelism and distributed architectures in new way. ... 잡음 감소와 2차 미분과 같은 방법은 수렴 속도를 향상하고, 심한 비선형과 불량 조건 같은 부정적 효과를 극복하며, 새로운 방식으로 병렬분산 계산 구조를 이용하는 길을 제시한다.” [Bottou2017]