

Ch6. 학습 관련 기술들

- Outline : 최적화 방법 (가중치 매개변수 탐색)
 - W 초기값, 하이퍼파라미터 설정방법
 - 오버피팅 대응책 - 가중치 감소, 드롭아웃 등 정규화 방법
 - 배치 정규화

⇒ 신경망(딥러닝) 학습의 효율과 정확도 ↑

1. 매개변수 갱신

- "최적화 (optimization)" : 손실함수 값을 최대한 낮추는 매개변수를 찾는 문제.

1) 확률적 경사 하강법 (SGD)

$$W \leftarrow W - \underbrace{\eta}_{\text{학습률 (learning rate)}} \underbrace{\frac{\partial L}{\partial W}}_{\text{W에 대한 손실함수(L)의 기울기}}$$

: 기울어진 방향으로 일정거리만 가는 방법.

코드 구현) `network = TwoLayerNet(...)`

* `optimizer = SGD()`

for `i` in range(10000):

`grads = network.gradient(x_batch, t_batch)`

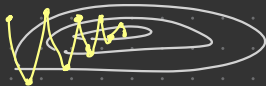
`params = network.params`

* `optimizer.update(params, grads)`

* 매개변수 갱신 (최적화)을 담당하는 클래스를 분리 ⇒ 모듈화

장점: 단순, 구현 easy.

단점



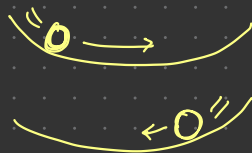
At 비등방성 함수(방향에 따라 성질, 기울기가 달라지는 함수),

탐색경로 비효율적!

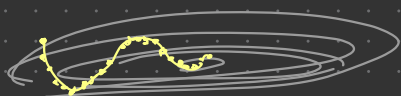
2) Momentum (Momentum)

$$\begin{cases} V \text{ (속도)} \leftarrow \alpha V - \eta \frac{\partial L}{\partial W} \\ W \leftarrow W + V \end{cases}$$

like 지면마찰, 공기저항



⇒ 기울기 방향으로 물체가 힘을 받아 가속된다는 "물리법칙" 적용.



SGD보다 α 쪽 방향으로 더 빠르게 다가감.

3) AdaGrad

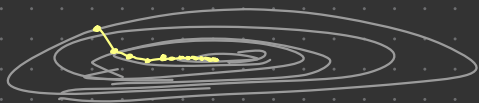
- 학습률(η)을 정하는 효과적인 기술 ⇒ "학습률 감소"
- AdaGrad : 개별 매개변수에 적응적으로 학습률(η)을 조정하면서 학습

$$\begin{cases} h \leftarrow h + \frac{\partial L}{\partial W} \\ W \leftarrow W - \eta \frac{1}{\sqrt{h}} \frac{\partial L}{\partial W} \end{cases}$$

$\circ \rightarrow$ 행렬의 원소별 곱셈
 \Rightarrow 학습률 감소가 매개변수 원소마다 다르게 적용됨.
 학습률 조정.

- 학습 진행할수록 갱신 강도가 약해지는 문제점

⇒ 개선책 : "RMSProp" (이러게 기울기 세히 있고, 새로운 기울기 정보 크게 반영하는 지수이동평균 (EMA) 방법 사용)



α 쪽 방향으로 갱신속도가 빠르게 약해짐.
(지그재그↓)

4) Adam \leftarrow 모멘텀 + AdaGrad
 구르는 움직임 적응적 갱신 정도 조정

결론: SGD, 모멘텀, AdaGrad, Adam 중에서 무엇을 사용할 지는

풀어야 할 문제와 하이퍼 파라미터 설정에 따라 적절히 선택한다!

2. 가중치의 초기값

- 가중치 감소 (Weight decay) 기법 - 가중치 매개변수 값이 작아지도록 학습.
 - 오버피팅을 억제해 범용성을 높이는 효과.

↳ 가중치 초기값을 이에 0으로 설정한다면?

⇒ bad! 가중치를 균일한 값으로 설정하면

오차역전파법에서 모든 가중치의 값이 똑같이 갱신되기 때문에.

결론) 가중치 초기값을 무작위로 설정해야 한다.

→ 적당히 과부하되야 신경망 학습이 효율적으로 이루어진다.

- 은닉층의 활성화값 분포 (코트-히스토그램 참고)

↳ 활성화함수 (시그모이드, ReLU 등)의 출력 레이어

- 1) 출력이 0에 가까워져야 미분 0이 다가감 ⇒ gradient vanishing 문제.
 (데이터가 0과 1에 치우쳐서 분포하는 경우)

- 2) 활성화 값이 0.5 부근에 치우쳐져 있음 ⇒ 포화력 제한
 (다수의 뉴런을 켜는 의미가 없음)

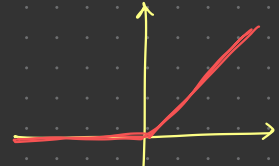
• Xavier 초기값

- 사비에르 글로르와 요슈아 벤지오의 논문에서 권장하는 가중치 초기값
Xavier Glorot, Yoshua Bengio
- 현재 일반적인 딥러닝 프레임워크들이 표준적으로 사용하고 있다.
- 초기값의 표준편차가 $\frac{1}{\sqrt{n}}$ 이 되도록 설정. (n 은 앞층의 노드 수)
⇒ 앞 층의 노드가 많을수록 대상 노드의 초기값으로 설정하는 가중치가 좁게 퍼진다.
- 중앙분포를 성형한 함수 (Sigmoid, tanh 함수) 에 적합.

• He 초기값

- 카이밍 헤의 이름을 딴 초기값
Kaiming He
- ReLU에 특화된 초기값.
- 초기값의 표준편차가 $\sqrt{\frac{2}{n}}$ 인 정규분포 사용

* ReLU



↑
음의 영역이 0이므로
더 쉽게 분포시키기 위해
그해의 계수가 필요한 것.

✓ 가중치 초기값은 신경망 학습에 아주 중요한 포인트!

3. 배치 정규화

Batch Normalization

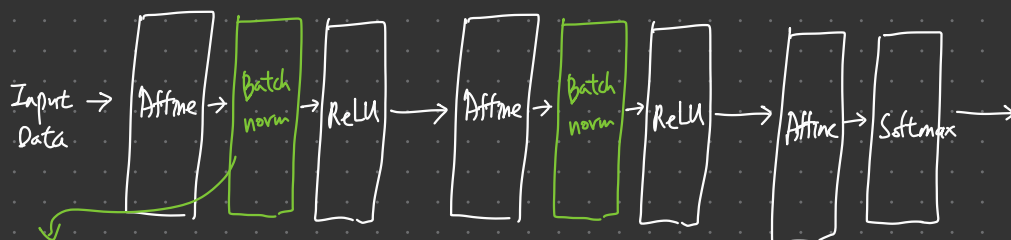
- 가중치 초기값 맞고, 각 층이 활성화층을 적당히 퍼뜨리도록 '강제'한다.

✓ 학습 속도 개선

✓ 초기값에 의존하지 않음. (초기값 선택 장애 안됨!)

✓ 오버피팅 억제 (드롭아웃 등 필요성 감소)

- 배치 정규화 계층 삽입 (데이터 분포를 정규화)



< 배치 정규화 알고리즘 >

1) 미니배치 $B = \{x_1, x_2, \dots, x_m\}$ 에 대해 평균 μ_B 과 분산 σ_B^2 구한다.
m개의 입력레이어 집합

$$\begin{cases} \mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \\ \sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \\ \hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \end{cases} \quad \begin{array}{l} 2) \text{ 그리고 평균} = 0, \text{ 분산} = 1 \text{ 이 되도록} \\ \text{정규화한다.} \end{array}$$

3) 정규화된 데이터: $\{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_m\}$

ϵ : epsilon (epsilon) : 아주 작은 값

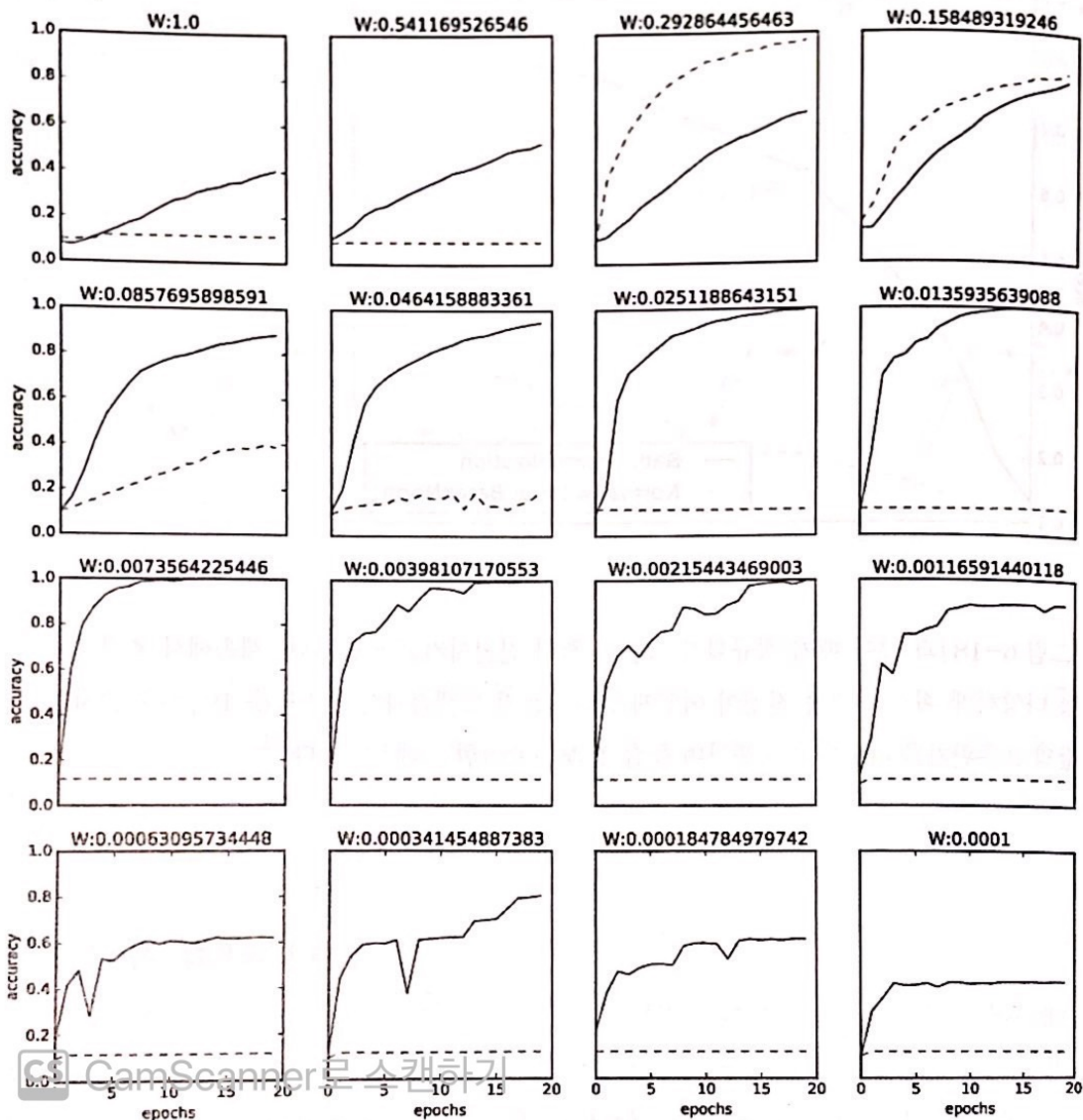
4) 정규화된 데이터에 고유한 확대 (Scale) 와 이동 (Shift) 변환

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

γ : 확대 β : 이동

$\gamma=1, \beta=0$ 부터 시작,
학습하면서 적합한 값으로 조정.

그림 6-19 실선이 배치 정규화를 사용한 경우, 점선이 사용하지 않은 경우 : 가중치 초기값의 표준편차는 각 그래프 위에 표기



학습 폭도가 훨씬 빨라지는군!

4. 바른 학습을 위해

4-1) 오버피팅

신경망이 훈련 데이터에만 지나치게 적응되어 그 외의 데이터에는 제대로 대응하지 못하는 상태

주로 \checkmark 매개변수가 많고 표현력이 높은 모델

\checkmark 훈련데이터가 적음

4-2) 가중치 감소

오버피팅은 가중치 매개변수의 값이 커서 발생하는 경우가 많다.

학습 과정에서 큰 가중치에 대해서는 그에 상응하는 큰 패널티를 부과하여 오버피팅을 억제하는 방법.

$$L2 \text{ norm} = \sum_{i=1}^n \sqrt{w_i^2}$$

$$L1 \text{ norm} = \sum_{i=1}^n |w_i|$$

$$L_{\infty} \text{ norm} = \text{Max}(|w_1|, |w_2|, \dots, |w_n|)$$

가중치 $W \rightarrow$ 가중치 감소 $\frac{1}{2} \lambda \|W\|^2$ 를 손실함수에 더한다.

정규화 세기 조절하는 하이퍼파라미터

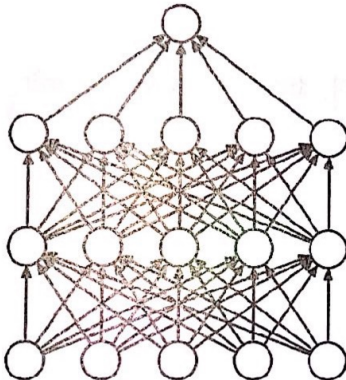
λ : 크게 설정할수록 큰 가중치에 대한 패널티가 커진다.

4-3) 드롭아웃 (Drop out)

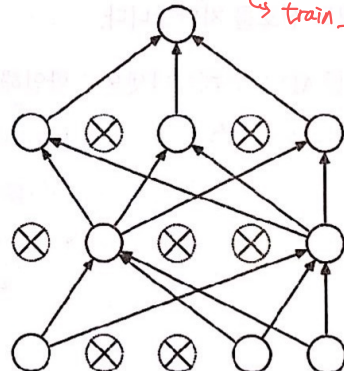
뉴런을 임의로 삭제하면서 학습하는 기법

훈련 때는 삭제한 뉴런 무작위 선택

시험 때는 모든 뉴런에 신호 전달
(훈련 때 삭제한 비율 곱하여 출력)



(a) 일반 신경망



(b) 드롭아웃을 적용한 신경망

$\hookrightarrow \text{train_flag} = \text{True}$ 일 때