

Regularization

overfitting = train data의 specific pattern 학습 (\approx noise)

이제 model의 capacity라 하는 밀접한 관계 (영향)을 갖게 된다.

"언제까지 general trend를 학습하는지 알수 없다"

이제 Regularization을 추가해준다. (성능개선 \downarrow 이던 w를 줄여가)

mode capacity control

$\min (Y - X\theta)^T (Y - X\theta) + \lambda \|\theta\|_2^2$ 으로 penalty term을 추가해준다.

L2 \Rightarrow w가 0이 되진 않는다.

λ 가 \uparrow 할수록 w의 값들이 더 작아져서 원의 점선으로 mapping된다.

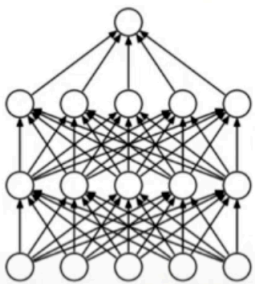
deep learning은 overfitting에 많이 취약하다.

(1) weight decay

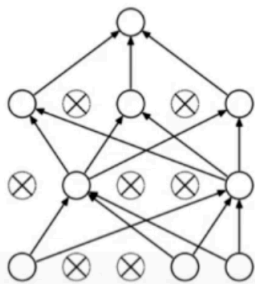
$$\Omega(w) = \sum \sum w_{ij}^2$$

or $\Omega(w) = \sum \sum |w_{ij}|$ 을 더해준다.

(2) Drop out



(a) Standard Neural Net



(b) After applying dropout.

모든 node 사용 X

특정 node 들은 확률적으로 사용한다.

\Rightarrow 이미지에 임의의 masking 처리한다. \Rightarrow 일반화 성능 \uparrow
Because 다양한 특징을 글로브 학습하고자 함.

node 1 2 ... n 각 node에 맞는 마스크
0.4 0.2 ... 0.8 샘플해서

가중치가 0.9 이상 node는 1 else 0 으로 만들어준다.

그때서 gradient가 1/2 개만 사용했기때문에 scale이 1/2 정도 낮게 학습한다.

그때서 test 시 가중치 scale을 2배 증가시켜준다.

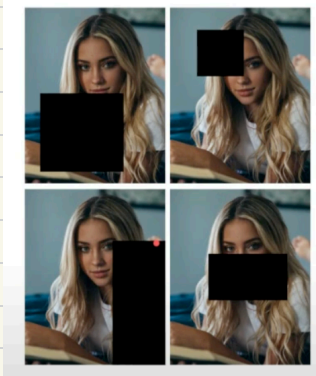
<test 시에 학습률 p를 부르기 전에>

test scale 무시 \Rightarrow 처음부터 train 시에 p를 주어서 gradient를 계산해준다.

(3) (v+0v)

Input data 자체에서 Masking 처리를 진행한다.

⇒ Input시 적용해야해서 train 시간 ↑의 단점.



(4) Early Stopping

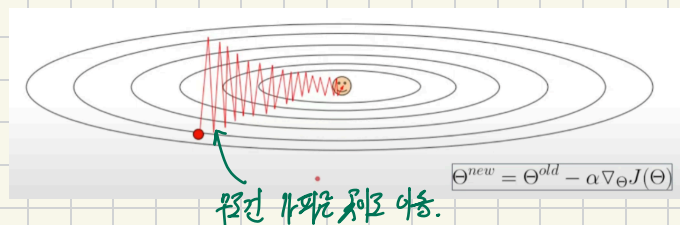
optimization

SGD 문제점 : (1) 느리다.

(2) saddle point 처리에 취약.

(3) mini-batch Big data set의 단점.

↳ Batch가 ↑ 할수록 무조건 성능이 올라간다.



(2) Momentum (관성)

$$SGD: x_{t+1} = x_t - \alpha \nabla f(x_t)$$

$$\text{Momentum: } v_{t+1} = \underbrace{\rho v_t}_{\text{이전에 사용된 gradient.}} - \alpha \nabla f(x)$$

$$x_{t+1} = x_t + v_{t+1}$$

(은 이전 관성을 평가하여 사용함.)

(3) Adagrad.

⇒ 이전까지 들어온 방향의 gradient의 sum을 나눠준다.

⇒ 이미 매우 작게 내려온 방향에 대해 보정을 잡아준다.

↳ 평국 다 0이되어 이동 거의 X

개변

(4) RMS props

⇒ 이전 모든 gradient가 아니라 window처럼 특정 구간에 대해서만 gradient를 sum.

(5) Adam

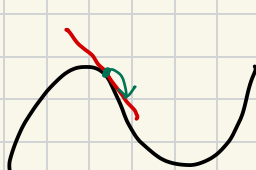
```

first_moment = 0.0
second_moment = 0.0
while True:
    dx = compute_gradient(x)
    (first_moment = beta1 * first_moment + (1 - beta1) * dx)
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    divider = np.sqrt(second_moment) + 1e-7
    x -= learning_rate * first_moment / divider

```

momentum + RMSprop이다.

- (1) first order optimization
- (2) Second order optimization



⇒ 기울기를 구하는 횟수↓ But 1개의 Hessian gradient를 구하는 시간이 매우 오래걸린다.

Batch Normalization

W 초기화 ⇒ zero mean unit variance 설정!
 Layer가 ↑ 많수록 설정 초기화 setting이 분리된다.

↳ 이를 해결 "Batch norm"

mini-batch 마다

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^k]}}$$

강제로 가운데를 변형
 ↳ 그냥 강제 처리해도 되는 No.

가중치 처리 후 다시 부피시켜준다. ⇒ 부피시키는 learnable parameter를 data로 부터 학습받게자 하는 Idea

$$W \Rightarrow \underbrace{\frac{y - \mu}{\sigma}}_{\text{Batch norm}} \Rightarrow \underbrace{\gamma \hat{y} + \rho}_{\text{부피}} \Rightarrow \text{Loss}$$

activation 에 들어가기 전에 넣게 된다.

But test 시에는 input이 1개이면 σ 가 1개에 대해서만 적용 0이 된다.

그래서 train시엔 사용안 쓰라 σ 들은 그대로 사용!

"Zero mean unit variance를 모든 layer에서 유지 가능하도록 하는 장치"