

[7-1. Hyperparameter Tuning]

< Tuning Process >

* 하이퍼 파라미터 우선 조정 순위

1순위) 학습률 (α)

2순위) 모멘텀 (momentum)의 $\beta \sim 0.9$

mini-batch size

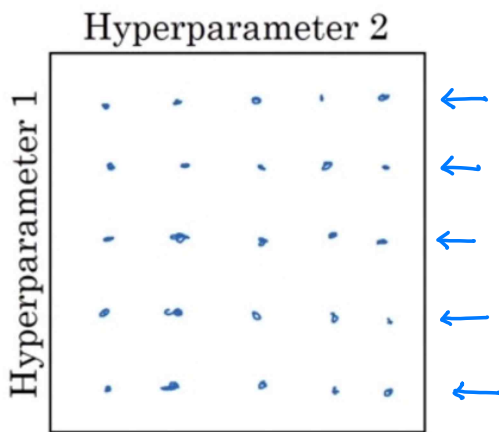
hidden units

3순위) # layers

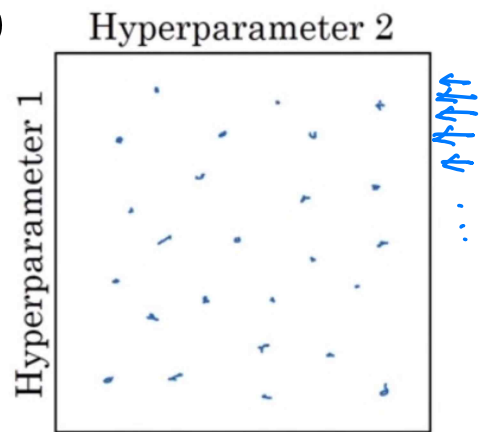
learning rate decay

Adam 알고리즘의 $\beta_1, \beta_2, \epsilon$
 $0.9, 0.999, 10^{-8}$

방법 1)



방법 2)



- 방법 1의 경우, 하이퍼파라미터의 개수가 비교적 많지 않을 때 잘 작동하는 편

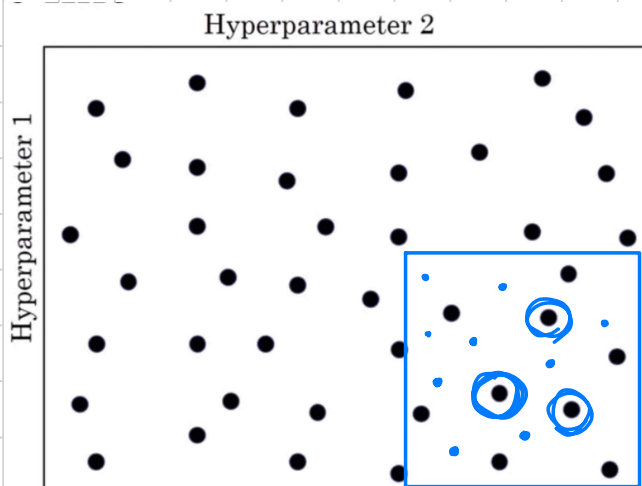
- 방법 2의 무작위 접근 방식이 더 좋다. 어떤 하이퍼파라미터가 문제 해결에 더 좋은지 알 수 있기 때문.

ex) hyperparameter 1: α , hyperparameter 2: ϵ 이라면,

즉, hyperparameter의 하나는 수치 조정에 중요한 값이고 하나는 애자면,

방법 1에서는 중요한 α 를 한번 밖에 세를 못하지만 방법 2는 많이 해서 더 나은 것 찾을

방법 3)



- 방법 3은 전체 하이퍼파라미터 공간에서 탐색하여 좋은 점을 찾은 후, 그 관점에서 더 정밀하게 탐색

< Using an Appropriate Scale to pick Hyperparameters >

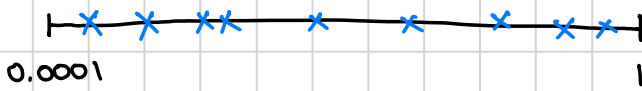
ex1) # hidden units : 50, ..., 100



layers : 2 ~ 4

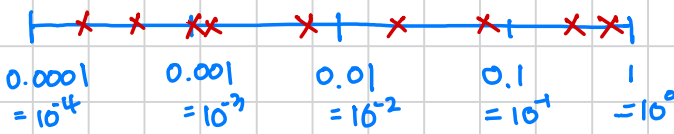
∴ #hidden units, # layers 는 신경망 크기를 이용하여 무작위로 뽑는 것이 합리적!

ex2) $\alpha = 0.0001, \dots, 1$



→ 90%의 값이 0.1 ~ 1 사이의 값

로그 척도



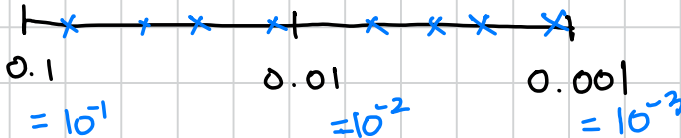
→ $r = -4 * \text{np.random.rand}()$

$\alpha = 10^r$

∴ 학습률(α)은 로그 척도를 이용하여 무작위로 뽑는 것이 합리적!

ex3) $\beta = 0.9, \dots, 0.999$

→ $1 - \beta = 0.1, \dots, 0.001$



→ $r \in [-3, -1]$

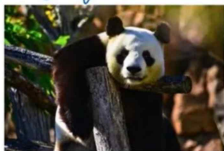
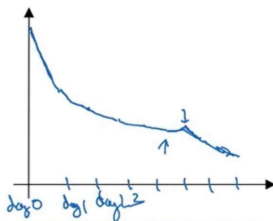
$1 - \beta = 10^r$

$\beta = 1 - 10^r$

∴ β 는 $1 - \beta$ 로 변환한 후, 로그 척도를 이용하여 무작위로 뽑는 것이 합리적!

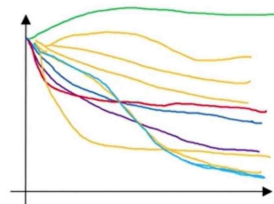
< Hyperparameter Tuning in Practice : Pandas vs Caviar >

Babysitting one model



Panda

Training many models in parallel



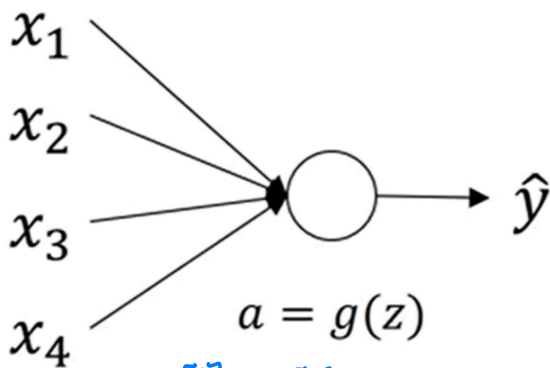
Caviar

Andrew Ng

- Babysitting one model (= 판다 접근) : 하나의 모델의 성능을 학습 도중을 지켜가면서 확인. CPU/GPU가 부족할 때
- Training many models in parallel : 여러 모델을 한번에 학습시킬 때 사용. CPU/GPU가 여유로울 때

[7-2. Batch Normalization]

<Normalizing Activation in a Network>



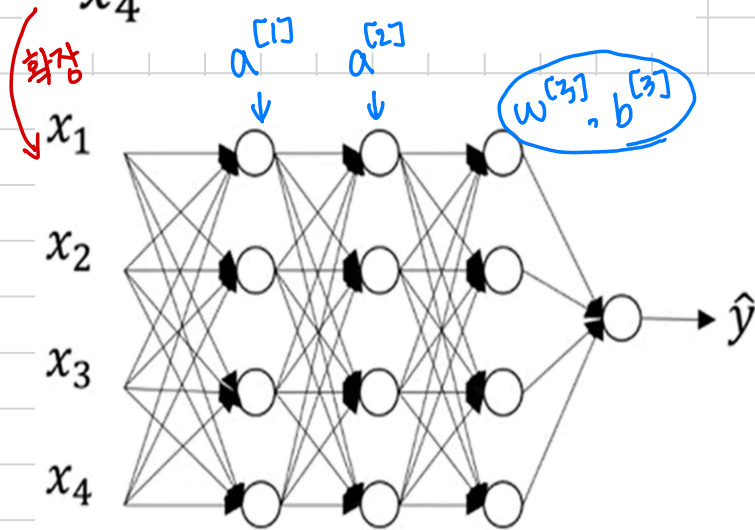
$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$X = X - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i x^{(i)2}$$

$$X = X / \sigma$$

element-wide



Can we normalize $a^{[2]}$ so as to train $w^{[3]}, b^{[3]}$ faster.

기존적으로는 $a^{[2]}$ 가 아니라 $z^{[2]}$ 를 normalize 함

* Implementing Batch Norm

- Given some intermediate values in NN

$$z^{(1)}, \dots, z^{(m)} \rightarrow z^{[2]}^{(i)}$$

* Batch Normalization

- 배치 정규화의 장점: 하이퍼파라미터의 탐색을 쉽게 만들어주며, 신경망과 하이퍼파라미터의 상관관계를 줄여줌

- 보통 활성화 함수 이전에 사용

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

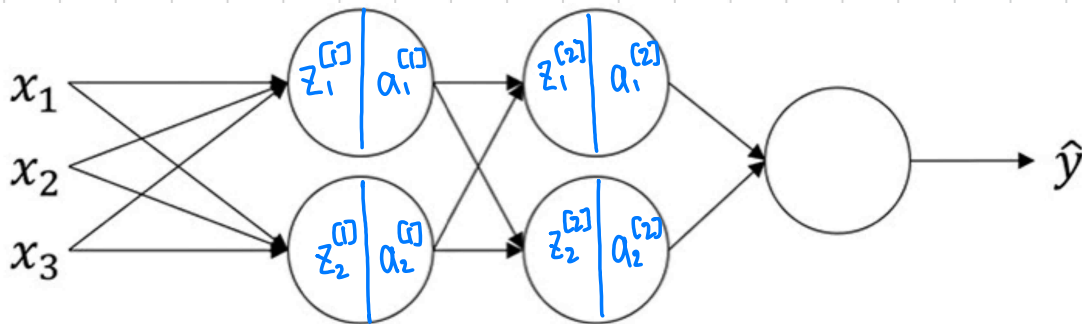
$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

learnable parameter in model

If $\gamma = \sqrt{\sigma^2 + \epsilon}$, $\beta = \mu$ 라면 $\tilde{z}^{(i)} = z^{(i)}$ 이므로 적절한 γ 와 β 를 선택하면 위의 나개의 식은 합동함수!
즉, γ 와 β 를 다르게 설정하면 hidden unit이 다른 평균과 분산을 갖도록 해야 함

\Rightarrow Use $\tilde{z}^{[2]}(i)$ instead of $z^{[2]}(i)$

<Fitting Batch Norm In to a Neural Network>



$$\begin{aligned}
 X &\xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow[\text{batch norm (=BN)}]{\beta^{(1)}, \gamma^{(1)}} \tilde{z}^{(1)} \xrightarrow[\text{activation function}]{\text{activation function}} a^{(1)} = g^{(1)}(\tilde{z}^{(1)}) \\
 &\xrightarrow{w^{(2)}, b^{(2)}} z^{(2)} \xrightarrow[\text{BN}]{\beta^{(2)}, \gamma^{(2)}} \tilde{z}^{(2)} \xrightarrow{\text{activation function}} a^{(2)} \rightarrow \dots
 \end{aligned}$$

: 선형결합 값을 계산하고, 이것을 BN시킨 후, 활성화함수를 거쳐 다음 계층

*Working with mini-batches

$$\begin{aligned}
 &\text{X}^{(1)} \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow[\text{batch norm (=BN)}]{\beta^{(1)}, \gamma^{(1)}} \tilde{z}^{(1)} \xrightarrow[\text{activation function}]{\text{activation function}} a^{(1)} = g^{(1)}(\tilde{z}^{(1)}) \rightarrow \dots \\
 &\text{X}^{(2)} \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow[\text{batch norm (=BN)}]{\beta^{(1)}, \gamma^{(1)}} \tilde{z}^{(1)} \xrightarrow[\text{activation function}]{\text{activation function}} a^{(1)} = g^{(1)}(\tilde{z}^{(1)}) \rightarrow \dots \\
 &\text{X}^{(3)} \rightarrow \dots
 \end{aligned}$$

⊕ Parameters : $w^{(l)}, b^{(l)}, \beta^{(l)}, \gamma^{(l)}$ 중에서 BN을 적용하면 z 의 평균을 바꾸기 때문에 상수항 $b^{(l)}$ 은 사용하지 않음

$$\therefore z^{(l)} = w^{(l)} a^{(l-1)}, \quad \tilde{z}^{(l)} = \gamma^{(l)} \cdot z_{\text{norm}}^{(l)} + \beta^{(l)}$$

*Implementing gradient descent

for $t=1$ to num Mini-Batches
compute forward prop on $x^{(t)}$

In each hidden layer, Use BN to replace $z^{(l)}$ with $\tilde{z}^{(l)}$

Use back prop to compute $dw^{(l)}, db^{(l)}, dr^{(l)}, d\beta^{(l)}$

Update parameters $w^{(l+1)} = w^{(l)} - \alpha dw^{(l)}$...

Work with momentum / RMSprop / Adam

<Why does Batch Norm work?>

- 입력특성 X 를 정규화하여 z 들을 비슷한 범위의 값들 (평균 0, 분산 1)로 만들어서 학습률 ↑
- 이전 층의 가중치 영향을 덜 받게 한다. 은닉층의 값의 분포의 변화를 줄여주어서, 입력값의 분포를 제한함. 즉, 입력값이 바뀌어서 발생하는 문제인 'covariance shift' 문제를 안정화 시킴.

∴ 앞층과 뒷층의 매개변수의 상관관계를 줄여서 $X \rightarrow Y$ 로 가는 mapping을 트rain한 경우,

학습속도를 향상시킴 (층이 깊을수록)

X 의 분포가 변했을 때 알고리즘을 변경하지 않고 ground true function에 적응하도록 하는 것

- 파라미터 정규화 (regularization)

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch. (∴ 평균과 분산이 고정됨)
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations. (∴ dropout은 은닉층의 확률에 따라 0또는 1을 곱하기때문에 공백 noise 존재)
- This has a slight regularization effect.

↳ batch normalization도 공백잡음 (x_{noise})과 드롭아웃 (-1)이 존재하여 regularization 효과 0

↳ 은닉층에 noise를 추가하는 것은

하나의 hidden unit에 너무 의존하지 않도록 하는 것

but) regularization의 효과가 있다는 것이지, 효과가 크지는 않으며 mini batch의 크기가 작으면 효과가 더 적어진다.

<Batch Norm at Test Time>

- BN은 한번에 하나의 mini-batch로 처리하지만, test 시에는 배치가 하나이기 때문에 평균과 분산을 계산할 수 없다.

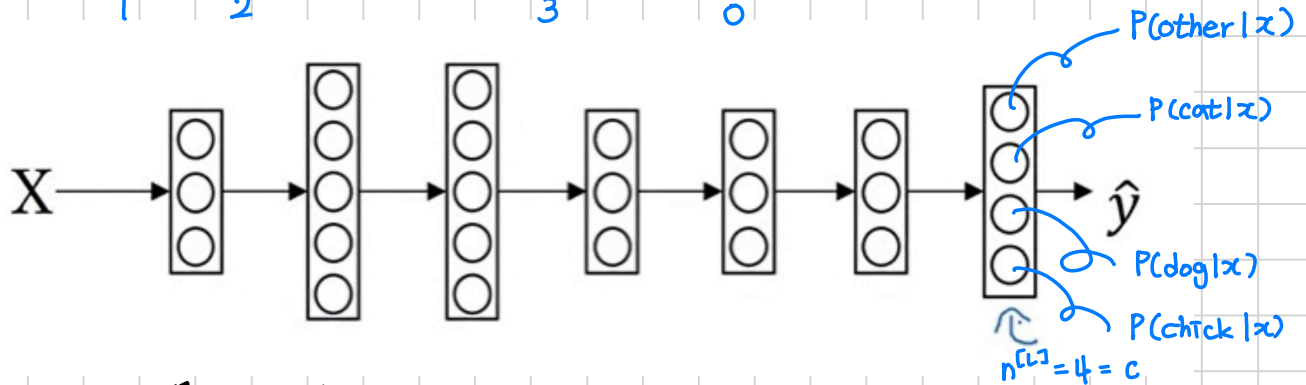
한번에 하나의 예제를 처리해야 할 수도 있음

∴ exponentially weighted average를 이용하여 각 μ 와 σ^2 를 산출 (across mini-batch)

[1-3 Multi-class Classification]

< Softmax Regression >

* Recognizing cats, dogs, and baby chicks, other



- softmax: 여러개의 클래스 분류시 사용하며 generalize logistic regression to C classes.

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \quad \text{activation function} \rightarrow$$

$$t = e^{z^{[L]}}$$

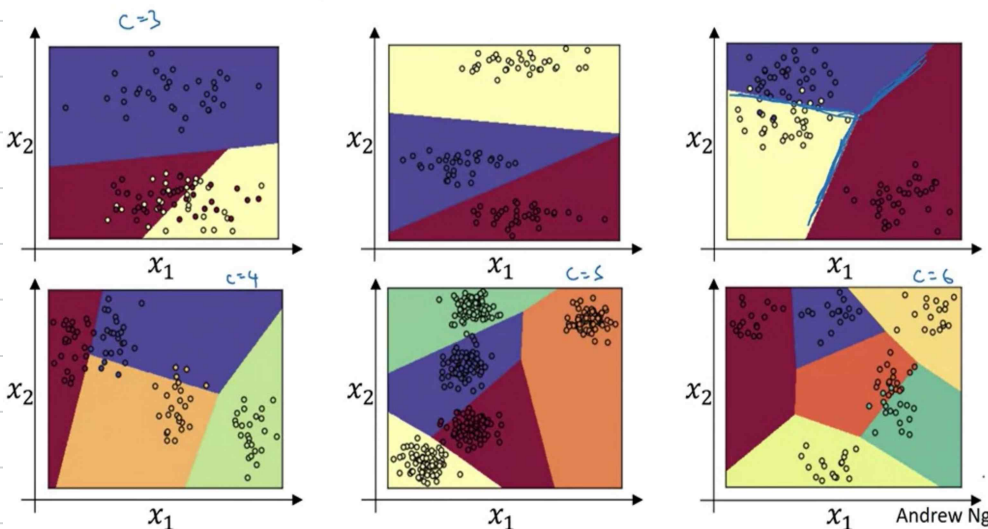
$$a_i = \frac{t_i}{\sum_{j=1}^C t_j}$$

Softmax examples

$$\begin{matrix} x_1 \\ x_2 \end{matrix} \rightarrow \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix} \rightarrow \hat{y}$$

$$z^{[L]} = (w^{[L]} \times t^{[L-1]})$$

$$a^{[L]} = \hat{y} = g(z^{[L]})$$



< Training a softmax Classifier >

ex)

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \quad \text{sum} = 1$$

* Loss function

$$- L(\hat{y}, y) = - \sum_{i=1}^m y_i \log \hat{y}_i$$

ex) $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $a^{(w)} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$ 의 loss function을 구해보면, $L(\hat{y}, y) = - \sum_{i=1}^m y_i \log \hat{y}_i = - \log \hat{y}_2$
즉, loss function을 최소화하는 것은 \hat{y}_2 를 최대화 하는 것!

$$- \text{Back prop : } dz^{[L]} = \hat{y} - y$$