

WEEK 3. Hyperparameter Tuning, Batch Normalization and Programming Frameworks

1. Hyperparameter tuning

• Hyperparameters (Andrew Ng's tips)

① $\alpha \rightarrow$ most important

② $\beta \sim 0.9$

③ mini-batch size

④ # of hidden units

⑤ # of layers & learning rate decay

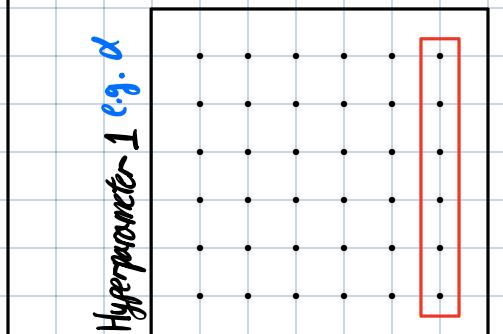
⑥ $\beta_1, \beta_2, \epsilon$ for Adam optimization \Rightarrow commonly $0.9, 0.999, 10^{-8}$

\Rightarrow How to explore?

- Try random values : Don't use a grid

- if # of hyperparameters \downarrow

\Rightarrow Search points in grid



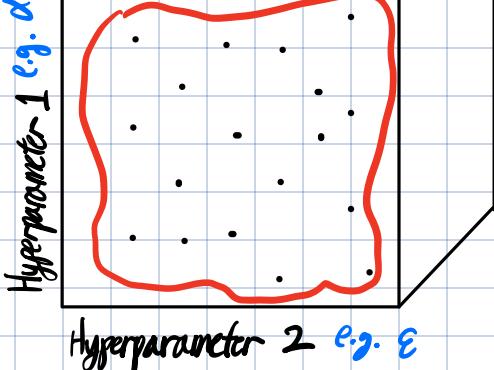
\Rightarrow 만약, 6 개의 서로 다른 α 를 선택할 때
그 모든 결과로 결과가 좋다면,
6개의 α 를 선택해서만 가능하지만 다른 것은

- if # of hyperparameters \uparrow (i.e. in DL)

\Rightarrow choose points randomly

\rightarrow sample of size hyperparameter

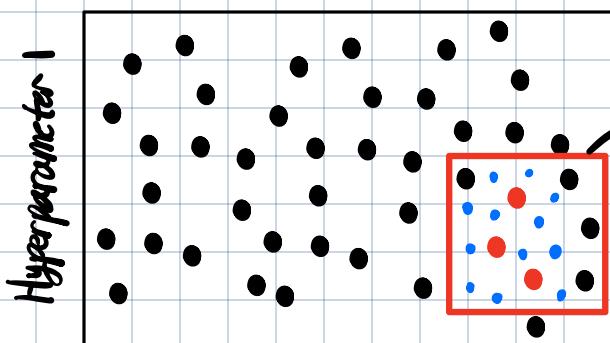




⇒ 25개의 서로 다른 값에 대해 학습 가능

- Course to find

Hyperparameter 2

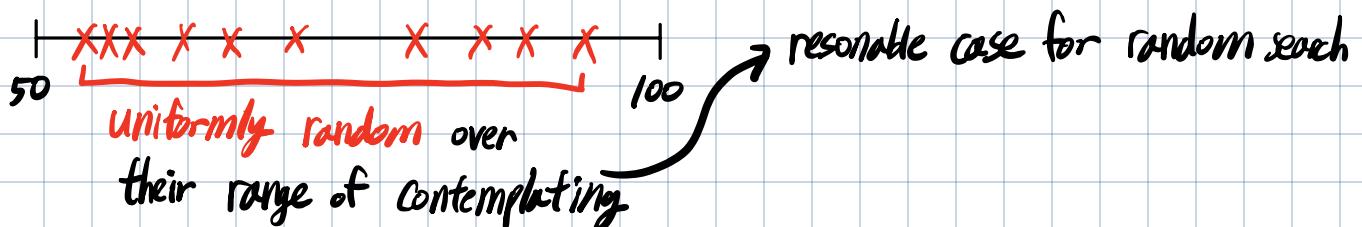


만약 빨간점 3개 이내로 model이 work well
→ 더 작은 양수로 함께 학대해서 더 조밀하게
점들을 터뜨려

- Using an appropriate scale to pick hyperparameters.

- Picking hyperparameters at random

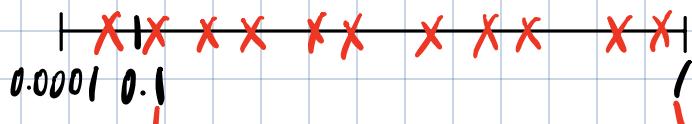
Let's say we explore best $n^{[1]}$ ($n^{[1]} = 50, \dots, 100$)



layers ($L = 2, \dots, 4$) \Rightarrow grid search is ok

- Appropriate scale for hyperparameters

Let's say if we explore best α ($\alpha = 0.0001, \dots, 1$)



90% of samples are over the range between 0.1 and 1

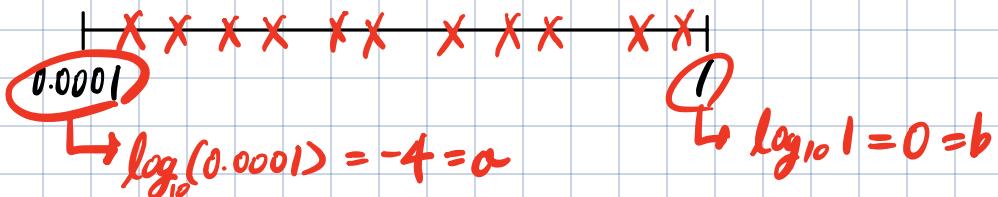
→ Unreasonable use for random search

$$\Rightarrow 0.0001 \approx 0.1 ?$$

\Rightarrow Change linear scale to **log scale** for random search.

① hyperparameter $\alpha \in \alpha = 10^r$ ($r \in [a, b]$) \Rightarrow We will define a and b .

② Effect of all that **log scale** define a and b



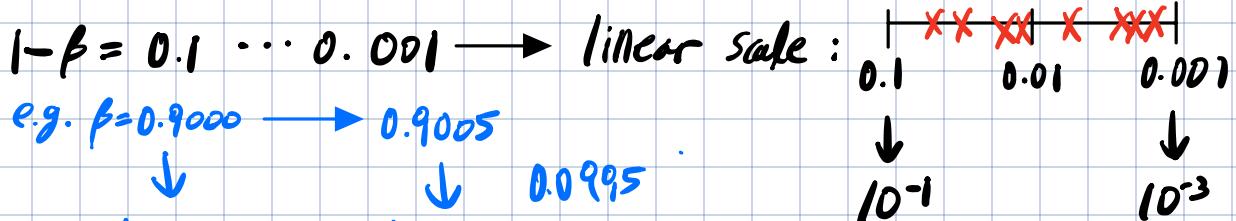
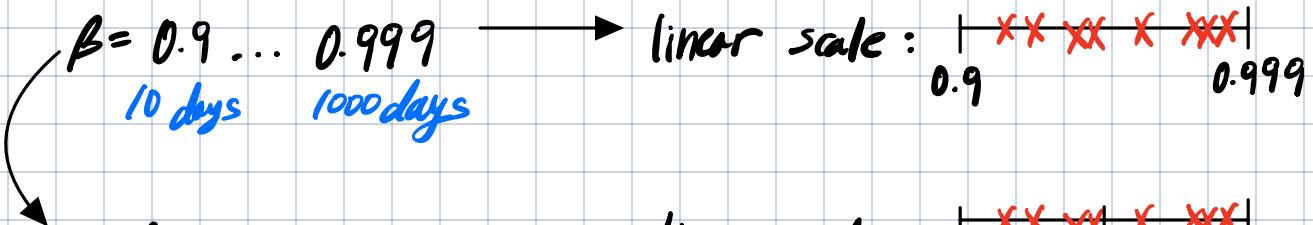
③ Overall we uniformly randomly all a and b to obtain $\alpha = 10^r$.

\Rightarrow hyperparameter ($\alpha = 10^r$) \equiv random search \Rightarrow $r \in [-4, 0]$.

$$r \in [-4, 0] \rightarrow r = -4 + np.random.rand()$$

$$\alpha = 10^r$$

- Hyperparameters for exponentially weighted averages.



$$\frac{1}{1-\beta} = 10 \quad \frac{1}{1-\beta} \approx 10$$

$$\beta = 0.999 \rightarrow 0.9995$$

$$\frac{1}{1-\beta} = 1000 \rightarrow \frac{1}{1-\beta} = 2000$$

$$r \in [-3, -1]$$

$$1 - \beta = 10^r$$

$$\beta = 1 - 10^r$$

\Rightarrow In linear scale, huge impact when $\beta \approx 1$ (sample of size small \Rightarrow huge impact)

• Hyperparameters tuning in practice: Pandas vs Caviar

- Re-test hyperparameters occasionally

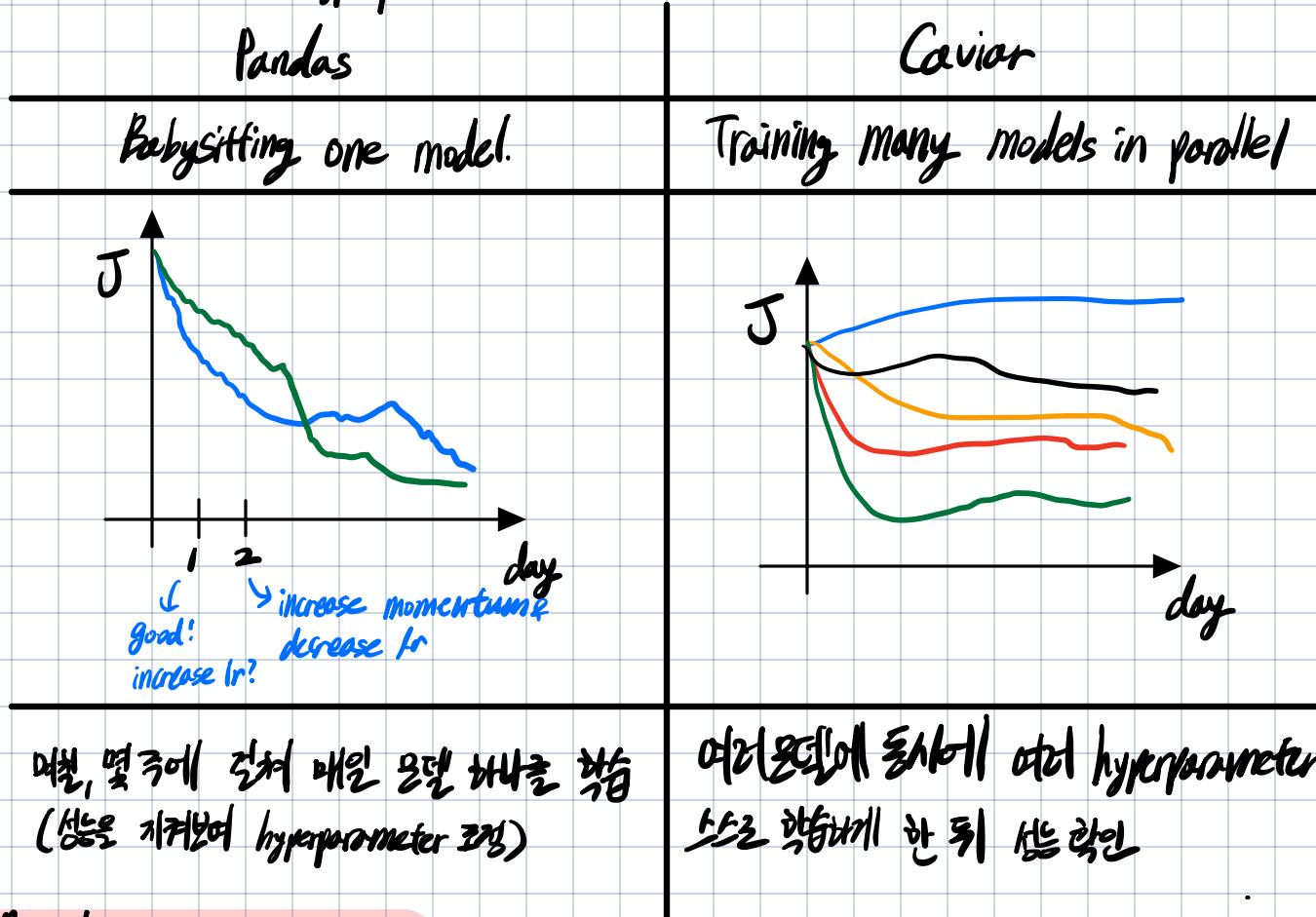
① hyperparameter를 다른 것들로 다른 영역에서 쓸 수도 있을 수도 있음.

e.g. NLP, Vision, Speech, Ads, logistics, ...

② 좋은 hyperparameter를 찾아도 아래 이유 때문에 가끔 재검증 필요

- [] 알고리즘의 뷰티
- [] 디버깅 간편
- [] 데이터 선택 시비 가능성도

- How to search hyperparameters?

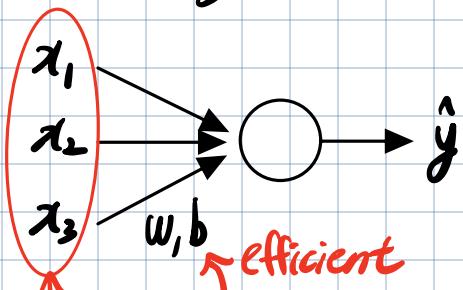


2. Batch Normalization

• Normalizing activations in a network.

⇒ Make the model robust to hyperparameters

- Normalizing inputs to speed up learning

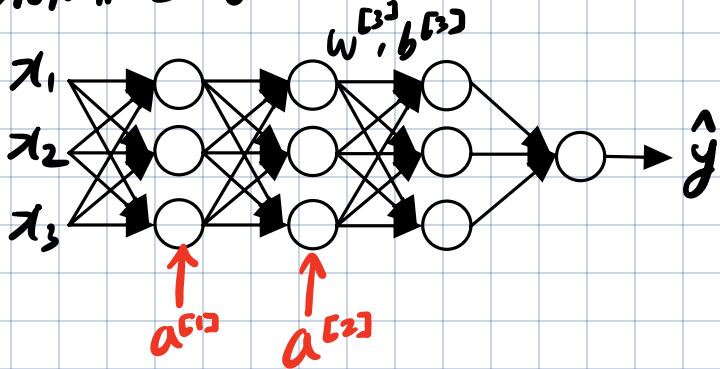


elongated $\rightarrow \begin{cases} \mu = \frac{1}{m} \sum_i x^{(i)} \\ X = X - \mu \\ \sigma^2 = \frac{1}{m} \sum_i (X^{(i)})^2 \\ X = X / \sigma \end{cases}$ easy to ...

Normalize) learning

Cost to
optimize

⇒ DNNNormalize는 입력을 뿐만 아니라 hidden layer의 값들도 정규화



⇒ Can we normalize $a^{[2]}$ for faster learning of $w^{[3]}$ & $b^{[3]}$?

A: Normalize $Z^{[2]}$ (commonly used)

- Implementing Batch Norm

Given some intermediate values in NN

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

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

$$Z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \rightarrow \text{make sure not to divide by 0}$$

$$\tilde{z}^{(i)} = f Z_{\text{norm}}^{(i)} + \beta \rightarrow \begin{array}{l} \text{learnable parameters} \\ \Rightarrow \text{gradient descent, ...} \end{array}$$

Use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

Why f & β ?

⇒ 입력을 가진 hidden unit의 μ & σ^2 가 항상 $(0, 1)$ 이 되는 건 bad
(hidden unit의 분포가 대칭이 아님과 비례하는 경우 가능)

⇒ f, β 를 이용해 hidden unit의 다른한 값을 가지게 조정

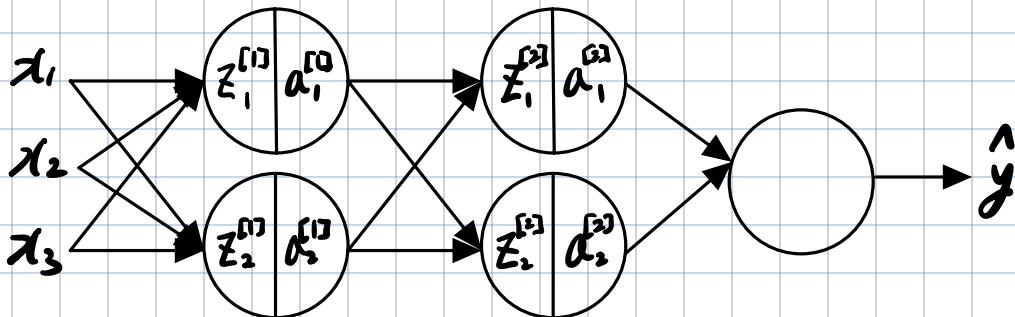
e.g.

학습할 때 sigmoid일 때 값이 0이 되면 학습정지 악화
∴ f와 β로써 문제 해결

$$\text{If } f = \sqrt{\sigma^2 + \epsilon}, \beta = \mu \Rightarrow \tilde{z}^{(i)} = z^{(i)}$$

• Fitting batch norm into a neural network

- Adding Batch Norm to a network



$$X \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow{\substack{\beta^{[1]}, \gamma^{[1]} \\ \text{Batch Norm(BN)}}} \tilde{z}^{[1]} \rightarrow \alpha^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} z^{[2]} \xrightarrow{\substack{\beta^{[2]}, \gamma^{[2]} \\ BN}} \tilde{z}^{[2]} \rightarrow \dots$$

different from
momentum / Adam

Learnable parameters :

$$\left\{ w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}, \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \right\}$$

$$\Rightarrow \text{Gradient descent : } \beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]}$$

↳ RMS prop, momentum etc. like this!

→ tf. nn. batch-normalization

- Working with mini-batches

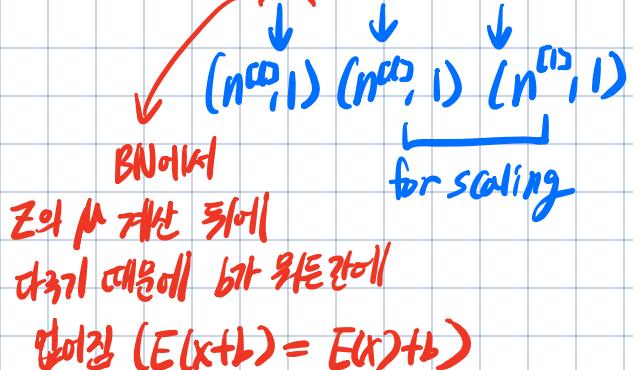
⇒ In practice, BN is often applied to mini-batches

$$X \xrightarrow{\substack{\text{313} \\ w^{[1]}, b^{[1]}}} z^{[1]} \xrightarrow{\substack{\beta^{[1]}, \gamma^{[1]} \\ \text{Batch Norm(BN)}}} \tilde{z}^{[1]} \rightarrow \alpha^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} z^{[2]} \rightarrow \dots$$

$$X \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]}$$

⋮ → 2000 mini batches use one BN

Parameters : $w^{[1]}, b^{[1]}, \cancel{\alpha^{[1]}}, \beta^{[1]}, \gamma^{[1]}$



e.g.

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$\bar{z}^{[l]} = w^{[l]} a^{[l-1]}$$

$$z_{\text{norm}}^{[l]} = \frac{z^{[l]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{[l]} = \gamma^{[l]} z_{\text{norm}}^{[l]} + \beta^{[l]}$$

$$\mu = \frac{1}{m} \sum_i x^{[l]}$$

β 가 차이로 평균을 계산

- Implementing gradient descent \rightarrow RMSprop, Adam, momentum, etc...

for $t = 1 \dots \text{num_minibatches}$

Compute forward prop on $X^{(t)}$

In each hidden layer, use BN to replace $z^{[l]}$ with $\tilde{z}^{[l]}$

Use backprop $\rightarrow dW^{[l]}, dB^{[l]}, df^{[l]}$

Update parameters:

$$w^{[l]} := w^{[l]} - \alpha dW^{[l]}$$

$$\beta^{[l]} := \beta^{[l]} - \alpha dB^{[l]}$$

$$f^{[l]} := f^{[l]} - \alpha df^{[l]}$$

- Why does Batch Norm work?

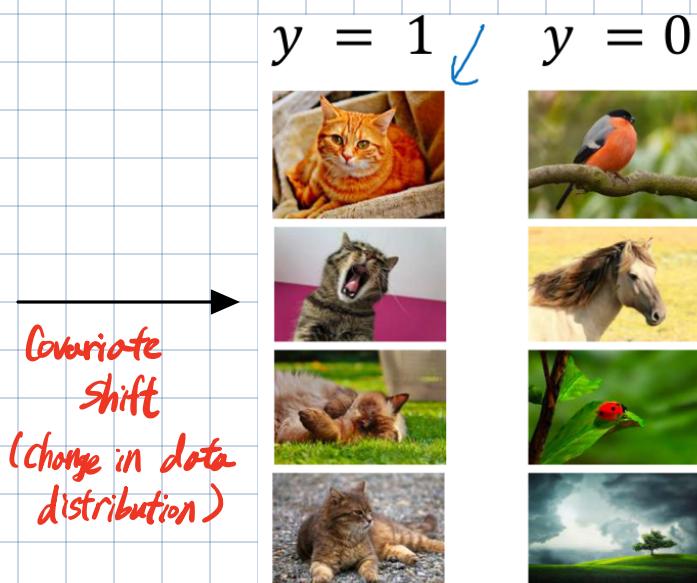
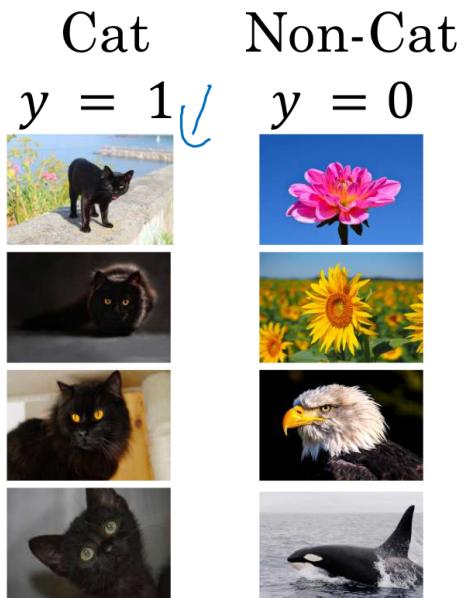
- Learning on shifting input distribution

NN에서 알맞은 가중치의 변화를 더하기 받는 정도

Shallow layer $>$ deeper layer

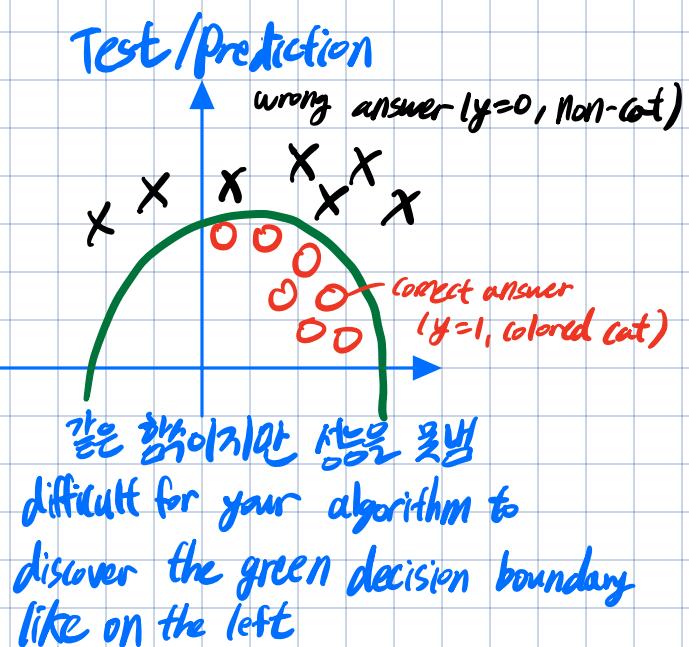
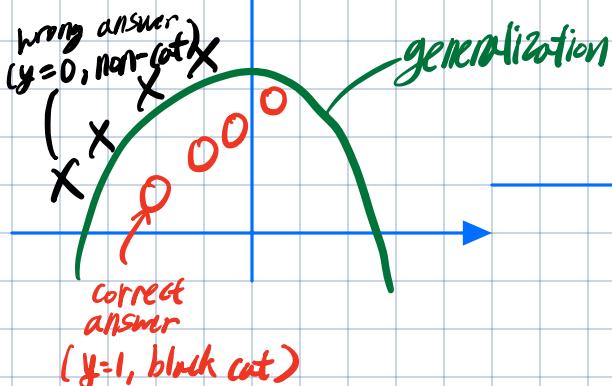
E.g. logistic regression (or deeper model)은 초기 모델이 무작위로 초기화 &

유기 고양이 vs. 비고양이 \Rightarrow low performance



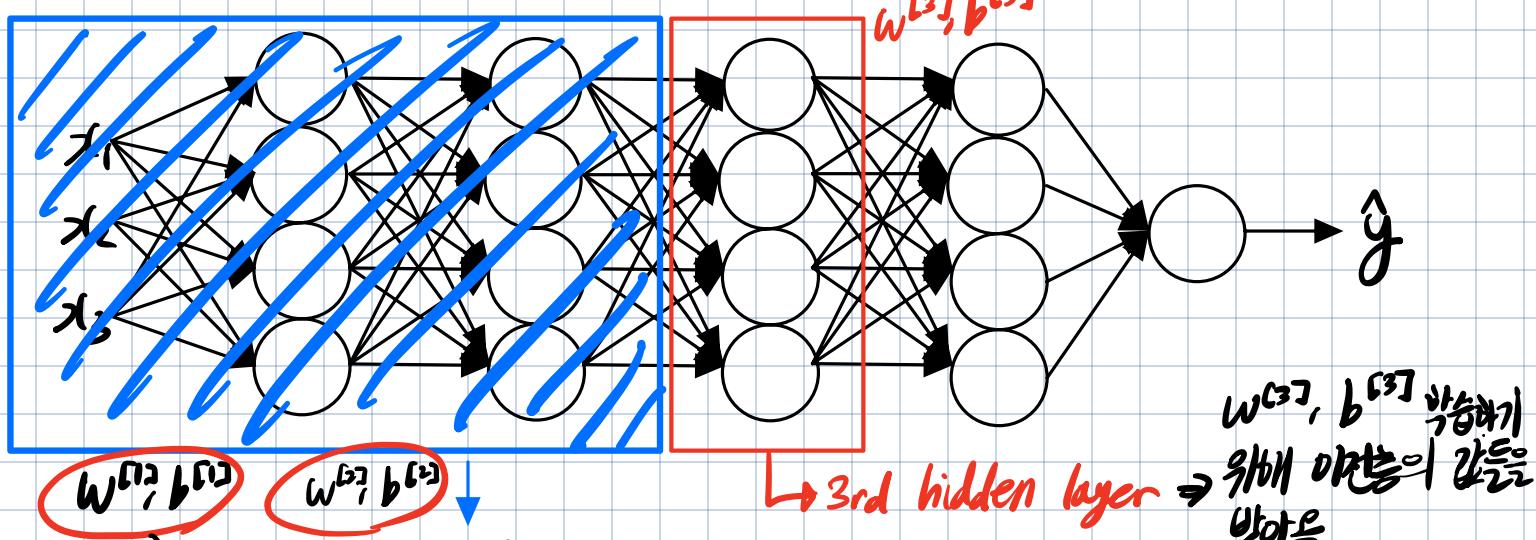
Covariate shift
(change in data distribution)

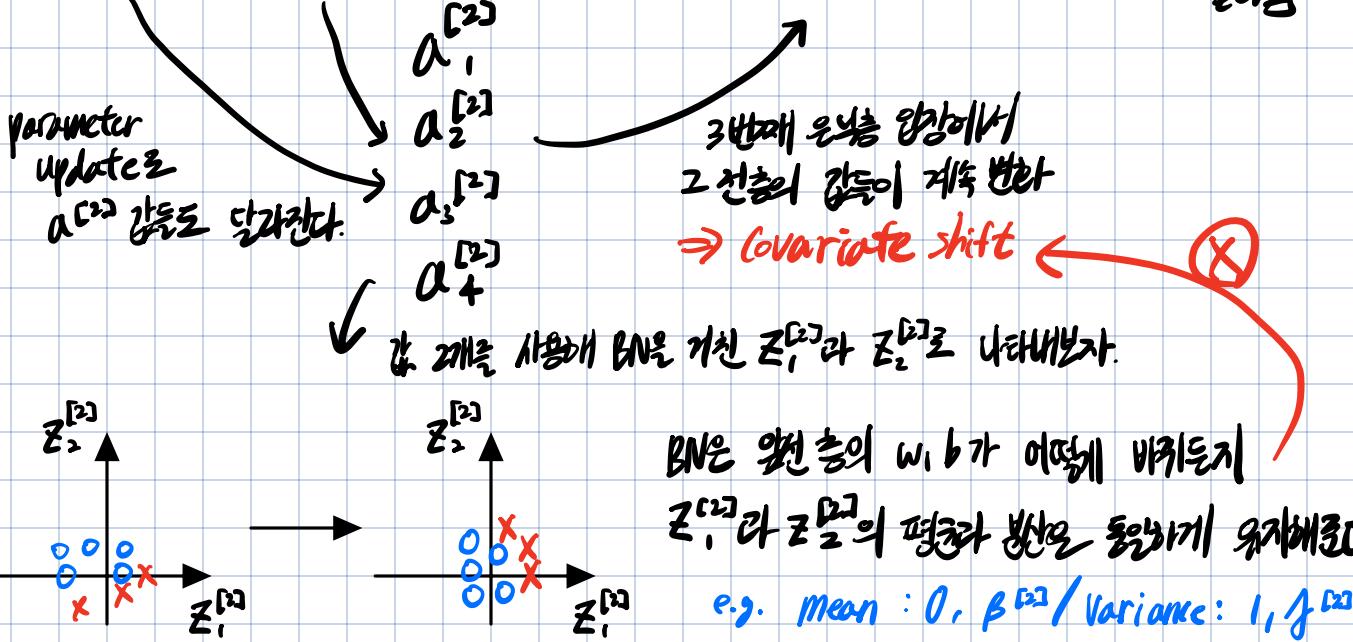
train & generalization



$\Rightarrow X \rightarrow Y$ 학습 시 X 의 분포가 바뀌면 (covariate shift) 예측을 해내는 문제

- Why this is a problem with neural networks?





\Rightarrow BN은 은닉 층의 분포가 변화하는 양을 줄여준다.

① 즉, 앞선 층에서 parameter가 바뀌었을 때 3rd hidden layer off
간 값의 분포를 재한해서 covariate shift를 안장한다

② 앞 층 parameter - 뒤 층 parameter 간 상관관계 ↓ \Rightarrow 뒤 층 활동에 용이

- Batch Norm as regularization

① Each mini-batch is scaled by the mean/Variance computed on just that mini-batch

\Rightarrow X 에 대해 $X_{\text{set}3}^n (n=64, 128, \dots)$ 가 더 작은 표준 deviation 존재
 \therefore mini-batch size ↑ \Rightarrow noise ↓ \Rightarrow regularization effect ↓

② 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

\Rightarrow noise가 있는 경우에 활동하기 때문에 $Z \rightarrow \tilde{Z}^{[l]}$ 과정에도 noise +
dropout을 학습에 따라 $\alpha^{[l]} \times 0$ or $\alpha^{[l]} \times 1 \Rightarrow$ 흔들림
BN은 흔들림(0) & 텁텁짐(1)

③ This has a light regularization effect

→ 훈련에 간접 추가하여 이전의 흐름이 현재의 흐름에 너무 의존하지 않도록 함
 (but, noise가 있어서 일관화 효과가 크지 X, 강한 흐름을 원하면 dropout과 함께 사용)

• Batch Norm at Test time

BN	Test
One mini-batch	a sample one at a time
$\mu = \frac{1}{m} \sum_i z^{(i)}$ $\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$ $z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$ $\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$ <p>Sum over zs in one mini-batch ($n=64, 128\dots$)</p> <p>There is no mini-batch in test set</p> <p>Come up with some separate estimates of μ & σ^2 from training set.</p>	μ, σ^2 : estimate using exponentially weighted average across the mini-batches <p>Layer 1:</p> $x^{(1)}, x^{(2)}, x^{(3)}$ $\downarrow \theta_1 \quad \downarrow \theta_2 \quad \downarrow \theta_3$ $\mu^{(1)[1]}, \mu^{(2)[1]}, \mu^{(3)[1]} \Rightarrow \mu$ $\sigma^{(1)[1]}, \sigma^{(2)[1]}, \sigma^{(3)[1]} \Rightarrow \sigma^2$ $z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$ $\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$ <p>track μ & σ^2 during training</p> <p>exponentially weighted averages</p>

3. Multi-class classification

• Softmax regression

- BN은 한 번에 하나의 mini-batch의 μ & σ^2 만 계산 \Rightarrow 테스트는 mini-batch X
 \Rightarrow 예측은 각각으로 다를 수 있어 필요
- Recognizing cats, dogs, and baby chicks



3

1

2

0

3

2

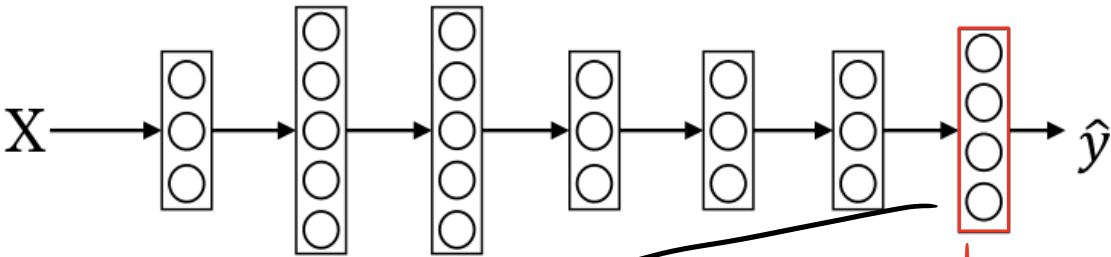
0

1

→ none of the classes

$C = \# \text{ of classes} = 4$ (label: 0 ~ C-1)

- Softmax layer



$$\hat{y} = (4, 1) = \begin{bmatrix} P(\text{other}|X) \\ P(\text{cat}|X) \\ P(\text{dog}|X) \\ P(\text{bc}|X) \end{bmatrix} \rightarrow \text{sum to 1}$$

C output units
 $N^{[L]} = C = 4$

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \quad (4, 1)$$

$$\text{e.g. } z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$t_{ij} \mid i=1 \dots 4$
normalize

Activation function for softmax:

$$t = e^{z^{[L]}} \text{ (element-wise)}$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}, \quad a_j^{[L]} = \frac{t_j}{\sum_{j=1}^4 t_j}$$

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

(4, 1)

(4, 1)

↳ ReLU, sigmoid-like or softmax-like
that, it's a vector

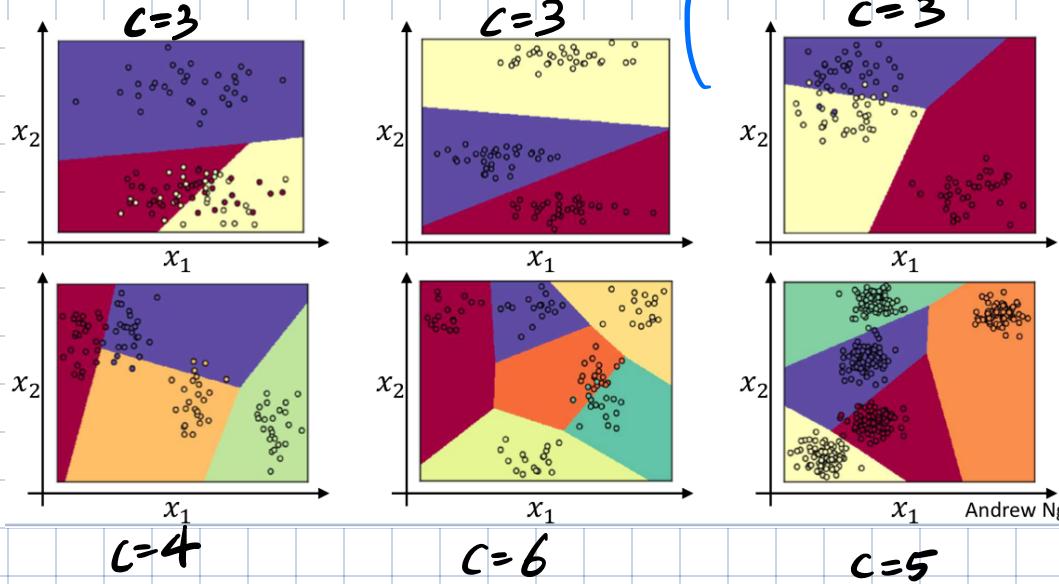
$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \cdot \sum_{j=1}^4 t_j = 196.3$$

$$a^{[L]} = \frac{t}{196.3} = \begin{bmatrix} \frac{e^5}{196.3} = 0.842 \\ \frac{e^2}{196.3} = 0.042 \\ \frac{e^{-1}}{196.3} = 0.002 \\ \frac{e^3}{196.3} = 0.114 \end{bmatrix}$$

- Softmax examples (w.o. hidden layer)

$$x_1 \quad x_2 \quad \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \hat{y} \quad (z^{[1]} = w^{[1]} x + b^{[1]} \\ a^{[1]} = \hat{y} = g(z^{[1]}))$$

hidden unit X
→ linear decision boundary



• Training a softmax classifier

- Understanding softmax

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \quad c=4 \quad g^{[L]}$$

→ 텐서의 1이 되도록 정규화

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

Softmax

hard max

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

⇒ Softmax regression generalizes logistic regression to c classes
(rather than just 2 classes)

(∴ If $c=2 \Rightarrow$ softmax reduces to logistic regression $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix} \rightarrow$ 예전에 0.842
→ 예전에 0.158)

- Loss function

$$y = \begin{bmatrix} y_1 & 0 \\ y_2 & 1 \\ y_3 & 0 \\ y_4 & 0 \end{bmatrix} \Rightarrow \text{cat} \quad \hat{y}^{(i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} = a^{[L](i)}$$

Our prediction

not working well on this example.

for one training sample: $\underline{\underline{J(\hat{y}, y)}} = -\sum_{j=1}^4 y_j \log \hat{y}_j$
minimize by GD ←

Since $y_1 = y_3 = y_4 = 0$, $J(\hat{y}, y) = -\sum_{j=1}^4 y_j \log \hat{y}_j = -\log \hat{y}_2 \Rightarrow$ should make \hat{y}_2 big ($0 \leq \hat{y}_2 < 1$)

⇒ Loss function은 훈련 세트에서 각각의 클래스가 놓인 것에 그 클래스의 확률을 선택한 것의 합이다.

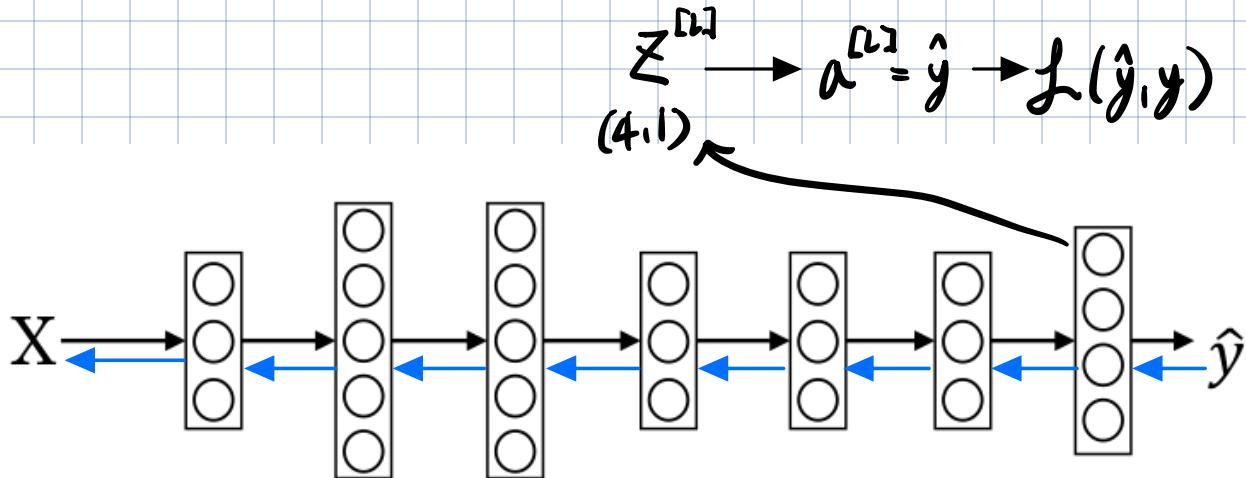
(= 2Klge 32)

for entire training set: $J(w^{(0)}, b^{(0)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (4, m)$$

$$\hat{Y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}] = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad (4, m)$$

- Gradient descent with softmax



$$\text{Backprop: } dZ^{[L]} = \hat{y} - y = \frac{\partial J}{\partial Z^{[L]}}$$

4. Programming Frameworks

• Deep Learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

Choosing deep learning frameworks

- Ease of programming (development and deployment)
- Running speed
- - Truly open (open source with good governance)