

manually decay after observation

Week-3

### Hyperparam Tuning Process

Hyper parameters rank

- |  |                            |
|--|----------------------------|
| □ $\alpha$   | ★ 3 <sup>rd</sup> priority |
| ○ $\beta$ NO.9   | □ 1 <sup>st</sup> priority |
| $\beta_1, \beta_2, \varepsilon \rightarrow 0.9, 0.99, 10^{-8}$ | ○ 2 <sup>nd</sup> pr       |

- \* # layers
- # hidden units
- \* learning rate decay
- mini batch size

Previously a grid was made ( $5 \times 5$ ). And all 25 params were used to test for best model.

But in case 1 hyperparam is more important than the other, then in such a case, you get to try only 5 values of that param

Now, we take 25 points only, but select randomly.

In case if we have more than 2 parameters, we use a cube. In such a case, even if we don't know which hyper parameter is important, we can still get the best result.

The area where optimal solution is obtained can be then more densely sampled.

But in order to choose, we need an appropriate scale to choose from.

For example  $\alpha = 0.0001, \dots, 1$

it may happen that 90% of samples are b/w 0.1 to 1. and only 10% are from 0.0001 to 0.1

In such a case, we might want to use a logarithmic scale.

$$0.0001 \quad 0.001 \quad 0.01 \quad 0.1 \quad 1$$

Python implementation {

$$\begin{aligned} r &= -4 \times \text{np.random.rand}() && \leftarrow r \in [0^{-4}, 1] \\ \alpha &= 10^r && \leftarrow 10^{-4} \dots 10^0 \end{aligned}$$

For exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

↓  
10 samples      ↓  
1000 samples  
avg.

$$1-\beta = 0.1 \dots 0.001$$

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

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

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

→ Retest hyperparameters occasionally

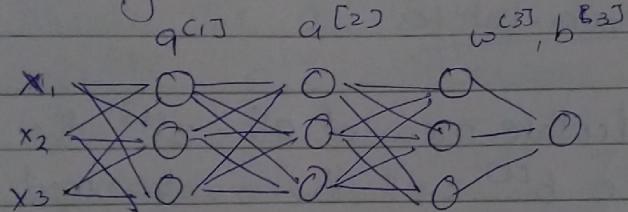
Two approaches

- Babysit a model [watch carefully the decrease in cost, make some changes, make more efficient]
- Training many models in parallel

Should choose based on your computational research

## Batch Normalization

Makes hyper parameter search problem easy, by increasing the scale



We can normalize  $w^{(3)}$ ,  $b^{(3)}$  so as to train  $w^{(3)}, b^{(3)}$  faster.

In this case we will normalize  $z^{(2)}$

Given some intermediate values in NN  $z^{(1)}, \dots, z^{(m)}$

$$\begin{aligned} \mu &= \frac{1}{m} \sum z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum (z^{(i)} - \mu)^2 \\ z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \end{aligned}$$

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

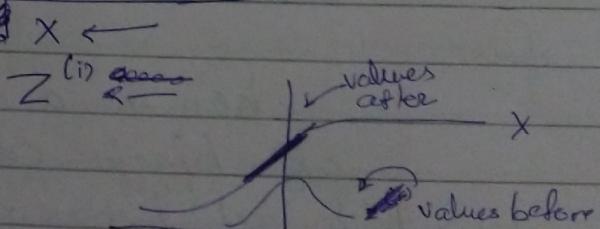
if  $\gamma = \sqrt{\sigma^2 + \epsilon}$  if  $\sigma^2 = 0$

$$\beta = \mu$$

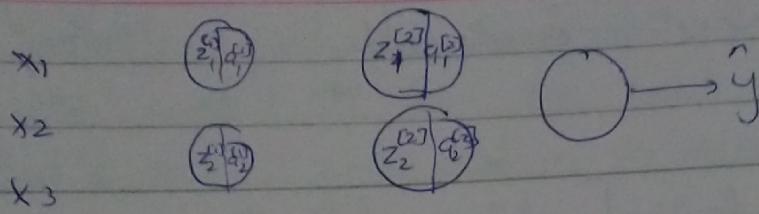
then  $\tilde{z}^{(i)} = z^{(i)}$

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

$$\begin{aligned} \mu &= \frac{1}{m} \sum x^{(i)} \\ x &= x - \mu \\ \sigma^2 &= \frac{1}{m} \sum (x^{(i)} - \mu)^2 \quad \text{Normalization} \\ x &= x / \sigma^2 \quad \text{Help in Logi. Regr.} \end{aligned}$$



So to adjust mean and variances deep in the NN we use  $\gamma$  &  $\beta$



$$X \xrightarrow{w^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow{\gamma^{[1]}, \beta^{[1]}, BN} \tilde{Z} \xrightarrow{\alpha^{[1]}} a^{[1]} = g^{[1]}(\tilde{Z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow{\gamma^{[2]}, \beta^{[2]}} \tilde{Z}^{[2]} \xrightarrow{\alpha^{[2]}} a^{[2]} \rightarrow \dots$$

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

To update  $\beta$ ,  
 $d\beta^{[l]}$  is used;  $\beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]}$   
→ Actually applied to mini-batches

When we actually calculate mean of all  $z^{[l]}$ , the constants are subtracted. Here,  $b^{[l]}$  will be subtracted from each entry and thus, they can be set to zero.  
In the normalized  $\tilde{Z}$ , beta is replaced by  $\beta^{[l]}$ .  
and is responsible for shift and bias.

$$Z^{[l]} = w^{[l]} a^{[l-1]}$$

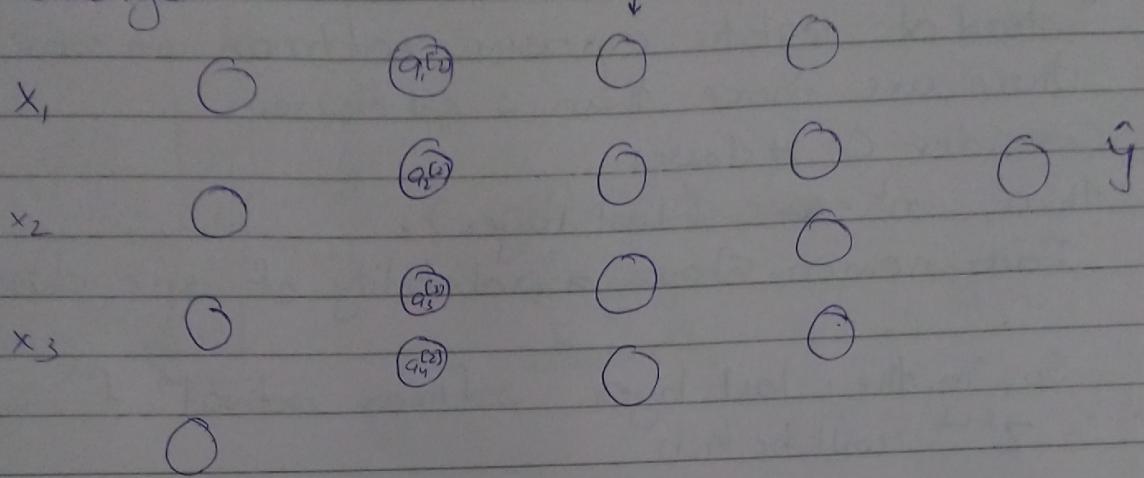
$$\tilde{Z}^{[l]} = \gamma^{[l]} Z_{norm} + \beta^{[l]}$$

dimension of  $Z^{[l]}$ ,  $\beta^{[l]}$  and  $\gamma^{[l]}$  is  $(n^{[l]}, 1)$

When in a case, if a model is trained on all black cats, and a coloured cat image is given to test, it may not understand.  
Such a case is called covariate shift.

Consider a The n/w if wants to learn about the new parameters there will be a lot of change in the weights.

What batch norm does, is make weights robust to changes



Consider the 3<sup>rd</sup> layer. Due to changes in the rest layer,  $a_1^{(2)}, \dots, a_n^{(2)}$  will change. batch norm makes sure that these changes do not affect the learning process of  $w^{(3)}, b^{(3)}$ . It makes sure that these values are more stable.

Thus, the early layers don't get to shift much, due to the constraints. This makes learning of later layers easier.

→ Batch Norm has a slight regularization effect.

Batch Norm at test time.

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

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

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

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

To estimate  $\bar{y}, \sigma^2$ , exponentially weighted avg (across mini-batches).

$$\begin{array}{c} x^{(1)}, x^{(2)}, x^{(3)} \\ \downarrow \\ \bar{y}^{(1)}, \bar{y}^{(2)}, \bar{y}^{(3)} \end{array}$$

$$\sigma^2 \approx \bar{\sigma}^2, \sigma^2 \approx \bar{\sigma}^2, \dots \rightarrow \bar{\sigma}^2$$

for each layer,  $\bar{y}, \sigma^2$  is taken of each mini batch, and averaged.

The finally averaged values are taken during test time and found to be robust.

## Softmax Regression

Instead of logistic regression, softmax is used, when there are more than 2 cat classes.

Consider  $C = \# \text{classes}$ .

then,  $n^2 = C$  (last layer).

Each neuron shows a probability of one class.

So, in the last layer softmax activation  $f^m$  is used.  
 $Z^{[L]}$  will be  $(4, 1)$

$$t = e^{(Z^{[L]})}$$

$$(4, 1) \quad a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^4 t_i} \quad a_i^{[L]} = \frac{t_i}{\sum t_i}$$

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

$$\begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^{-3} \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 0.1 \end{bmatrix} \quad \sum_{j=1}^4 t_j = 176.3$$

$$\begin{array}{l} \text{---} \rightarrow 0.842 \left( \frac{e^5}{176.3} \right) \\ \text{---} \rightarrow \frac{e^2}{176.3} = 0.042 \\ \text{---} \rightarrow \frac{e^{-1}}{176.3} = 0.002 \\ \text{---} \rightarrow \frac{e^{-3}}{176.3} = 0.001 \end{array}$$

Without hidden layer it can work only for linear decision boundary.

## Training Softmax Classifier

### Loss function

$$y^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{cat}$$

$y_1 = y_3 = y_4 = 0$

~~$\frac{\partial L}{\partial y^{(1)}} = -1$~~

$$L(\hat{y}, y) = - \sum_{j=1}^m y_j \log \hat{y}_j$$

↑  
make small

$$= -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

Make  $\hat{y}_2$  big

$$J(\omega^{(0)}, b^{(0)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$Y = [y^{(1)} \ y^{(2)} \dots \ y^{(m)}]$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}_{(4, m)}$$

$$\hat{y} = [\hat{y}^{(1)}, \dots, \hat{y}^{(m)}]$$

Backpropagation:

$$\underbrace{\frac{\partial z^{(l)}}{\partial z^{(l)}}}_{(4, 4)} = \hat{y} - y$$

### Deep Learning Frameworks

Caffe / Caffe2

CNTK

DL4J

Keras

Lasagne

mxnet

PaddlePaddle

TensorFlow

Theano

Torch

Tensorflow