

Course 2

Improving Deep Neural Networks: Hyperparameters Tuning, Regularization and Optimization

WEEK 1

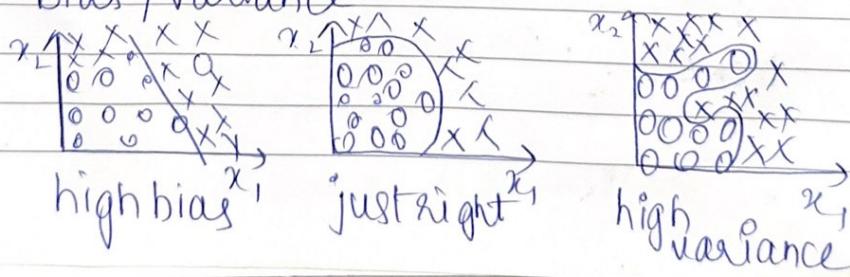
→ Train / Dev / Test sets

Data divided →

- ① Training set
- ② Development set
- ③ Test set

- Above sets are used to evaluate the performance of machine learning algorithm
- Traditional ratios for Train / Dev / Test sets are 98% / 1% / 1% where there is a million examples.

→ Bias / Variance



$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(h_i(x_i), y_i) + \frac{\lambda}{2} \|w\|^2$$

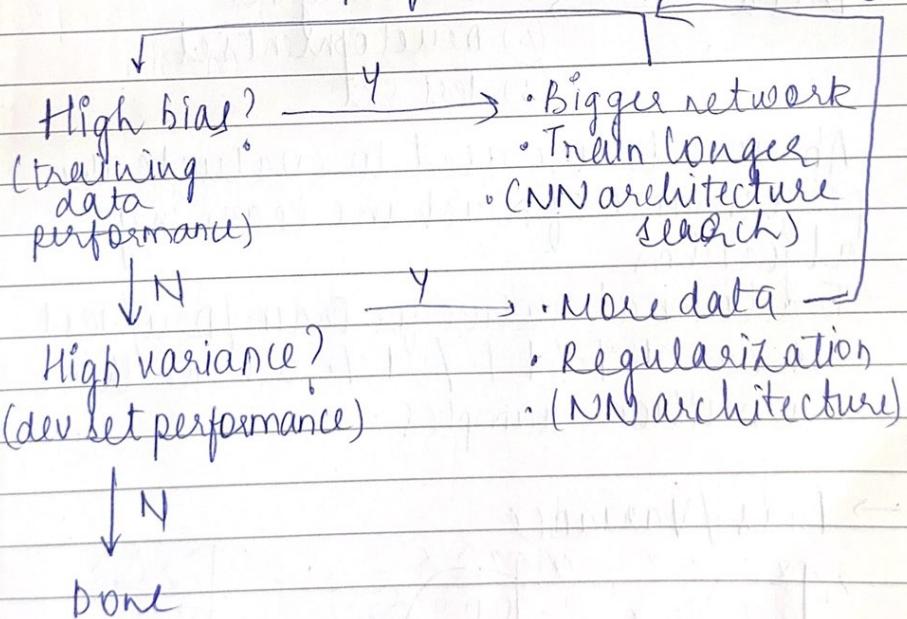
$w \in \mathbb{R}^n, b \in \mathbb{R}$

← Logistic Regression

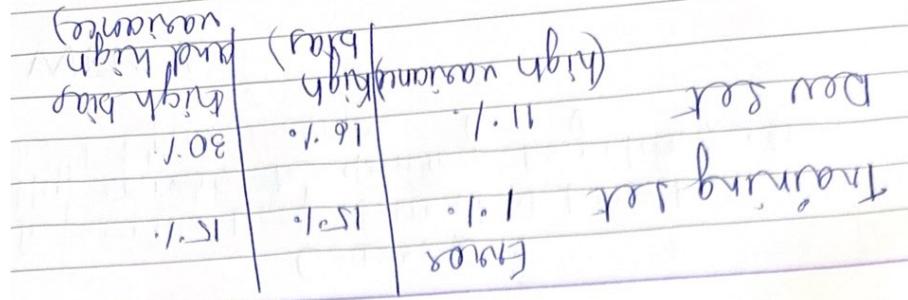
← Regularization

	Error		
Training Set	1%	15%	15%
Dev Set	11% (high variance)	16% (high bias)	30% (high bias and high variance)

→ Basic Recipe for Machine learning



← Basic Recipe for Machine Learning



→ Regularization

→ Logistic Regression

$$\min_{w,b} J(w, b) \quad w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$+ \frac{\lambda}{2m} b^2$$

↳ this term is omitted

→ L₂ Regularization

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

→ L₁ Regularization

$$\frac{1}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{1}{2m} \|w\|_1$$

$\lambda \rightarrow$ Regularization Parameter

This auxiliary function to update the network
 by calculating linear
 sum, and then activation function with
 in loops, the weight will be updated
 to that of the activation parameter
 function for every activation helps
 building too large
 neural networks the weight matrix from
 each term to the cost function that
 regularization is done by adding an
 reducing the variance of the model
 regularization helps with overfitting by
 ← noisy regularization reduces overfitting



→ Neural Network

$$\begin{aligned}
 J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) \\
 &\quad + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|^2 \\
 \|w^{[l]}\|^2 &= \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2
 \end{aligned}$$

"Frobenius
Norm"

$w \rightarrow$ matrix
of dimension
 $(n^{[0]} \times n^{[1]})$

$$dw^{[l]} = (\text{from backprop}) + \frac{1}{m} w^{[l]}$$

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

"Weight decay"

$$w^{[l]} = w^{[l]} - \alpha \left[(\text{from backprop}) + \frac{1}{m} w^{[l]} \right]$$

$$= w^{[l]} - \alpha \frac{1}{m} w^{[l]} - \alpha (\text{from back-prop})$$

"using up to
what margin"

"WON
margin"

$$\begin{aligned} & \left(\sum_{i=1}^m \sum_{j=1}^n w_j h_j(x_i) - b \right)^2 \\ & + \lambda \left(\sum_{j=1}^n w_j^2 \right) \end{aligned}$$

Normal Network



→ Why Regularization Reduces Overfitting

- Regularization helps with overfitting by reducing the variance of the model
- Regularization is done by adding an extra term to the cost functⁿ that penalizes the weight matrices from being too large
- Intuition for why regularization helps is that if the regularization parameter is large, the weights will be relatively small, and the activation function will be relatively linear
- This results in a simpler network that is less prone to overfitting
- When implementing regularization, make sure to plot the new definition of the cost functⁿ that includes the regularization term.

→ Dropout Regularization

- It works by randomly eliminating nodes in the network and all of the

outgoing connections from those nodes

- This results in a much smaller network which is then trained using backpropagation
- Implementation of Dropout

"Inverted Dropout"

$$l = 3, \text{keep-prob} = 0.8$$

$$\begin{aligned} d3 &= \text{np.random.rand}(a3.shape[0], \\ &\quad a3.shape[1]) < \text{keep-prob} \\ f[a3] &= \text{np.multiply}(a3, d3) \\ [a3] / &= \text{keep-prob} \end{aligned}$$

→ these steps ensure that the expected value of the activations remains the same

- At test time, dropout is not used as it would add noise to the predictions

→ Understanding Dropout

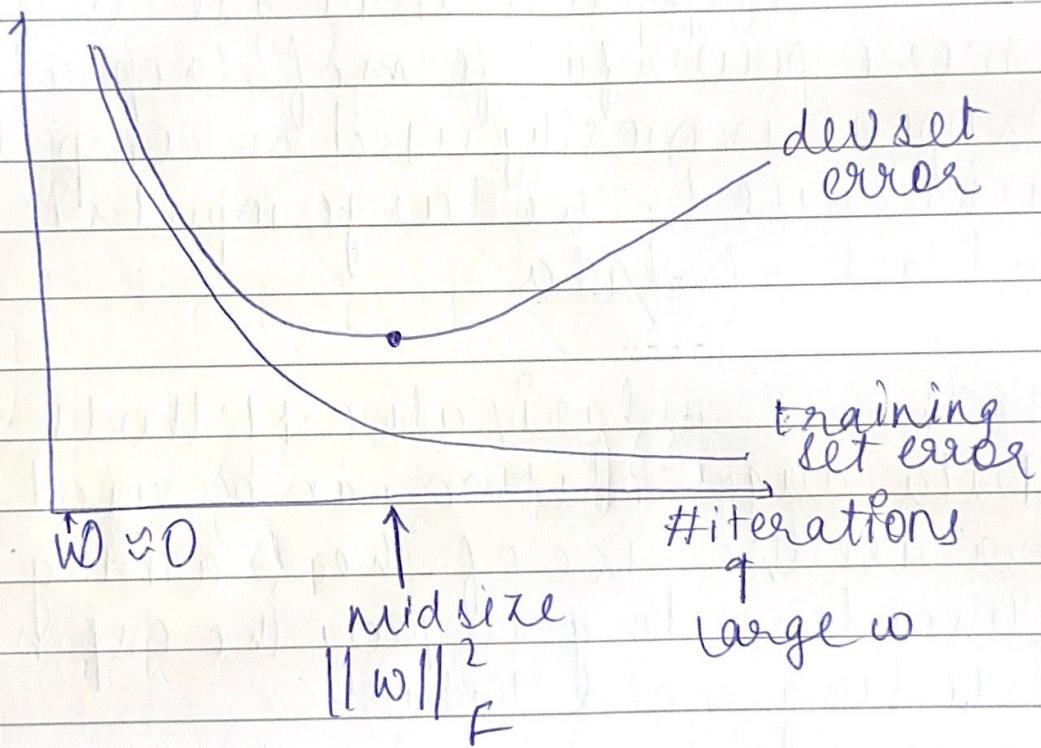
- Due to randomness, we cannot rely on any feature.

- This encourages the unit to spread out its weights, which has similar effect to L2 regularization
- It is possible to vary the keep-prob by layer
- For layers with a lot of parameters, a lower keep-prob can be used to apply a more powerful form of dropout
- Dropout is mostly used in computer vision due to the large inputsizes and lack of data

→ Other Regularization Methods

- Data augmentation can be used to increase the size of the training set without needing to pay the expense of collecting more data
 - Eg: flipping images horizontally, taking random crops of the image, and imposing random rotations and distortions
 - Early stopping is another technique used to reduce overfitting, which involves

plotting the training and dev set errors and stopping the training of the neural network when the dev set error and stopping the training of the neural network when the dev set error is at its lowest.



- Early stopping couples the tasks of optimizing the cost functⁿ and reducing variance, which can make the search space of hyperparameters more complex

SETTING UP YOUR OWN OPTIMIZING OR OPTIMIZATION PROBLEM

→ Normalizing Inputs

- It is a technique used to speed up training of a neural network
- It corresponds to two steps:
 - ① subtract out / zero out the mean (μ)
 - ② Normalize the variances (σ)

Subtract

mean

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

$$x' := x - \mu$$

Normalize

variance

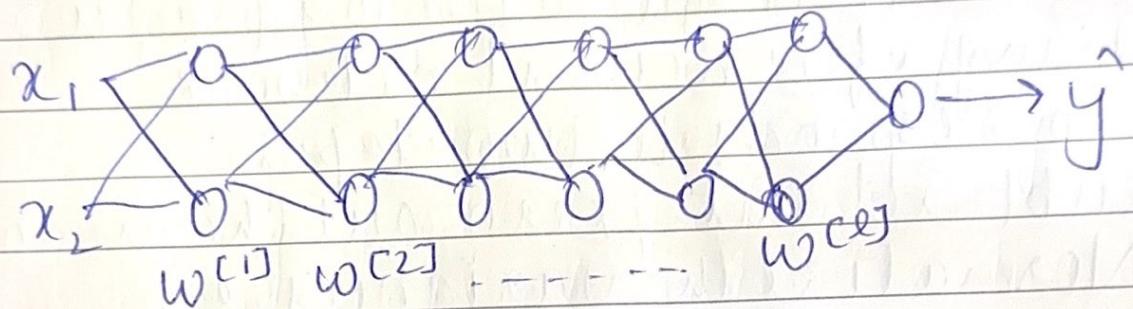
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} * 2$$

$$x' / \sigma$$

- Normalizing inputs helps to make the cost function more symmetric and easier to optimize
- If input features come from different scales, then it is important to normalize the features.
- If input features come from similar scales, normalizing is less imp. but beneficial

→ Vanishing / Exploding Gradients

- These problems occur when training very deep neural networks



$$g(z) = z$$

$$b^{[l]} = 0$$

$$y = w^{[1]} w^{[2]} w^{[3]} \dots$$

$$w^{[1]} x = z_1 = a_1 = g(z_1)$$

$$w^{[2]} a_1 - z_2 = a_2 = g(z_2)$$

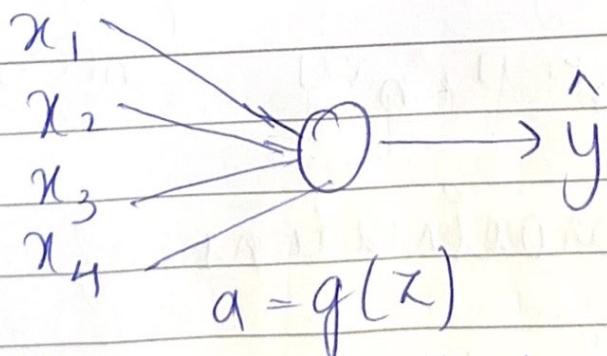
- If each weight matrix $w^{[l]}$ is just a bit larger than one times the identity \hat{y} will be very large and grow exponentially with no. of layers
- Similarly if weight matrix $w^{[l]}$ is less than one then \hat{y} will decrease exponentially.

- The same argument can be used to show that the derivatives or gradients will also increase/decrease exponentially as a function of the no. of layers

→ Weight Initialization for Deep Networks

- Vanishing and exploding gradients are problems that can occur in very deep neural networks
- To solve the above problem, we must be more careful for the choice of random initialization for the neural network

.



$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

- To make z not blow up/become too

small, the variance is set to $1/n$, where $n = \text{no. of input features}$

- If using ReLU functⁿ, the variance of w should be set to $2/n$
- This helps reduce the vanishing and exploding gradients problem
- Other variants of initialization used are:

① Xavier Initialization

$$\text{variance} = \sqrt{\frac{2}{n^{(l-1)} + n^{(l)}}}$$

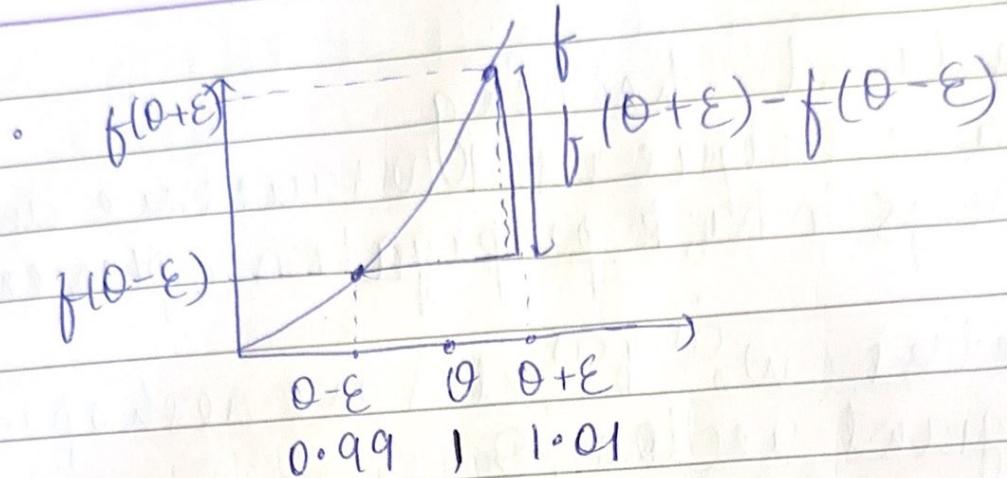
$$\tanh = \sqrt{\frac{1}{n^{(l-1)}}} \rightarrow w^{(l)} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}$$

②

$$\text{variance} = \sqrt{\frac{2}{n^{(l-1)} + n^{(l)}}} \left(\frac{2}{n^{(l-1)}} \right)$$

→ Numerical Approximations of Gradients

- Numerical approximations of gradients can be used to verify the correctness of a backpropagation implementation



$$\frac{f(0+E) - f(0-E)}{2E} \approx g(0)$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.001$$

$$g(0) = 30^2 = 3$$

approx error = 0.0001

$$f'(0) = \lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon} \rightarrow O(\epsilon^2)$$

$$f' = \frac{f(0+\epsilon) - f(0)}{\epsilon} \rightarrow O(\epsilon)$$

This two sided difference formula is what is used when doing gradient checking.

→ Gradient Checking

• It is a technique used to save time and find bugs in back propagation implementation

- Parameters ($w^{(1)}, b^{(1)}$, etc) are reshaped into a giant vector, θ
- Cost functⁿ (J) is now a functⁿ of theta
- Implementation:
for each i :

$$d\theta_{\text{approx}}^{(i)} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\hat{=} d\theta^{(i)} = \frac{\partial J}{\partial \theta^i}$$

$$d\theta_{\text{approx}} \hat{=} d\theta$$

- $d\theta_{\text{approx}}$ and $d\theta$ are compared by calculating Euclidean distance b/w them

$$\| d\theta_{\text{approx}} - d\theta \|_2$$

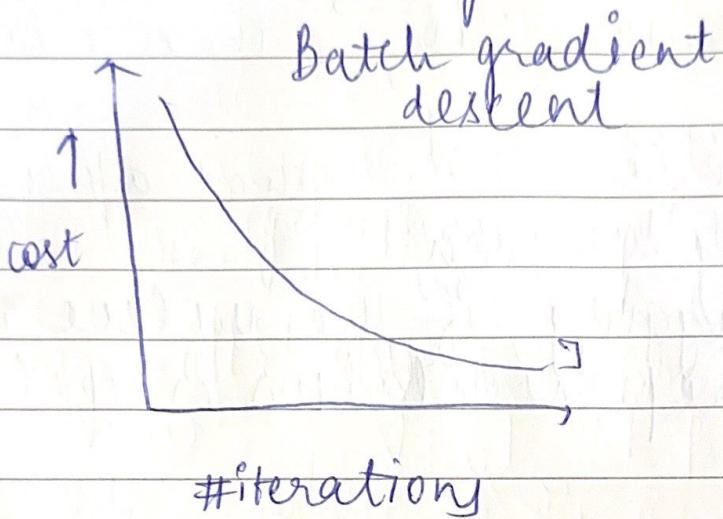
- If distance is of order 10^{-7} or smaller, the derivative approximation is correct
- If distance is $> 10^{-5}$, then a bug is suspected
- If distance is $> 10^{-3}$, then a bug is definitely present
- After debugging, if the value is small, the implementation is likely correct

→ Gradient Checking Implementation Notes

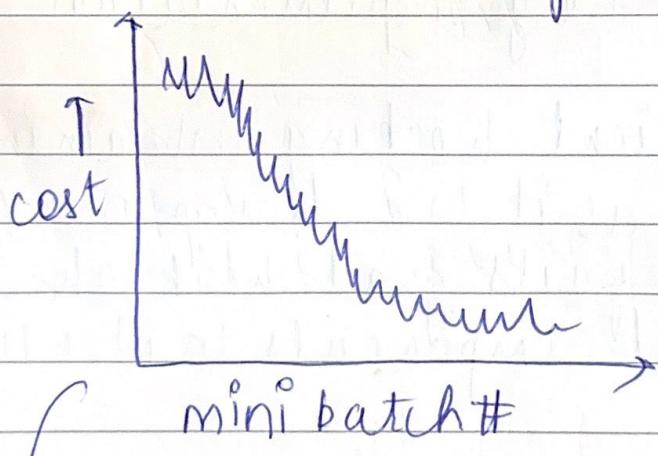
- Don't use gradient checking in training, only to debug as it is a slow process
- If algorithm fails grad check, check the individual components to identify the bug
- Remember regularization
- Grad check doesn't work with dropout so turn off dropout and use grad check to double check the algorithm without dropout
- Run at random initialization; perhaps again after some training

WEEK 2 : OPTIMIZATION ALGORITHM

→ Mini-batch gradient descent



Batch gradient
descent



Mini batch gradient descent

- The cost may not decrease on every iteration, but should trend downwards.
- Size of minibatch should be b/w 1 to the size of training set.

- Batch gradient takes longer time per iteration

- Stochastic gradient is noisy and never converges

- Mini batch gradient descent is more efficient method

- Stochastic gradient is noisy and can be reduced by using smaller learning rate

- If training set is small, then use gradient descent

- Typical mini batch sizes are below 64 and 512, it is faster if size is a power of 2

- Make sure the mini-batch fits in CPU/GPU memory

- Minibatch is a hyperparameter that can be searched to find most efficient value

→ Exponentially Weighted Average

- It is a faster optimization algorithm

- General formula EWA $\Rightarrow V_t = 0.9 V_{t-1} + (0.1) \delta_t$

- $v_t = \beta v_{t-1} + (1-\beta) \theta_t$

→ Gradient descent with momentum

- The basic idea is to compute an exponentially weighted average of gradients and use that to update weights
- This algorithm helps to dampen oscillation in the vertical direction and move quickly in the horizontal direction
- Compute d_w and d_b on the mini batch

$$v_{dw} = \beta v_{dw} + (1-\beta) d_w$$

$$v_{db} = \beta v_{db} + (1-\beta) d_b$$

$$w = w - \alpha v_{dw}$$

$$b = b - \alpha v_{db}$$

→ RMS prop

- Root Mean square prop
- It works by keeping an exponentially weighted average of the squares of the derivatives
- Compute d_w, d_b on current mini batch

$$s_{dw} = \beta s_{dw} + (1-\beta) d_w^2$$

$$s_{db} = \beta s_{db} + (1-\beta) d_b^2$$

- $w_t = w - \alpha \frac{dw}{\sqrt{sdw}}$

- $b_t = b - \alpha \frac{db}{\sqrt{sdb}}$

→ Adam Optimization Algorithm

- Implementation:

$$v_{dw} = 0, s_{dw} = 0, v_{db} = 0, s_{db} = 0$$

- On iteration t , compute dw and db using current mini batch

- $v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw$ momentum

- $v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$

- $s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2$ RMSprop

- $s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2$

- corrected

$$\hat{v}_{dw} = v_{dw} / (1 - \beta_1^t); \hat{v}_{db} = v_{db} / (1 - \beta_1^t)$$

$$\hat{s}_{dw} = s_{dw} / (1 - \beta_2^t); \hat{s}_{db} = s_{db} / (1 - \beta_2^t)$$

→ Learning Rate Decay

- Reduces learning rate for faster algorithm
- This helps to reduce noise in the mini batches and allows the algorithm to converge to the minimum more quickly
-

$$\alpha = \frac{1}{1 + \text{decay rate} \times \text{epoch number}} \times \alpha_0$$

$$\alpha = 0.95^{\text{epoch number}} \cdot \alpha_0$$

$$\alpha = \frac{k}{\sqrt{\text{epoch number}}} \cdot \alpha_0 \text{ or } \frac{k}{\sqrt{t}} \cdot \alpha_0$$

$$w_{ab} = P_{ab} + C P_{ab} w_{ab}$$

WEEK 3

→ Normalizing Activations in a Network

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

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

$$\cdot x = x / \sigma$$

• Implementing Batch Norm
Given some values in NN
 $z^{(1)}, \dots, z^{(m)}$

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

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

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

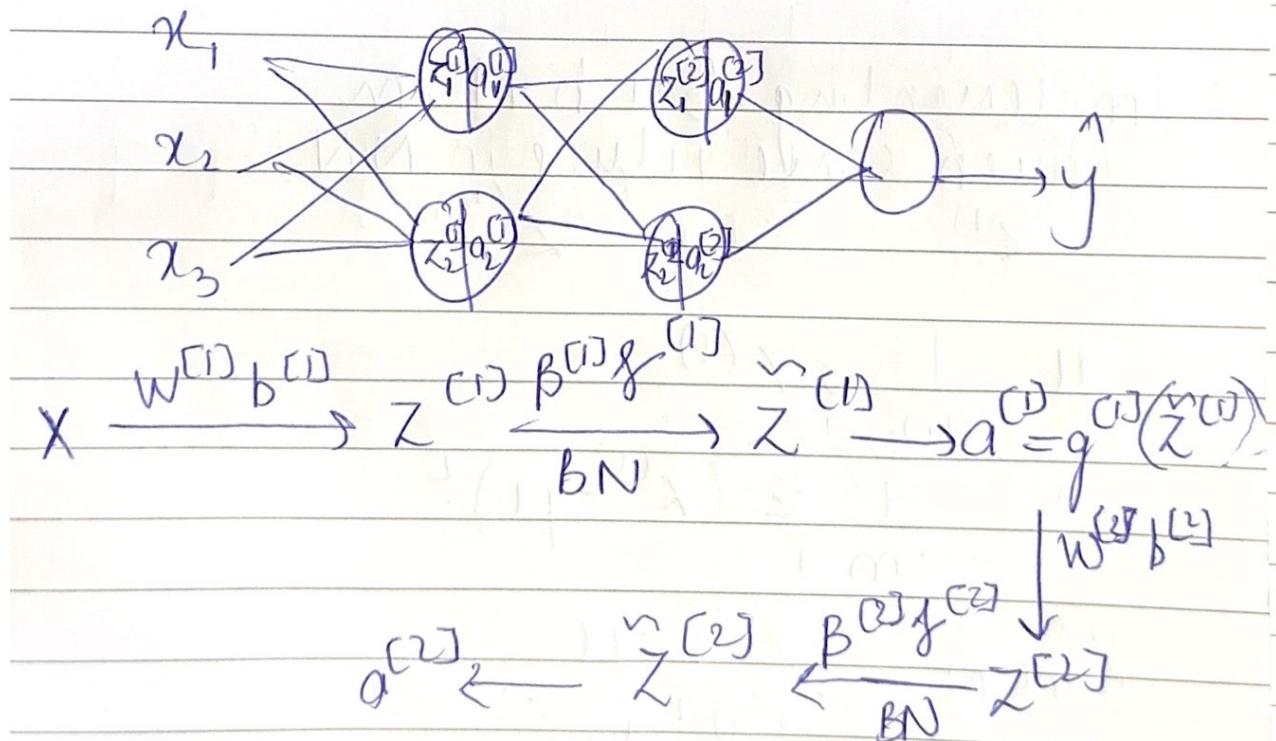
$$\tilde{z}^{(i)} = f \underset{\text{norm}}{\overset{\downarrow}{z}}^{(i)} + \beta$$

learnable
parameters
of model

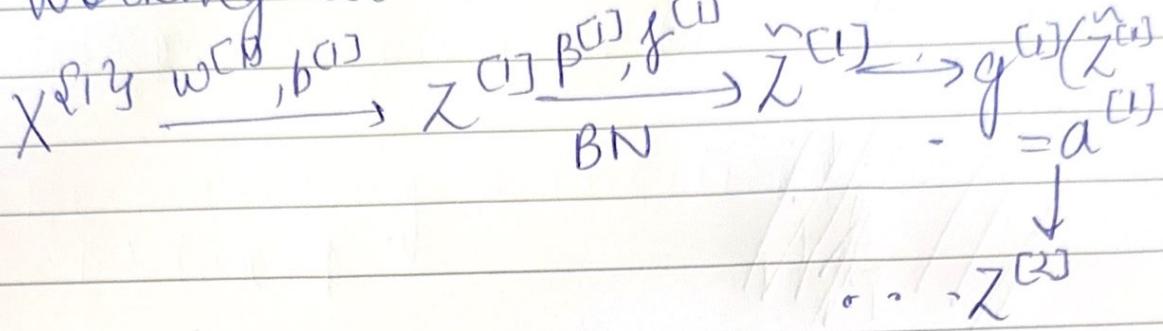
if $f = \sqrt{\sigma^2 + \epsilon}$
and $\beta = \mu$
then

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

→ Fitting Batch Norm into a Neural Network



- Working with mini batches



- Implementing gradient descent
for $t=1 \dots n$ of mini batches
compute forward prop on $X^{(t)}$
in each hidden layer, use BN for
 $Z^{(t)}$ with $Z^{(t)}$

- use backprop and compute $dW^{[l]}$, $d\beta^{[l]}$,
and $df^{[l]}$

update parameters

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

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

$$f^{[l]} := \dots$$

→ Softmax Regression

- It is a type of logistic regression used for multi-class classification
- C : no. of classes
- \hat{y} : 4 -dimensional vector, containing the probabilities of each of the four classes

$$\bullet t = e(z^{\text{eff}})$$

$$a^{\text{eff}} = \frac{z^{\text{eff}}}{\sum t_i}$$

$$a^{\text{eff}} = \frac{e^z}{\sum t_i}$$

$$\sum t_i \\ y_i = 1$$