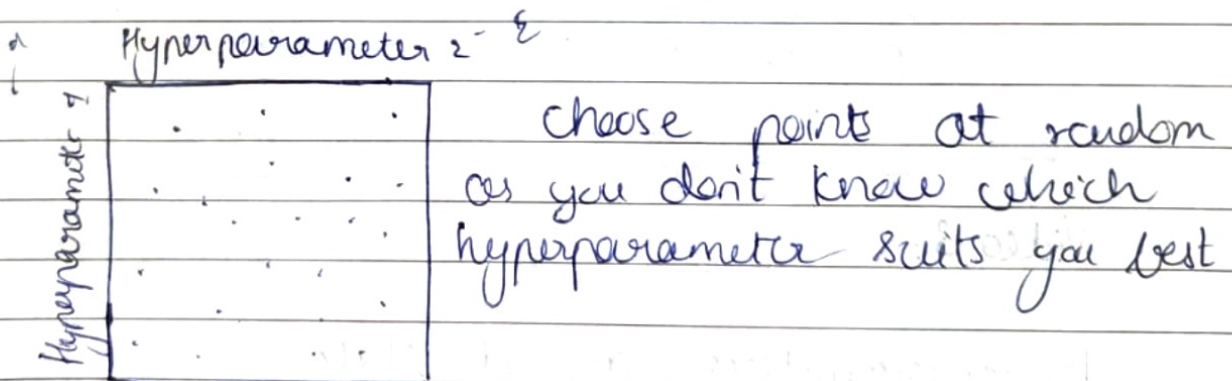
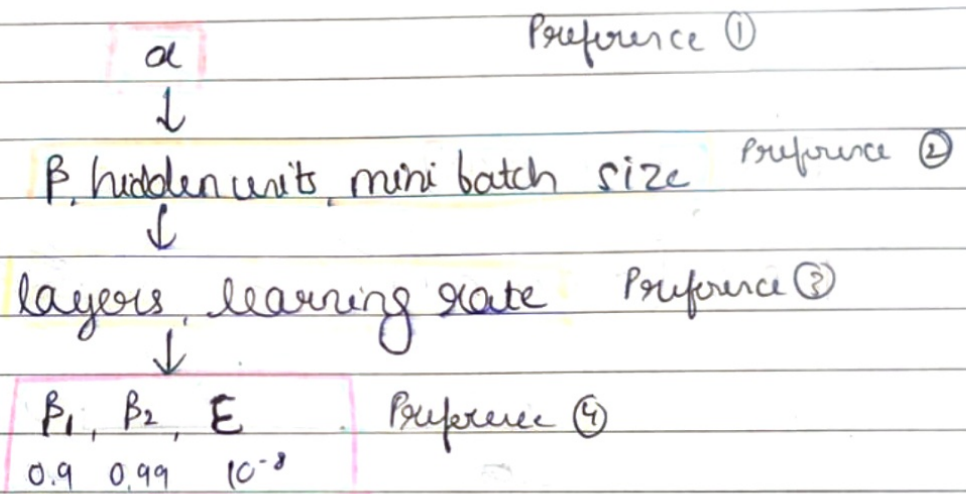


Week 3

Video 1: Tuning Process

Hyperparameters

Preference for tuning & importance



coarse to fine sampling - when you find some hyperparameters values that give you a better performance - zoom into a smaller region around these values and sample more densely in this space

each hidden layers activation

This has regularisation effect.

Batch norm causes the input values to become more stable, so that the later layers of the neural network has more firm ground to stand on. And even though the input distribution changes a bit, it changes less and it allows each layer to learn by itself a little independently of the other layers and this speeds up learning.

Video 4: Batch Norm at Test time

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

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

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

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

- ① We compute mean & Variance of mini batch
- ② In testing we process one example at a time. The mean & Variance of one example won't make sense
- ③ Thus, we compute an estimated value of mean & Variance to use it in testing time

Working with mini batches

In batch normalisation, the mean of $z^{(l)}$ values in the layer gets zeroed out and is replaced by $\beta^{(l)}$ that affects the biased terms.

Algorithm:

- ① For $t = 1$... number of mini batches
- ② compute forward propagation on $x^{(t)}$
- ③ In each hidden layer use BN to replace $z^{(l)}$ with $\tilde{z}^{(l)}$
- ④ use backprop to compute $dw^{(l)}$, $df^{(l)}$, $d\gamma^{(l)}$
- ⑤ update parameters:

$$\begin{aligned} w^{(l)} &= w^{(l)} - \alpha dw^{(l)} \\ \beta^{(l)} &= \beta^{(l)} - \alpha df^{(l)} \\ \gamma^{(l)} &= \gamma^{(l)} - \alpha d\gamma^{(l)} \end{aligned}$$

Video 5: Why does batch norm work

Covariate shift: characteristics of the input data set change from training to test dataset.

Batch norm as regularisation

Each minibatch is scaled by the mean/variance computed on just that minibatch.

This adds some noise to the values $z^{(l)}$ within the minibatch. So similar to dropout it adds noise to

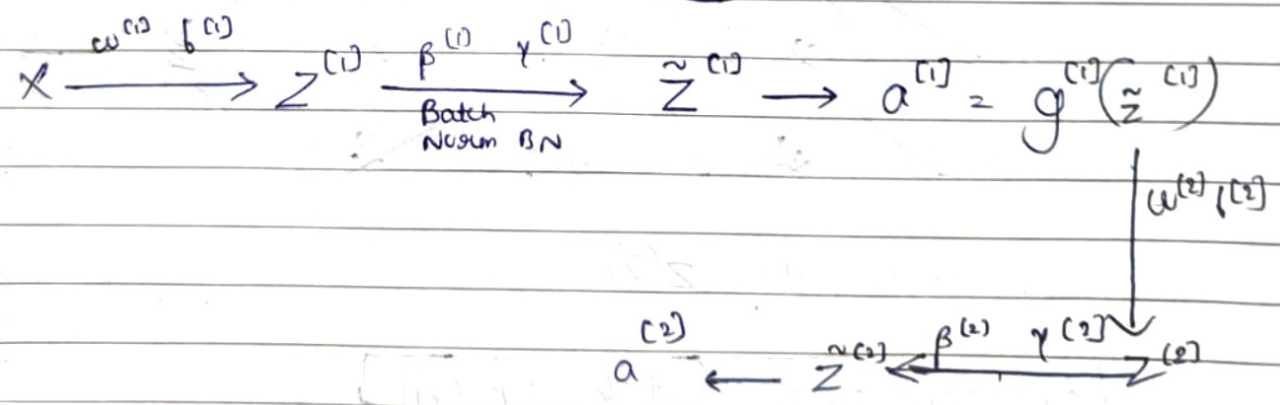
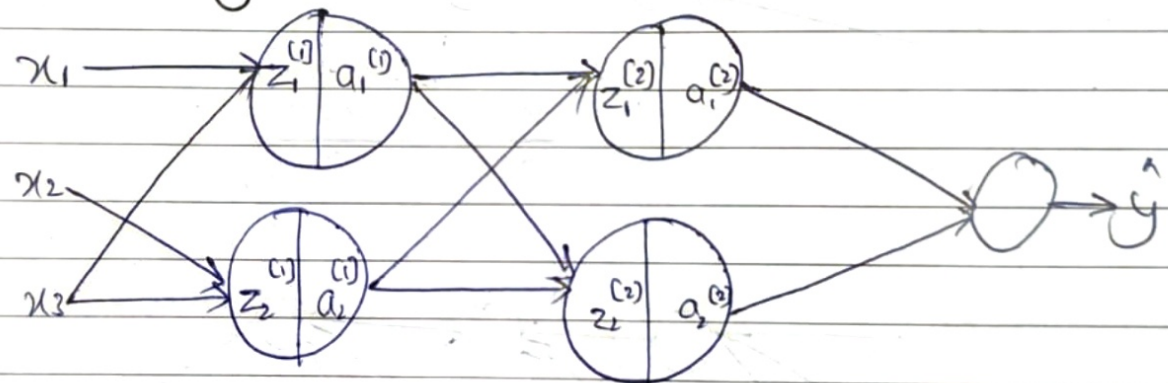
$$z^{(i)}_{\text{norm}} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

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

γ & β would update weights - learnable parameters

If $\beta = \sqrt{\sigma^2 + \epsilon}$ & if $\gamma = 1$ then $\tilde{z}^{(i)} = z^{(i)}_{\text{norm}}$

Video 2: Fitting Batch Norm Into a Neural Network



Parameters:

$$\left\{ \begin{array}{l} \omega^{(1)} \quad b^{(1)} \\ \beta^{(1)} \quad \gamma^{(1)} \end{array} \right\}, \left\{ \begin{array}{l} \omega^{(2)} \quad b^{(2)} \\ \beta^{(2)} \quad \gamma^{(2)} \end{array} \right\}, \dots, \left\{ \begin{array}{l} \omega^{(L)} \quad b^{(L)} \\ \beta^{(L)} \quad \gamma^{(L)} \end{array} \right\}$$

Video 2: Using an appropriate scale to pick hyperparameters

Use logarithmic scale:

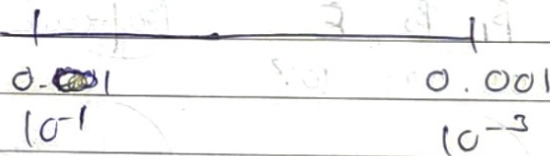
$$\gamma = -4 * \text{np.random.rand}() \leftarrow \gamma \in [-4, 0]$$

$$\alpha = 10^{\gamma} \leftarrow \alpha \in [10^{-4}, 10^0]$$

$$10^{\alpha} \dots 10^6$$

$$\beta \in [0.9, 0.99]$$

$$1-\beta \in [0.001, 0.1]$$



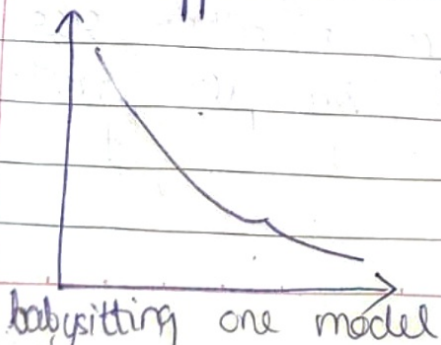
$$\gamma \in [-3, -1]$$

$$\beta = 1 - 10^{\gamma}$$

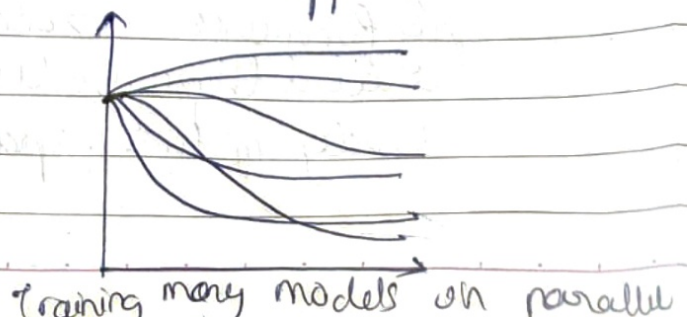
Video 3:

Hyperparameters can get stale
(thus, reevaluate)

Panda approach



Colvia approach

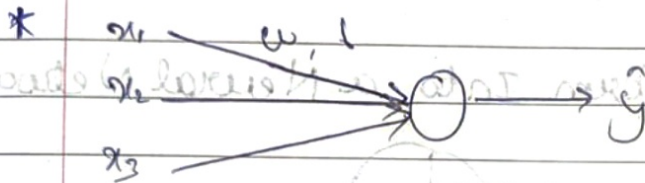


Batch Normalisation

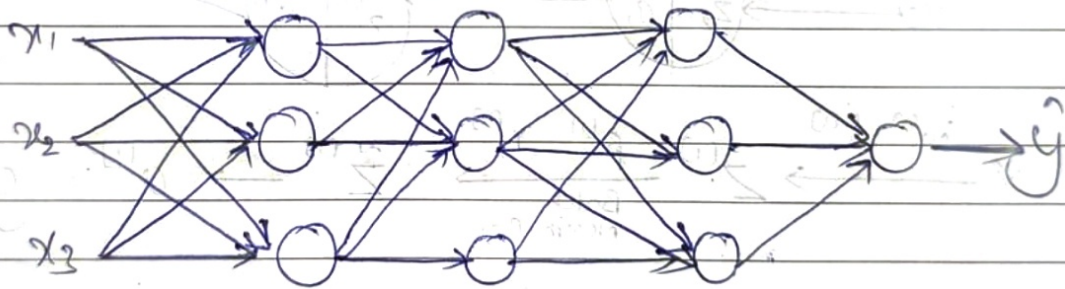
Video 1:

* Advantages of batch normalisation -

- (i) It makes your hyperparameter search problem easier.
- (ii) neural network becomes more robust
- (iii) You have a lot of suitable options for hyperparameters
- (iv) Helps train deep networks.



$$x' = \frac{x - \mu}{\sigma}$$



We normalise $z^{(2)}$

Implementing batch norm

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

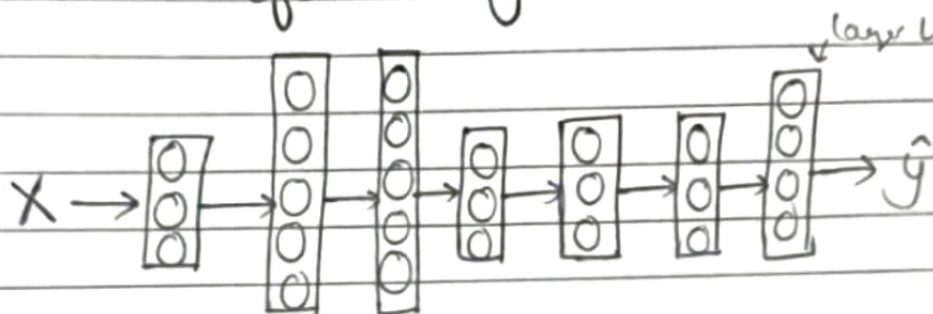
$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

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

- We can use weighted average across mini batches.
- we will use estimated values of mean & variance to Test.

Multi Class Classification

Video 1: Softmax Regression



$$\hat{g} \text{ is } (4, 1)$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

Activation function for softmax

$$t = e^{z^{(L)}} \quad (4, 1)$$

$$a^{(L)} = \frac{e^{z^{(L)}}}{\sum_j t_j}$$

Video 2: Training a Softmax Classifier

Softmax is a generalisation of logistic ^{activation} regression function to C classes. If $C=2$, softmax reduces to logistic regression.

$$L(y, y\text{-hat}) = -\sum_j (y[j] * \log(y\text{-hat}[j]))$$