

* Hyperparameters (Week 3)

$\alpha \approx 0.9$, $B_1, B_2 \approx 0.999$, $\varepsilon = 10^{-8}$, [layers], hidden unit,
learning rate decay, mini-batch size.

- ①
- ② — }
- ③ } Priority order
of hyperparameters
- ④ — }

- It's difficult to know the priority & range of hyperparameters beforehand, thus we'll have to explore.
- By using random search grid, we're ensuring that we're exploring widely as there's more values to try in random.

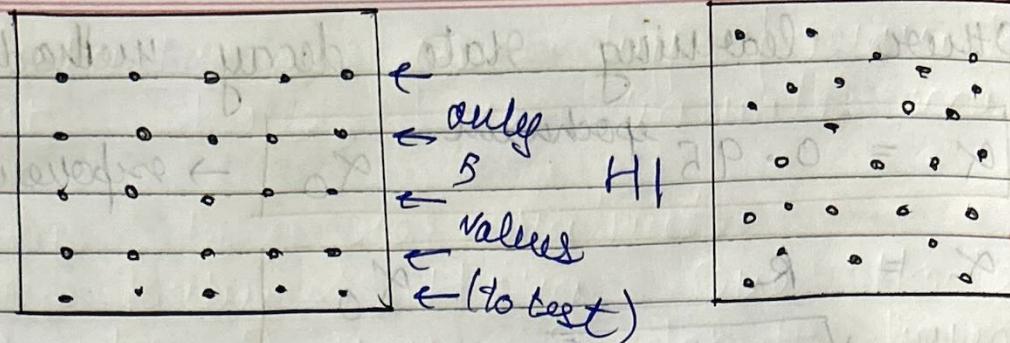
H2

H2

classmate

Date _____
Page _____

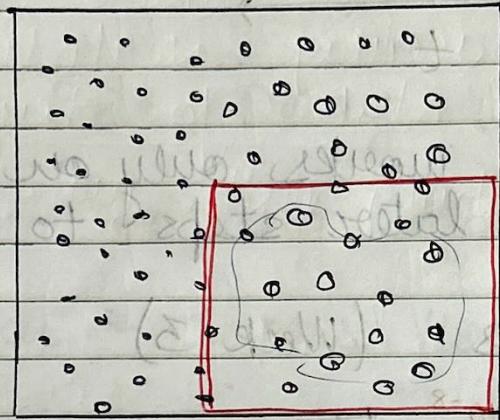
H1



Using Grid

Using Random

Coarse to fine :



after we realize

→ that H1 & H2's optimum values are concised in that area.

* Appropriate Scale for Hyperparameters

$$\alpha = 0.0001 \quad 1$$

number line scale
 $0.0001 \quad 0.1 \quad 1$

↑ only 10% resources here; unfair.

~~log scale~~
 $0.0001 \quad 0.001 \quad 0.01 \quad 0.1 \quad 1$

thus, dividing resources this way is giving uniform random noise

$$\eta = -4 * \text{np.random.rand}()$$

$$\eta \in [-4, 0]$$

$$\alpha = 10^{\eta} \leftarrow 10^{-4}, \dots, 10^0$$

Logic :

$$10^a \dots 10^b$$

$$\text{here } a = \log_{10} 0.0001 \\ = -4$$

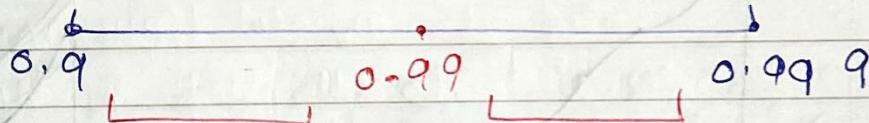
$$b = \log_{10} 1 \\ = 0$$

* Hyperparameters for exponentially weighted averages:

$$\beta = 0.9 \dots 0.999$$

\downarrow

\downarrow



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

allows
dense
exploring
to tune
properly

0.1 0.01 0.001

10^{-1} 10^{-2} 10^{-3}

resources are equally divided

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

$$1 - \beta = 10^\eta \Rightarrow \beta = 1 - 10^\eta.$$

- We're using logarithmic scale over linear scale to distribute resources more efficiently, & because w.r.t. 0.1 to 0.001 $1 - \beta$ is highly sensitive to small changes as well.

- Intuitions for hyperparameters do get stale, re-evaluate occasionally.

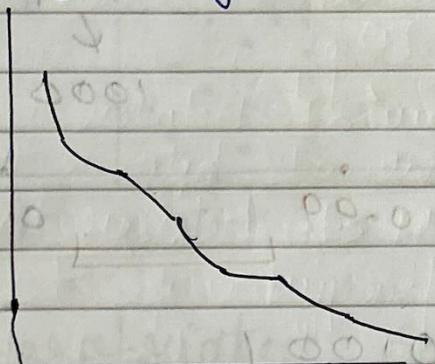
Ways to re-evaluate:

1. Baby-sitting One Model
→ when you don't have GPUs
high computation

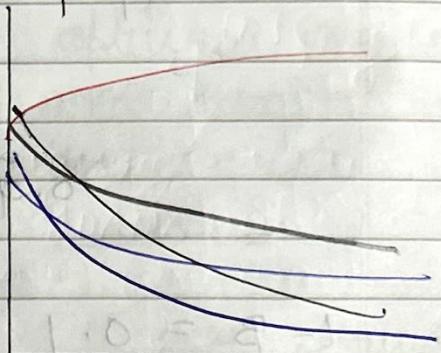
2. Training many models in parallel.

which one to choose depends on computation power available & the amount of data. Cavier.

Huge, huge, then panda.



One model
(Panda)



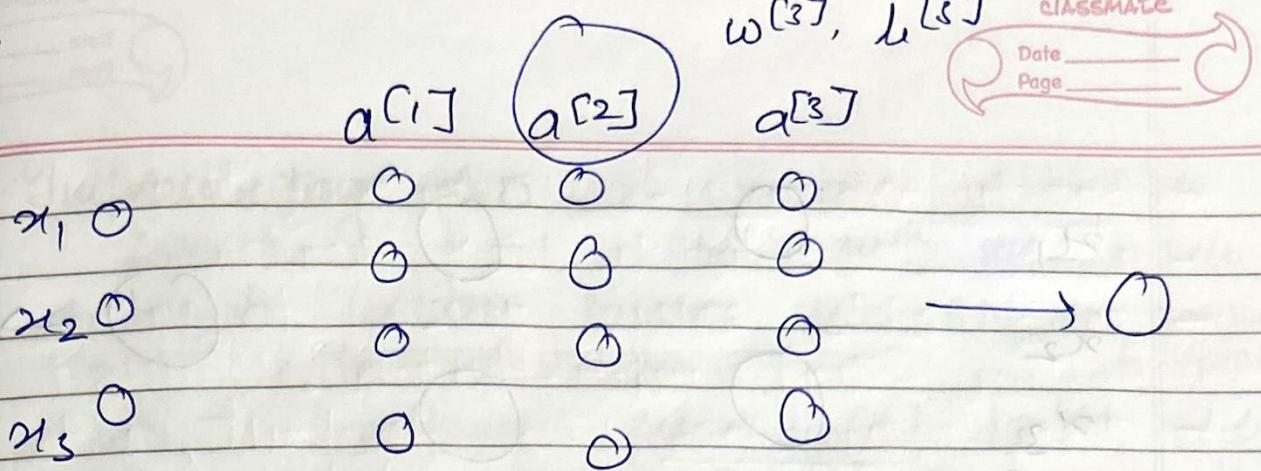
Many Models
(cavier)

* Batch Normalisation

$$\mu = \frac{1}{m} \sum x^{(i)} \quad x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum (x^{(i)})^2 \quad x = x / \sigma^2$$





here, if we normalise $a^{[2]}$ so as to compute $w^{[3]}, b^{[3]}$ faster

Further, normalising $z^{[2]}$ is an even better option.

• Implementing Batch Norm:
applies normalisation to deeper units

Given some intermediate value in NN:

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

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

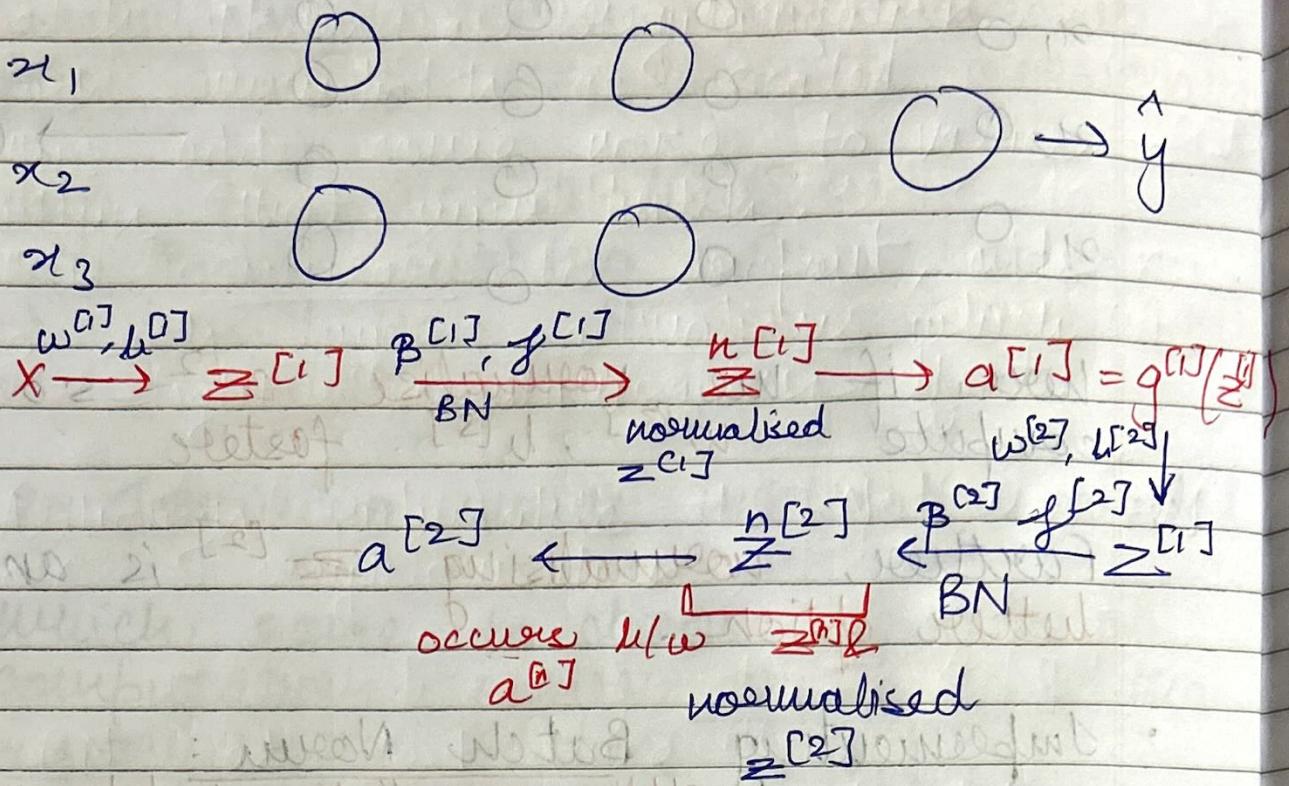
$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

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

$$\gamma = \sqrt{\sigma^2 + \epsilon} \quad \beta = \mu$$

$$z^{[\alpha](i)} \gg z^{[\alpha](i)}$$

mean & variance, are controlled by γ & β normalised throughout the network.



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

- not related to momentum
- not related to weighted averages.
- not related to RMSprop & ADAMs
- not a hyperparameter

Now, can apply OGP :

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

- BN is applied on mini-batches (n_{avg}, m) , (n_{avg}, m) , (n_{avg}, m)
- Parameters : $w^{[l]}, b^{[l]}, \beta^{[l]}, \gamma^{[l]}$
 can be eliminated, as in the mean subtraction step, BN zeroes out the mean of $b^{[l]}$

• Implementing GD:

for $t = 1$ to n^{th} mini batch
 In each hidden layer, use BN to normalise & replace
 $z^{[l]}$ with $\hat{z}^{[l]}$
 Use backprop for $d\omega^{[l]}$, $d\beta^{[l]}$, $df^{[l]}$

Update parameters:

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

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

:

* Learning on Shifting Input Distribution:

- BN reduces the problem of input batches changing, it lowers the impact of previous layers on the upcoming one. Allows each layer to learn by itself & thus optimise the code.
- Due to this, the previous layers do not get switched around as they're constrained to mean & variance.

Eg: shifting Input

Training Data: All black cats.

When tested on another-colour cat as test set, will not work well.

Similarly in a NN of 4 layers, let's consider only 3 & 4 layer, then 3rd is dependent on $a^{[2]}$, thus change in input affects the output

* Batch Norm as regularization:

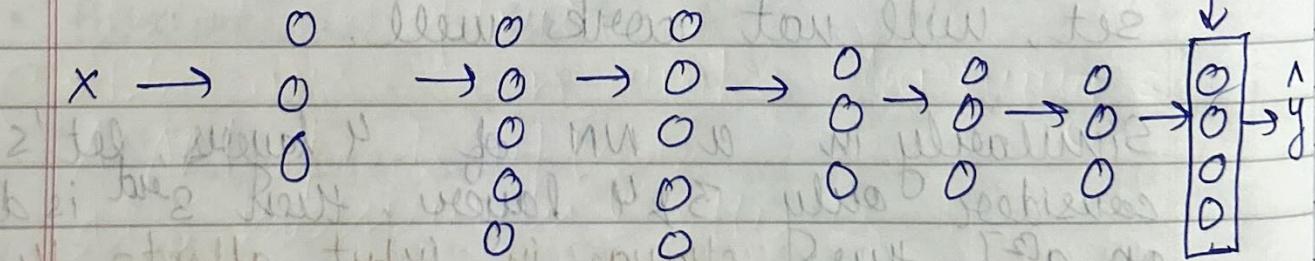
- Each mini batch is scaled
- This adds some noise to values ≥ 0 within that mini-batch.
- Thus, similar to dropout, adds noise to each layer's activations which has a slight regularization effect.
- Although advisable to Not use BN for regularisation.

* BN at test time

- At test time, we do not take an entire mini-batch, rather every example individually
- Thus μ & σ^2 for single value is useless. Thus people usually implement an exponentially weighted average to keep track of μ & σ^2 obtained from train set.

* Multi Class Classification

Softmax Regression layer L



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

Activation funcⁿ:

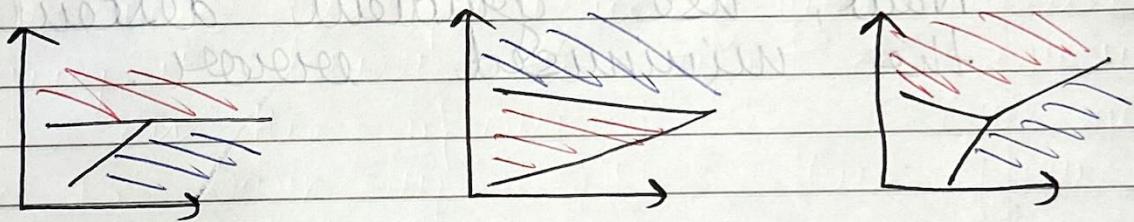
$$t = e^{z^{[L]}} \rightarrow \text{temporary variable}$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum t_i}, \quad a_i^{[L]} = \frac{t_i}{\sum t_i}$$

term to normalise

- Difference from ReLU & tanh as it uses vector input & output.

- More than 2 classes.



- Can give many linear & non linear decision boundary.

Eg:

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

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

$$= \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

Software Regression generalises logistic regression to C classes.

If $C = 2$, Software reduces to logistic regression

* Loss funcⁿ in Software:

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

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

Now, use gradient descent calculate the minimised error.

* The problem of local optima;

- GD does not run well.
- Not all points that appear as local optima are minima, they're actually saddle points
- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow.
- ADAM speeds up w/ also momentum & others