

## • Hyperparameters Tuning

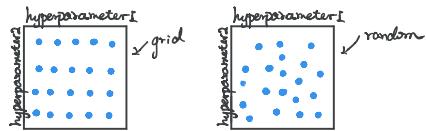
The following will show how to systematically organize hyperparameters tuning process. Depending on the algorithms, different sets of hyperparameters need to be tuned:  $d, \beta_1, \beta_2, \epsilon, \# \text{ layers}, \# \text{ hidden units}, \text{learning rate decay}, \text{mini-batch size etc.}$

most important second important third important rarely change, 0.9, 0.999,  $10^{-8}$

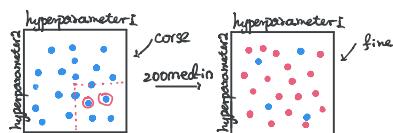
- How to search different values of a range of hyperparameters:

Random sampling:

Suppose we only tune 2 hyperparameters,



Coarse to fine sampling:



Traditionally, we create grid for hyperparameters to try out value systematically, but if we don't know which hyperparameter is more important and we can only try a few different values (e.g. 4 or 5 in 20 models), it's not very cost-effective. If we randomly sample values, we can try more different values in the same amount of models.

If some points in a specific area perform better than the others, we can zoom in that area and sample (randomly) more density.

- How to pick the scale of hyperparameters:

At linear scale:

- # hidden units  $n^{[L]}: 50, \dots, 100$

- # layer L: 2, 3, 4, ...

However, for other hyperparameters, they are more sensitive in one ranges than in other ranges (e.g.  $0.999 \rightarrow 0.9995$  in  $\beta$  has way more impact on the model than  $0.9 \rightarrow 0.905$  in  $\beta$ ). So we can't scale linearly anymore.

At log scale:

- learning rate  $\alpha: 0.001, \dots, 1 \Rightarrow \log \alpha$  ranges from  $\log 10$  to  $\log 10$

then  $r = -b + np.random.rand()$   $\Rightarrow$  sample  $\alpha$  in log scale  
 $\alpha = 10^r$

-  $\beta: 0.9 \dots 0.999$

instead of sampling  $\beta$  directly, we sample  $1-\beta$

$1-\beta: 0.1 \dots 0.001 \Rightarrow$  similar to  $\alpha$ ,  $\log(1-\beta)$  ranges from  $\log 10$  to  $\log 10$ .

then  $r = -b + np.random.rand()$

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

## • Batch Normalisation

Batch normalisation enable hyperparameter search easier and makes the NN much more robust. What batch normalisation does is to normalise  $z^{[l]}$  for any hidden layer to train  $w^{[l]}$  and  $b^{[l]}$  much faster (There are debates about whether we should normalise  $z$  or  $A$  but normalising  $z$  is done much more frequently so here I will present normalising  $z$ ).

How to implement:

Given some intermediate values in NN in layer  $l$ :  $z^{[l+1]}, \dots, z^{[L]}$

$$\left\{ \begin{array}{l} \mu = \frac{1}{m} \sum z^{[l]} \\ \sigma^2 = \frac{1}{m} \sum (z^{[l]} - \mu)^2 \\ z_{\text{norm}}^{[l]} = \frac{z^{[l]} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \hat{z}^{[l]} = \gamma z_{\text{norm}}^{[l]} + \beta \end{array} \right.$$

element-wise multiplication and addition, so different values for different  $z$ s in layer  $l$ .  
trainable parameters:  $\gamma$  and  $\beta$  enable  $z$  to have different mean and variance (other than 0 and 1 in the normalised version), i.e. they control the range of value of  $z$  in the algorithm.  
note: If  $\gamma = \sigma^{-1}$   $\Rightarrow \hat{z}^{[l]} = z^{[l]}$  (inverting the normalisation)  
 $\beta = \mu$

Adding Batch Norm to NN:

$$X \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{BN}]{\gamma^{[1]}, \beta^{[1]}} \hat{z}^{[1]} \xrightarrow{A^{[1]}, g^{[1]}(\hat{z}^{[1]})} z^{[2]} \xrightarrow[\text{BN}]{\gamma^{[2]}, \beta^{[2]}} \hat{z}^{[2]} \xrightarrow{A^{[2]}, g^{[2]}(\hat{z}^{[2]})} \dots \xrightarrow{A^{[L]}, g^{[L]}(\hat{z}^{[L]})} A^{[L]}$$

same dimension as  $z^{[L]}$

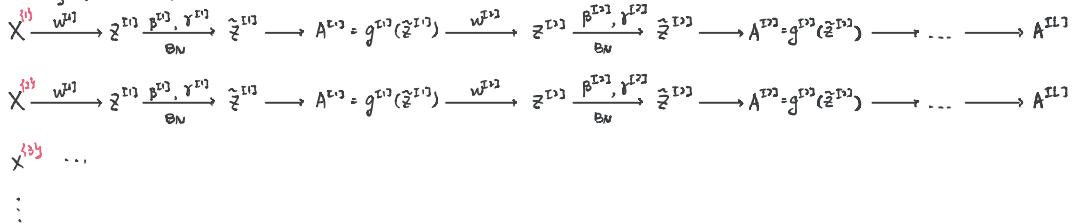
then for each iteration:

- compute forward prop on  $X$  (in each hidden layer, replace  $z^{[l]}$  with  $\hat{z}^{[l]}$ )
- use back prop to compute d $w$ , d $b$ , d $\sigma$
- update  $w$ ,  $b$ ,  $\sigma$  (using gradient descent, momentum, RMS, Adam etc.)

but in practise we will use programming framework so we won't need to implement these details ourselves. (tf, nn, batch-normalization)

note: since in the normalization step, the mean is subtracted and then later added back in, so effectively the parameter  $b$  is cancelled out. Therefore we can omit  $b$ , i.e.  $\hat{z}^{[l]} = w^{[l]} A^{[l]}$ .

Working with mini-batches:



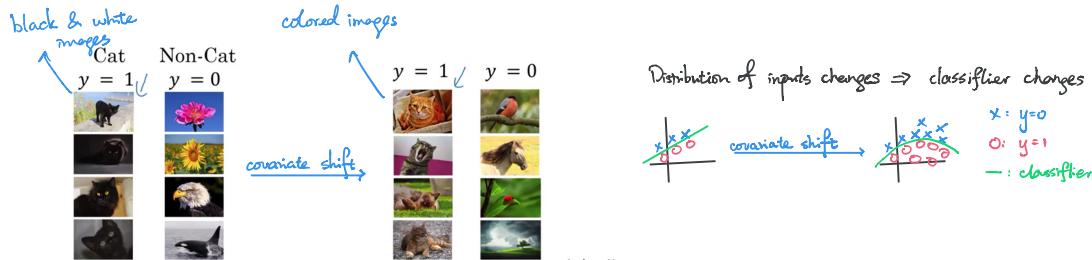
for  $t=1 \dots \text{minibatch\_num}$ :

- compute forward prop on  $X$  (in each hidden layer, replace  $z^{[l]}$  with  $\hat{z}^{[l]}$ )
- use back prop to compute d $w$ , d $b$ , d $\sigma$
- update  $w$ ,  $b$ ,  $\sigma$  (using gradient descent, momentum, RMS, Adam etc.)

Why does Batch Norm work:

- It normalizes inputs and makes the cost function bowl-like, which helps speed up the training.
- It also makes weights in deeper layer of the NN, more robust to changes to weights in earlier layers of the NN, by make each layer learn independently. The following will explain more about it:

If the model has learned how  $X$  maps to  $Y$ , then when the distribution of  $X$  changes, we might need to retrain the algorithm to make correct predictions.



What it means for NN is that, if the input change, parameters in earlier layers change, which also leads to changes in later layers, i.e. from the perspective of the  $l^{\text{th}}$  layer, the hidden unit values change constantly and thus suffering from the problem of covariate shift.

What batch norm does is to reduce the amount that the distribution of these hidden unit values shift around, because parameters  $\beta$  and  $\sigma$  control the mean and variance. So even if the values of hidden units in earlier layers change when parameters update, the mean and variance will remain the same. Thus, these values will become more stable and the later layers won't have to adapt as much if the early layers change, i.e. each layer can learn independently of other layers. This has the effect of speeding up of learning in the whole network.

Batch Norm at test time

Since at test time, we need to process the examples one at a time, so we won't have a mini batch of examples the same size as the training process. So we need a different way of coming up with  $\mu$  and  $\sigma^2$ . In typical implementations of batch norm, what we usually do is to estimate these using an exponentially weighted average where the average is across the mini batches:

suppose we have mini-batches  $X^{(1)}, X^{(2)}, X^{(3)}, \dots$   
for each batch, we compute the mean and variance:  $\mu^{(1)}, \sigma^{(1)}, \mu^{(2)}, \sigma^{(2)}, \dots$

then the exponentially weighted average of these means and variances become the estimates for the mean and variance of  $\tilde{z}^l$  for layer  $l$ , i.e.  $\hat{\mu}^{(l)}$ ,  $\hat{\sigma}^{(l)}$   
 finally,  $\tilde{z}_{\text{norm}} = \frac{\tilde{z} - \hat{\mu}^{(l)}}{\sqrt{\hat{\sigma}^{(l)} + \epsilon}}$  and  $\tilde{z}^{(l)} = T\tilde{z}_{\text{norm}} + \beta$  learned in training process for the test sample.

#### • Multi-class classification

The generalization of logistic regression called Softmax Regression, which makes predictions to recognize one of multiple classes. The Softmax regression has different number of units in the output layer to that of binary regression, which is equal to the number of classes:

$$\text{let } C = \# \text{ classes}, n^{(l)} = C$$

$$\text{Then } A^{(l)} = \begin{bmatrix} a_1^{(l)(m)} \\ a_2^{(l)(m)} \\ \vdots \\ a_C^{(l)(m)} \end{bmatrix}, \text{ which output the probability of being each class, i.e. } A^{(l)} = \begin{bmatrix} \Pr(\text{class 1}) \\ \Pr(\text{class 2}) \\ \vdots \\ \Pr(\text{class } C) \end{bmatrix}$$

$$\text{Since these are probabilities, } \sum_{i=1}^{n^{(l)}} a_i^{(l)(m)} = 1.$$

The Softmax regression also has different activation function for the output layer:

$$\text{Suppose we have } \tilde{z}^{(l)}, \text{ which is } n^{(l)} \times 1$$

$$\text{Then } A^{(l)} = \frac{e^{\tilde{z}^{(l)}}}{e^{\tilde{z}^{(l)}} \text{ summed up vertically}} \rightarrow \text{element-wise division}$$

$$\text{i.e. suppose } \tilde{z}^{(l)(m)} = \begin{bmatrix} \tilde{z}_1^{(l)} \\ \tilde{z}_2^{(l)} \\ \vdots \\ \tilde{z}_{n^{(l)}}^{(l)} \end{bmatrix}, \text{ then } e^{\tilde{z}^{(l)(m)}} = \begin{bmatrix} e^{\tilde{z}_1^{(l)}} \\ e^{\tilde{z}_2^{(l)}} \\ \vdots \\ e^{\tilde{z}_{n^{(l)}}^{(l)}} \end{bmatrix}, \text{ and finally } a^{(l)(m)} = \frac{e^{\tilde{z}^{(l)}}}{\sum_{i=1}^{n^{(l)}} e^{\tilde{z}_i^{(l)}}} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \checkmark \text{ normalise across different outputs so they sum to 1}$$

The name softmax is in contrast to hardmax, which maps probabilities to only 0 and 1, e.g.  $\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ . If  $C = 2$ , softmax regression is reduced to binary regression.

Training the Softmax classifier:

$$\text{loss function: } L(\hat{y}, y) = - \sum_{j=1}^{n^{(l)}} y_j \log \hat{y}_j \rightarrow \text{turns out to be a number because } y_j \log \hat{y}_j \text{ will become 0 in other } n^{(l)-1} \text{ classes, e.g. if } y^{(0)} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, L(\hat{y}^{(0)}, y^{(0)}) = -y_1 \log \hat{y}_1. \text{ Minimising it is equivalent to maximising } \hat{y}.$$

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

$$\text{gradient descent: } dL^{(l)} = \frac{\partial L}{\partial z^{(l)}} = \hat{y} - y \rightarrow \text{in practise, we only need to specify forward prop and the programming framework will compute backprop.}$$

#### • Deep learning frameworks

There are many deep learning frameworks that can help us apply deep learning more easily.

- Caffe/Caffe2
  - CNTK
  - DL4J
  - Keras
  - Lasagne
  - mxnet
  - PaddlePaddle
  - TensorFlow → we will show how to use tensorflow in another notebook.
  - Theano
  - Torch
- Choosing deep learning frameworks
- Ease of programming (development and deployment)
  - Running speed
  - Truly open (open source with good governance)