

Week 3: Hyperparameter Tuning, Batch Normalization

Midclass Programming Frameworks

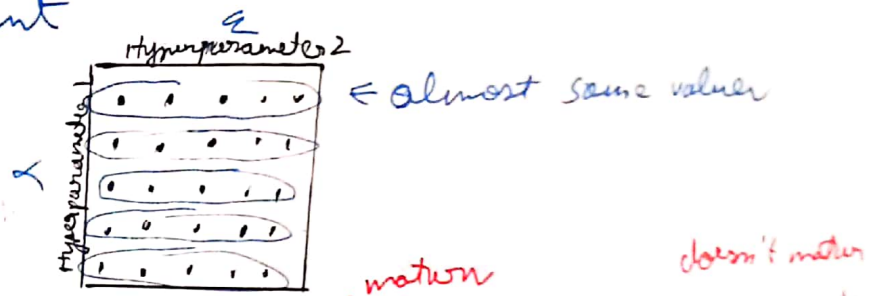
Hyperparameter Tuning

Tuning Process:

- We need to choose α , β , $(\beta_1, \beta_2, \epsilon)$ -adam, # layers, # hidden units, learning rate decay, mini-batch size

- Most important
- Second important
- Third important

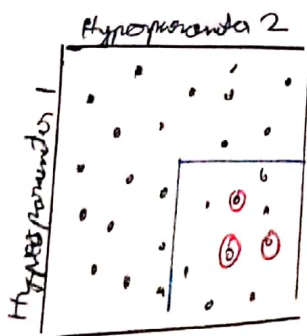
• People used to use a grid, but this is inefficient



If hyperparameter 1 is α and hyperparameter 2 is ϵ , we would try 25 different ~~values~~ of examples, out of which only 5 values of α (even though this is more important). Hence every row will have almost same answer since ϵ doesn't matter much.

Instead using random values is better since that way we ~~check~~ ^{try} 25 values of α .

- Another method is coarse to fine,

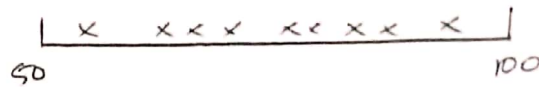


If we are getting good results in a particular area, we can focus our search there.

Using an appropriate scale to pick hyperparameter:

Suppose we want to pick values for

$$h^{[L]} = 50, \dots, 100 \quad - \text{between 50 and 100}$$



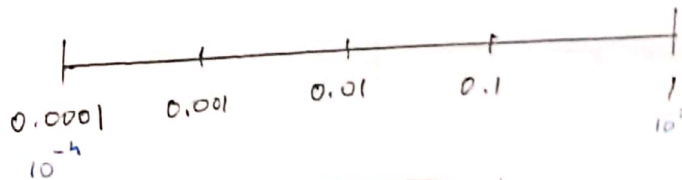
Then we can randomly pick between the range 50 to 100

However consider finding values for an exponential hyperparameter like:

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

If we use the previous method, then 90% of the values will be between 0.1 - 1

\therefore we use log



$$r = -4 + \text{np.random.rand()}$$
$$\alpha = 10^r$$

$$\leftarrow r \in [-4, 0]$$

$$\leftarrow 10^{-4} \dots 10^0$$

For exponentially weighted averages:

$$\beta = 0.9 \dots 0.9999$$

$$\text{So we take } 1 - \beta = 0.1 \dots 0.0001$$

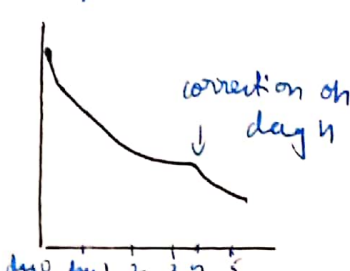
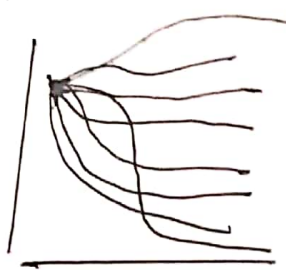
$$\therefore r \in [-1, 1]$$

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

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

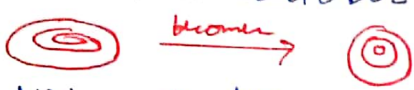
Hyperparameters tuning in practice:

Pandas VS Carver

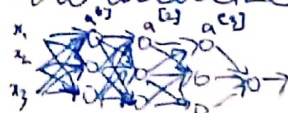
<u>Panda</u>	<u>Carver</u>
<ul style="list-style-type: none">• Pandas only have one (or two) babies and give them full attention  <ul style="list-style-type: none">• Babysitting one model (setting a hyperparameter value. If it doesn't work change during training)• Used when dataset is very large	<ul style="list-style-type: none">• Carver lay 1000s of eggs and let see which survive  <ul style="list-style-type: none">• Training many models in parallel• Used when good CPU/GPU is available for multiprocessing

Batch Normalization

Normalizing activations in a network

- We previously saw that normalising the input ~~values~~ ^{features} x helps gradient descent learn faster since the contours become more uniform \rightarrow 

- However in a NN, we have other layers, and we can normalize them too!



- We normalise z and not a

Let the intermediate values in NN be,

$$z^{(1)}, \dots, z^{(m)} \\ \text{or } z^{[i]}$$

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

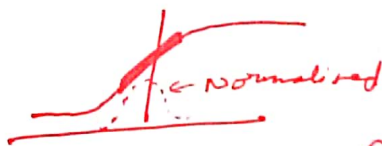
$$2. \sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

$$3. z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$4. \tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

we pass this to activation function

Here γ and β are learnable parameters
Suppose we want to use sigmoid function
that can't use normalized values



In such a case we use $\tilde{z}^{[i]}$ instead of $z^{[i]}$

If $\gamma = \sqrt{\sigma^2 + \epsilon}$ & $\beta = \mu$, it learns this

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

$$\therefore \tilde{z}_{\text{norm}}^{(i)} = z^{(i)}$$

Fitting Batch Norm into a NN

$$x \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{Batch Norm}(\gamma, \beta)]{\gamma^{[1]}, \beta^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \dots$$

Place in the middle, before passing z through the activation function, normalise it.

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

Usually all this is done by inbuilt code in tensorflow.

• working with mini-batches

$$x^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow g^{[1]}(\tilde{z}^{[1]}) = a^{[1]} \rightarrow \dots$$

$$x^{[2]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow \dots$$



here we compute over the batch size and not over n

• b doesn't matter, you can set it to 0.

Since $z^{[L]} = w^{[L]} a^{[L-1]} + \underbrace{b^{[L]}}_{\leftarrow \text{this is automatically removed when we zero out the mean } z^{(i)} - \mu}$

Implementing gradient descent:

for $t = 1 \dots \text{num-min-batches}$

compute forward prop on $x^{[t]}$

In each hidden layer, use BN to replace $z^{[L]}$ with $\tilde{z}^{[L]}$

Use backprop to compute $dw^{[L]}, d\beta^{[L]}, d\gamma^{[L]}, df^{[L]}$

Update parameters

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

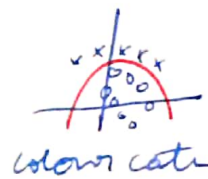
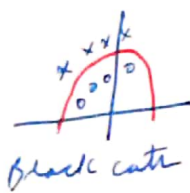
$$\beta^{[L]} := \beta^{[L]} - \alpha d\beta^{[L]}$$

$$\gamma^{[L]} := \gamma^{[L]} - \alpha d\gamma^{[L]}$$

(works with moments, RMSprop or Adam).

Why does Batch Norm work?

1. Like how normalising the input speeds up learning, normalising hidden features also speeds up learning
2. Covariate shift:
Let's say we train on black cat images but test on other colour cats, then it'll not recognise.



even though function is same it won't work because the distribution of x has changed. But if we use Batch Norm, this doesn't happen because it standardizes the data at every layer (the mean and variance remain same)

Batch Norm also has a regularization effect (slight):

- Happens since each mini-batch is scaled by the mean / variance computed on just that mini-batch.
- This adds some noise to $z^{[L]}$ within that mini-batch. So similar to dropout, it adds some noise to each hidden layer's activations.
- Don't use this for regularization - its effect is very slight. As the mini-batch size increases, this effect decreases.

Batch Norm at Test Time:

During testing we process an individual training example, and so we don't have a batch size. As a result, how do we find μ and σ^2 .

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

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

What we do is, during training we use exponentially weighted average to get a μ and σ^2 till as how much ever examples have been trained

$$\begin{array}{ccccc} x^{[1]} & , & x^{[2]} & , & x^{[3]} \dots \\ \uparrow & & \uparrow & & \uparrow \\ \mu^{[1]} & & \mu^{[2]} & & \mu^{[3]} \end{array}$$

This is similar to $\theta_1, \theta_2, \theta_3 \dots$

$$1. \mu_1 = \beta \mu_0 + (1 - \beta) \theta_1$$

$$2. \mu_2 = \beta \mu_1 + (1 - \beta) \theta_2 \dots$$

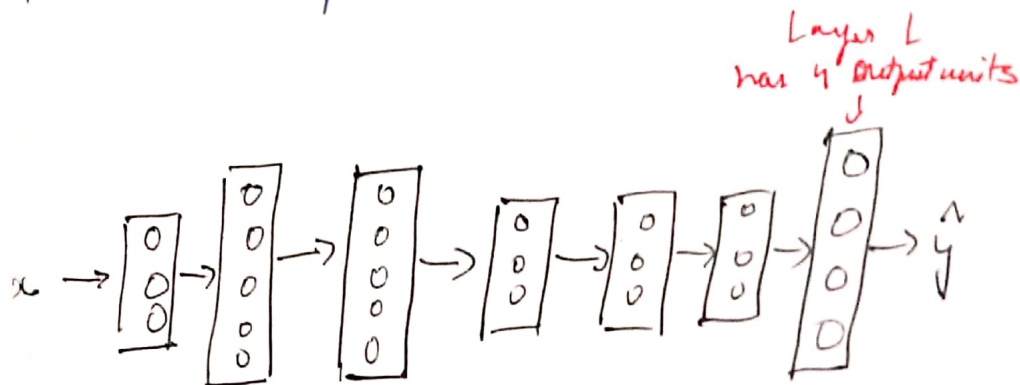
If we ~~test~~ now we take μ_2 . So we keep a track of μ and σ^2 while training.

Multiclass classification - softmax

Before we were doing binary classification (either 0 or 1), now we will use the softmax function for multiclass classification

no. of classes (ex: cat, dog, baby chick, horse)
2
3
4

For example, let $L = 4$



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

Activation function

$$\rightarrow t = e^{(z^{[L]})} \leftarrow \text{element wise power} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \quad \text{if } z = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$\rightarrow a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^n t_i} \quad \therefore a_i = \frac{t_i}{\sum_{i=1}^n t_i}$$

Here $a^{[L]}$ will be (4,1) dimension

unlike ReLU or sigmoid where it is (1,1)

Ex: $z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$, $t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$ $\sum_{i=1}^4 t_i = 176.3$ \rightarrow adding all up

$$\therefore a^{[L]} = \frac{t}{176.3} \quad a = \begin{bmatrix} \frac{e^5}{176.3} \\ \frac{e^2}{176.2} \\ \frac{e^{-1}}{176.2} \\ \frac{e^3}{176.2} \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \leftarrow \begin{matrix} 84.2\% \\ 4.2\% \\ 0.2\% \\ 11.4\% \end{matrix}$$

Training a softmax classifier:

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

Hardmax - it takes the highest element and makes it one while others are 0. $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

So softmax is a 'softer' version that maps to probabilities. $\begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$

→ Softmax regression generalises logistic regression to C classes.

i.e. if $C=2$, softmax reduces to logistic

regression ($a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$ ← logistic regression output & ignore this)

Loss function:

Suppose $y^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ ← cat $a^{[L]}(1) = \hat{y}^{(1)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$ ← \hat{y}_2 needs to be big

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^4 y_i \log \hat{y}_i$$

← here since $y_1 = y_3 = y_4 = 0$, we get $-y_2 \log \hat{y}_2$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 1 & 0 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix}$$

(n,m)

$$\hat{Y} = \begin{bmatrix} \hat{y}^{(1)} & \hat{y}^{(2)} & \dots & \hat{y}^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} 0.3 & \dots & \dots & \dots \\ 0.2 & \dots & \dots & \dots \\ 0.1 & \dots & \dots & \dots \\ 0.4 & \dots & \dots & \dots \end{bmatrix}$$

(n,m)

Hence if loss should be small, \hat{y}_2 needs to be big. Since \hat{y}_2 will be learnt to be large, it'll predict cat properly as it trains.

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

backprop,

$$dz^{[L]} = \hat{y} - y \quad \therefore dz^{[0]} \text{ will be } (y, 1)$$

But we only need to do forward prop properly, the framework (tensorflow) will figure out backprop.

Introduction to Programming Frameworks

Deep Learning Frameworks:

Rather than building everything from scratch, it is more practical and efficient to use libraries:

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- Paddle Paddle
- TensorFlow
- Theano
- Torch

Choosing a deep learning framework:

- Ease of programming (development & deployment)
- Running speed
- Truly open (open source with good governance)

Tensorflow

Program to find w for $J(w) = w^2 - 10w + 25$
loss
(Answer should be 5):

```
import numpy as np
import tensorflow as tf
```

```
coefficients = np.array([[1], [-10], [25]])
```

```
w = tf.Variable([0], dtype=tf.float32)
```

```
x = tf.placeholder(tf.float32, [3, 1])
```

```
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]
```

```
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```
init = tf.global_variables_initializer()
```

```
with tf.Session() as session:
```

```
    session.run(init)
```

```
    print(session.run(w))
```

```
    for i in range(1000):
```

```
        session.run(train, feed_dict={x: coefficients})
```

```
    print(session.run(w))
```

