

Playlist-2

Divide Data in

Train	Dev	Test
-------	-----	------

~~Di~~ Bias \propto train data
variance \propto dev/test data

High Bias \rightarrow Big network } try
(training data problem) Train longer
(try diff NN architecture)

High variance \rightarrow More data } try
(dev set performance) Regularisation
(try diff NN architecture)

\Rightarrow Regularisation

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

Regularisation parameter λ
doesn't contribute much

\mathcal{L}_2 regularisation $\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$ euclidean norm

\mathcal{L}_1 regularisation we use, $\frac{\lambda}{2n} \|w\|_1 = \frac{\lambda}{2n} \sum_{j=1}^n |w_j|$

This makes w more sparse, it may take less memory.

For NN,

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{2m} \sum_{i=1}^m d(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l)}} (w_{ij}^{(l)})^2, \quad w: (n^{(l-1)} \times n^{(l)})$$

↑ "Frobenius Norm" ← *error diff*
(weight diff)

Now,

$$dw^{(l)} = dJ \leftarrow \text{error diff}$$

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

$$w^{(l)} \rightarrow 0$$

weight decay $w^{(l)} := w^{(l)} - \frac{\lambda}{m} dw^{(l)} = \alpha (\text{back prop})$

weight decay is also called L2 regularisation

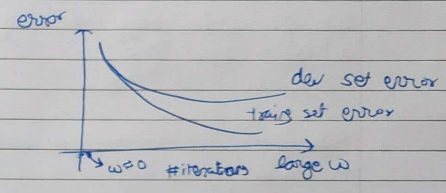
% $\uparrow w^{(l)} \downarrow$
 $\therefore z^{(l)} \downarrow \Rightarrow z^{(l)} = w^{(l)} x + b^{(l)}$
So ~~there~~ it will become linear
and not have good boundaries

Dropout Regularisation

Train small-small parts of NN, based on probability we train some & choose to not train some.

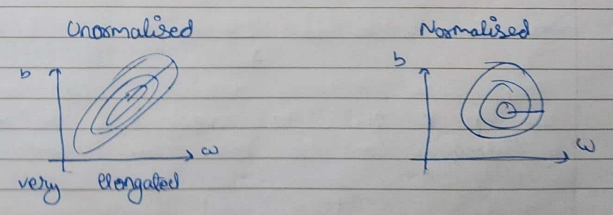
Downside: J is not well defined.

Early Stopping



Normalization

It is scaling the data to increase the variance if the scale of the inputs is not similar



For NN,

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{2m} \sum_{i=1}^m \alpha (\hat{y}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l)}} (w_{ij}^{(l)})^2, \quad w: (n^{(l-1)} \times n^{(l)})$$

"Frobenius Norm" ← $\|w\|_F$
 (adding each element)
 (addition of each row)

Now,

$$dw^{(l)} = \frac{dJ}{dw^{(l)}} \leftarrow \text{error diff}$$

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

$$w^{(l)} \rightarrow w^{(l)} - \alpha dw^{(l)}$$

weight decay $w^{(l)} := w^{(l)} - \frac{\lambda}{m} dw^{(l)} = \alpha (\text{back Prop})$

weight decay is also called L2 regularisation

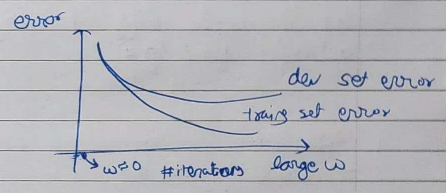
If $\alpha \uparrow$ $w^{(l)} \downarrow$
 $z^{(l)} \downarrow$ as $z^{(l)} = w^{(l)} x + b^{(l)}$
 So $\alpha \uparrow$ it will become linear
 and not have good boundaries

Dropout Regularisation

Train small-small parts of NN, based on probability we train some & choose to not train some.

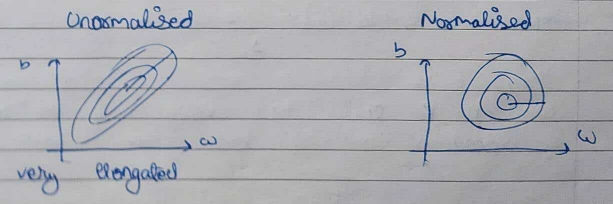
Downside: J is not well defined.

Early Stopping



Normalization

It is scaling the data to increase the variance if the scale of the inputs is not similar.



→ weight initialization

$$Z = w_1 x_1 + w_2 x_2 \dots + w_n x_n + b$$

↑ we want 'w' to be small so that z is smaller.

$$w^{(0)} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}\left(\frac{1}{n^{(l-1)}}\right)$$

$$\text{Var}(w) = \frac{1}{n} \quad \text{for ReLU} \quad \text{Var}(w) = \frac{2}{n}$$

$$\text{for } b \quad \text{np.sqrt}\left(\frac{2}{n^{(l-1)}}\right)$$

So we need diff. ini. of weight for diff. activation func. $g(z)$

→ Gradient Checking

$$J(\theta) = J(\theta_1, \theta_2, \dots)$$

for each i :

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

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

Keep $\epsilon \approx 10^{-7}$

we check this for every i , if the magnitude of $d\theta_{\text{approx}}$ and $d\theta$ have same dimensions and are almost same, we are correct.

To check, $\frac{|d\theta_{\text{approx}} - d\theta|^2}{(|d\theta_{\text{approx}}|^2 + |d\theta|^2)} \approx \frac{10^{-7}}{10^{-5} + 10^{-5}} \approx 10^{-3}$ great

→ Gradient Checking implementation notes

- don't use this in training - only to debug if it is very slow
- If algo fails grad check, look at all components (w, b, \dots) to try to identify bug.
- Remember Regularisation

$$\text{eg: if } J(\theta) = \frac{1}{n} \sum L(y, \hat{y}) + \frac{\lambda}{2m} \sum ||w^{(i)}||_2^2$$

$d\theta$ = gradient of J wrt θ , also include this term

- It doesn't work with dropout
- Run random ini., perhaps after training, just to be sure that, initially $w, b \neq 0$

WEEK-2

→ Mini-Batch Grad. descent / stochastic gradient descent

Good for $m \geq 2000$

It is very helpful for large data set, so here we divide data into mini-batches, for each step of grad. descent we need to process data, so it makes sense if we already start grad. descent of mini-batches of processed data.

eg:

$$(n \times m) X = \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \end{bmatrix}$$

$$(1 \times m) y = [y^{(1)}, \dots, y^{(m)}]$$

for $m \uparrow$ eg: 5000000, it makes sense to make a batch of let say 1000

$$x^{(1)} \dots x^{(1000)} = X^{(1)}$$

$X^{(1)}$ is for mini batch

∴ for $t=1$ to 5000 :

Forward prop on $x^{(t)}$

$$z^{(1)} = w^{(1)} x^{(t)} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

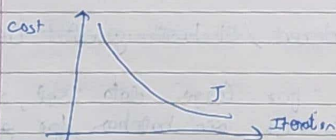
vertical implementation
(1000 g)

compute cost $J^{(t)} = \frac{1}{1000} \sum \frac{1}{2} (y - \hat{y})^2 + \frac{\lambda}{2C(1000)} \sum ||w^{(a)}||_F^2$ (i.e. $\lambda \neq 0$)

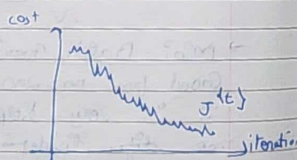
Backprop to compute gradients cost $J^{(t)}$ using $(x^{(t)}, y^{(t)})$

$$w^{(2)} := w^{(2)} - \eta d w^{(2)}$$

Batch grad. des.



mini Batch



If mini batch size = 1000 then it is batch,

If size = 1, it is called stochastic grad. descent

Recommend to keep size in 2nd net

→ Exponentially weighted average

eg: $V_t = \beta V_{t-1} + (1-\beta) \theta_t \rightarrow$ It is an exponential function

~~is a dynamic~~

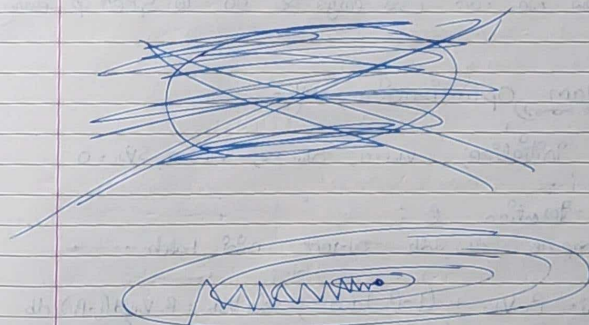
turns out $V_t \approx \frac{1}{1-\beta}$ days avg temp

for $\beta = 0.98$

$V_t = 50$ days avg temp

Bias Correction is important as in early period the function doesn't work as expected, but although later it becomes correct.

→ Gradient Descent with momentum



more ↓ slower the learning rate
more ↔ faster

→ Implementation

On iteration t:

Compute dw, db on current mini batch

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

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

$$w = w - \alpha V_{dw}, \quad b = b - \alpha V_{db}$$

(default $\beta = 0.9$)

Remember

→ RMSprop (great mean sq. propagation)
For it assume $b \downarrow$ & $w \leftrightarrow$
On iteration t :

Compute dw, db on current mini batch

$$S_{dw} = \beta S_{dw} + (1-\beta) \frac{dw^2}{\text{small}}, \quad S_{db} = \beta S_{db} + (1-\beta) \frac{db^2}{\text{large}}$$

element wise

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}, \quad b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

just to make sure denominator $\neq 0$ ϵ is very small i.e. 10^{-8}

∴ Now we'll get large w (i.e. \leftrightarrow) & small b (i.e. \rightarrow)
Now we can use large α so to speed up learning

→ Adam Optimization Algo

We initialize, $V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$

On iteration t :

Compute dw, db current mini batch

momentum

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

RMSprop

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

Corrected

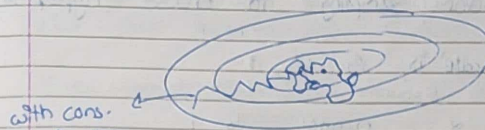
$$V_{dw} = V_{dw} / (1-\beta_1^t), \quad V_{db} = V_{db} / (1-\beta_1^t)$$

$$S_{dw} = S_{dw} / (1-\beta_2^t), \quad S_{db} = S_{db} / (1-\beta_2^t)$$

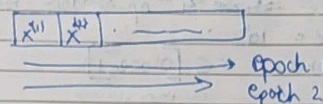
$$w := w - \alpha \frac{V_{dw}}{\sqrt{S_{dw} + \epsilon}}, \quad d := w - \alpha \frac{V_{db}}{\sqrt{S_{db} + \epsilon}}$$

α - needs to be tune, β_1 - generally 0.9, β_2 - 0.999 (dw^2), ϵ - 10^{-8}

→ Learning rate decay



with cons. & it may never converge



1 epoch = 1 pass through data

$$\alpha = \frac{\alpha_0}{1 + \text{decay rate} * \text{epoch num}}$$

eg. epoch	α	Let $\alpha_0 = 0.2$ decay rate = 1
1	1	
2	0.67	
3	0.5	

Using learning rate decay it will initially take large step then smaller to converge

Other methods

$\alpha = 0.95^{\text{epoch num}}$ (exponential)
or $\alpha = \frac{k}{\sqrt{\text{epoch num}}}$
or discrete staircase α
or decrease α manually

WEEK 3

→ Use random sampling to choose hyperparameters

→ But If the scale is too big,

eg: for α 0.0001 to 1

then we use logarithm for better samples

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

$\alpha = 10^r$

where $r \in [0.0001, 1]$

+ For exponentially weighted values

eg: $\beta = 0.9 \dots$ to 0.999

$1-\beta = 0.01 \dots$ to 0.001

10^{-1} to 10^{-3}

$r \in (-3, -1)$

$1-\beta = 10^r$

$\beta = 1 - 10^r$

Because we use $\frac{1}{1-\beta}$ so when $\beta \rightarrow 1$ it's not

There can be 2 approaches

Babysitting one model \rightarrow like single child full attention making sure it survives

Training many parallel models \rightarrow multiple child, getting the fastest one (more compute)

Batch Normalization

If this applies to your algo, it helps you reduce time for hyperparameters search

Gives some intermediate values in NN $z^{(1)} \dots z^{(n)}$

~~Search for~~

$z^{(1)}, \dots, z^{(n)}$

$$\mu = \frac{1}{n} \sum z^{(i)}$$

$$\sigma^2 = \frac{1}{n} \sum (z_i - \mu)^2$$

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

we normalize to $z^{(i)}$ to have mean = 0, & var = 1

But we don't want all z to have it,

so,

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

learnable parameters of model

Now we'll use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

we use to normalize x to $z^{(i)}$

Now we can regulate mean & var by regulating γ & β so we are in control.

$$x \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\text{Batch Norm}(\beta, \gamma)} \tilde{z}^{(1)} \rightarrow a^{(1)} = g^{(1)}(\tilde{z}^{(1)}) \rightarrow z^{(2)}$$

Parameters: $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$
 $\beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}$

In practise BN is usually performed in mini-batches

As mean of $z^{(k)} = 0$, so in its eqⁿ

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

$b^{(k)}$ is eliminated as it gets cancelled to make mean = 0

Parameters: $w^{(k)}$, $\beta^{(k)}$, $\gamma^{(k)}$

Implementing Grad descent.

for $t=1$ to no. of minibatches:

Compute forward prop on $x^{(t)}$

In each hidden layer use BN to replace $z^{(k)}$ with $\tilde{z}^{(k)}$

Use backprop to compute $w^{(k)}$, $\beta^{(k)}$, $\gamma^{(k)}$

Update them:

→ As we normalize x so that there is not much variance in input so it fastens the process, BN does the same this to $z^{(k)}$ as $z^{(k)}$ is kinda input to next hidden layer so it fastens the process

→ Also we know that the output of a layer depends on the output of all layers behind it so changing values of any layer before can make great change in current one. So what batch norm does is it ensure no matter the values change, their mean & variance doesn't change that much making learning process easy.

→ Each minibatch is scaled by mean/variance computed on just that mini batch

→ This adds some noise to values $z^{(k)}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations

→ This has a slight regularisation effect.

By reducing batch size noise increases & so does regularisation effect.

→ Softmax Regression

$C = \# \text{ classes}$

It is used if we want multiple of not just 0 or 1

So here instead of one output unit at end there will be C units

β will be $(C \times 1)$ instead (1×1)

eg: let $C=4$,

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

(4×1)

Activation function:

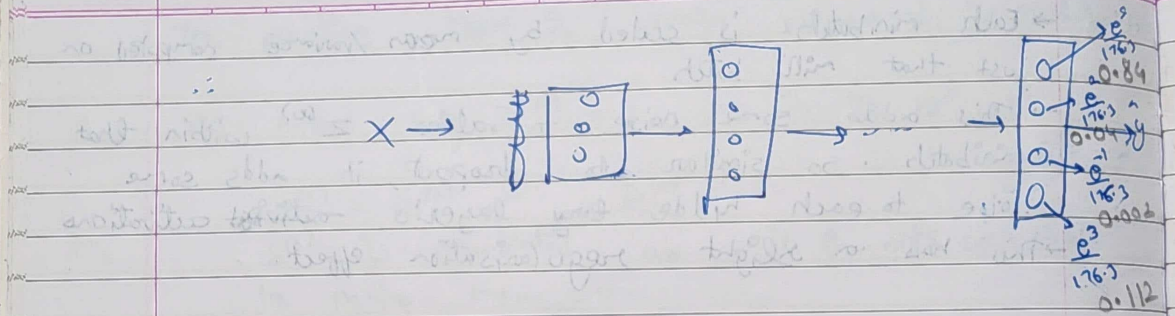
$$\text{temp} = e^{z^{(k)}}$$

$$a^{(k)} = \frac{e^{z^{(k)}}}{\sum_{j=1}^C \text{temp}_j}$$

(4×1)

$$a_i^{(k)} = \frac{t_i}{\sum_{j=1}^C t_j}$$

$$\text{ex, } z^{(k)} = \begin{bmatrix} 3 \\ 2 \\ 1 \\ 0 \end{bmatrix}, t = \begin{bmatrix} 0.5 \\ 0.2 \\ 0.1 \\ 0.2 \end{bmatrix}, \sum_{j=1}^C t_j = 1.0, a_i^{(k)} = \frac{t_i}{1.0}$$



→ $\hat{y} = 1$ (positive region)
 $\hat{y} = 0$ (negative region)

(1×1) matrix $\begin{pmatrix} 1 & -1 \end{pmatrix}$ (1×2) \rightarrow New $\hat{\mu}$