



Neural Networks & Deep Learning

ReLU → Rectified Linear Unit

{(a function that is either 0 or +ve linear)
↳ eliminates -ve values



Densely Connected Layers

all the o/p's from the previous layer are connected to each and every node in the present layer.

Supervised Learning

when the i/p and o/p of an application are clearly stated and the model is trained using i/p and expected o/p's as training data

Structured & **Unstructured Data** → images, text, speech etc.

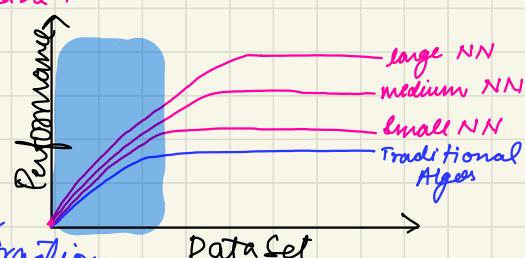
Traditionally computers were unable to process unstructured data but that has changed with the introduction of Neural Networks.
↳ databases with tables

* Neural Nets have also improved structural data understanding

Performance vs Dataset Size.

* NN have become more and more prominent due to the rise in amounts of data as only they can give performance gains

For low amounts of data the relative ordering depends on feature extraction

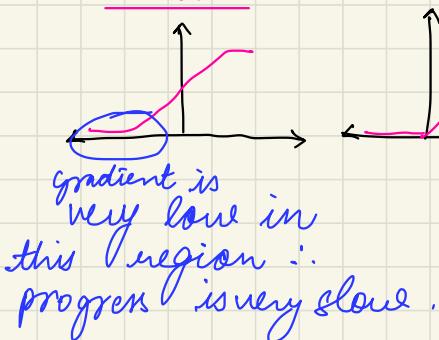


* earlier the advancements in NN were based on Data & Computation advancements but now algorithmic transformations are all the rage.

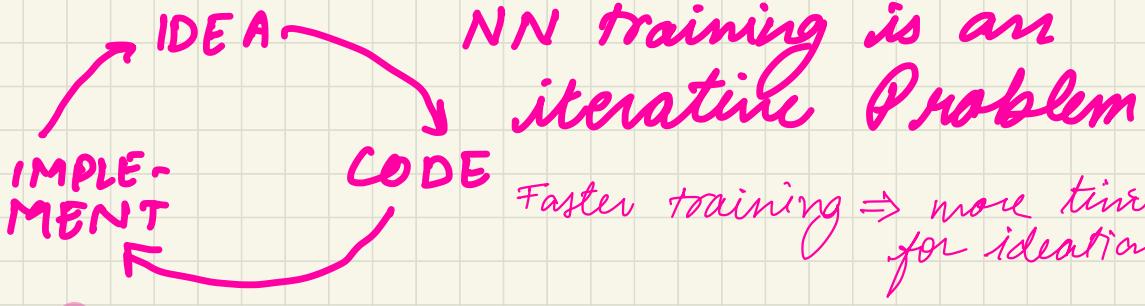
→ latest advancement

Replacement of Sigmoid with ReLU

Reason



gradient is 1
∴ gradient descent works faster



Binary Classification

Image Representation

They are stacked in a column vector X for NN input.

3 channels for pixel values for Red, Green and Blue.

For an image of size $h \times w$, the dimensions are $h \times w \times 3$.

$$X = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{n \times 1}$$

For m training examples, X becomes

$$X = \begin{bmatrix} 1 & 1 & & 1 \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ 1 & 1 & & 1 \end{bmatrix}_{n \times m}$$

$\xrightarrow{\text{i/p to a neural network}}$

O/p is represented as.

$$Y = [y^{(1)} \ y^{(2)} \ y^{(2)} \cdots \ y^{(m)}]_{1 \times m}$$

Logistic Regression

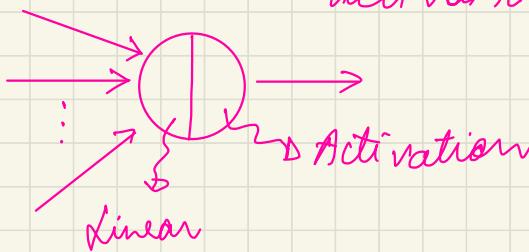
Let $\hat{y} = P(x=1|x)$ (probability of $y=1$)

Let W be a weight matrix of size n where n is the number of features

$$W = \mathbb{R}^{n \times 1}$$

Represents one node of a neural network. $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$ sigmoid function linear relationship

NODE = linear Relationship + activation function



Now, Maximum Likelihood estimation is

$$= \prod_{i=0}^C P(y=i|x) \quad (C = \# \text{ of classes})$$

For Binary Classification

$$\begin{aligned} MLE &= P(y=0|x) \cdot P(y=1|x) \\ &= (1-\hat{y})^{1-y} \cdot (\hat{y})^y \end{aligned}$$

We need to maximize MLE for accurate results.

Taking log on both ends.

$$\log(MLE) = (1-y) \log(1-\hat{y}) + y \log(\hat{y})$$

We can also achieve max likelihood by maximizing $\log(MLE)$ or minimizing

$$[-\log(MLE)]. \text{ So, } -\log(MLE) \text{ can be our cost}$$

$$\alpha(y^{(i)}, \hat{y}^{(i)}) = -[(1-y) \log(1-\hat{y}) + y \log(\hat{y})]$$

over the entire data set

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

Here, $J(w, b)$ is a convex function and hence can be reduced by gradient descent.

Gradient Descent

Aim: Take small steps to reach global min

Repeat { learning rate

$$w_j := w_j - \alpha \frac{\partial J(w, b)}{\partial w_j}$$

→ slope

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

?

$$\frac{\partial J(w, b)}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \ell(g_i^{(i)}, y_i^{(i)})}{\partial w_j}$$

$$\frac{\partial \ell(g_i^{(i)}, y_i^{(i)})}{\partial w_j} = \frac{\partial \ell(g_i^{(i)}, y_i^{(i)})}{\partial \hat{y}_i^{(i)}} \cdot \frac{\partial \hat{y}_i^{(i)}}{\partial w_j}$$

$$\begin{aligned} \frac{\partial \ell(g_i^{(i)}, y_i^{(i)})}{\partial \hat{y}_i^{(i)}} &= -y_i^{(i)} \frac{1}{\hat{y}_i^{(i)}} + (1-y_i^{(i)}) \frac{1}{1-\hat{y}_i^{(i)}} \quad | z = w^T x + b \\ &= \frac{y_i^{(i)} - \hat{y}_i^{(i)}}{\hat{y}_i^{(i)}(1-\hat{y}_i^{(i)})} \end{aligned}$$

$$\frac{\partial \hat{y}_i^{(i)}}{\partial w_j} = \frac{\partial g_i^{(i)}}{\partial z_i^{(i)}} \cdot \frac{\partial z_i^{(i)}}{\partial w_j}$$

$$= \frac{-1}{(1+e^{-z_i^{(i)}})^2} \times e^{-z_i^{(i)}} \times -1 \cdot x_j^{(i)}$$

$$\frac{\partial \hat{y}^{(i)}}{\partial w_j} = \frac{e^{-z^{(i)}}}{(1+e^{-z^{(i)}})^2} x_j = \frac{(1+e^{-z^{(i)}}) - 1}{(1+e^{-z^{(i)}})^2} x_j$$

$$= \frac{1}{1+e^{-z^{(i)}}} \cdot \frac{(1+e^{-z^{(i)}}) - 1}{1+e^{-z^{(i)}}} \cdot x_j$$

$$= \hat{y}^{(i)} \cdot (1 - \hat{y}^{(i)}) x_j$$

$$\frac{\partial \alpha(\hat{y}, y^{(i)})}{\partial w_j} = \frac{\hat{y}^{(i)} - y^{(i)}}{\hat{y}^{(i)}(1 - \hat{y}^{(i)})} \cdot x_j \cdot \cancel{\hat{y}^{(i)}(1 - \hat{y}^{(i)})}$$

$$= [\hat{y}^{(i)} - y^{(i)}] x_j^{(i)}$$

$$\frac{\partial J(w, b)}{\partial w_j} = \underbrace{\sum_{i=1}^m [\hat{y}^{(i)} - y^{(i)}] \cdot x_j^{(i)}}_m$$

VECTORIZATION

$$X = \begin{bmatrix} 1 & | & 1 & | & \dots & | & 1 \\ x^{(1)} & | & x^{(2)} & | & \dots & | & x^{(m)} \\ | & | & | & | & & | & | \end{bmatrix}_{n \times m} \quad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1}$$

$$X = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]_{1 \times m}$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\hat{y} = \sigma(W^T X + b)$$

$$\frac{\Delta J}{W} = X(\hat{y} - y)^T \quad \frac{\Delta J}{B} = \sum (\hat{y} - y)$$

$$W = W - \frac{\alpha}{m} X (\sigma(W^T x + b) - y)^T$$

$$b = b - \frac{\alpha}{m} \sum (\sigma(W^T x + b) - y)$$

Shallow Neural Networks

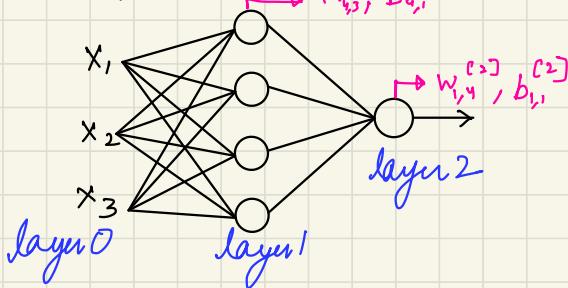
Superscript $[i] \rightarrow$ layer i

example (2 layer NN)

$$W^{[1]}, b^{[1]}$$

$$a^{[0]} = X_{(n,m)}$$

* dimensions of $W^{[i]}$
 $= (i_0/p, i_1/p)$
 for any layer i



* dimensions of $b^{[i]}$
 $= (i_1/p, 1)$
 for any layer i

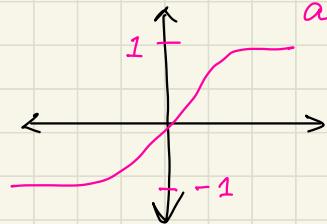
Let $g(z)$ be the activation function for the neural net so,

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} \\ a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = g(z^{[2]})$$

Activation Functions

$$a = \text{Tan}(z)$$

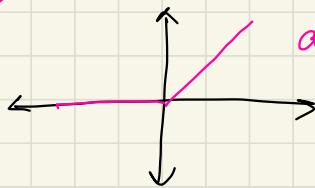


$$a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

(shifted Sigmoid)

can be better for hidden layers as data is equally spread and mean is closer to 0.

$$a = \text{ReLU}(z)$$

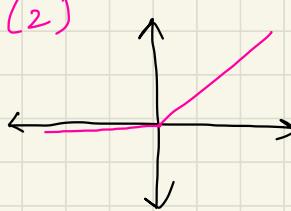


$$a = \max(0, z)$$

learns faster

use when unsure about what to choose or want regression output

$$a = \text{Leaky ReLU}(z)$$



$$a = \max(0.01z, z)$$

slightly better performance than ReLU

* use Sigmoid for last layer only in case of Binary Classification.

Why Linear Functions are not used

Because multiple linear layers can be composed to a single linear layer.

Can be used as o/p layer with regression problems.

Derivatives

$$\# g(z) = \frac{1}{1+e^{-z}} \quad (\text{Sigmoid})$$

$$g'(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right) = g(z)(1-g(z))$$

$$\# g(z) = \tanh(z)$$

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= 1 - (\tanh(z))^2$$

$$g'(z) = 1 - [g(z)]^2$$

$$\# g(z) = \text{ReLU}(z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

(should be undefined at 0, but doesn't matter in software implementation)

Gradient Descent

Repeat {

$$dW^{[1]} = \frac{\partial J}{\partial w^{[1]}} , \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

$$dW^{[2]} = \frac{\partial J}{\partial w^{[2]}} , \quad db^{[2]} = \frac{\partial J}{\partial b^{[2]}}$$

$$w^{[1]} := w^{[1]} - \alpha dW^{[1]}$$

$$b^{[1]} := b^{[1]} - \alpha db^{[1]}$$

$$w^{[2]} := w^{[2]} - \alpha dW^{[2]}$$

$$b^{[2]} := b^{[2]} - \alpha db^{[2]}$$

y

CODE Forward Propagation (Binary Classification)

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = g(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g(Z^{[2]}) = \sigma(Z^{[2]})$$

Back Ward Propagation

$$dZ^{[2]} = A^{[2]} - Y$$

$$\frac{1}{m} dZ^{[2]} A^{[1]} {}^T$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdim=True})$$

$$dZ^{[1]} = W^{[2]} {}^T dZ^{[2]} * g^{[1]}'(Z^{[1]})$$

↑ element wise product

$$\frac{1}{m} dZ^{[1]} X {}^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdim=True})$$

why weights should not be 0?

if all weights are zero
hidden units in one layer
are symmetric.

∴ initialize with small random numbers.

if W is too large, the activation is saturated.

Deep Neural Networks

$L = \# \text{ layers}$

$n^{[0]} = n_x$ $n^{[L]} = \# \text{ units in layer } L$
 $n^{[0]} = 1$

Generalized F.P.

$$\begin{aligned} z^{[l]} &= w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

dimensions

$$w^{[l]} = (n^{[l]}, n^{[l-1]})$$

$$b^{[l]} = (n^{[l]}, 1)$$

$$A^{[l]} = (n^{[l]}, m)$$

$$z^{[l]} = (n^{[l]}, m)$$

Why deep neural networks work?

example: For face detection earlier layers may detect edges, followed by edge combinations to detect facial features, followed by part of face detections and so on.

Deep Learning → Builds a system from primitive level to upwards toward complex

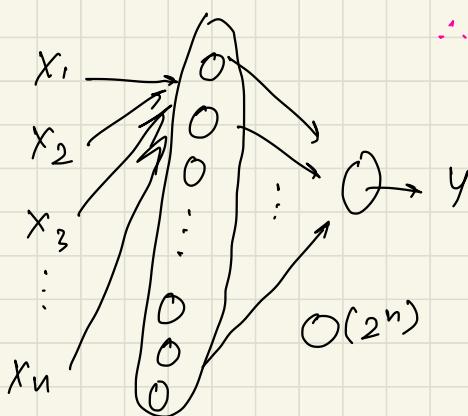
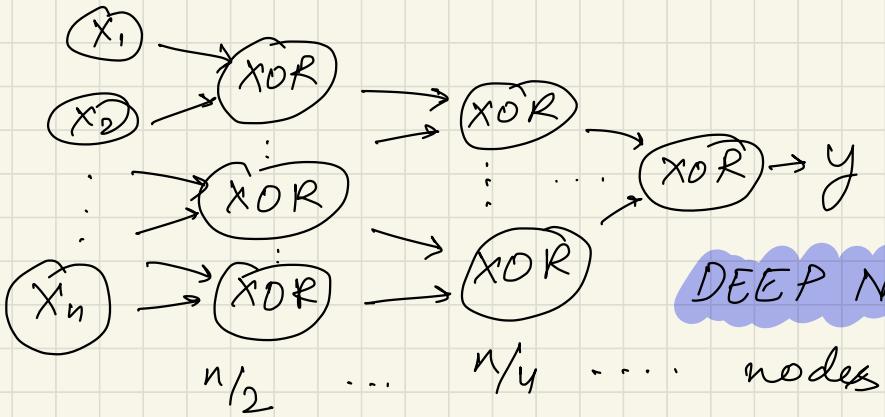
* DL are modelled against brain as brain does the same (starting from base to building up)

Circuit Theory

if a digital circuit is modelled using large NN, it will require less number of nodes per layer

But if the NN is shallow, the number of nodes increases exponentially

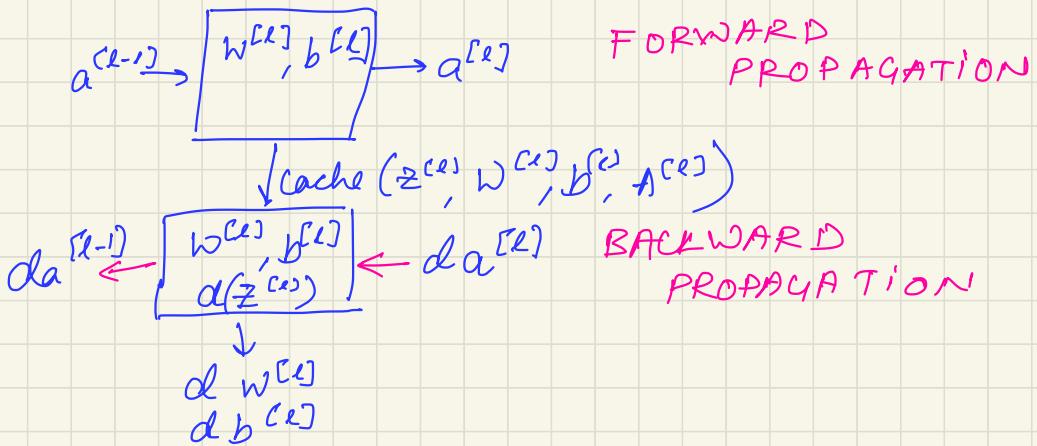
e.g. $y = x_1 \oplus x_2 \oplus x_3 \dots \oplus x_n$



∴ DNN can learn more complex systems & functions easily

Propagation

For any layer l



Forward

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

Backward

$$d_z^{[l]} = da^{[l]} * g'(z^{[l]})$$

$$d_w^{[l]} = [d_z^{[l]} \cdot (a^{[l-1]})^T] / m$$

$$d_b^{[l]} = \frac{1}{m} \text{ np.sum } (d_z^{[l]}, \text{axis}=1, \text{Keepdims=True})$$

$$da^{[l-1]} = (w^{[l]})^T \cdot d_z^{[l]}$$

element wise
product

FOR LAST LAYER WITH SIGMOID

$$\partial A^{[L]} = -Y/A^{[L]} + (1-Y)/(1-A^{[L]})$$

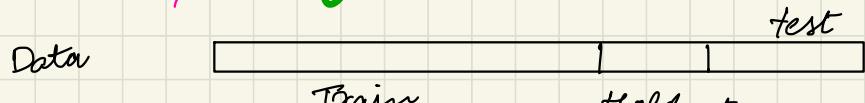
Hyperparameters

- * learning rate
- * # iterations
- * # layers (L)
- * hidden units n_0, n_1, \dots, n_{L-1}
- * choice of activation function

Hyperparameters are the parameters that influence the model architecture and learning of actual computable parameters

Improving Deep Neural Nets

Train, Dev and Test Sets \rightarrow



Train

1

test

1

In case of a smaller dataset

- Holdout covers validation
- Development set
- "Dev"

60-20-20 is a good breakup

But as the datasets become larger we use smaller %.

e.g. For example in a dataset with 1,000,000 data points dev and test sets can just be 10,000 points as that is enough to get an idea of the system.

98-1-1

Dev, testsets should be from the same distribution as the training set, or else, the results may be inaccurate.

If a non-biased evaluation of the model is not required it's okay to not have a test set.

Train set \rightarrow training the parameter
Dev set \rightarrow Fine tuning the hyper parameters
Test set \rightarrow gives an un-biased evaluation of the model.

Bias & Variance

High Bias ~ the model fit to the data is extremely simple.
[underfitting] (Model is biased towards simplicity)

High Variance ~ the model fit to the data is extremely complicated.
[overfitting] (Model surrenders to the variance in data)

{
Train Error \downarrow & Dev Error \uparrow \rightarrow high Variance
Train Error \uparrow & Dev Error \uparrow \rightarrow high Bias
Train Error \uparrow & Dev Error \uparrow \rightarrow high Bias & Variance
Train Error \downarrow & Dev Error \downarrow \rightarrow low Bias & Variance
wrt to Bayes Error.

Basic Steps for DL optimization

- ① High Bias (visualize training data performance.)
Sol: Bigger Network
Learn longer
Change NN architecture

② After Bias Reduction, Check for High Variance

Sol: More data

Regularization

Change NN architecture

Regularization

L2 Regularization :- Adding $\frac{\lambda}{2m} \|W\|_2^2$ to cost

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

$$\|W\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

L1 Regularization :- Adding $\frac{\lambda}{m} \|W\|_1$ to cost

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{m} \|W\|_1$$

$$\|W\|_1 = \sum_{j=1}^{n_x} |w_j|$$

* L2 regularization can lead to W being sparse.
⇒ Better in compressing a model.

λ → regularization parameter

Neural Network Regularization

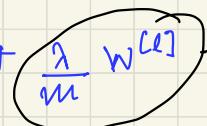
$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$$\|W^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l-1)}} (w_{ij}^{(l)})^2$$

→ FROBENIUS NORM

Now,

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



Reason why
L2 is also
known as
weight decay

$$W^{[l]} = W^{[l]} - \alpha \left[\text{derivative from back prop} + \frac{\lambda}{m} W^{[l]} \right]$$

$$W^{[l]} = \left(1 - \frac{\alpha \lambda}{m}\right) W^{[l]} - \alpha (\text{derivative})$$

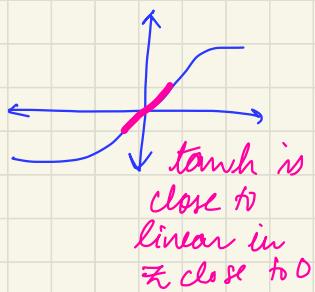
Why Regularization Reduces Overfitting?

with large enough λ , w_{so} some hidden units have less impact making model more general.

$$\lambda \uparrow \Rightarrow W^{[l]} \downarrow \Rightarrow Z^{[l]} \downarrow \Rightarrow a^{[l]} \approx z^{[l]}$$

$$(z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]})$$

∴ layers are roughly linear
hence, even deep models become somewhat simpler.



Dropout Regularization

Removing some nodes from each layer based on a dropout probability.

if keep-probability for a layer $l = kp$

$dl = np.random.rand(a_l.shape[0], a_l.shape[1]) < kp$

$a_l = np.multiply(a_l, dl)$ sets sum values to 0.

a_l / kp increases the impact by chance of node removal.

e.g. if $kp = 0.8$ a_l is increased by 20%.

↳ removes out the node removal

Different units should be removed at different iterations of gradient descent.

during test \rightarrow dropout should be eliminated to avoid noise in results

division by K_P makes sure test results aren't required to be scaled as dropout is missing in test.

Why dropout works?

Since features to a NN node are eliminated at random, weight doesn't depend on any one feature
 \therefore weights are well spread out.
Has similar impact as L2.

Keep probability can be different for diff layers.

$K_P = 1$ (for last layer)

K_P can also be used for inputs but the number should be close to 1.

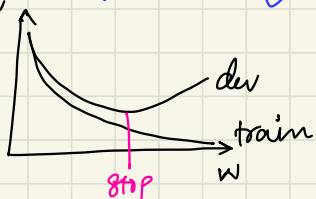
Other ways to Regularize Data

Data Augmentation \rightarrow modify data set slightly

e.g. Flip images horizontally or zoom them a little
to create more data points

Distorting a little helps too

Early Stopping \rightarrow Plot Train & Dev set Costs. and stop training half way.



DISADVANTAGE \rightarrow We compromise on Cost Optimization.
L2 Regularization is a better fit.

Normalization

- ① Modify for mean to be zero
- ② Modify for variance normalization

$$\textcircled{1} \quad \bar{X} = \frac{1}{m} \sum_{i=1}^m X^{(i)}$$

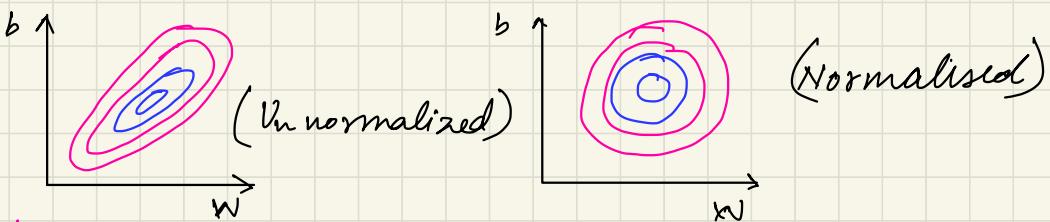
$$x = X - \bar{X}$$

$$\textcircled{2} \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m X^{(i)} \times 2$$

$$x = X / \sigma$$

INTUITION

For normalized data the weights for each feature are within the same range which ensures that the cost function is symmetric to some extent. This means, larger learning rates can be chosen to train the system.



Vanishing / Exploding - Gradients

The gradients either become too big or too small.

if $w > 1$ then over the layers the gradient of y_j will explode to a large number for multiple layers

if $w < 1$ the gradient will become quite small.

To rectify the vanishing / exploding gradient
 $\text{Var} = \frac{1}{n}$ for DNN and $\frac{2}{n}$ for ReLU layers weights

For non-ReLU layers

$$W^{[L]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{[L-1]}}\right)$$

For ReLU layers

$$W^{[L]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{[L-1]}}\right)$$

Variance $\frac{2}{n^{[L-1]} + n}$ can also be used for tanh

Gradient Checking

- ① Concatenate $w^{[1]}, b^{[1]}, w^{[2]}, \dots, w^{[L]}, b^{[L]}$ into a vector Θ .
- ② Concatenate $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ into a vector $d\Theta$.

$$J(w^{[1]}, \dots, b^{[L]}) = J(\Theta)$$

for each i :

$$d\Theta_{\text{app}}^{[i]} = \frac{J(\Theta, \Theta_2, \dots, \Theta_i + \epsilon, \dots) - J(\Theta, \dots, \Theta_i - \epsilon, \dots)}{2\epsilon}$$

$$d\Theta_{\text{app}}^{[i]} \approx d\Theta^{[i]}$$

$$\text{Check } \frac{\|d\Theta_{\text{app}} - d\Theta\|_2}{\|d\Theta_{\text{app}}\|_2 + \|d\Theta\|_2} \approx 10^{-7} \quad \text{for } \epsilon = 10^{-7}$$

$(10^{-5}$ is also good)

- * Should not be used in training (only debugging)
- * Remember to use regularization
- * Does not work with dropout

Optimization Algorithms

mini-batch gradient descent

$$y \quad x_{(n \times m)} = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} \\ & & & & x^{(1001)} & \dots & x^{(2000)} \\ & & & & & \ddots & x^{(m)} \\ x^{\{1\}} & & & & & x^{\{2\}} & \\ & & & & & & & x^{\{2000\}} \end{bmatrix}$$

$$M = 5,000,000$$

batch_size = 1000

for $t=1$ to 5000

Forward Propagation ($x^{E_{CY}}$)

Backward Propagation (x^{t+1}, y^{t+1})

update parameters.

We basically run GP on smaller batches of dataset.

epoch → 1 pass through the entire dataset

batch size = m

batch size = 1

Gradient Descent

Scholastic G.D.

(Can be extremely noisy)

(we loose the progress by vectorization)

Optimal batch size
is somewhat noisy
as well but
Vectorization can
be used to reduce
noise

→ Typically a power of 2
 $2^6, 2^7, 2^8, \dots$

Exponentially Weighted Average

moving average of data over a time period

$$V_t = \beta V_{t-1} + (1-\beta) O_t$$

if $\beta = .9$ (averaging over last 10 values)

$\beta = .98$ (averaging over last 50 values)

Bias Correction

we start from $V_0 = 0$

due to which, the initial values are biased towards 0.

∴ for initial values

$$\text{Corrected Weighted avg} = \frac{V_t}{1 - \beta^t}$$

Gradient Descent with momentum

① Compute dW, dB on current mini batch

$$VdW = \beta VdW + (1-\beta) dW$$

$$VdB = \beta VdB + (1-\beta) dB$$

$$③ W = W - \alpha VdW ; B = B - \alpha VdB$$

* $\beta = 0.9$ works best mostly but can also be tuned

Root Mean Square Prop (RMSProp)

① Compute dW, dB on the current mini-batch

$$② SdW = \beta SdW + (1-\beta) dW^2 \quad \leftarrow \text{element wise}$$

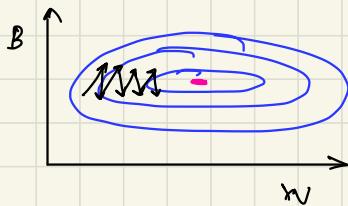
$$SdB = \beta SdB + (1-\beta) dB^2 \quad \leftarrow \text{element wise}$$

$$③ W = W - \alpha \frac{dW}{\sqrt{SdW}} ; B = B - \alpha \frac{dB}{\sqrt{SdB}}$$

Intuition

(Horizontal) $dW^2 \leftarrow \text{small}$ \rightsquigarrow dividing by a large number increases horizontal shifting of the gradient

(Vertical) $dB^2 \leftarrow \text{large}$ } decreases vertical movement



α can be as large as the chances of movement getting skewed are decreased.

* element with larger sq. in general gets damped and that with smaller values are scaled to ensure equal movement in both directions

We add $\epsilon = 10^{-8}$ to the denominators to ensure non-zero denominators

$$W = W - \frac{\alpha dW}{\sqrt{sdw} + \epsilon}; B = B - \frac{\alpha dB}{\sqrt{sdb} + \epsilon}$$

Adam Optimization (Momentum + RMS Prop)

On iteration t :

- ① Compute dW, dB using current mini-batch
- ② $VdW = \beta_1 VdW + (1-\beta_1) dW; VdB = \beta_1 VdB + (1-\beta_1) dB$
- ③ $SdW = \beta_2 SdW + (1-\beta_2) dW^2; SdB = \beta_2 SdB + (1-\beta_2) dB^2$
- ④ $v_{dW}^{\text{corrected}} = VdW / (1-\beta_1^t); v_{dB}^{\text{corrected}} = VdB / (1-\beta_1^t)$
- ⑤ $SdW_{\text{corrected}} = SdW / (1-\beta_2^t); SdB_{\text{corrected}} = SdB / (1-\beta_2^t)$

(9)

$$\omega = w - \alpha \frac{v_{dw}^{corr}}{\sqrt{s_{dw}^{corr}} + \epsilon}; \quad b = b - \alpha \frac{v_{db}^{corr}}{\sqrt{s_{db}^{corr}} + \epsilon}$$

Hyperparameters (by Adam Paper)

$\alpha \rightarrow$ tuned

$\beta_1 \rightarrow 0.9$

$\beta_2 \rightarrow 0.999$

$\epsilon \rightarrow 10^{-8}$

Adam \rightarrow Adaptive Moment Estimation

Learning Rate Decay

as learning reaches closer to minimum, smaller learning rate helps us achieve results faster rather than oscillating about the solution.

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

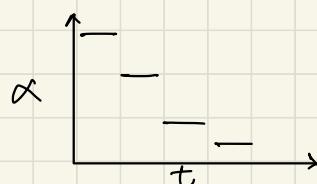
Exponential decay

$$\alpha = (0.95)^{\text{epoch}} \alpha_0$$

Sqroot decay

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0$$

discrete staircase



Local Optima

in lower dimensional spaces, local minima have 0 derivatives.

But for high dimensional spaces, saddle points may have 0 derivatives

Plateau

learning(deri) ≈ 0 for a large time



→ RMS prop, Momentum and Adam can help.

Hyper Parameter Tuning

$\alpha \rightarrow$ most important

β

hidden layer units } \rightarrow some what important.
mini-batch size

layer

learning rate decay } \rightarrow least important
among the lot.

How to:

- ① Choose a region
- ② Pick random points in that region
- ③ Create models for those points
- ④ If some points in a sub-region are better performers, choose that region and increase point density.
- ⑤ Repeat 3.

Parameter 1	Parameter 2	Parameter 3
-	.	.
.	-	.
.	.	-
-	.	.
.	-	.
.	.	-

Scale

linear Scale $\rightarrow \alpha$, # hidden units

logarithmic scale $\rightarrow \alpha$

(0.0001 \rightarrow 1)

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

$$\alpha = 10^x$$

Sample $(1-\beta)$ von logarithmic range. (0.9 \rightarrow 0.999)

$$\tau = (-3, -1)$$

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

PANDA APPROACH \rightarrow Train 1 model and modify hyperparameters each day

CAVIAR APPROACH \rightarrow Train multiple models in parallel and then select the best model

Batch Normalization

normalizing $Z^{(l)}$ in hidden layers to make the model train faster for $w^{(l)}, b^{(l)}$

$$\mu = \frac{1}{m} \sum_i Z^{(l)}$$

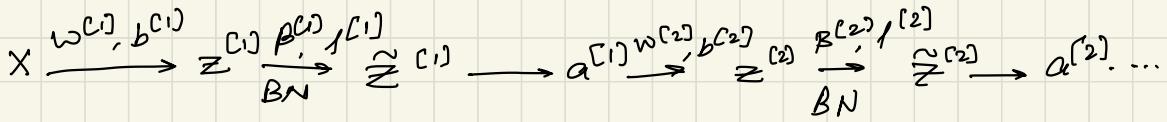
$$\sigma^2 = \frac{1}{m} \sum_i (Z^{(l)} - \mu)^2$$

$$Z_{\text{norm}}^{(l)} = \frac{Z^{(l)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{Z}^{(l)} = \gamma Z_{\text{norm}}^{(l)} + \beta$$

↑ learnable parameters

Use $\tilde{Z}^{(l)}$ instead of $Z^{(l)}$



parameters

$$w^{(1)}, \beta^{(1)}, w^{(2)}, \beta^{(2)} \dots$$

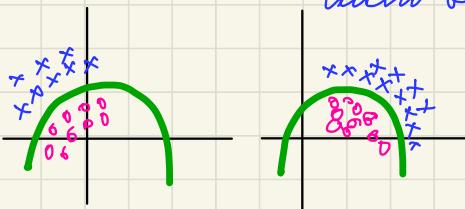
$$\beta^{(0)}, \gamma^{(0)}, \beta^{(2)}, \gamma^{(2)} \dots$$

We can eliminate $b^{[l]}$ as batch normalization zeros out the mean

$$z^{[l]} = \gamma^{[l]} a^{[l-1]}$$

β and γ can also be updated using GD, GD with momentum or Adam.

Covariate shift \Rightarrow the ground truth stays the same but the distribution of data-set changes



Batch Normalization stabilizes the layer later's variance and mean which ensures that even if the distribution changes for the first layer, the model still adapts to it.

Batch norm solves Covariate Shift

* Batch norm has slight regularization effects as mean & variance are computed per mini-batch.

Testing

- ① Compute μ and σ^2 using exponentially weighted avg during the training process for all minibatches.
- ② Use this avg value during test as values of σ^2 and μ .

Multiclass Classification

$$C \rightarrow \# \text{ of classes}$$
$$n^{[C]} = C$$

Last layer uses softmax to implement this

Activation Function

$$t = e^{(z^{[C]})}$$
$$a^{[C]} = \frac{e^{z^{[C]}}}{\sum_{i=1}^C t_i}$$

$$a_i^{[C]} = \frac{t_i}{\sum_{j=1}^n t_j}$$

Kerr function

$$\chi(\hat{y}, y) = \sum_{j=1}^c y_j \log \hat{y}_j$$

Tensor Flow:

- ① Define variables `tf.variable(0, dtype)`
- ② Initialize gradient tape (records forward prop)
- ③ Set trainable variables list
- ④ Compute gradients `tape.gradient(cost, trainable-va)`
- ⑤ `optimizer.apply_gradients(zip(grad, trainable-va))` for one step of training.

OR

- ① Define variables
- ② Define optimizer `tf.keras.optimizers.Adam(0.1)`
- ③ define cost-functions()
- ④ call `optimizer.minimize([cost-functions, [variables]])`
for one step of training.