

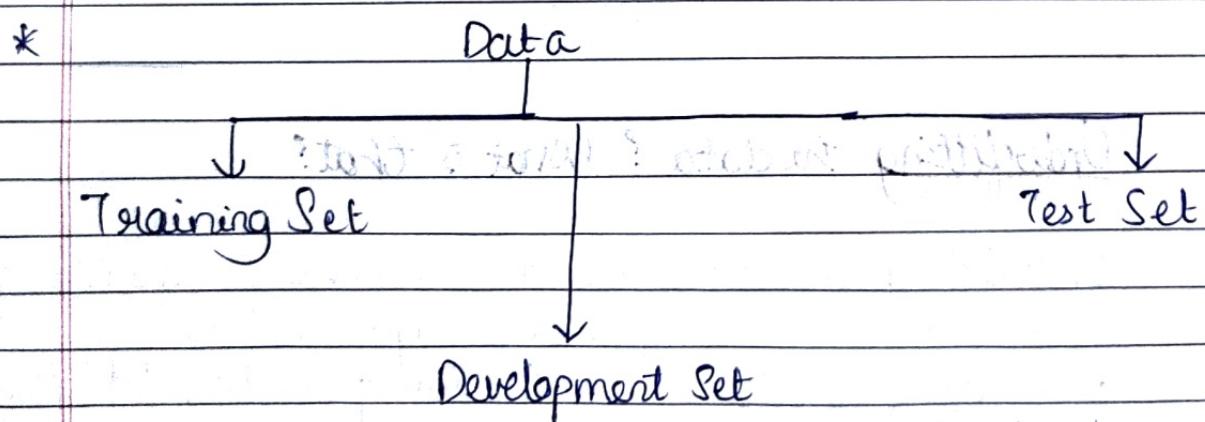
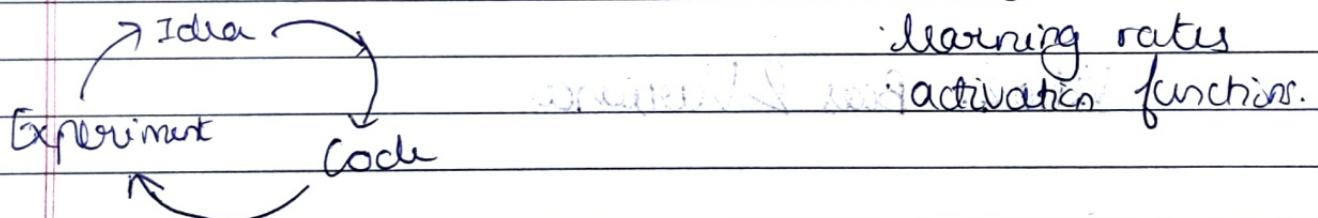
IMPROVING DEEP NEURAL NETWORK:

HYPERPARAMETER TUNING, REGULARISATION

OPTIMISATION

Video 1: Train / Dev / Test Sets

- * Deep learning is a highly iterative process. You must decide various factors like:
 - layers
 - hidden units
 - learning rates
 - activation functions.



Eg. of splitting data

70% / 30%
↓ ↓
train test

60% / 20% / 20%
Train Dev Test

- * The goal of the development set is to see which algorithm works best.
- * When there is ~~is~~ ^{very huge} a huge dataset, the dev & test sets are a small % age of the total data
- * The goal of test set is that it gives you an estimate of accuracy
- * Dev and test set should come from same distribution
- * Since test set gives you an unbiased estimate of performance, its not necessary.

Video 2: Bias & Variance

(Bias variance trade off?)

Underfitting the data? What is that?

Underfitting occurs when the machine learning model is too simple to capture the underlying patterns in training data. It fails to learn relationship between input features & target variable effectively.

The model's ^{performance} performs well on training data ^{but} and on unseen data is poor. It over simplifies problem and doesn't capture complex relations.

Page No.
Date

Cause: using a very simple model
insufficient training data
inadequate feature selection

Solution: use a more complex NN

Overfitting the data! What???

Overfitting occurs when the machine learning model is too complex / flexible and it fits the training data closely, capturing unnecessary stuff like noise & random fluctuations rather than the ~~nature~~ ~~patterns~~ patterns.

The model performs well on training data but poorly on unseen data.

It has low bias and high variance

It is sensitive to small changes in training data.

Bias: Bias refers to the error introduced by approximating a real world problem, which may be complex, by a simplified model. The model makes strong assumptions about the data resulting in underfitting

Variance: Variance refers to the model's sensitivity to small fluctuations or noise in the training data. The models are very complex and tend to fit the noise in the data rather than underlying patterns

how well or fitting
bias

Train Set	1%	15%	15%	0.5%
Dev Set	11%	16%	30%	1%
Variance problem	high variance	high bias	high bias	low bias low variance

humane \approx 0% error

Bayes error - lowest possible error achievable

Analysis is relative to human error

Video 3: Basic Recipe for Machine Learning

Q1: high bias? (bad training data performance)

- Solution:
- (i) Bigger network
 - (ii) Train longer
 - (iii) NN architecture that suits your problem
 - (iv) test

Q2: high variance? (bad on dev set performance)

- Solution:
- (i) more data
 - (ii) regularization
 - (iii) test NN architecture

Bias Variance Tradeoff:

In the pre-deep learning era, we didn't have as many tools to just reduce or just reduce variance.

But today,

train a big network — reduces bias



get more data



regularise

} reduce variance



Regularising Your Neural Network

Video 1: Logistic Regression

$w \in \mathbb{R}^n, b \in \mathbb{R}$

$$\min_{w, b} J(w, b)$$

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

↑
regularisation parameter

↳ regularisation:

$$\|w\|_2^2 = \sum_{j=1}^n w_j^2 \stackrel{?}{=} w^T w$$

$$\text{L1 regularisation: } \frac{\lambda}{2m} \sum_{i=1}^m |\omega_i|$$

$$\frac{\lambda}{2m} \|\omega\|_1$$

λ - hyperparameter that requires tuning

Neural Network:

$$J(\omega^{(l)}, b^{(l)}, \dots, \omega^{(0)}, b^{(0)}) = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\omega^{(l)}\|^2$$

Frobenius norm $\rightarrow \|\omega^{(l)}\|_F^2 = \sum_{i=1}^{n(l+1)} \sum_{j=1}^{n(l)} (\omega_{ij}^{(l)})^2$

- i \rightarrow no. of neurons in current layer $n^{(l)}$
- j \rightarrow columns of weight matrix \rightarrow no. of neurons in previous layer $n^{(l-1)}$

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

$$\omega^{(l)} = \omega^{(l)} - \alpha d\omega^{(l)}$$

$$\omega^{(l)} = \omega^{(l)} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} \omega^{(l)} \right]$$

$$= \omega^{(l)} \left(1 - \frac{\alpha \lambda}{m} \right) - \alpha$$

↑
weight decay

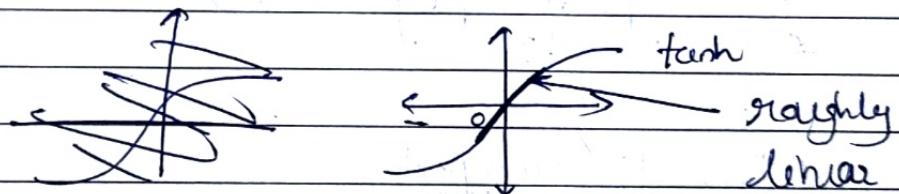
Video 3: Why Regularisation Reduces Overfitting

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

- * If λ is too large a lot of 'w's will be close to 0 which will make the NN simpler
- * If λ is good enough it will just reduce some neural networks overfit

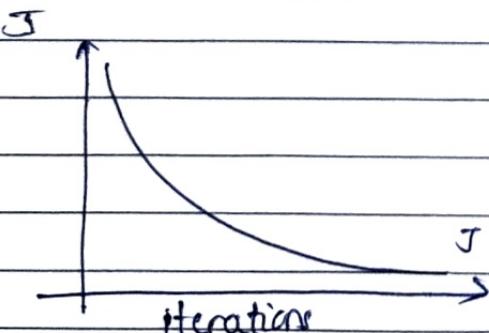
With tanh function:

- * If λ is too large, w's will $\rightarrow 0$



- * If λ is just right only some tanh functions will become roughly linear

Implementation tip:



Video 4: Dropout Regularisation

* Dropout Regularisation: It eliminates some neurons/weights at each iteration based on probability

* Inverted Dropout:

$$\text{Keep prob} = 0.8$$

$$l = 3$$

d_3^* = dropout vector for layer 3

$d_3 = \text{np.random.rand}(a[1].\text{shape}[0], a[l].\text{shape}[1]) < \text{keep prob}$

→ if number generated is less than 0.8, then it will be dropped

$$a_3 = \text{np.multiply}(a_3, d_3)$$

$$\# a_3^* = d_3$$

$$a_3 = a_3 / \text{keep prob}$$

This is so that expected value of a_3 does not change

* Do not use dropout at test time: it adds noise to predictions

Video 5: Understanding Dropout

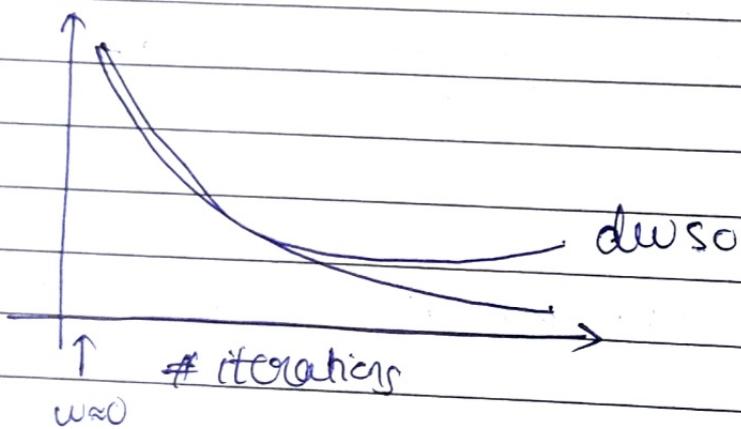
- * One cannot rely on just one feature so you should spread out weights
- * Input layer should have 1 or near 1 keep prob because you don't want to eliminate a lot of features.
- * Best way - have some layers where you apply dropout and ~~some~~ some where you don't & then just one hyperparameter which is a keep-prob for the layers where you do not apply them
- * Dropout is a regularization technique to prevent overfitting.
- * It is used in Computer Vision where overfitting occurs when a ML model performs well on training data but poorly on unseen data
- * Drawback of dropout - cost function is not well defined.
Solution: apply dropout during ~~test~~ training only.

Video 6: Other Regularisation methods

1. Data Augmentation:

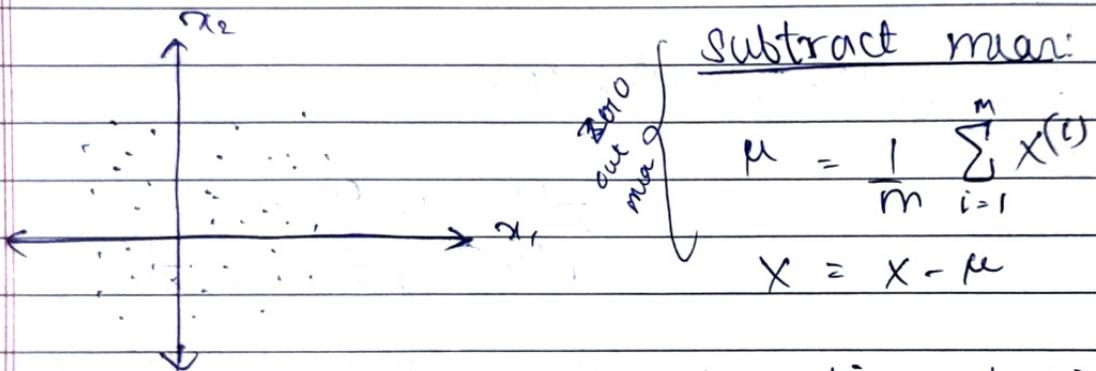
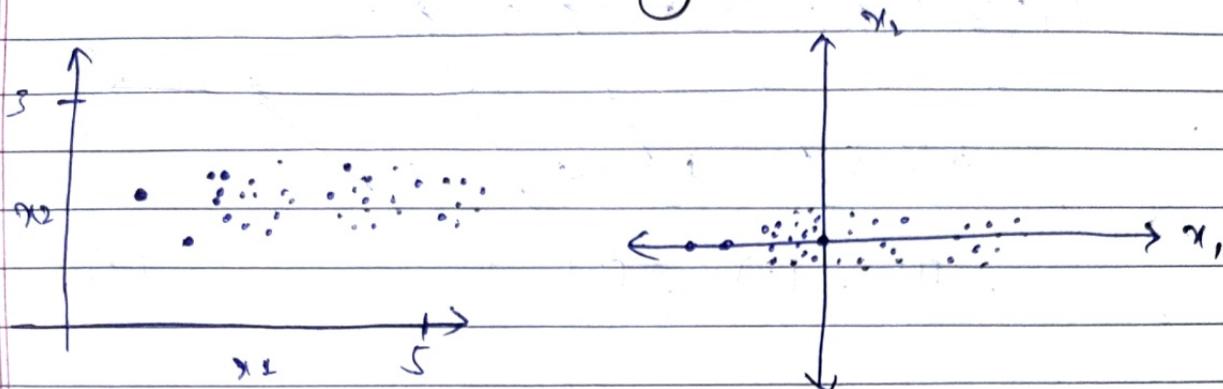
- * Flip pictures horizontally
- * Apply random position / rotation
- * In OCR - introduce rotations / distortions

2. Early Stopping



Week 1: Setting Up Optimisation Problem

Video 1: Normalise training sets

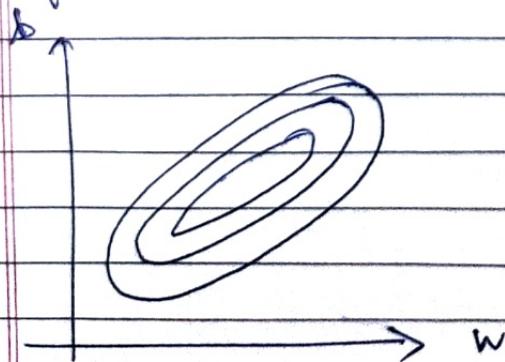


Normalise variance:

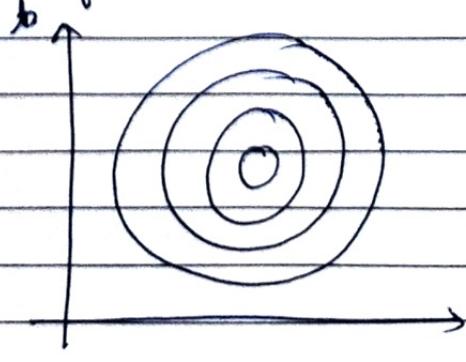
vector with variance of each feature $\rightarrow a^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \cdot x^{(i)}$

$$x / a^2$$

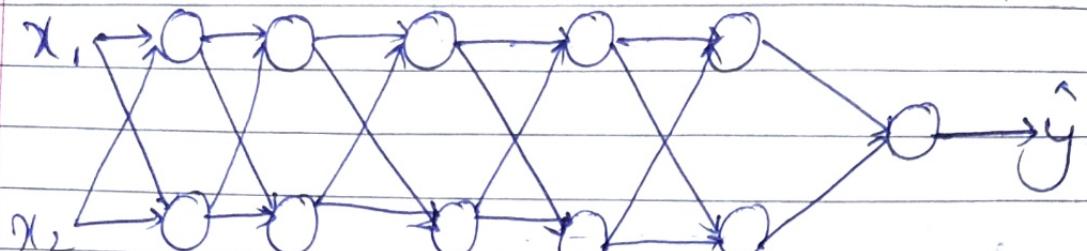
If we don't normalise:



If we normalise



Video 2: Vanishing / Exploding Gradients



$$w^{[1]} \quad w^{[2]} \quad w^{[3]} \quad \dots \quad w^{[L]}$$

$$g(z) = z \quad b^{[L]} = 0$$

$$\hat{y} = w^{[L]} w^{[L-1]} w^{[L-2]} \dots w^{[3]} w^{[2]} w^{[1]} x$$

$$z^{[1]} = w^{[1]} x$$

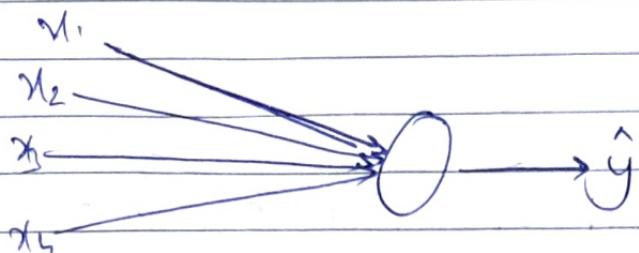
$$a^{[1]} = g(z^{[1]}) = z^{[1]}$$

$$a^{[2]} = g(z^{[2]}) = g(w^{[2]} a^{[1]})$$

$w^{[1]} > 1 \rightarrow$ activations can explode

$w^{[1]} < 1 \rightarrow$ activations decrease exponentially

Video 3: Weight Initialisation for Deep Networks



$$a = g(z)$$

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

large $n \rightarrow$ smaller w_i

$$\text{var}(w) = \frac{1}{n} \frac{2}{n}$$

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

ReLU

$$g^{(l)}(z) = \text{ReLU}(z)$$

Video 4: Numerical Approximation of Gradients

Video 5: Gradient Checking

* Gradient checking approximates the gradients and is used to find errors in backpropagation but is slower than gradient descent

* Take $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$ → reshape to a vector of size(Θ)

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = J(G)$$

Take $d\theta^{(0)}, d\theta^{(1)}, \dots, d\theta^{(L)}$ & reshape to a vector $d\theta$.

For each i :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{\text{approx}} \approx d\theta$$

Check:

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx 10^{-7}, \text{ great}$$

$$\epsilon = 10^{-7}, \quad \downarrow \quad \text{worry}$$

Videos: Gradient Checking Implementation Notes

- Don't use in training - only to debug [really slow process]
- If algorithm fails grad check, look at components to identify bug: look at the different values of i to see which are the values of $d\theta_{\text{approx}}$ really different from $d\theta$
- Regularisation

$$J(\theta) = \frac{1}{m} \sum_i L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_j \|w^{(j)}\|^2$$

- Doesn't work with dropout
- run at random initialisation; perhaps again after training