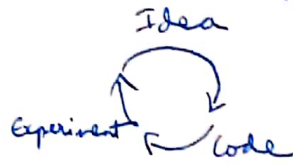


Course 2: Improving Deep Neural Networks

Week 1: Practical Aspects of Deep Learning

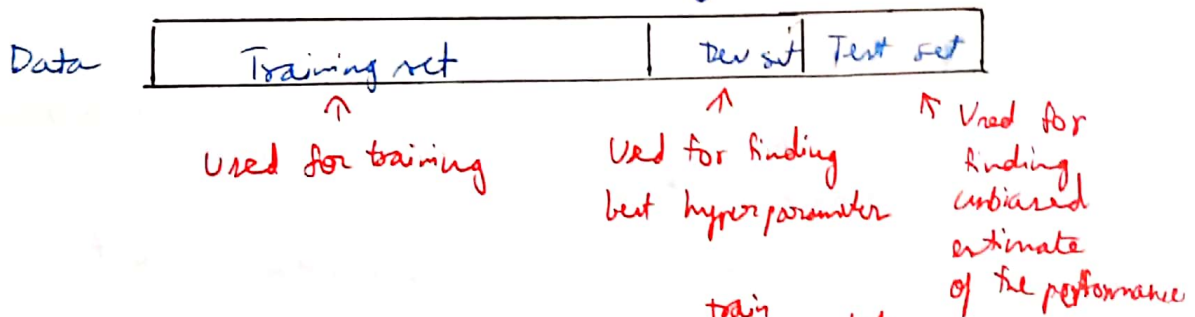
- Applied ML is a highly iterative process.



- To find the best hyperparameters (#layers, #hidden units, etc), we start off with an idea of what might be the best, write the code, experiment it and then decide a new idea that can be better....
- Intuition from one field won't work in another (ex: a person who is in MLP doing computer vision)

Train/dev/test sets

also called hold out,
cross validation



- Usually people split is as 70%/30% or 60%/20%/20%.
And this is fine if you have less training examples (100 - 10,000)
But if you have, let's say a million training examples then it's better to divide as 98%/1%/1%, even if it's 1%, it's still 10K training examples & serves the purpose

Mismatched train/test distribution

Make sure dev and test set come from same distribution

Training set:

Cut pics from
webpage

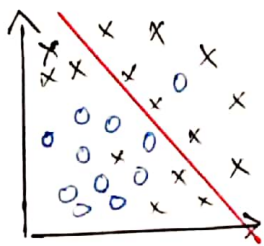
Dev/test sets:

Cut pics from users
using your app

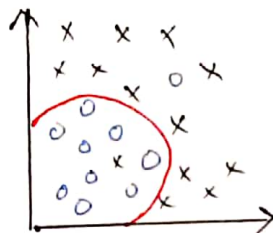
Here pics maybe of better quality than this

It's fine to train from webpage pics. (like most people do with scraping), but make sure dev & test set come from same distribution

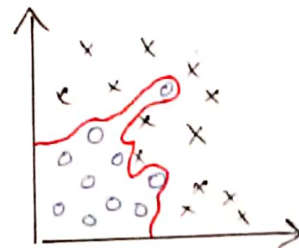
Bias / Variance



High bias
Underfitting



"Just Right"



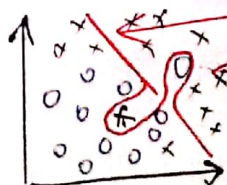
High Variance
Overfitting

Take a classification example:

| | | | | |
|------------------|---------------|-----------|---------------------------|-------------------------|
| | 1% | 15% | 15% | 0.5% |
| Train set error: | 1% | 15% | 30% | 1% |
| Dev set error: | 11% | 16% | | |
| | High Variance | High Bias | High Bias & High Variance | Low Bias & Low Variance |

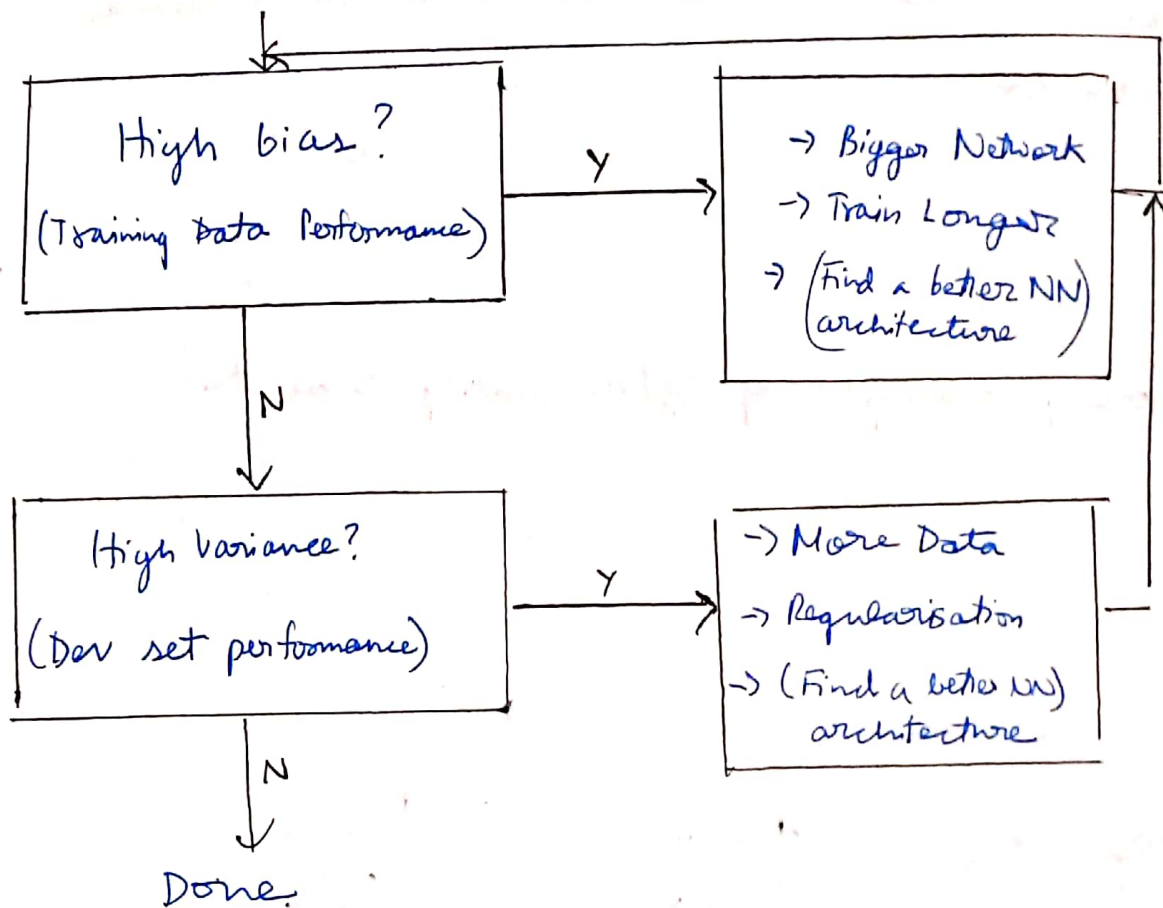
Usually humans have ~0% error - Optimal/Bayes error is 0%
If it were 15%, then this would have low bias & variance

What does high bias and high variance look like?



It'll have high bias in some places
& high variance in some places
This occurs a lot in high dimension models

Basic Recipe for Machine Learning



Regularization

We add the regularization term to the cost function to prevent overfitting.

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

← regularization parameter

L_2 regularisation: $\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$ ← used more often

L_1 regularisation: $\frac{\lambda}{2m} \sum_{i=1}^m |w_1| = \frac{\lambda}{2m} \|w\|_1$

→ In logistic regression

In neural network:

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

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ji}^{[l]})^2$$

If $\|w^{[l]}\|_2^2$ it is called "Frobenius norm"

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

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

→ This is also called weight decay

$$\text{since } w^{[l]} := w^{[l]} - \alpha \left[(\text{from bp}) + \frac{\lambda}{m} w^{[l]} \right]$$

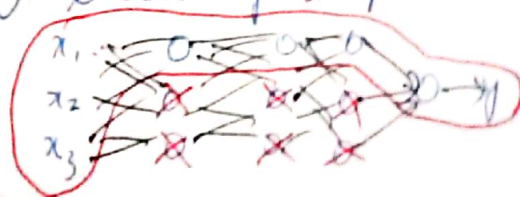
$$w^{[l]} := w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]} - \alpha (\text{from bp})$$

$$w^{[l]} := w^{[l]} \left(1 - \frac{\alpha \lambda}{m} \right) - \alpha (\text{from bp})$$

→ $w^{[l]}$ is slightly reducing (decaying)

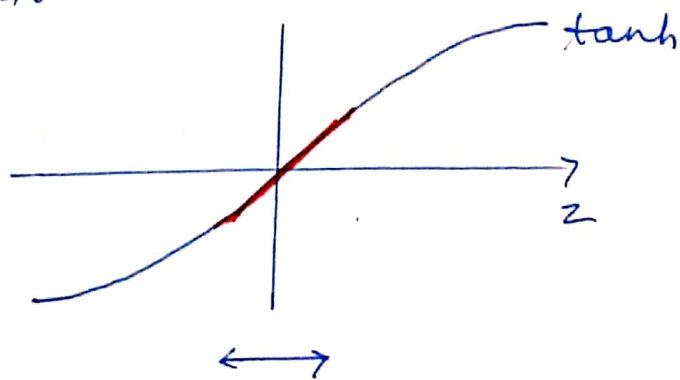
How does regularization prevent overfitting?

By penalising the w terms, $w^{[l]} \approx 0 \leftarrow$ it starts reducing to almost 0. Therefore the neuron's start to weaken and give low outputs. This is equivalent to the NN becoming simpler network (like ~~linear~~ logistic regression)



Therefore the variance starts to reduce since the network is simpler.

Another intuition is as follows,



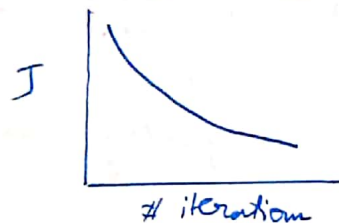
As $\lambda \uparrow$, $w^{[L]} \downarrow \therefore z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \downarrow$

Therefore ~~$w^{[L]}$~~ $z^{[L]}$ will be confined in the red zone where it is linear.

\therefore It starts to behave more linearly (more bias)

How to debug:

$$J(\dots) = \sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_e \|w^{[e]}\|_F^2$$



When the regularisation term is added, J will reduce monotonically - if it doesn't, then something is wrong. If the regularisation term isn't there, the J may or may not reduce monotonically.

Dropout Regularization

Here you 'drop-out' a certain percentage of neurons at every iteration \Leftarrow (forward + backward pass) of gradient descent.

Ex: a:

If ~~keep-prob~~ Keep-prob = 0.8

This means 20% of the neurons will be dropped

$d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1])$
(create a matrix with 20% as 0 - for this
80% as 1 - example keep-prob
element wise)

$a3 = \text{np.multiply}(a3, d3)$
The 0 elements make ~~the~~ 20% of $a3$ 0 (since multiply)

$a3 /= \text{keep-prob}$

\uparrow by dividing by 0.8 $\frac{a3}{0.8} = \frac{0.3}{0.8} = 0.3 \times \frac{10}{8}$

inverted. We pump up $a3$ so when we do
dropout $z^{[n]} = W^{[n]} \cdot a^{[3]} + b^{[n]}$

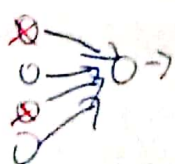
The value of $z^{[n]}$ doesn't reduce
(pumping up $a3$ compensates the 20% reduction)

Don't use drop down at test time

Why does it work?

- At every iteration, the NN is being reduced in size so only 80% of the neurons are active (in the above example)

- A neuron can't rely on any one feature, so it has to spread out weights

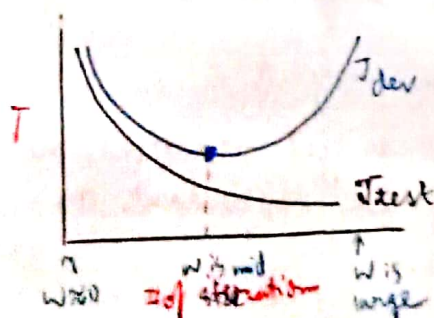


\Leftarrow Here any feature can get removed, so it will equally be dependant on all the features

- You can set different keep-prob for every layer. So this way you can choose which layer to penalise more
- ~~Set~~ Set keep-prob for input layer to 1.0 (Don't drop any inputs - not ideal)
However this is done in computer vision (set to 0.9) since there are so many inputs (pixel values).
- Disadvantage:
Harder to debug. You can't plot J with #iterations and expect a monotonic reduction \searrow .
Instead while testing, set all the layers' keep-prob to 1.0 (Don't drop any neuron).

Other regularisation methods

- Data augmentation. you can flip or distort or crop existing images to get ~~more~~ more data. This is done if you have less data and can't get more. But it isn't as good as getting new images.
- Early stopping:



We stop gradient descent early so it doesn't overfit the data

Disadvantage of Early stopping:

- Orthogonalisation - splitting the tasks

- Optimise cost function }
- Gradient descent ...

- Prevent overfitting }
- Regularisation ...

Usually we first focus on optimising cost function and then ~~use~~ use tools to prevent overfitting, however in early stopping we are combining these 2 tasks. We don't find the optimal cost function on one hand and don't prevent overfitting properly on the other. Rather it's better to separate these 2 tasks and use appropriate ~~tools~~ ^{tools} for that task.

Normalizing inputs

This is used to make the data more symmetrical

→ Subtract out / zero out the mean

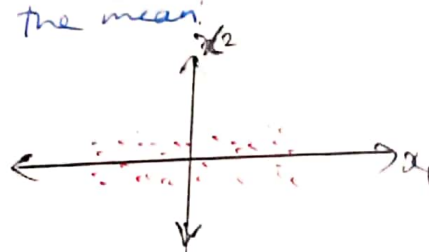
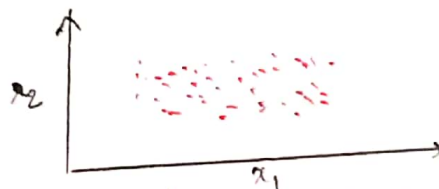
$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \mu$$

→ Normalize the variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$$

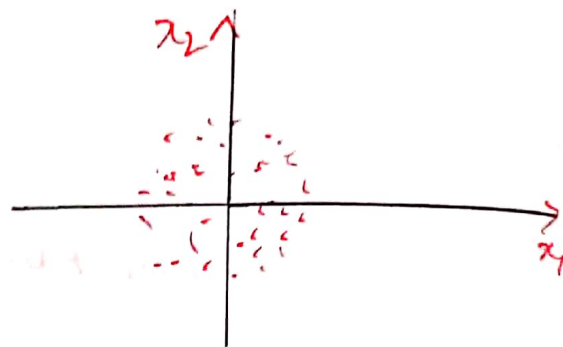
$$x := \frac{x}{\sigma}$$



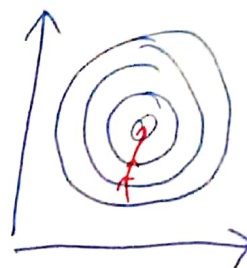
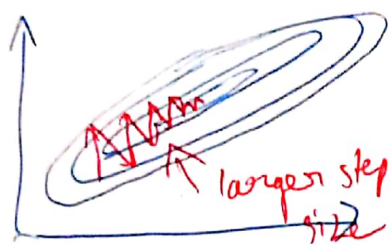
Use the same μ & σ^2 to normalise test set

∴ Combining both,

$$x := \frac{x - \mu}{\sigma}$$

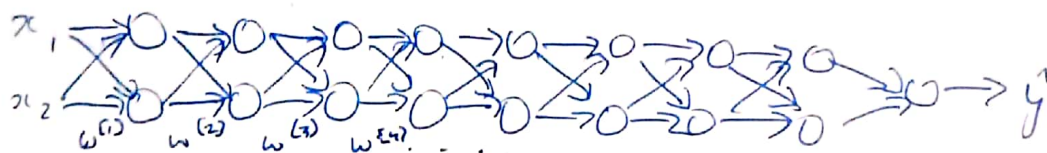


By doing this, we can use a smaller gradient descent step



Vanishing/Exploding gradients

- One problem with deep NN is that the data explodes (increases exponentially) or vanishes (decreases exponentially).



Let $g(z) = z$ - linear function and $b^{[L]} = 0$

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

① Suppose $w^{[2]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

$$\therefore \hat{y} = 1.5^L x$$

← Here if L is more, then data will explode

② Suppose $w^{[2]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$

$$\therefore \hat{y} = 0.5^L x$$

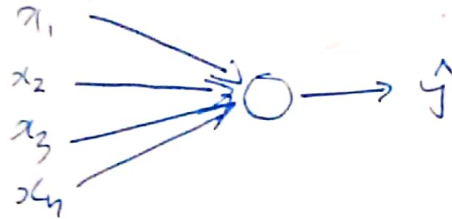
← Here if L is more, then data will vanish

$\therefore W^{[l]} > I$: Data exploder
 $\therefore W^{[l]} < I$: Data vanisher

There is a partial solution for this:

Weight initialization

take a single neuron for example,



let $b=0$ for this example

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

Now we can see that if we don't want z to be too big,
 for large $n \rightarrow$ smaller w_i

Setting variance of w_i , $\text{var}(w_i) = \frac{1}{n}$ or $\frac{2}{n}$ for ReLU
 solves the problem

we do this by,

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

Other variants,

for tanh = $\sqrt{\frac{1}{n^{[l-1]}}}$ we use 1
xavier initialization

other, $\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$

But for ReLU use this

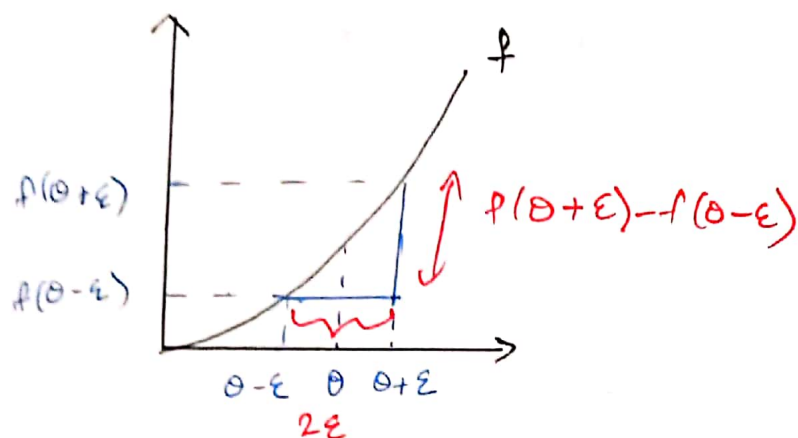
He initialization $\rightarrow \sqrt{\frac{2}{n^{[l-1]}}}$



recommended

Numerical Approximation of gradients

$$\text{Let } f(\theta) = \theta^3$$



$$\therefore \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

Example, for $\theta = 1$ and $\epsilon = 0.01$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

Gradient Checking

1) Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ

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

2) Take $dw^{[1]}, db^{[1]}, \dots, dw^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$

Now we use this to check if $d\theta$ is the gradient of $J(\theta)$

3) Grad Check:

for each i :

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

check $\|d\theta_{\text{approx}} - d\theta\|_2$

$$\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2$$

If we get of the order, 10^{-7} great!!

for $\epsilon = 10^{-7}$ $\rightarrow 10^{-5}$
 $10^{-3} \leftarrow$ bug is there

Gradient Checking implementation notes

1. Don't use in training - only to debug
2. If algorithm fails grad check, look at components to try to identify bugs
3. Remember regularization
4. Doesn't work with dropout
5. Run at random initialization, perhaps again after some training.