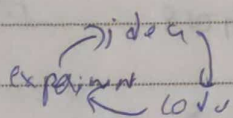# COURSE 2: PRACTICAL ASPECTS OF DEEP LEARNING

## WEEK 1

### TRAIN/DEV/TESTS

Applied ML is a highly iterative process
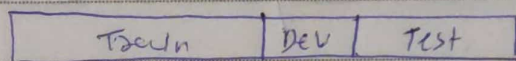
we start with ~appeox hyperparameters then do

experiment ) idea
             ) code

to improve them

Total
Data is usually divided into 3 parts

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

↓
Cross-check

In early stages when we had less data
~ some tens thousands — we used 70/30 %.
train / test

But now when we have millions of units we
reduce the % of test, dev set

NN with diff hyper periods
↓
dev, Test set — used to verify which algo does better
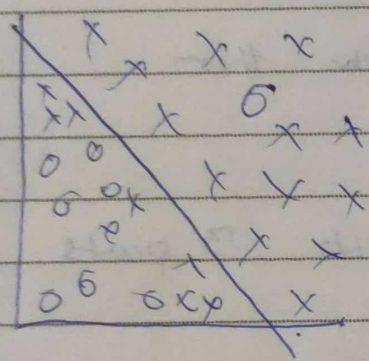prediction

Since we have millions of datasets even a small %
of them can be used to verify which algo is best

→ Test/dev and Train data should come from same distribution !! No mismatch should take place
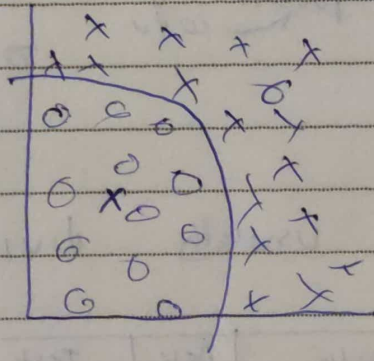
→ It is ok if you don't have a Test set
→ you evaluate them on dev set
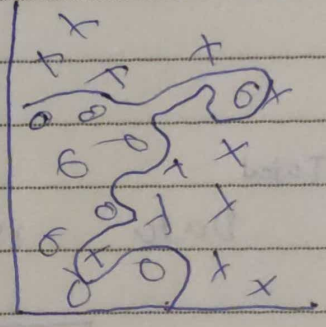
# BIAS AND VARIANCE

↗ ↘ → our predictions for a classification problem



Highly biased (under fitting) — Just right — High variance (over fitting)

To understand better

Let the
Train set error: 1%  } → it means you've
dev set error: 11%  }  over fitted the model
↓
high variance

Let the
Train set error: 15%  } → it means you've
dev set error: 16%  }  under fitted the model
↓
High bias

Let the.
train set error : 15% } → high bias and
dev set error : 30% } high variance
// bad model

let the
train set error : 0.5% → low bias and low
dev set error : 1% variance.

eg: of high bias and high variance



# BASIC RECIPE FOR ML

1st ask

you have high bias ? ──yes→ → Try bigger network
↓ more importance train longer
depends on train set on Grad dient use better optimize'n algo
descent. //then look to a better
↓No. suited NN architecture

then ask

high variance ? ──yes→ Get more data.
↓ Regularization
depends on dev set Search for other better NN architecture
↓No.

Earlier there was a tradeoff → if you try to decrease bias
variance went up and vv
But now, set of things to do to increase bias/variance are diff
so no tradeoff

# REGULARIZATION

↳ use it your model is overfitting → high variance

lets d understand this using Logistic regression

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

W ∈ Rⁿ

b ∈ R

$$\|w\|^2 = \sum_{j=1}^{m} w_j^2 = w^T w$$

called L2 norm

we add this term
to regularize w

[L2 Regularization]

why regularizing only for w?

Actually we can add the term $\frac{\lambda}{2m} b^2$ but

it won't matter that much.

you can also use L1 regularization

here you add $\frac{\lambda}{m} \sum_{i=1}^{n} |w| = \frac{\lambda}{2m} \|w\|_1$

↳ if you use this then w will be sparse,

they will be lots of zeros in w.

L2 regularization is used very often.

$\lambda$ = regularization parameter
↓
you've to tune it

L2 regularization in Neural Networks

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

$$\|w^{(l)}\|^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{(l)})^2 \qquad w \text{ is a } n^{[l]} \times n^{[l-1]} \text{ matrix}$$

"Frobenius norm"

then how do you calculate apply Gradient Descent here
we did it using backprop which gave dw.

$$dw^{(l)} \rightarrow \frac{\partial J}{\partial w^{(l)}}$$

✓ $dw^{(l)} \rightarrow$ wartean backprop gave $+ \frac{\lambda}{m} w^{(l)}$

↓
correct dw ←

when using regularization

then that w update becomes

$$w^{(l)} = w^{(l)} - \alpha \cdot dw^l = w^l - \alpha \left[ \text{from backpop} + \frac{\lambda w^{(l)}}{m} \right]$$

$$= w^l - w^{(l)} \frac{\alpha\lambda}{m} - \alpha(\text{from backpop})$$

$$\underbrace{w^{(l)} \left( 1 - \frac{\alpha\lambda}{m} \right)}_{} - \qquad ↓$$

∴ L2 regularization is
also called "weight decay" ←

↓

b/z in add" to subtracting what we were prusly doing we
are also subtracting some multiple of w $\left[ \frac{w\alpha\lambda}{m} \right]$

# Why does Regularization prevent overfitting?

graphical answer

with large $\lambda$ → basically we reduce down $w$ → reduces $z$

↓                                    makes sure our
we penalize                          lies
when $w$     sigmoid/tanh result times around zero //
$J$ for
is large

// also when you're doing regularization
make sure you use the updated value
for $J$


# Dropout Regularization

we randomly remove some nodes from layers

then we remove all the weights and incoming/
outgoing edges → basically remove those node's existence
↓
→ makes our network somewhat ~~bigger~~ small

But we removed nodes at random right? then
why does it work?
It just does

How do we implement this?

Iterate each layer   let's say keep prob → 0.8

$$d^{(3)} = np.random.rand(a^{(3)}.shape[0], a^{(3)}.shape[1]) < keep.prob$$
mask
we'll keep 80% of nodes.

then our a3 becomes

$$a3 = np \; multiply \; (a3, d3) \longrightarrow a3^* = d3$$

$$a3 \; /= keep-prob$$

// dont implement dropout directly during test time

But again why does dropout (which is random) work?
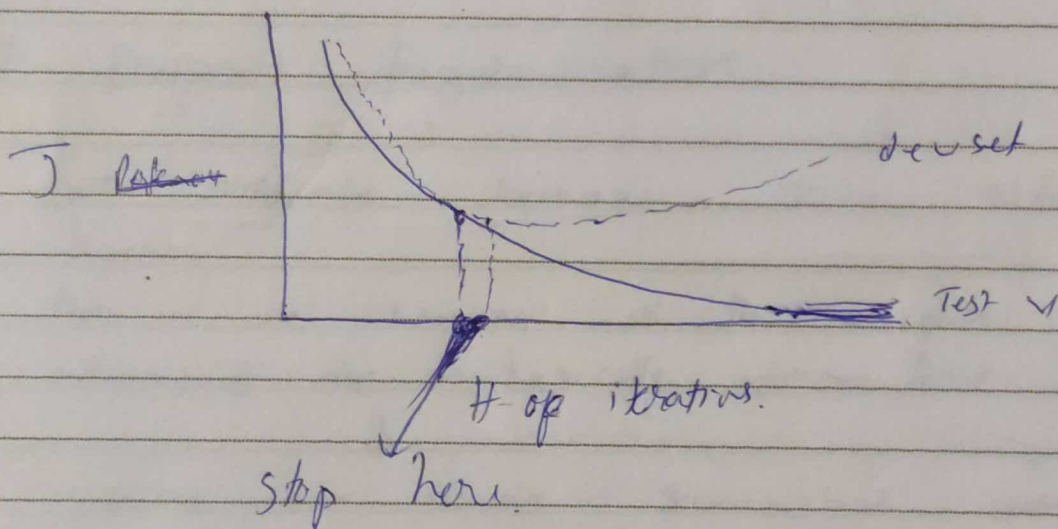
# Understanding dropout.

# Other Regularization Techniques

→ Bigger Network (always helps)
↓

#1
data augmentation { → get more data,
→ if you dont have access to more data
generate more data using the data you already
have

## #2 Early stopping    (orthogonalization)

plot but test and dev set simultaneous



J Rate

dev set

Test ↓
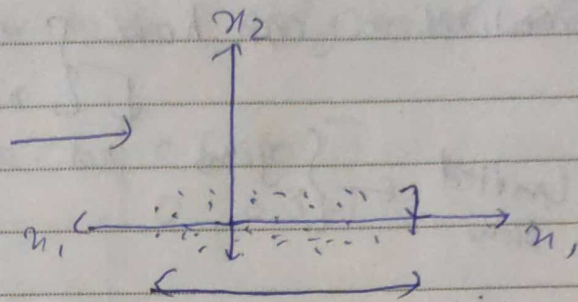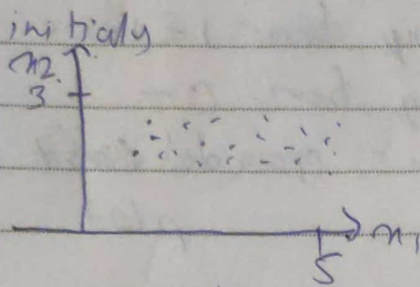
# of iterations.

stop here.

// you can just use L2 regularization instead of
Early stopping
↓
computationally somewhat
expensive bcs you've to try
different values of λ

# Normalizing input
↳ helps in speeding up training

say $x = \begin{bmatrix} n_1 \\ n_2 \end{bmatrix}$

initially



subtract mean:

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

$$x := n - \mu$$

// $n_1$ variance
is larger than
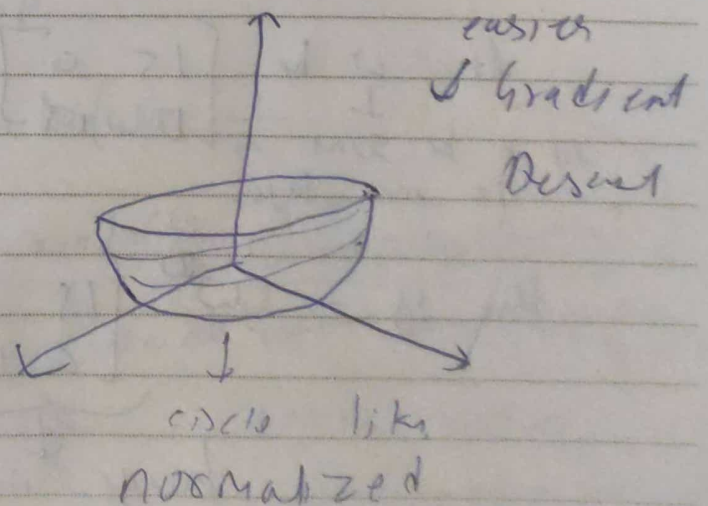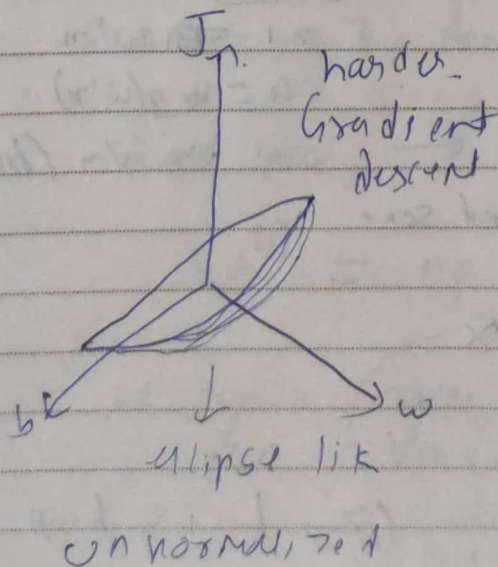$n_2$ //

normalize variance

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

$n/ = 6$



it you use μ and σ to normalize
train set, use same μ and σ to normalize
test set.

Normalization also helps in getting a good
cost function



harder
Gradient
descent

easier
Gradient
Descent

ellipse like
unnormalized

circle like
normalized

unnormalized can have $\begin{cases} n_1 \text{ range from } 1-1000 \\ n_2 \text{ range from } 0-1 \end{cases}$

makes Gradient $\leftarrow \begin{cases} \text{gives not such a good cost } f^n \\ \qquad\qquad\qquad\qquad\qquad plot \end{cases}$
descent slow

Normalization makes the range of all inputs similar

# Vanishing / Exploding Gradients

$\begin{cases} \text{Sometimes your gradient becomes vv small} \\ \text{or vv big} \rightarrow \text{makes training harder} \end{cases}$

specially in deep nns.

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \overset{big}{\overset{\downarrow}{L}}$
consider a deep nn with $L$ layers

let $b = 0$ for now and say we're using
a linear activation $f^n$.

then $y = W^{[L]} W^{[L-1]} W^{[L-2]} \cdots W^{[3]} W^{[2]} W^{[1]} \underline{x}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad z1 \leftarrow W'x.$

let $w$ be $\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$
$\qquad\qquad\qquad\qquad\qquad\qquad a = g(w'x)$
let it be same
in each layer.                and so.               $x$ w'x (linear)

then $y = W^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$

$\underbrace{\qquad\qquad\qquad\qquad}$

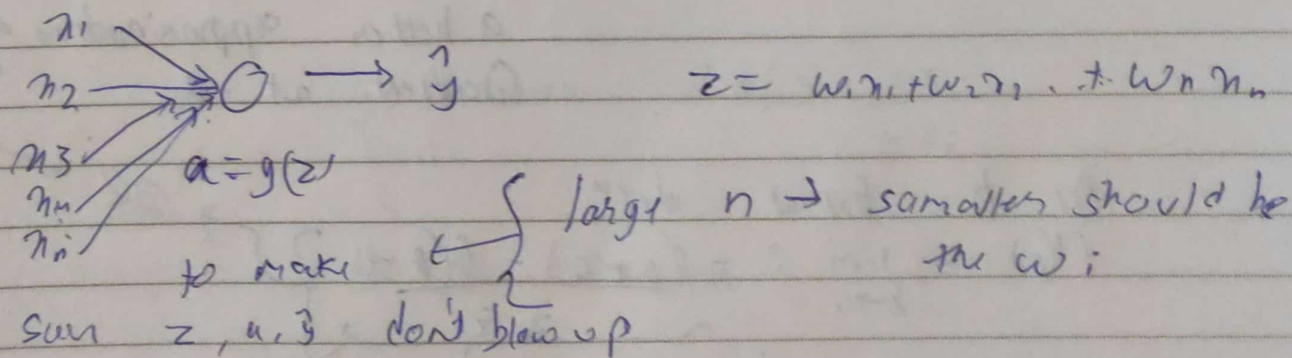huge num ber (z L is high

this makes y v big and gradient also big

Similarly it $W = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$

y will becom v small and hence
gradient will also becom v small

→ takeaway : don't take weight matrix of large
↓ or small no.
be careful while initializing matrix
randomly.

# weight initialization for Deep Networks

Single neuron example



$z = W_1 n_1 + W_2 n_2 \ldots + W_n n_n$

$a = g(z)$

{ large n → smaller should be
the $w_i$

to make $z, a, \hat{y}$ don't blow up
sum

Good idea → set $var(w_i) = \frac{1}{n}$

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

it your using ReLU          internal random generation makes
Then   Var $(w_i) = \frac{2}{n}$          works better
                                                          it equal
↓
* np.sqrt($\frac{2}{\text{DIMS}}$)  // gaussian random nd.
generation ??

If you're using tanh →
$$\sqrt{\frac{1}{n_{i}}} \times \text{ this guy}$$

np.random.randn(shape-1) × this guy

Xavier initialize ↗

another version   $\sqrt{\frac{2}{n^{[L-1]} \times n^{[L]}}}$

# Numerical approximation of Gradient

say you're  a $t^n$ f

$f(a) = \theta^3$



instead of the usual $n_1$
$y_1$

$f'(\theta) \approx \frac{n_1 + n_2}{y_1 + y_2}$  gives

a better approximation of Gradient at a

Ac to formal def$^n$.

$$f'(\theta) = \lim_{\varepsilon \to 0} = \frac{f(\theta + \varepsilon) - f(\theta - \varepsilon)}{2\varepsilon} \approx$$

and the difference b/w actual gradient and gradient calculated using   is  $O(\varepsilon^2)$

small $\varepsilon$ → small error

if we use ↓  the error is $O(\varepsilon)$

↓
relatively bad

# Gradient checking

1) Take $w^{(1)}, b^{(1)}, \ldots, w^{[L]}, b^{(L)}$ and reshape into a big vector $\theta$.

then $J(w^{(1)}, b^{(1)}, \ldots, w^{(L)}, b^{(L)}) = J(\theta)$

2) Take $dw^{(1)}, db^{(1)} \ldots dw^{(L)}, db^{(L)}$ and reshape into big vector $d\theta$. → same dimension as $\theta$

{ Is $d\theta$ same as gradient of cost function?

this is called

↳ Gradient check. (grad check.)

↓

calculate $d\theta_{appros}$

for each $i$

$$d\theta_{appror}[i] = \frac{J(\theta_1, \theta_2 \ldots \theta_i + \varepsilon, \ldots) - J(\theta_1, \theta_2 \ldots \theta_i - \varepsilon \ldots)}{2\varepsilon}$$

and you have $d\theta[i] = \frac{dJ}{\partial \theta_i}$

Now is

$$d\theta \approx d\theta_{approx} ?$$

how to decide whether they're approximately equal or not?

get euclidean distance → $d\theta_{appror} - d\theta = \frac{\|d\theta_{appror} - d\theta\|_2}$

check $\dfrac{\|d\theta_{appror} - d\theta\|_2}{\|d\theta_{appror}\|_2 + \|d\theta\|_2} \approx 10^{-7}$ ? yes → Great

$10^{-5} \sim$ okayish

$10^{-3} \to$ hard

↓ you decide (usually v small)

↓ think your derivatives off

# Gradient Checking implementation Notes

- Don't use Grad check in or training (~~it~~ is slow)
  ( use only to debug )

- If algo fails grad check, look at components of $d\theta_{approx}$ and $d\theta$ to identify bug
  - check which values of higher diff
    figure out which $d\theta$ and hence which $dw$
    and $db$ are causing it.

- Remember about regularization terms

- Grad check doesn't work with dropout
  $\downarrow$
  cz v random.

  ( you can fix which nodes to drop once
    then do Grad check )

- Run at random initialization.
  $\downarrow$

  maybe your implementation of backprop, gradient
  is only correct when w, b are close to 0

  so start with random initialization, run grad check
  & it train for some time then run grad
  check again