

# TensorFlow

## Lecture 1 - TensorFlow Basics & Gradient Descent

pip install TensorFlow

rank n tensor ( $\Rightarrow$  n dimensional array)

TensorFlow datatype  $\neq$  numpy datatype

$a = \text{tf.constant}(2, \text{tf.float32})$

`print(a.numpy())` casts as a normal float.

Two types of matrix multiplication!

$A^{n \times m}, B^{m \times l}$

Hadamard product

maths

$$c_{ij} = a_{ij} b_{ij}$$

python

`tf.multiply(a, b)`

Standard  $C^{n \times l}$

$$c_{ij} = \sum_{k=1}^m a_{ik} b_{kj}$$

`tf.matmul(a, b)`

Obviously, this is

$$A^{n \times m} \times B^{n \times m} \rightarrow C^{n \times m}$$

E.g.

so different matrix dimensions required.

$a = \text{tf.constant}([[1, 2], [3, 4]])$

$b = \text{tf.constant}([[1, 1], [1, 1]])$

$c = \text{tf.matmul}(a, b)$

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 3 \\ 7 \end{pmatrix}$$

can't add datatypes  
that don't match  
 $\Rightarrow$  you must cast

Datatypes

int32, float32

Should always specify datatype when you create a tensor.

`t = cast(a, tf.int32)`

to cast from one to another.

True/False is boolean

$\xrightarrow{\text{cast}}$

1, 0 in `tf.int32`

## Tensor Shape

can't add matrices of different sizes...

But there is 'broadcasting'

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} + 1 = \begin{pmatrix} 2 & 3 \\ 4 & 5 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} + (10, 20) = \begin{pmatrix} 11 & 22 \\ 13 & 24 \end{pmatrix}$$

## Elementwise Tensor operations

easy. `tj. FUNCTION(a). numpy()`

comparisons too `tj. operator ([1, 2, 3], [2])`  
 $\Rightarrow [\text{false}, \text{false}, \text{true}]$

Wedge operator `tfp.Brauel`

$$a \star b = \text{tf. multiply}(a, b)$$

[Hadamard product]

## Variables

`tj. Variable(—)` can be reassigned

$$x = \text{tj. constant}([0.3], \text{tf. float32})$$

$$x = \text{tj. constant}([-1], \text{tf. float32})$$

also reassess...

but need to declare variable type..

`w. assign-add` / `w. assign-sub`

analogous to `+=` or `-=`

## Aggregation Functions

functions	actions
<code>tj(). reduce - sum(a)</code>	sums elements
<code>tj(). reduce - mean(a)</code>	mean of elements
<code>tj(). reduce - max(a)</code>	max of elements
<code>tj(). argmax(a)</code>	counts index where max appears

a could  
be a  
matrix  
etc...

### Example

$$a = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}, \text{ float } \geq$$

$$tj(). reduce\_sum(tj(). cast(a > 1, tj(). Float(3) \geq)). numpy() = \boxed{5}$$

*if entry bigger than 1, set as true*

true  $\rightarrow$  cast 1  
 false  $\rightarrow$  cast 0

Sum overall result.

### Algebra on arrays

$$a = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

Coordinate	Ax is	Output
<code>tj(). reduce - sum(a, axis=0)</code>	$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$	[4 6]
<code>tj(). reduce - sum(a, axis=1)</code>	$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$	[3 7]

## Automatic Differentiation

- NOT numerical or symbolic  
 → give well to node  $f(x)$ . It applies chain rule all the way down.

E.g -

```
x = tf.Variable(3.0, tf.float32)
with tf.GradientTape() as tape:
    y=tf.pow(x,2.0)
dydx=tape.gradient(y, x)
print(dydx.numpy())
6.0
```

$$\begin{aligned}y &= x^2 \\ \frac{dy}{dx} &= 2x \\ \left. \frac{dy}{dx} \right|_{x=3} &= 6\end{aligned}$$

```
x=tf.Variable(4.0,tf.float32)
y=tf.Variable(2.0,tf.float32)
with tf.GradientTape() as tape:
    f=tf.pow(x,2.0)*3.0+y
[dfdx, dfdy]=tape.gradient(f, [x,y])
```

```
print(dfdfx.numpy(), dfdy.numpy())
24.0 1.0
```

$$\begin{aligned}f(x, y) &= 3x^2 + y \\ \frac{\partial f}{\partial x} &= 6x \\ \frac{\partial f}{\partial y} &= 1\end{aligned}$$

Difg not considered here..

```
x=tf.constant(4.0,tf.float32)
with tf.GradientTape() as tape:
    tape.watch(x)
    f=tf.pow(x,3.0)
dfdx=tape.gradient(f, x)
print(dfdfx.numpy())
48.0
```

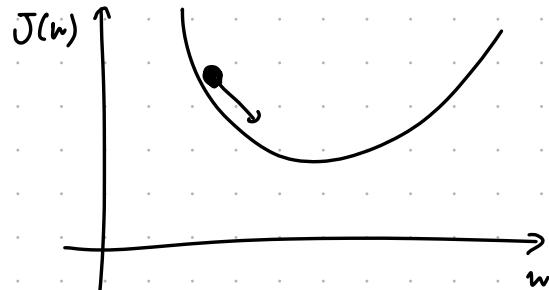
Can use auto diff on

$$J: \mathbb{R}^{2 \times 2} \times \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$J(w, x, y) = \| \tanh(wx) - y \|^2$$

### Introduction to Gradient Descent

$$w_{t+1} = w_t - \gamma \frac{dJ}{dw}$$



### Example

Find minimum of  $y = x^2 - 4x + 4$

```
import tensorflow as tf

eta = 0.1 # learning rate
x = tf.Variable(10.0, tf.float32) # arbitrary initial value

for i in range(50):
    with tf.GradientTape() as tape:
        y=#TODO put in formula for y in terms of x here
        dydx=tape.gradient(# TODO finish this line
        x.assign(#TODO finish x_(t+1)=x_t-eta*dydx
        print("iteration:",i, "x:", x.numpy(), "y:", y.numpy())
```

These are also built in optimizers!

We can speed this up wrt by more derivative calculations outside of the loop & using @tf.function.

def do\_update(x): & instead, use x.assign(→)

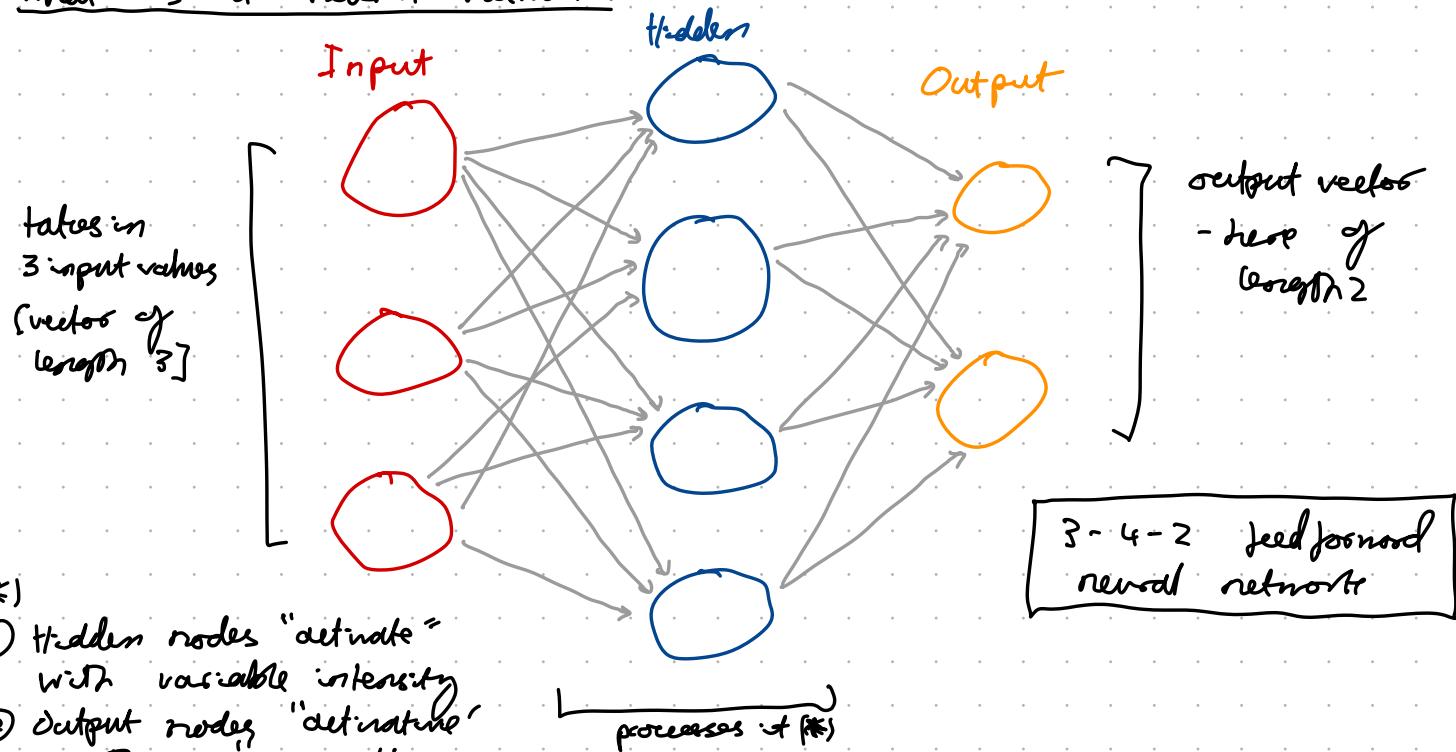
## Lecture 2 - Neural Networks

Note: deep learning is deep neural networks. Deep learning is a subset because neural networks started nothing. Why?

- ① More training data
- ② More CPU/GPU power
- ③ Clever neural net architecture
  - CNNs
  - Dropout
  - LSTM nodes
  - Transformers
- ④ Better activation functions & weight initializations.
  - ReLU
  - Xavier initialization
- ⑤ Better learning algorithms (e.g. Adam)

breakthroughs  
that made  
machine  
learning  
today

### What is a neural network?



A feedforward neural network:  $y = f(w, x)$

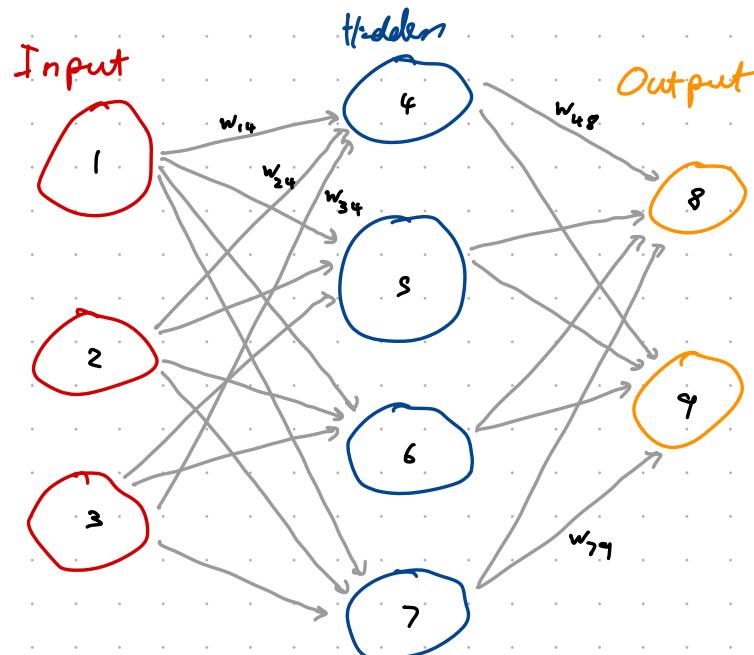
- $w$  - weights (a list of tensors)
- Give enough appropriately set weights & neural networks can represent ANY function. [Universal function approximators]

Example	Input	Output
Vision system	grid of pixel intensities over an image	classification of image (car / house)
Language translation	vector of words in English sentence	vector of words in French sentence

# Feed forward AlexNet

- Number the nodes
  - Label the weights  
-  $w_{ij}$  connects node  $i$  to node  $j$ .
  - Input vector to this network must be dimension 3  
 $x = (x_1, x_2, x_3)$
  - Die "electrical" signals passing along the wires to node 4  
Therefore total to
- $$S_4 = x_1 w_{14} + x_2 w_{24} + x_3 w_{34}$$

$$S_k = \sum_{k=1}^3 x_k w_{ek}$$



- Activation of node 4: Node 4 will fire when  $S_4$  is bigger than some threshold  $b_4$

$$\left( a_4 = \begin{cases} 1 & \text{if } S_4 > b_4 \\ 0 & \text{if } S_4 \leq b_4 \end{cases} \right) \Leftrightarrow \left( a_4 = g(x, w_{14} + x_2 w_{24} + x_3 w_{34} - b_4) \right)$$

where  $g(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$

activation of node 4

In total:

$$\begin{aligned} a_4 &= g(x, w_{14} + x_2 w_{24} + x_3 w_{34} - b_4) \\ a_5 &= g(x, w_{15} + x_2 w_{25} + x_3 w_{35} - b_5) \\ a_6 &= g(x, w_{16} + x_2 w_{26} + x_3 w_{36} - b_6) \\ a_7 &= g(x, w_{17} + x_2 w_{27} + x_3 w_{37} - b_7) \end{aligned}$$

or with matrices:

$$(a_4 \ a_5 \ a_6 \ a_7) = g\left(\begin{pmatrix} x_1 & x_2 & x_3 \end{pmatrix} \begin{pmatrix} w_{14} & w_{15} & w_{16} & w_{17} \\ w_{24} & w_{25} & w_{26} & w_{27} \\ w_{34} & w_{35} & w_{36} & w_{37} \end{pmatrix} - \begin{pmatrix} b_4 & b_5 & b_6 & b_7 \end{pmatrix}\right)$$

And for the output layer:

$$(a_8 \ a_9) = g\left((a_1 \ a_2 \ a_3 \ a_4) \begin{pmatrix} w_{48} & w_{49} \\ w_{58} & w_{59} \\ w_{68} & w_{69} \\ w_{78} & w_{79} \end{pmatrix} - (b_8 \ b_9)\right)$$

By "training" the neural network, we aim to find values of the  $w_{ij}$  and  $b_i$  so the neural network behaves as we want.

More correctly:

$$h_1 = (a_1 \ a_2 \ a_3 \ a_4) \quad b_1 = -(b_4 \ b_5 \ b_6 \ b_7)$$

Hidden layer:  $h_1 = g(x, w_1 + b_1)$

$$w_1 = \begin{pmatrix} w_{14} & w_{15} & w_{16} & w_{17} \\ w_{24} & w_{25} & w_{26} & w_{27} \\ w_{34} & w_{35} & w_{36} & w_{37} \end{pmatrix}$$

Output layer:  $y = g(h_1, w_2 + b_2)$

$$x = (x_1 \ x_2 \ x_3)$$

rate  
minus  
(forwards for  
 $w_2, b_2$ )  
 $y = (a_8 \ a_9)$

Note: for gradient descent to work, the activation functions  $g(x)$  needs to be differentiable so  $g(x) = 1_{x>0}$  doesn't work!

$$\frac{e^x - e^{-x}}{e^x + e^{-x}} \rightarrow 1$$

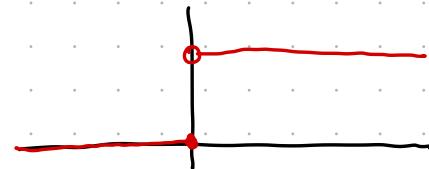
Change it to

$$g(x)$$

Picture

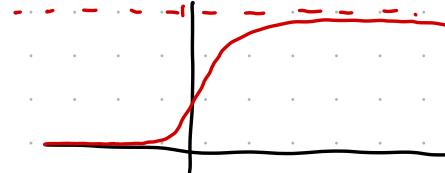
now

$$g(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$



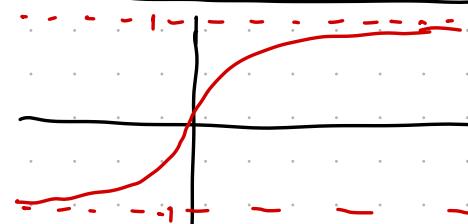
step functions

$$g(x) = \frac{1}{1 + e^{-x}}$$



logistic sigmoid function

$$g(x) = \tanh(x)$$



Note: Anything non-linear, smooth, increasing works...

$$g(x) = \frac{1}{1 + e^{-x}}$$

in maths

in TensorFlow code

for  
input  
vector  
 $x$

$$h_1 = g(x w_1 + b_1)$$

`h1=tf.sigmoid(tf.matmul(x,W1) + b1)`

$$y = g(h_1 w_2 + b_2)$$

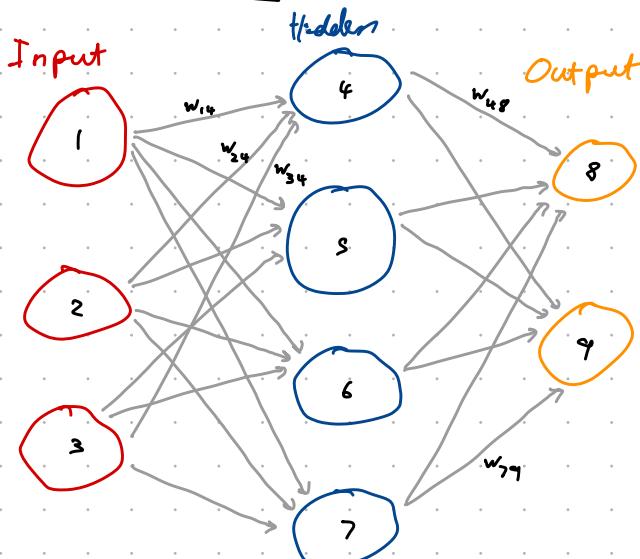
`y=tf.sigmoid(tf.matmul(h1,W2) + b2)`

Can easily add a second hidden layer:

```
h1=tf.sigmoid(tf.matmul(x,W1) + b1)
h2=tf.sigmoid(tf.matmul(h1,W2) + b2)
y=tf.sigmoid(tf.matmul(h2,W3) + b3)
```

(as only matrix multiply if both arguments have equal ranks.  $\Rightarrow \text{rank}(x) = 2$ )

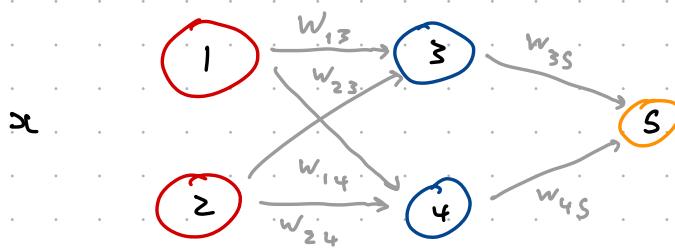
Shape cheat:



Variable	Shape
$x$	[1, 3]
$w_1$	[4, 3]
$w_2$	[4, 2]
$b_1$	[1, 4]
$b_2$	[1, 2]
$h_1$	[1, 4]
$y$	[1, 2]

*weights matrices "kernels"*

Exercise: Build a 2-2-1 feedforward network



$$h_1 = (a_3 \ a_4)$$

$$h_1 = g(xw_1 + b_1)$$

$$y = g(h_1 w_2 + b_2)$$

In full:

$$(a_3 \ a_4) = g \left( (x_1 \ x_2) \begin{pmatrix} w_{13} & w_{14} \\ w_{23} & w_{24} \end{pmatrix} - (b_3 \ b_4) \right)$$

$$a_S = g \left( (a_3 \ a_4) \begin{pmatrix} w_{3S} \\ w_{4S} \end{pmatrix} - b_S \right)$$

```
import tensorflow as tf

x=tf.constant([[1,1]], tf.float32) # our input vector
tf.random.set_seed(1) # this is just for teaching
# Build our random weight and bias matrices, of appropriate shapes
W1 = tf.Variable(tf.random.truncated_normal([#, #], stddev=0.1))
b1 = tf.Variable(tf.random.truncated_normal([#, #], stddev=0.1))
W2 = tf.Variable(tf.random.truncated_normal([#, #], stddev=0.1))
b2 = tf.Variable(tf.random.truncated_normal([#, #], stddev=0.1))

# define our feed-forward neural network here:
def run_network(x):
    h1=#TODO
    y=#TODO
    return y

print(run_network(x).numpy())
```

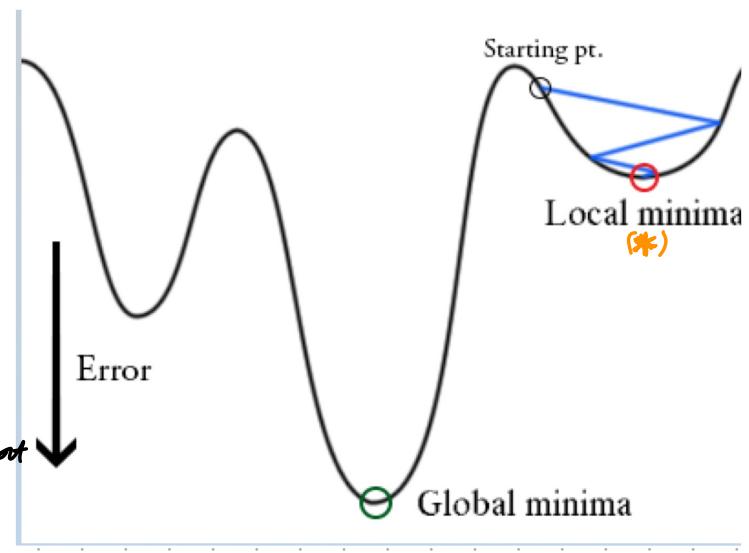
easy now!  
Just put in the numbers...

This script is in lecture2\_notebook.ipynb under "Exercise 1". Complete it now.

### Randomizing initial weights

Must randomize weights in a neural network before training.

- Symmetry breaking for gradient descent
- Randomized start point for gradient descent.
- Need multiple starts from different random weights to avoid (\*)
- magnitude **VERY** important...



## Learning the XOR function

Need to choose the weight & bias tensors to make the network behave as desired.

Now need to modify our feedforward network to pass through all 4 input vectors.

Input		Output
A	B	
0	0	0
1	0	1
0	1	1
1	1	0

XOR Truth Table

- Change `x=tf.constant([[1,1]], tf.float32)`  
To `x=tf.constant([[0,0],[0,1],[1,0],[1,1]], tf.float32)`

Now, instead of passing in  $x = (x_1, x_2)$  we pass in a vector of vectors. 4 different  $1 \times 2$  vectors, so our output (transposed to broadcast with  $b_1$  &  $b_2$ ) is  $1 \times 4$ .

We output  $y = \begin{pmatrix} 0.531... \\ 0.531... \\ 0.531... \\ 0.531... \end{pmatrix}$  but we want  $t = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$

The squared error between output & target is

$$L = \|y - t\|^2$$

We train the network by doing gradient descent on this.

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

TensorFlow will calculate this for us using Autodiff (Backpropagation)

where  $w = (w_1, w_2, b_1, b_2)$  ← all weights & all biases.

## Defining our loss Function

$$y\_labels = tf.constant([0, 1, 1, 0])$$

```
def calc_loss():
```

$$L = \|y - t\|^2$$

`y=run_network(x)`

$$\delta_i = y_i - t_i$$

`deltas=tf.subtract(y, y_labels)`

$$L = \sum_i \delta_i^2$$

`squared_deltas=tf.square(deltas)`
`loss=tf.reduce_sum(squared_deltas)`
`return loss`

Now, training our network

This lets us implement gradient descent.

`optimizer = tf.keras.optimizers.SGD(0.5)`
`for i in range(20000):`

variables to minimize

`optimizer.minimize(calc_loss, [W1,b1,W2,b2])`

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

`if (i%1000)==0:`

input function

`print("iteration ",i," loss", calc_loss().numpy())`
`print(run_network(x).numpy())`

These are going on in these three lines \*

optimizer = tf.keras.optimizers.SGD(eta)

for i in range(20000):

    optimizer.minimize(calc\_loss, [W1,b1,W2,b2])

Is short for

$$w \leftarrow w - \eta \frac{dL}{dw}$$

$$L = \|y - t\|^2$$

for i in range(20000):

    with tf.GradientTape() as tape:

        L=calc\_loss()

$$\# \text{define our } L = \|y - t\|^2$$

[dLdW1,dLdW2,dLdb1,dLdb2]=tape.gradient(L, [W1,W2,b1,b2]) calculate partial derivatives

W1.assign(W1-dLdW1\*eta)

$$\frac{dL}{dw} = \left( \frac{dL}{dw_1}, \frac{dL}{dw_2}, \frac{dL}{db_1}, \frac{dL}{db_2} \right)$$

W2.assign(W2-dLdW2\*eta)

short hand for backpropagation auto diff calculation!

b1.assign(b1-dLdb1\*eta)

b2.assign(b2-dLdb2\*eta)

These lines implement gradient descent:

$$w_1 \leftarrow w_1 - \eta \frac{dL}{dw_1}$$

$$w_2 \leftarrow w_2 - \eta \frac{dL}{dw_2}$$

$$b_1 \leftarrow b_1 - \eta \frac{dL}{db_1}$$

$$b_2 \leftarrow b_2 - \eta \frac{dL}{db_2}$$

A LOT going on & getting condensed into just three lines!

These give new & improved data  
for  $(w_1, w_2, b_1, b_2) = w$ .

A "pattern" is a pairing of input & output vectors. Our neural network has learnt 4 patterns.

$$x = \begin{pmatrix} (0,0) \\ (1,0) \\ (0,1) \\ (1,1) \end{pmatrix} \rightarrow y = \begin{pmatrix} 0.01094\dots \\ 0.9874\dots \\ 0.9896\dots \\ 0.0098\dots \end{pmatrix}$$

## Running Neural Networks on New Data

Now, lets evaluate our network on a new pattern (make a prediction based on our training data).

print(run\_network([[0.5,0.5]]).numpy())

Now, restored the random seed function, will make a table to see if we always converge to a valid solutions:

Attempt #	Loss at iteration			x value	
	0	10,000	20,000	$((0,0), (1,0), (0,1), (1,1))$	$(0.5, 0.5)$
1	1.0...	0.0007	0.0003	$(0.009, 0.99, 0.99, 0.008)$	0.9918
2	1.0...	0.9999...	0.999	$(0.06, 0.93, 0.92, 0.05)$	0.062
3	1.0...	0.9999...	0.00124	$0.015, 0.985, 0.98, 0.014$	0.0148
4	1.0...	0.0007	0.0003	$0.007, 0.99, 0.99, 0.009$	0.9919

doesn't seem to always converge to zero C

## Limitations of XOR

To train deeper networks than XOR, we need

- better weight initializations
- better activation functions
- more training data
- more CPU time

You can access the weights and biases for a layer as

- "layer1.trainable\_variables", which evaluates as [W1,b1]
- Note: You must run the network with an input tensor x before accessing trainable\_variables though

M.Fairbank, University of F...

## Introducing keras layers

### Functions

Create two  
network  
layers

### Before

```
W1 = tf.Variable(tf.random.truncated_normal([2,2], stddev=0.1))
b1 = tf.Variable(tf.random.truncated_normal([1,2], stddev=0.1))
W2 = tf.Variable(tf.random.truncated_normal([2,1], stddev=0.1))
b2 = tf.Variable(tf.random.truncated_normal([1,1], stddev=0.1))
def run_network(x):
    h1=tf.sigmoid(tf.matmul(x,W1) + b1)
    y=tf.sigmoid(tf.matmul(h1,W2) + b2)
    return y
```

creates  
new  
for us  
random  
tensor  
us!

### After

```
layer1=tf.keras.layers.Dense(2, activation=tf.nn.sigmoid)
layer2=tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
def run_network(x):
    h1=layer1(x)
    y=layer2(h1)
    return y
```

layer 1 & layer 2  
act as callable  
functions

Even  
easier:

```
layer1=tf.keras.layers.Dense(2, activation=tf.nn.sigmoid)
layer2=tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
def run_network(x):
    h1=layer1(x)
    y=layer2(h1)
    return y
```

```
layer1=tf.keras.layers.Dense(2, activation=tf.nn.sigmoid)
layer2=tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
model = tf.keras.Sequential([layer1,layer2])
def run_network(x):
    return model(x)
```

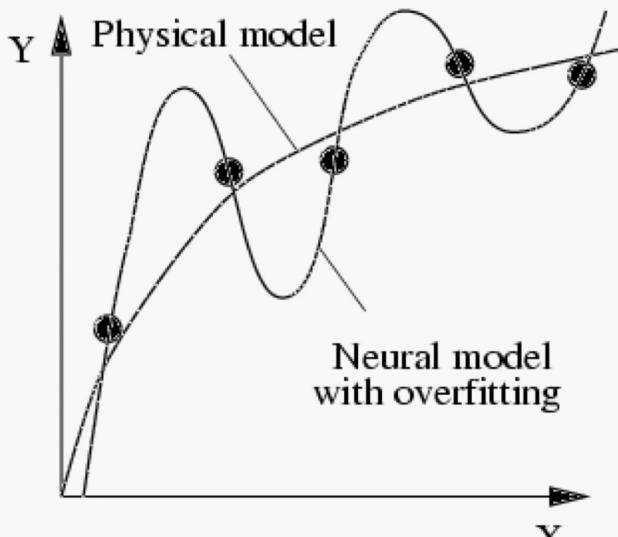
Fit loop

```
optimizer = tf.keras.optimizers.Adam()
for i in range(10000):
    optimizer.minimize(calc_training_loss, trainable_variables)
    if (i%100)==0:
        [train_loss, train_acc, train_cross_entropy]=calc_loss_and_accuracy(inputs_train, labels_train)
        [test_loss, test_acc, test_cross_entropy]=calc_loss_and_accuracy(inputs_test, labels_test)
        print("iteration : ", iteration, " loss: ", train_loss.numpy(), " accuracy: ", train_acc.numpy(), " testLoss: ", test_loss.numpy(), " test_accuracy: ", test_acc.numpy())
```

```
optimizer = tf.keras.optimizers.Adam()
model.compile(optimizer=optimizer, loss='mse') # This uses the "mean squared error" as the loss function
history=model.fit(x, y, labels, epochs=10000, batch_size=x.shape[0])
```

model.trainable\_variables will give us  
[W1,b1,W2,b2]

This is how most people  
train tensorflow + keras  
neural networks



"Fit" because trying to fit our neural network's behaviour to match the data.

BUT, really its running the gradient descent training loop in full. This is what we've done by hand above.

## Understanding Neural Networks video

- Neurons can either fire or not, a 0 or 1. one of (0) functions
- The input of one neuron is the output of many others. Each can be given a weight. These can excite or inhibit.

- You train to adjust weights.
- The bigger the network:
  - more layers
  - more nodes per layer
- The better the accuracy will be.
- Activation could be step/cat/etc.

## 2-2-1 FOR neural network training

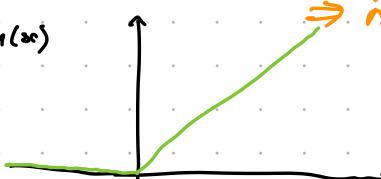
```
import tensorflow as tf

layer1=tf.keras.layers.Dense(2, activation=tf.nn.sigmoid)
layer2=tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
model = tf.keras.Sequential([layer1,layer2])

print("Untrained outputs", model(x))
print("trainable vars",len(model.trainable_variables))
print("trainable weights",len(model.trainable_weights))
optimizer = tf.keras.optimizers.SGD(0.5)
model.compile(optimizer=optimizer,
              loss='mse') # This uses the "mean squared error" as the loss function
print("Running training loop...")
history=model.fit(x, y_labels, epochs=10000, batch_size=x.shape[0],
                   verbose=1 # change this to 1 to show output
                  )
print("Finished training")
print("Trained outputs",model(x).numpy())
Epoch 9994/10000
1/1 [=====] - 0s 28ms/step - loss: 6.1407e-04
Epoch 9995/10000
1/1 [=====] - 0s 24ms/step - loss: 6.1400e-04
Epoch 9996/10000
1/1 [=====] - 0s 18ms/step - loss: 6.1392e-04
Epoch 9997/10000
1/1 [=====] - 0s 15ms/step - loss: 6.1385e-04
Epoch 9998/10000
1/1 [=====] - 0s 16ms/step - loss: 6.1378e-04
Epoch 9999/10000
1/1 [=====] - 0s 18ms/step - loss: 6.1370e-04
Epoch 10000/10000
1/1 [=====] - 0s 21ms/step - loss: 6.1363e-04
Finished training
Trained outputs [[0.02736123]
 [0.97644806]
 [0.9764228 ]
 [0.02439287]]
```

Tricks for Improving Network Training

## ① Better activation functions

- $g(x) = \frac{1}{1+e^{-x}}$  (logistic sigmoid) is bad. Vanishing gradient possible, then weights don't change because  $w \leftarrow w - \eta \frac{\partial L}{\partial w}$  ← If this is small, weights don't change → network doesn't learn.
  - $g(x) = \tanh(x)$
  - $g(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases}$  is best. e.g. nn.relu(x)
- Relu (Rectified Linear Function)**
- 
- $$g(x) = \max(0, x)$$

## ② Better weight initialisations

- Need random initialised weights to help gradient descent & break symmetry. Need to avoid vanishing / exploding learning gradients ( $\frac{\partial L}{\partial w_i}$ ) so magnitude of initial weights is important.
- Glorot/Xavier initialisation. Set the standard deviation for a weight matrix of dimension  $n \times m$

$$\sigma = \sqrt{\frac{2}{n+m}}$$

This formula varies depending on:

- Activation functions used [how much input before neuron fires?]
- Weight distributions used ( $\geq$  surely always normal?)
- What exact research you follow.

Easy to implement when using higher level functions to create your network layers...

```
layer1=tf.keras.layers.Dense(2, activation=tf.nn.sigmoid, kernel_initializer='glorot_uniform')
```

Preparing Your Training DataInput Pre-Processing

- Neural Networks work best when input vectors normalized to have variance 1. For each 12 columns of your dataset input tensor, subtract off column means & divide each column by the column s.d.

```
import numpy as np
inputs=[[0.0,0],[0,10],[10,0],[10,1]]
column_means=np.mean(inputs, axis=0)
inputs-=column_means
column_stds=np.std(inputs, axis=0)
inputs/=column_stds
```

$$\text{inputs} = \begin{pmatrix} 0 & 0 \\ 0 & 10 \\ 10 & 0 \\ 10 & 1 \end{pmatrix}$$

Column means:  $S_{11} = 5, S_{12} = 5, S_{21} = 5, S_{22} = 5$

$$\text{Subtracted inputs} = \begin{pmatrix} -5 & -5 \\ -5 & 5 \\ 5 & -5 \\ 5 & -5 \end{pmatrix}$$

$$\sigma = (\sigma_1, \sigma_2)$$

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

- Can also do "whitening", scales very well: if the input data is multivariate gaussian, then the whitened data will be gaussian w/ zero mean & identity covariance matrix → remove undesirable correlations, eigenvectors vectors scoreless ...

$$\text{inputs normalized} = \begin{pmatrix} -\frac{5}{\sqrt{50}} & -\frac{5}{\sqrt{50}} \\ -\frac{5}{\sqrt{50}} & \frac{5}{\sqrt{50}} \\ \frac{5}{\sqrt{50}} & -\frac{5}{\sqrt{50}} \\ \frac{5}{\sqrt{50}} & -\frac{5}{\sqrt{50}} \end{pmatrix}$$

Beginners often forget this step.

**VERY IMPORTANT  
TO NORMALISE  
YOUR INPUTS!!!**

## Categorical Inputs of One-hot encoding

When there is no logical order to data numerically, don't give as a single numerical input into the neural network (NN). Instead use 4 different inputs (one-hot encoding):

Teacher

Teacher	Doctor	Student	Unemployed
1	2	3	4

NOT THIS  
(single numerical input)

$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$	Doctor
$\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$	Student
$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	Unemployed

THIS

(4 numerical inputs)

### Example

- Age of a child: 1 input
- Colour of eyes: multiple inputs (Brown, Blue, Green etc)
- Position of tennis ball: 2 inputs ( $x$  &  $y$ )
- Alphabet vision system:  $26 \times 2$  characters + punctuations.

Use tf.one-hot functions to do this.

Seems like you can reverse it with argmax ???

## Training Classification Networks

For say a vision system (classification), we want the NN to output probabilities. Say we have chair, table, door, bed. We want a output vector  $y = (-0.2, 0.5, 0.1, 0.1) \rightarrow ?$

$$p = (0.175, 0.382, 0.237, 0.287)$$

$$p_i = \frac{e^{y_i}}{\sum_k e^{y_k}}$$

softmax function

makes positive

normalizes so they sum to 1.

### Cross-entropy loss function

Our NN now outputs a p vector with  $p_i = \frac{e^{y_i}}{\sum_k e^{y_k}}$ . We want these to become TRUE probabilities. To do this, "train" the NN w/ a "cross entropy loss function". (Physics???) For a single pattern, this is defined to be

$$L = - \sum_i t_i \log(p_i)$$

CROSS ENTROPY LOSS FUNCTION

where  $t$  (vector) is the one-hot encoded true output classification &  $i$  is the category index.

### Example

Consider a "table",  $t = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$ , actual NN output  $y = \begin{pmatrix} -0.2 \\ 0.5 \\ 0.1 \\ 0.1 \end{pmatrix}$

Category index of "table" is 1 so

	chair	table	door	bed
$t_i$	$t_0$	$t_1$	$t_2$	$t_3$
value	0	1	0	0

so here,  $L = -\log(p_1)$

THIS MEANS: when trying to minimise  $L$ , it has to minimise  $-\log(p_1)$ .  $\Rightarrow$  maximise  $p_1 \Rightarrow$  maximise  $y_1$  compared to  $y_0, y_2, y_3$ .



## Where does the formula come from?

Can show (proof??) that:

$$\left( \text{minimizing } L = - \sum_i t_i \log(p_i) \text{ over your whole dataset} \right) \Leftrightarrow$$

Example

Finding the weights which maximise your NN's estimation of the total probability that your given training dataset members each have their given (fixed) category label.

- ① Train NN w/ 60 cat, 40 dog images. Then NN should learn to output  $P(\text{cat}) \approx 1$  for each cat image,  $P(\text{dog}) \approx 1$  for every dog image. Then hope w/ new dog image, get  $P(\text{dog}) \approx 1$  for the new data.
- ② 100 identical photos, 40 labelled "1", 60 labelled "2"; then NN

[see previous reading... function is clever!]

## In Tensorflow

SoftMax in TensorFlow:

```
y=tf.keras.activations.softmax(y, axis=1)
```

Cross entropy in TensorFlow:

```
cce=tf.keras.losses.SparseCategoricalCrossentropy()  
cross_entropy = tf.reduce_mean(cce(y_true=train_labels, y_pred=y))
```

- This expects `y_pred` to have had softmax applied to it on entry.
- Also it expects `train_labels` to be an integer for the category number (so the label is NOT one-hot encoded)

The variable `cross_entropy` is returned by our "calc\_loss()" function, to optimise through gradient descent

```
optimizer.minimize(calc_loss, trainable_variables)
```

## MNIST Problem

Breakdown of code

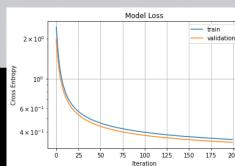
- Load data in two, (`train_images, train_labels`) & save for test.
- Pick some pictures
- Reshape the image from matrix to vector. Rescale from greyscale to float.
- Create a 1 layer deep 10 node network.
- SGD to optimise, cross entropy loss function.
- Plot Loss vs iteration #
- Test against new data.

Floating between 70 - 95% accurate.

## Using Matplotlib

To plot the graph just use:

```
import matplotlib.pyplot as plt  
plt.plot(history.history['loss'], label="train")  
plt.plot(history.history['val_loss'], label="validation")  
plt.title('Model Loss')  
plt.yscale('log')  
plt.ylabel('Cross Entropy')  
plt.xlabel('Iteration')  
plt.grid()  
plt.legend()  
plt.show()
```



## Using Matplotlib

To log results for plotting, note that the keras "fit" method returns a "history" variable:

```
history = keras_model.fit(..)
```

This contains a dictionary:

```
history.history = {"loss": [...], "val_loss": [...], "sparse_categorical_accuracy": [...], "val_sparse_categorical_accuracy": [...]}
```

Each of the values of this dictionary is a numeric array, which we can plot.

## Summary: Classification vs Regression Problems



- Use softmax to obtain probabilities

$$p_i = \frac{e^{y_i}}{\sum_k e^{y_k}}$$

- Use cross entropy loss function to train

$$L = - \sum_i t_i \log(p_i)$$

- Use max to pick NN's favoured category

Example: Iris Dataset

4 inputs, 3 outputs. See notebook for code...

### Fighting Overfitting

whole point of machine learning is to learn a model that's useful on unseen data. To stop it, we can

- ① Early stopping: stop training when orange curve  $T$ .
- ② More training data: impossible to learn it all
- ③ Use a simpler model: fewer hidden layers/nodes, constrain the weights to be smaller (Regularisation,  $L^2$ ). Think Occam's Razor.
- ④ Use Dropout (major breakthrough!)
- ⑤ Combine multiple neural networks (ensemble learning)

### $L^2$ Regularisation

$$\text{Total Loss} = \text{cross entropy} + k_2 \sum_i w_i^2$$

choose the constant  
 $k_2 = (0.01)$

We aim to constrain the weights to be smaller. How?

Add a  $\sum_i w_i^2$  term to the loss function! [for all neural weights  $w_i$ ]

Gradient descent minimises "Loss" so this'll force most weights to  $0$  in  $L^2$ .

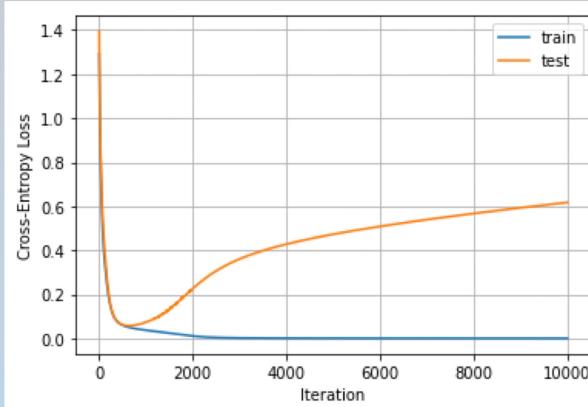
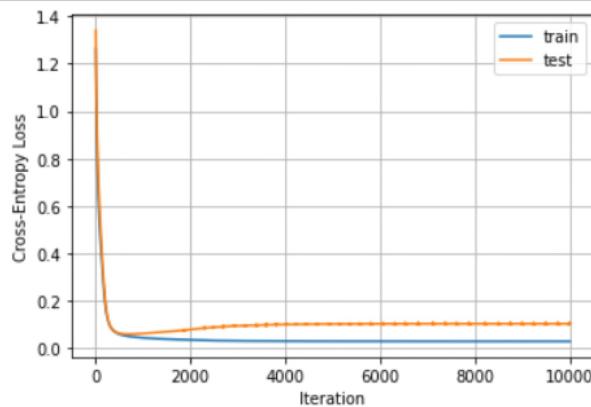
If all weights smaller  $\Rightarrow$  "simpler" model  $\xrightarrow[\text{razor}]{\text{occam's}}$  The simpler model that explains the data is more likely to be correct.

## TensorFlow code:

modifies loss functions to be composite.

```
hids=[4,20,20,3]
k_l2=0.001
layer1=tf.keras.layers.Dense(hids[1], activation='tanh', kernel_regularizer=keras.regularizers.l2(k_l2))
layer2=tf.keras.layers.Dense(hids[2], activation='tanh', kernel_regularizer=keras.regularizers.l2(k_l2))
layer3=tf.keras.layers.Dense(hids[3], activation='softmax', kernel_regularizer=keras.regularizers.l2(k_l2))
keras_model = tf.keras.Sequential([layer1,layer2,layer3])
```

## Results

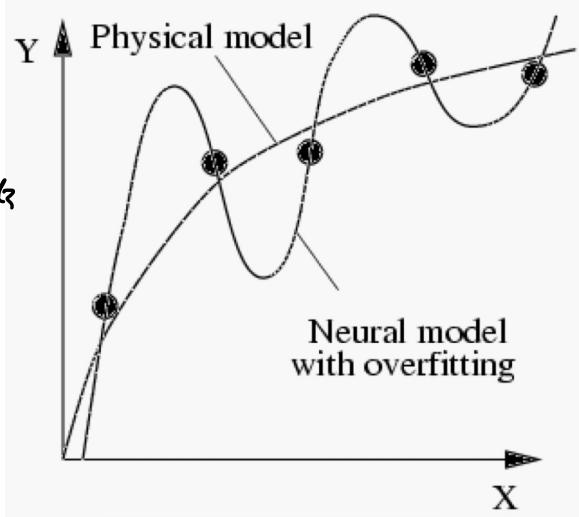


with  $L^2$  regularisation  
[by  $k_{L^2}$ , overfitting ↓ more]

without  $L^2$  regularisation

## Problems with $L^2$ regularisation

- When you apply  $L^2$  regularisation, we stiffen the curve on right, prevents it wiggling too much (stops overfitting). It also makes the curve less flexible.
- $L^2$  regularisation:
  - Increases "bias" (stiffness)
  - Reduces "variance" (flexibility)

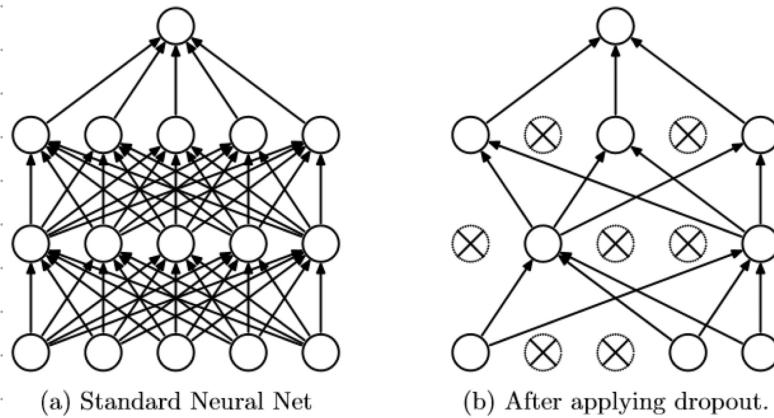


## Fighting Overfitting: Dropout

Modify hidden layers so nodes randomly switch off 50% of the time.

Radical approach, paper in 2014.

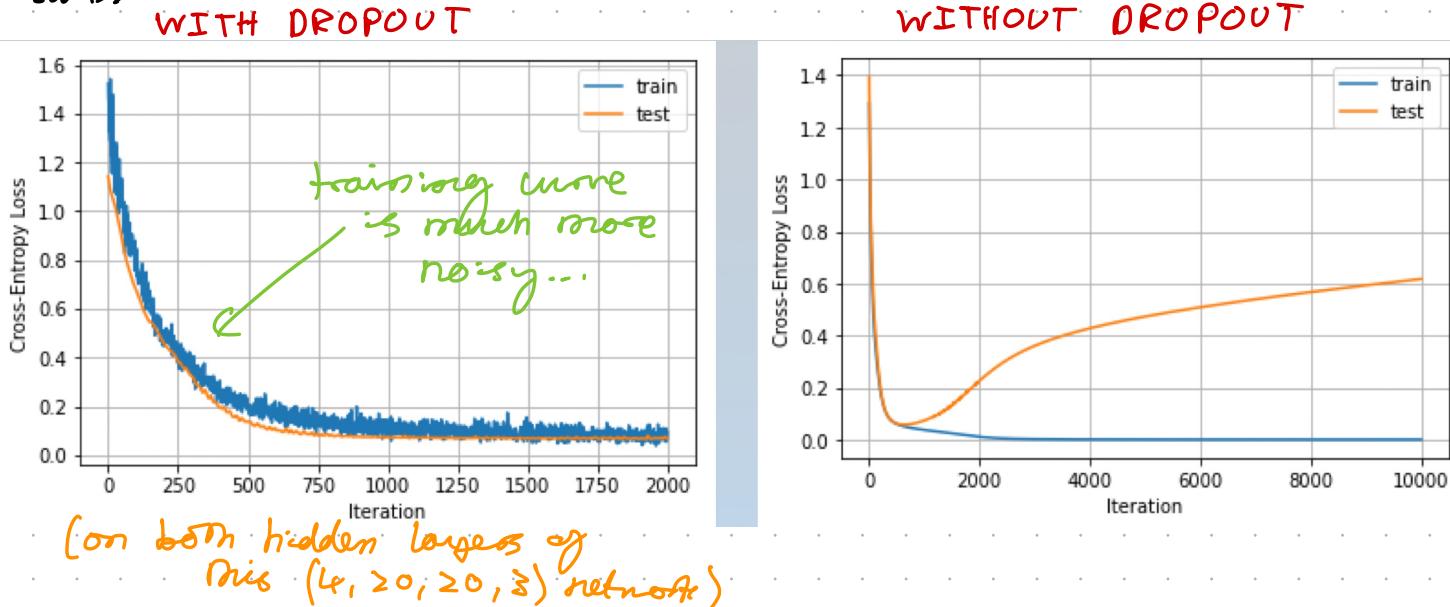
→ makes it much harder for neural network to simply memorise all the data.



## In TensorFlow code:

```
hids=[4,20,20,3]
layer1=tf.keras.layers.Dense(hids[1], activation=tf.tanh)
layer1do=tf.keras.layers.Dropout(rate=0.5)
layer2=tf.keras.layers.Dense(hids[2], activation=tf.tanh)
layer2do=tf.keras.layers.Dropout(rate=0.5)
layer3=tf.keras.layers.Dense(hids[3], activation=tf.keras.activations.softmax)
model = tf.keras.Sequential([layer1,layer1do,layer2,layer2do,layer3])
```

## Results:



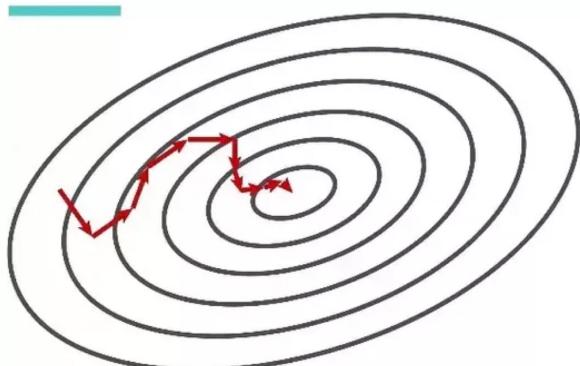
## Takeaways:

- Dropout >  $L^2$  regularisation (Sometimes)
- Battle overfitting w/ combo of  $L^2$  &/or dropout
- Dropout is a Big breakthrough in deep learning

## Using Mini-Batches

On massive datasets, you need miniBatches. MiniBatches changes "Gradient Descent" to "Stochastic Gradient Descent".

### Stochastic Gradient Descent



### Why better?

- Can shake gradient descent out of local minima.
- Improves generalisation.
- SGD can be the best learning algorithm (despite being slow)

We don't count # training iterations, we count # training epochs

$$\text{epoch \#} = \frac{\text{iterations} \times \text{minibatch size}}{\text{training set size}}$$

The fit loop will automatically shuffle mini-batches of size 20 for us...

```
history = keras_model.fit(  
    inputs_train,  
    labels_train,  
    batch_size=20,  
    epochs=2000,  
    validation_data=(inputs_test, labels_test),  
)
```

We can save weights too

This to load a model

# Save the final model:

```
keras_model.save("IrisModel")
```

```
model2 = keras.models.load_model('IrisModel')
```