

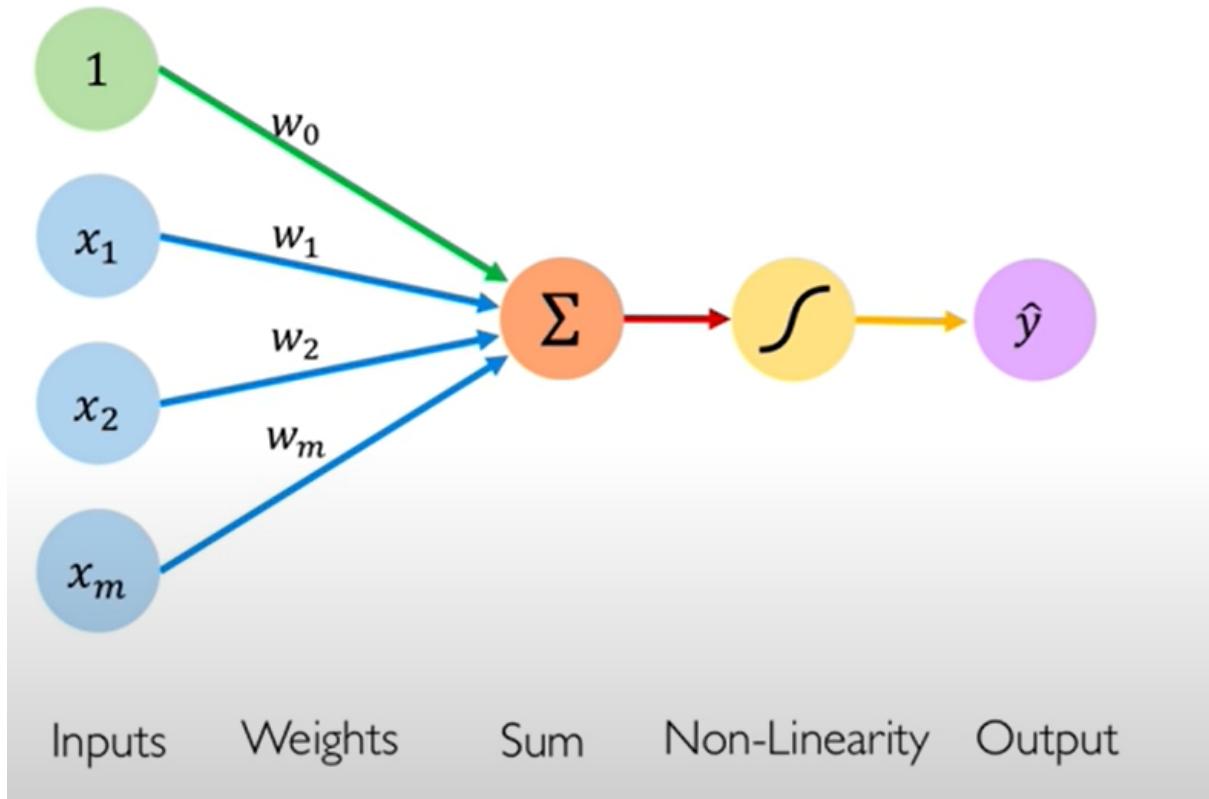


[1] Introduction to Deep Learning

What is Deep Learning

- Artificial intelligence - Any technique that enables computers to mimic human behavior
 - Machine Learning - Ability to learn without explicitly programmed
 - Deep Learning - Extract from data using neural networks

The Perceptron: Forward Propagation



- multiply the input with the weight and bias and take the sum of all of them
- Take the single number and pass it through a non-linear activation function and you get the output

The equation diagram shows the mathematical formulation of the neuron's computation. The output \hat{y} is given by the formula:

$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

Annotations explain the components: a purple arrow labeled 'Output' points to \hat{y} ; a red arrow labeled 'Linear combination of inputs' points to the summation term; a green arrow labeled 'Bias' points to w_0 ; and a yellow arrow labeled 'Non-linear activation function' points to the function g .

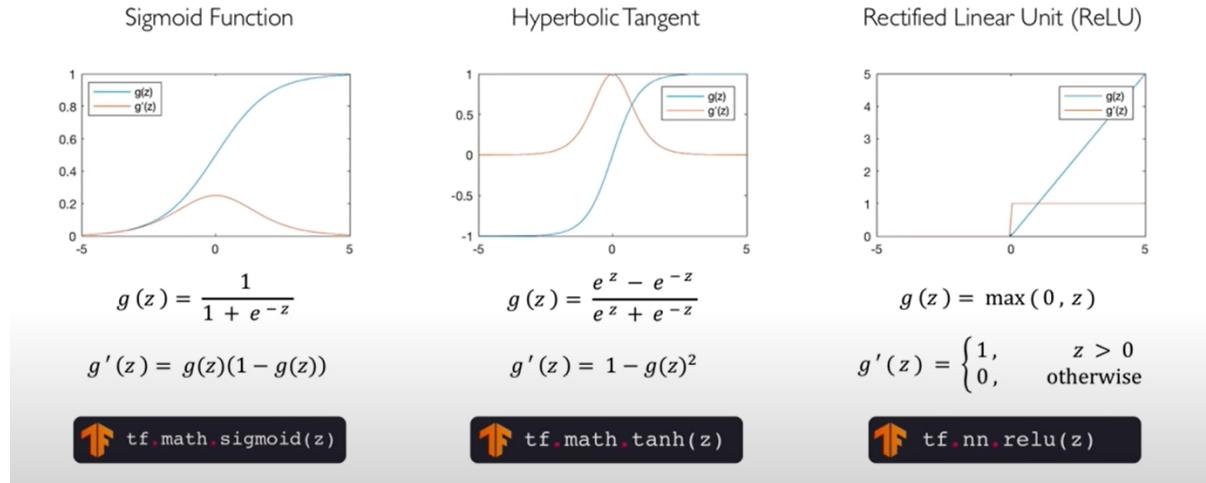
Below the main equation, another form is shown:

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

where: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

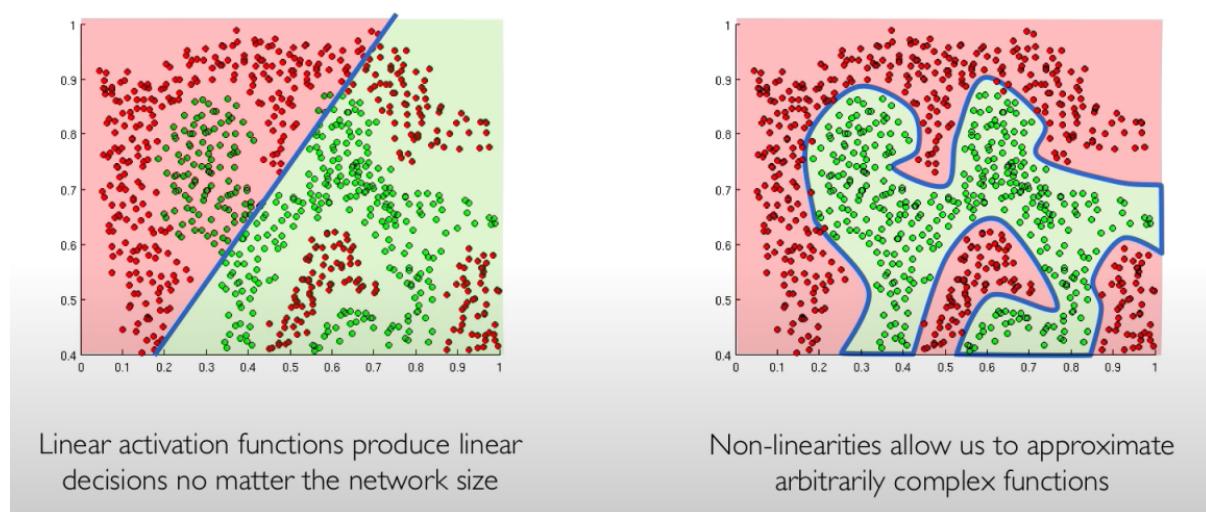
Activation Functions

- Sigmoid
 - It takes transfers the output into a scalable output between 0 and 1
- tanh
 - number between -1 to 1
- ReLU

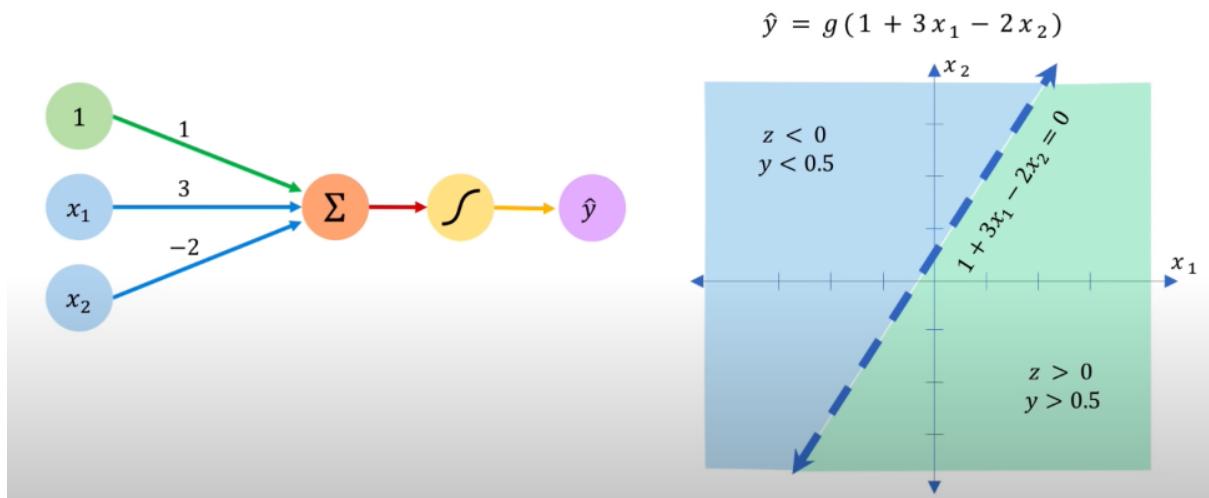


- Why activation functions?
 - Activation function introduces a complexity to your neural network

Example (separate red and green dots):



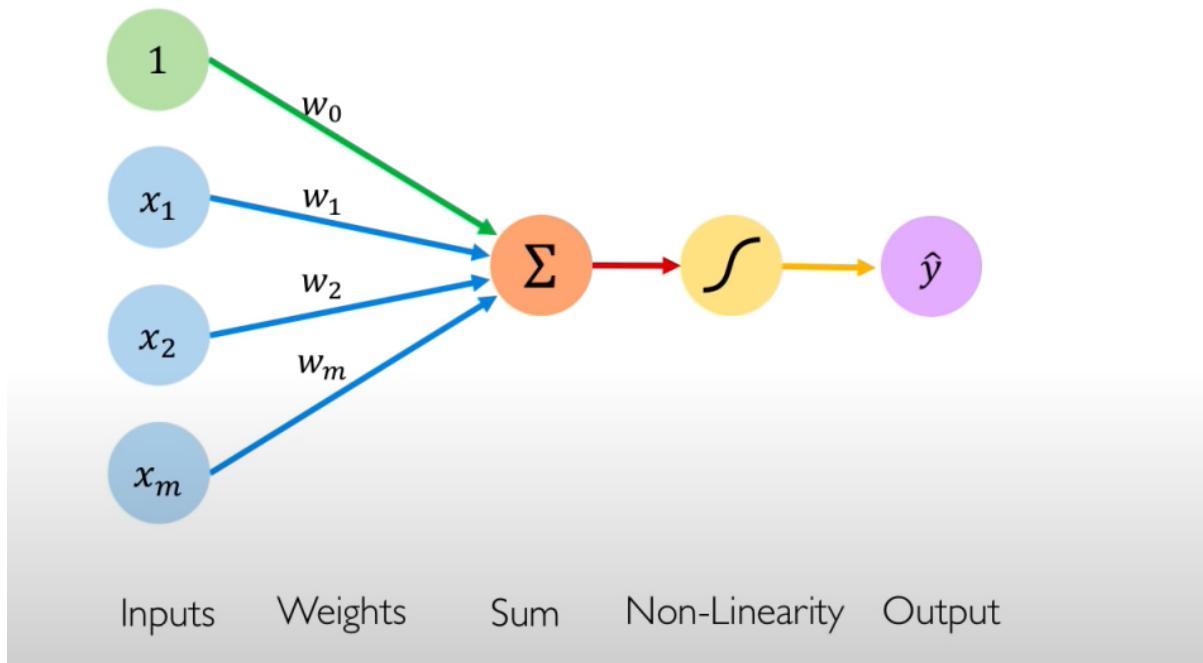
Example:



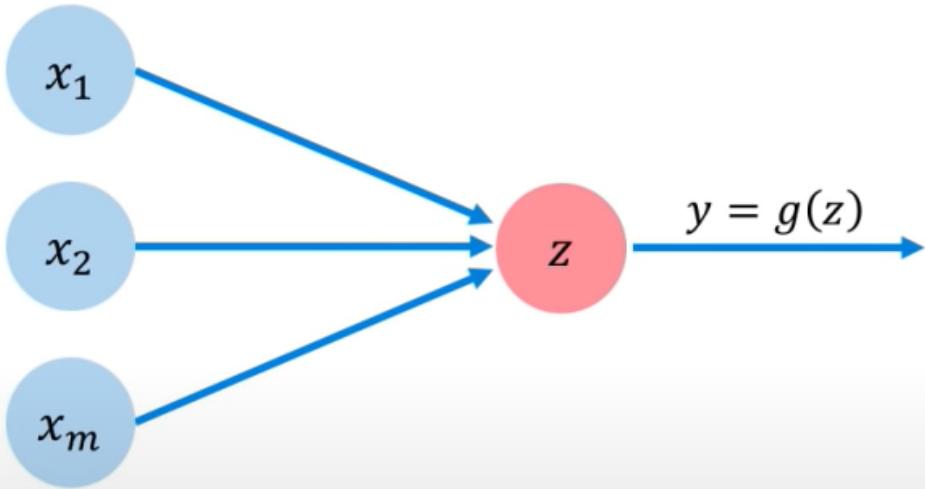
Building Neural Networks with Perceptions

The Perceptron: Simplified

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

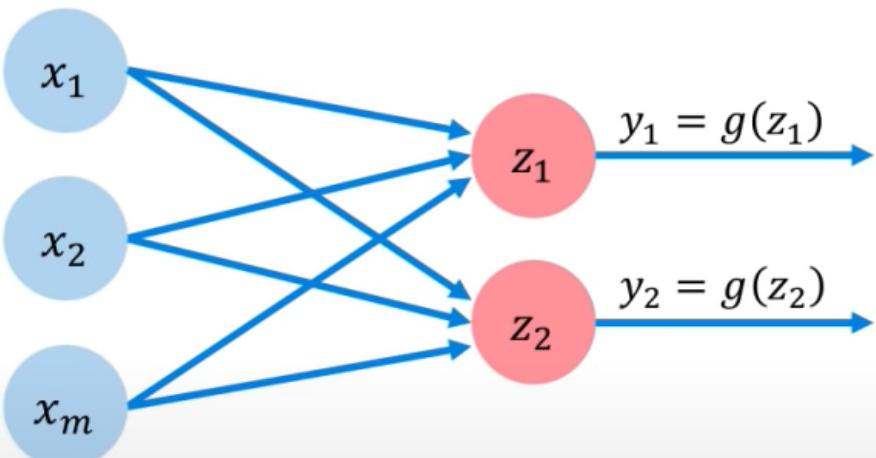


- How it works
- take inputs and apply dot product to the weight, add the bias and apply a non-linearity



$$z = w_0 + \sum_{j=1}^m x_j w_j$$

- The diagram above has no bias and no weight label

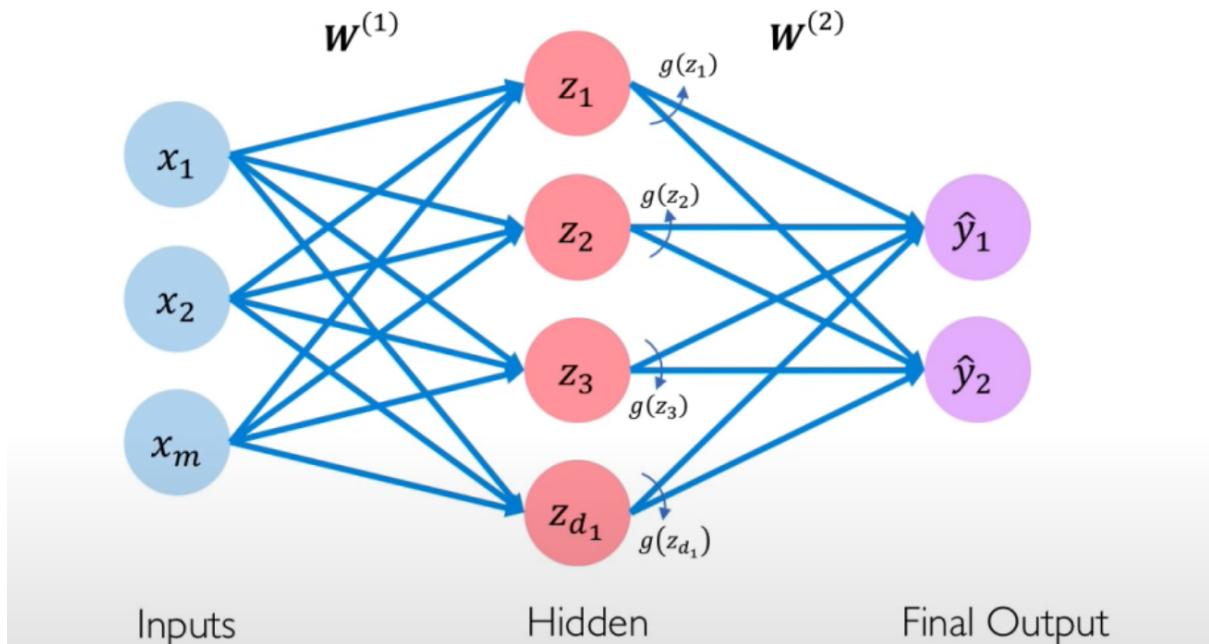


$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

- Diagram above is a Multi-Level neural network
- The only difference between this and single level is an extra perceptron in the network
 - y1 is output one
 - y2 is output two
- Because all outputs are connected to all inputs so those layers are called **dense** layers
- How to build the a Dense Layer
 - TensorFlow have implemented a dense layer for you so you don't have to do it manually

```
 import tensorflow as tf
layer = tf.keras.layers.Dense(
    units=2)
```

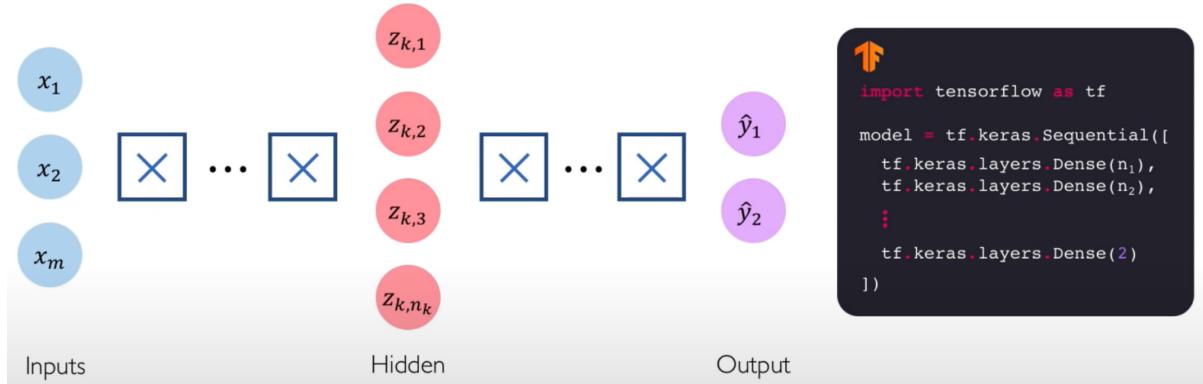
- The Diagram below is a single layer neural network



- The code below Implements a single neural network in TensorFlow

```
model = tf.keras.Sequential([
    tf.keras.layers.Dense(n),
    tf.keras.layers.Dense(2)
])
```

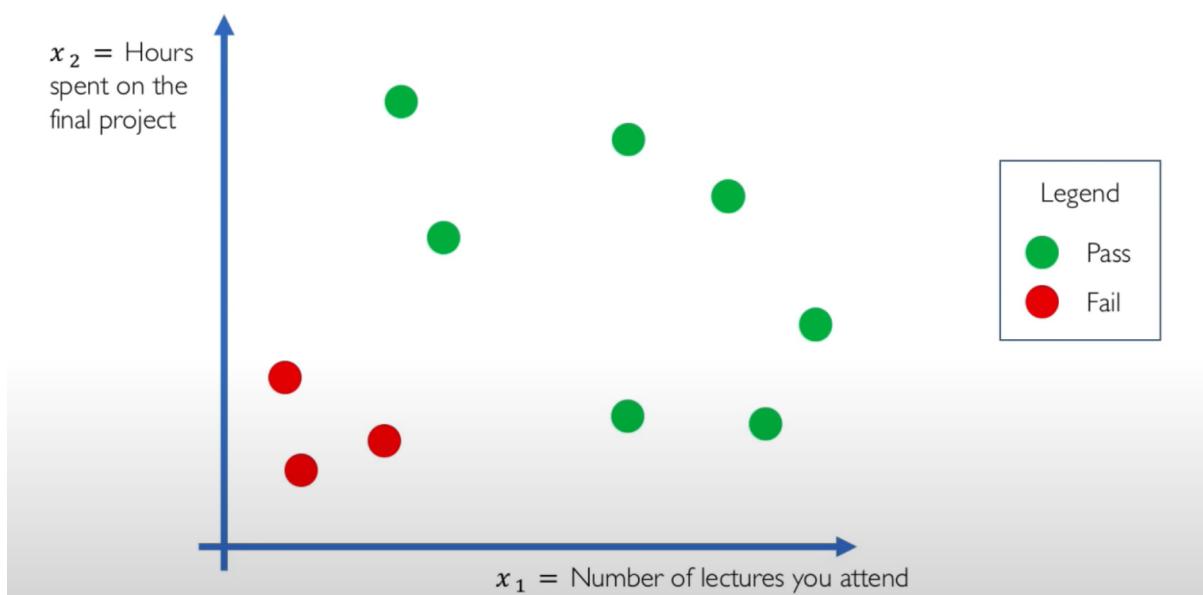
- The code below implements a deep neural network in TensorFlow



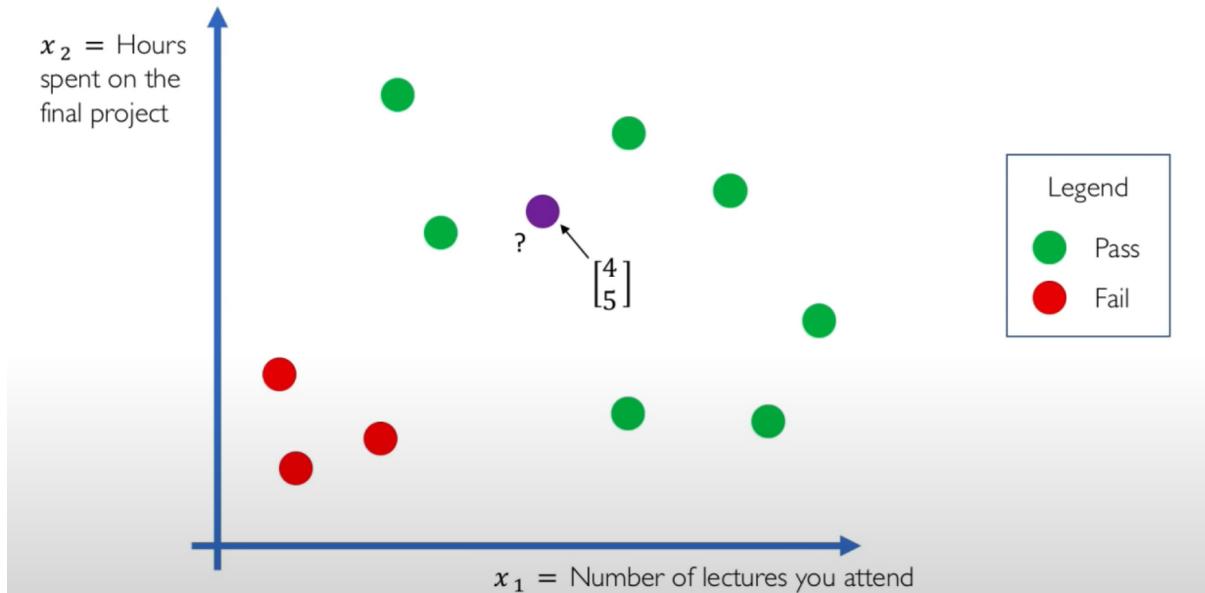
Applying Neural Network

Example Problem

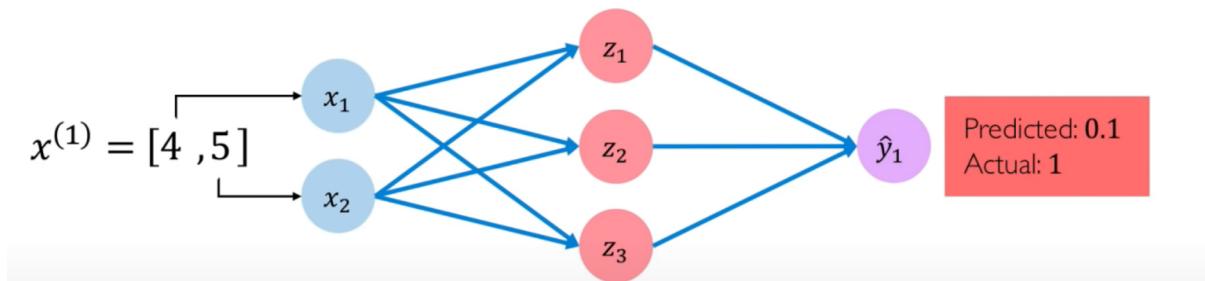
- Will I pass this class
 - x_1 = Number of lecture attened
 - x_2 = Number of hours spent on final project



- The purple dot is you, will you pass the class?



- If we just pass in your results into the network, the result you get is very wrong, you are expecting 1 but you got 0.1 which is 10 percent change for you to pass



- The reason is because the network was not trained therefore it does not know what is right and what is wrong it is just guessing the answer with a really low accuracy
- To train the network we use the loss function to determine how wrong the network is
- The loss function we will be using is softmax
- softmax compares how different the distribution are

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)}}_{\text{Actual}} \log \left(f(x^{(i)}; \mathbf{W}) \right) + (1 - \underbrace{y^{(i)}}_{\text{Actual}}) \log \left(1 - f(x^{(i)}; \mathbf{W}) \right)$$



```
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```

- Now, let's say we want to predict your grade instead of whether you are going to pass[
- The loss function we will be using is called mean square error
- this could be used for regression models that output continuous real numbers

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{(y^{(i)} - f(x^{(i)}; \mathbf{W}))^2}_{\text{Actual} \quad \text{Predicted}}$$

- Take the prediction of the network and subtract it from the actual number and square it, which would be the loss and which the model would try and change the weight to minimize the loss

Training The Neural Network

- Training the network basically tries to find the weight that achieves the lowest loss

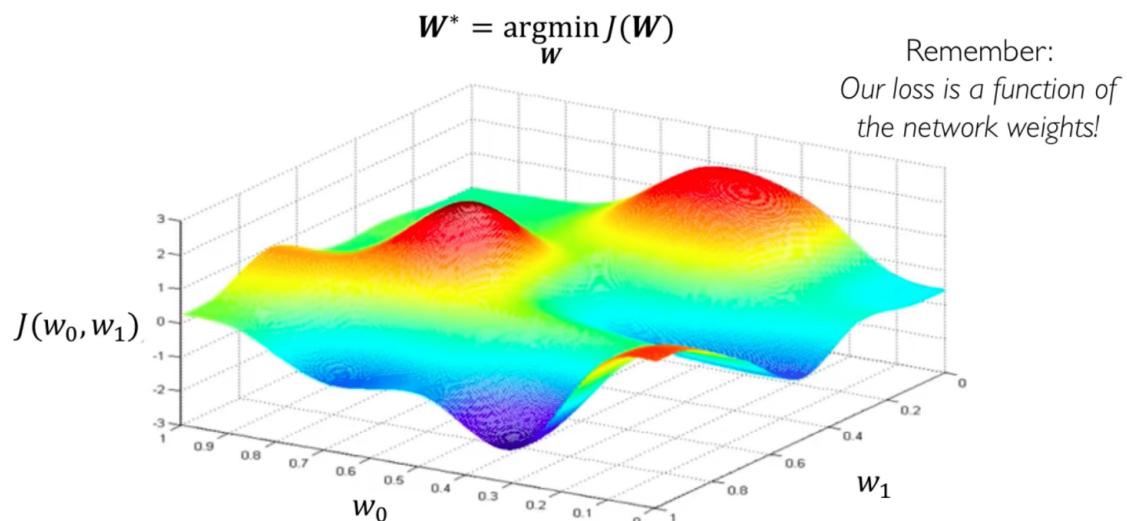
$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:
 $\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$

- This graph plots the weights and the loss



- This is algorithm gradient descend\

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

- Learning rate is how strong you want each gradient to be (ADA)
- Sudo code in actual code that is hard coded

```

import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

while True:    # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
        gradient = g.gradient(loss, weights)

    weights = weights - lr * gradient

```

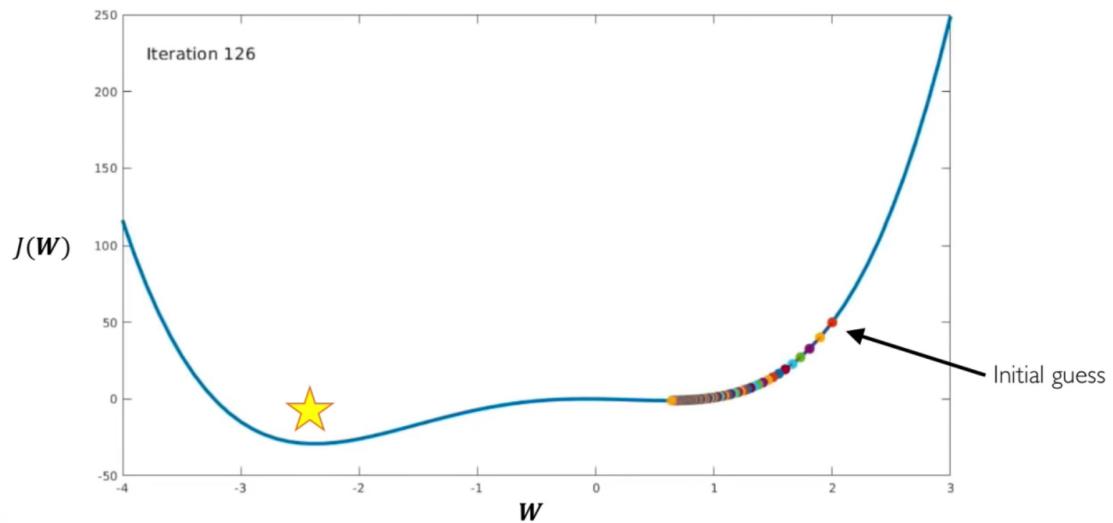
Optimization

- function for optimization

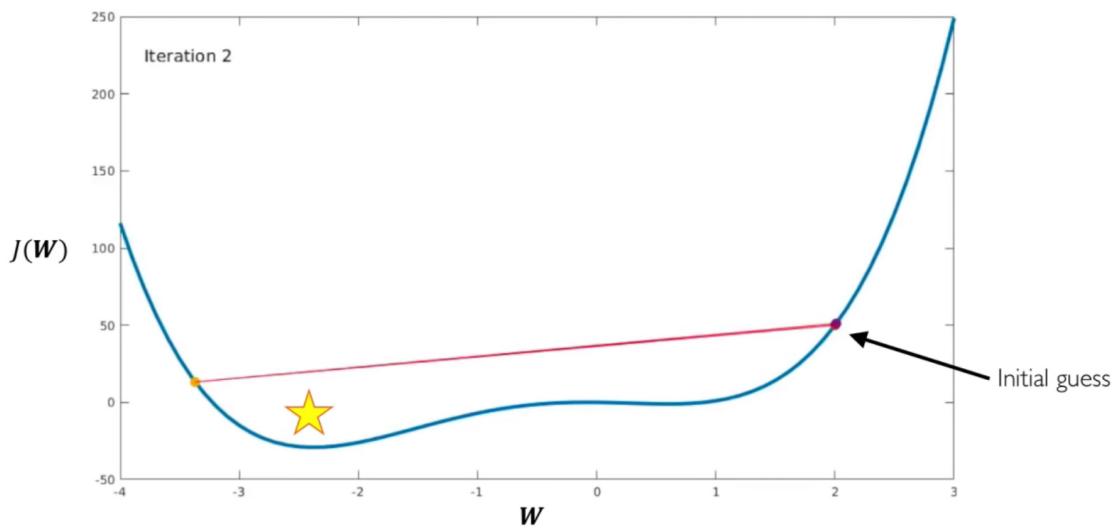
$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

- Setting learning rate is important
 - If the learning rate is too small then the model could possibly be stuck at local minima
 - If the learning rate is too high then the model could overshoot

Small learning rate converges slowly and gets stuck in false local minima



Large learning rates overshoot, become unstable and diverge



- How to set a learning rate?
 - Try different ones and see which one is "just right"
 - Adaptive learning rate that changes with different landscapes
- Adaptive learning rate
 - Learning rate is no longer fixed
 - Can be larger or smaller depending on:
 - how large the gradient is
 - how fast learning is happening

- size of particular weight
- ETC....
- The image below shows different types of learning rates

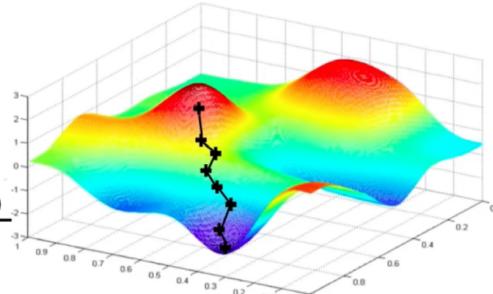
Algorithm	TF Implementation	Reference
• SGD	 <code>tf.keras.optimizers.SGD</code>	Kiefer & Wolfowitz. "Stochastic Estimation of the Maximum of a Regression Function." 1952.
• Adam	 <code>tf.keras.optimizers.Adam</code>	Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.
• Adadelta	 <code>tf.keras.optimizers.Adadelta</code>	Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.
• Adagrad	 <code>tf.keras.optimizers.Adagrad</code>	Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.
• RMSProp	 <code>tf.keras.optimizers.RMSProp</code>	

Mini Batches

- It is basically a approximation of the gradient

Algorithm

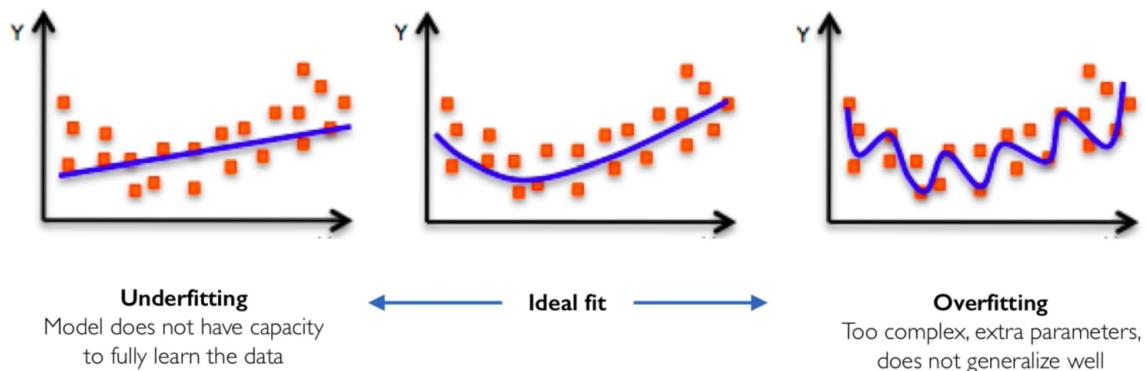
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



- Mini batches makes estimating to gradient much quicker since it takes less to compute and have multiple points
- Could increase learning rate
- Faster training because we could split it up to different GPU or computers

Overfitting

- We would need to test the model with unseen test data so the model don't overfit



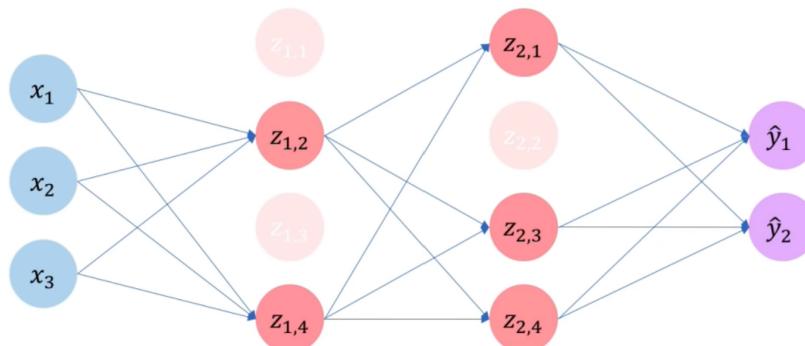
Regularization

- Technique to constrain our optimization problem to discourage complex models
- Improve accuracy for unseen data

Dropouts

- Drops out a percentage of data on during each training
- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

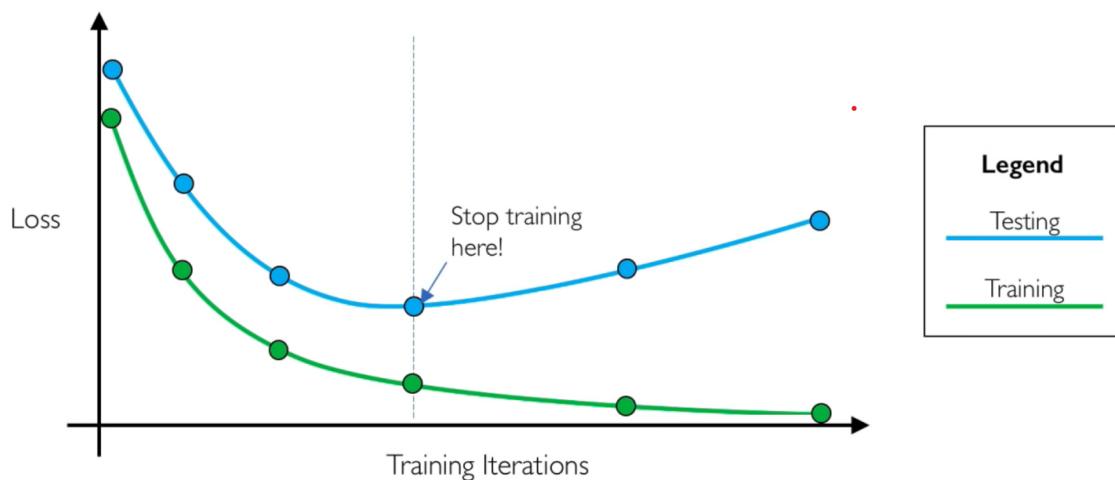
`tf.keras.layers.Dropout(p=0.5)`



- Forced to generalize better and build a better model and prevent overfitting

Early stopping

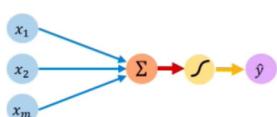
- Stop the training before it has the chance to overfit



Summary:

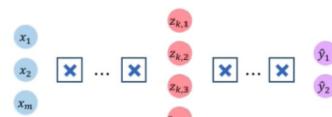
The Perceptron

- Structural building blocks
- Nonlinear activation functions



Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

- Adaptive learning
- Batching
- Regularization

