

# Types of Regularization Techniques

There are two main types of regularization techniques: L1(Lasso) and L2( Ridge) regularization

## 1) Lasso Regularization (L1 Regularization)

In L1 you add information to model equation to be the absolute sum of theta vector ( $\theta$ ) multiply by the regularization parameter ( $\lambda$ ) which could be any large number over size of data (m), where (n) is the number of features.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(h_\theta(x^{(i)}), y^{(i)}) + \frac{\lambda}{m} \sum_{j=1}^n |\theta_j|$$

## 2) Ridge Regularization (L2 Regularization)

In L2, you add the information to model equation to be the sum of vector ( $\theta$ ) squared multiplied by the regularization parameter ( $\lambda$ ) which can be any big number over size of data (m), which (n) is a number of features.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(h_\theta(x^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

## Advanced Learning Algorithms - Course 2

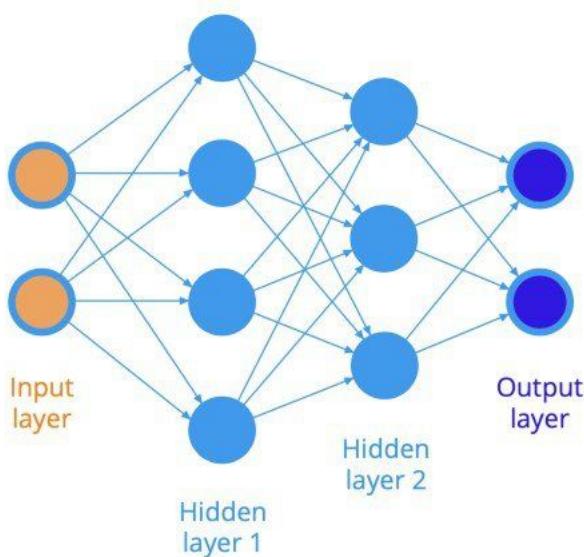
### NEURAL NETWORKS

- Learning Algorithms that mimic the Human Brain.
- Speech Recognition, Computer Vision, NLP, Climate Change, Medical Image Detection, Product Recommendation.
- Just like human neuron, neural networks use many neurons to receive input and gives an output to next layer of network.

- As amount of data increased the performance of neural networks also increased.
- Layer is a grouping of neurons, which take input of similar features and output a few numbers together.
- Layer can have single or multiple neurons.
- If the layer is first, it is called Input Layer
- If the layer is last it is called Output Layer
- If the layer in middle, it is called Hidden Layer
- Each neurons have access to all other neurons in next layer.

## Activations (a)

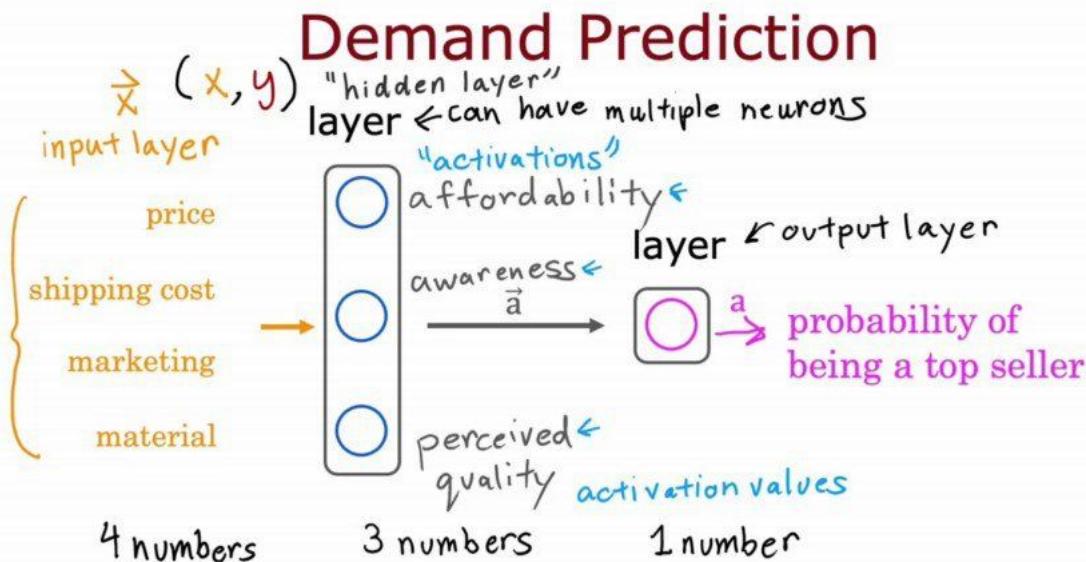
- Refers to the degree of high output value of neuron, given to the next neurons.
- Neural Networks learn its own features. No need of manual Feature Engineering.



## Demand Prediction of a Shirt?

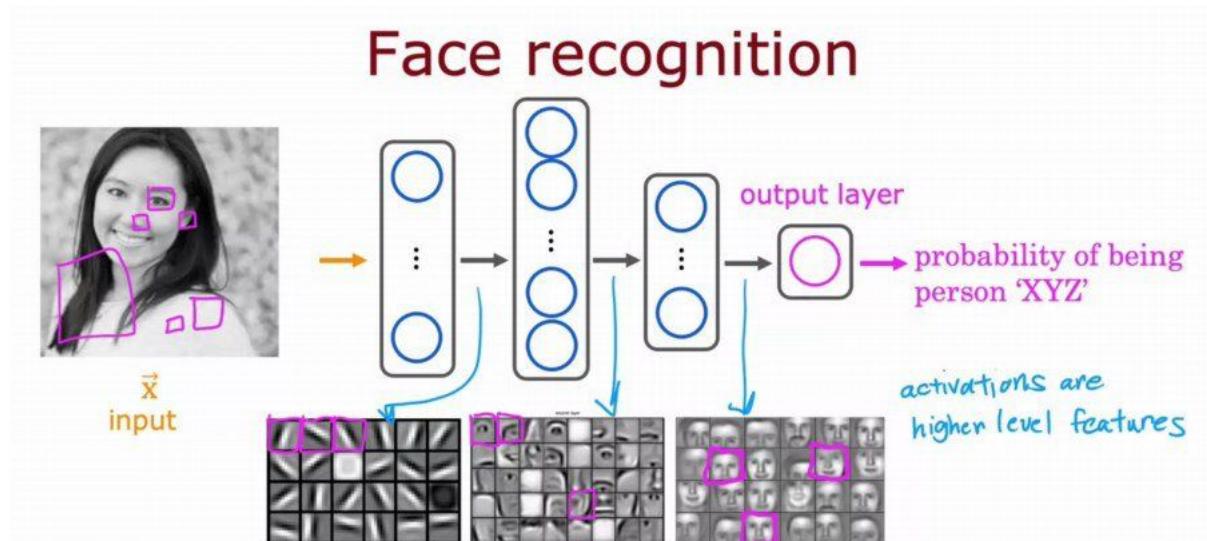
- affordability, awareness, perceived quality are activations of each neurons.
- It is basically a Logistic Regression, trying to learn much by its own.
- activation,  $a = g(z) = 1 / (1 + e^{-z})$
- Number of layers and neurons are called Neural Network Architecture.

- Multi-layer perception is also known as MLP. It is fully connected dense layers.



## Image Recognition

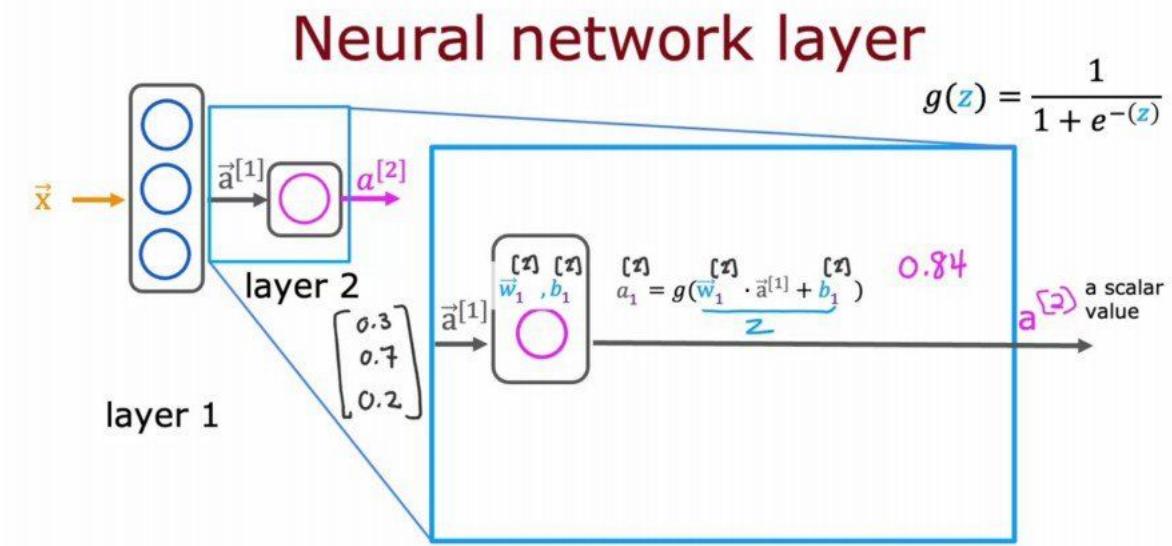
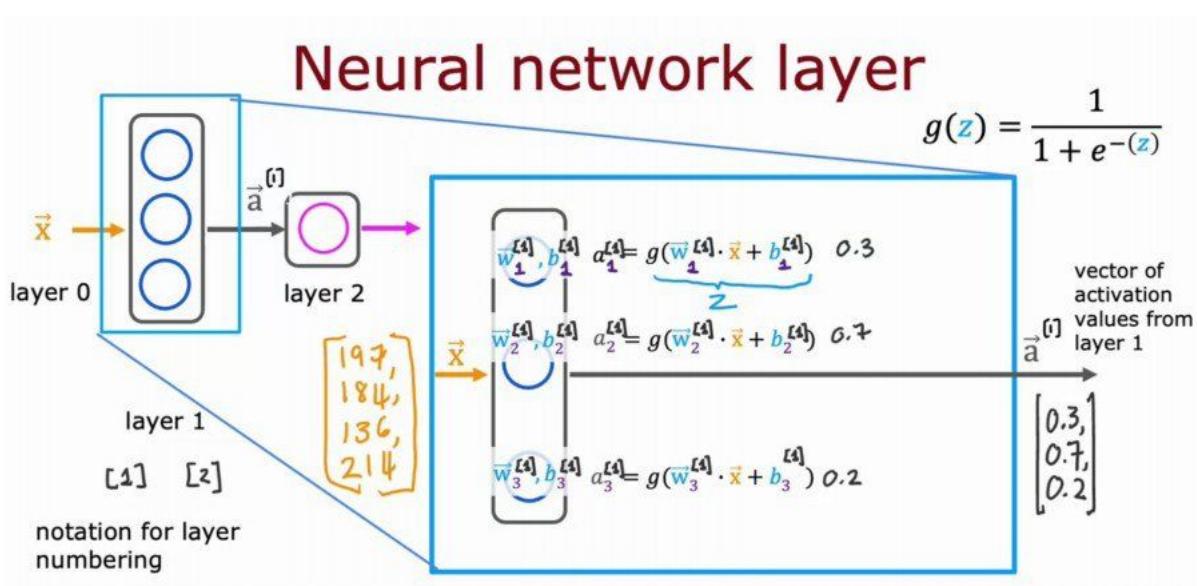
- For a  $1000 * 1000$  pixel image we have a vector of 1 million elements.
- We build a neural network by providing that as an input.
- In each hidden layer they are improving the recognition all by itself.



## Neural Network Layer

- Each neuron is implementing the activation of Logistic Regression.

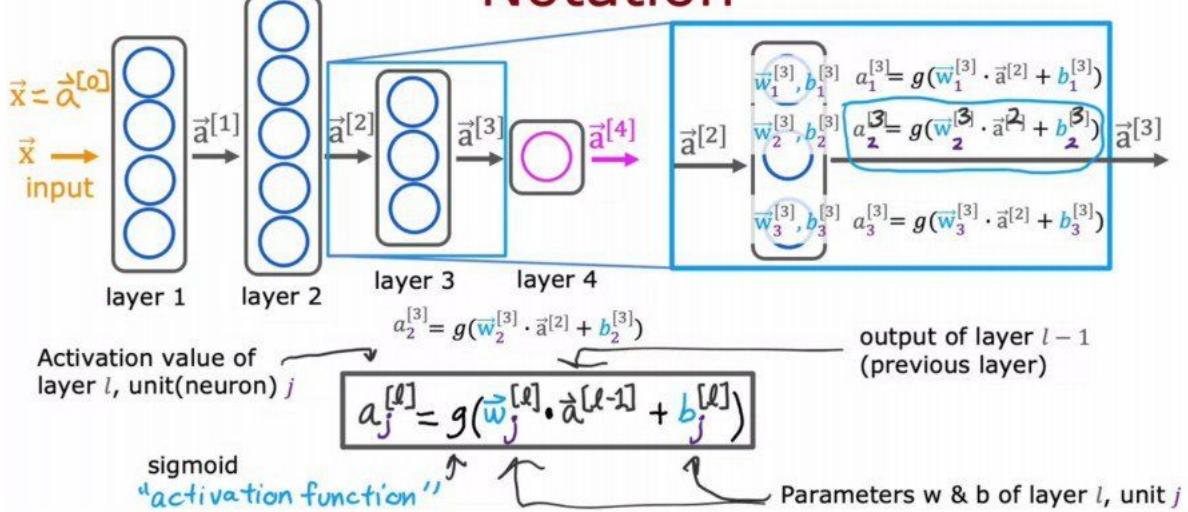
- Super script square bracket 1 means quantity that is associated with Neural Network 1 and so on.
- Input of layer 2 is the output of layer 1. It will be a set of vectors.
- Every layer input a vector of numbers and output a vector of numbers



## Complex Neural Network

- $a_j[l]$  where  $j$  is the neuron/unit and  $l$  is the layer
- $g(w \cdot a + b)$  is the activation function (sigmoid function etc can be used here)
- input layer is  $a[0]$

## Notation

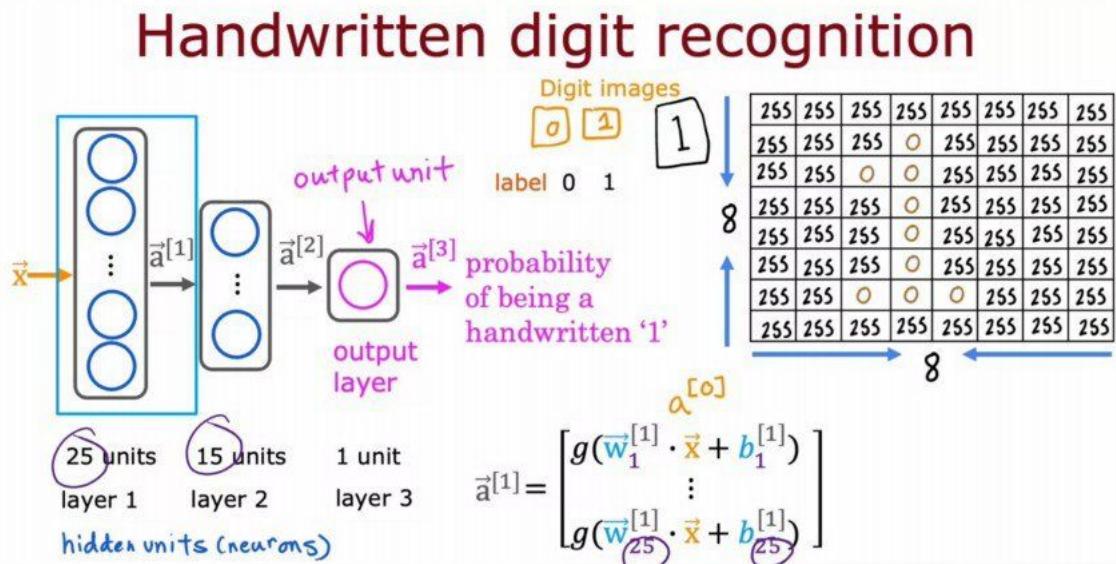


## Inference/Make Prediction (Handwritten Digit Recognition)

### Forward Propagation

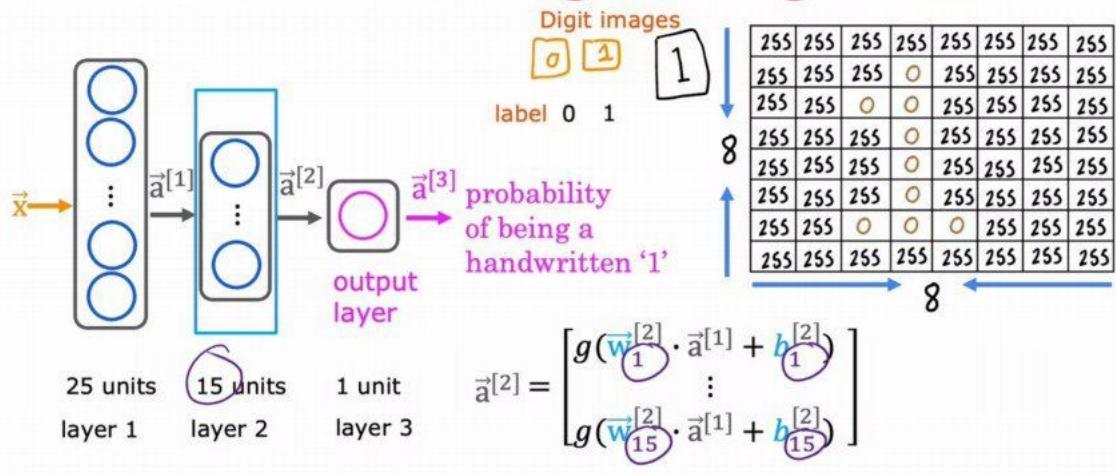
- If activation computation goes from left to right it is called, Forward Propagation.

Calculation of first layer vector

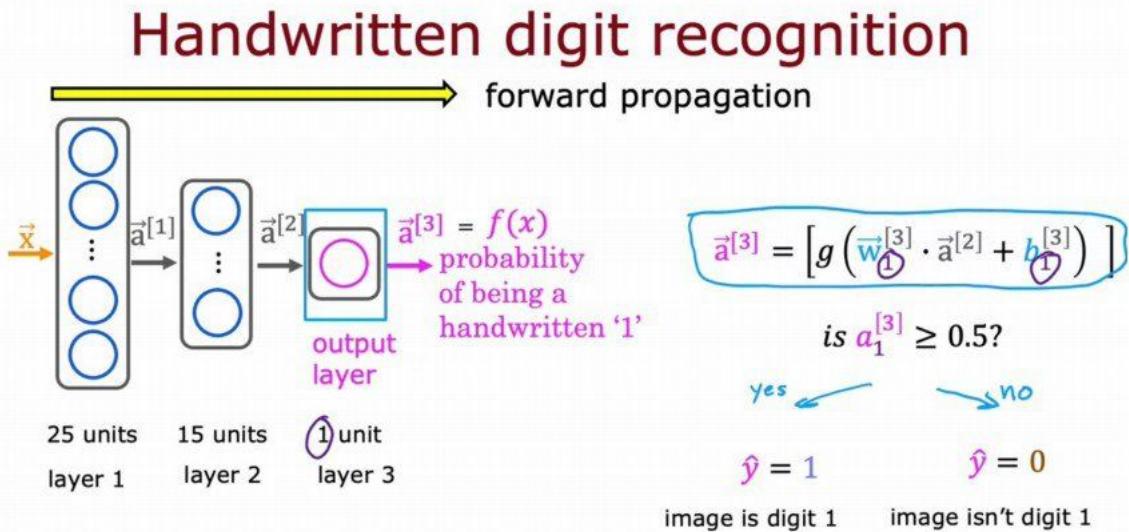


Calculation of second layer vector

## Handwritten digit recognition



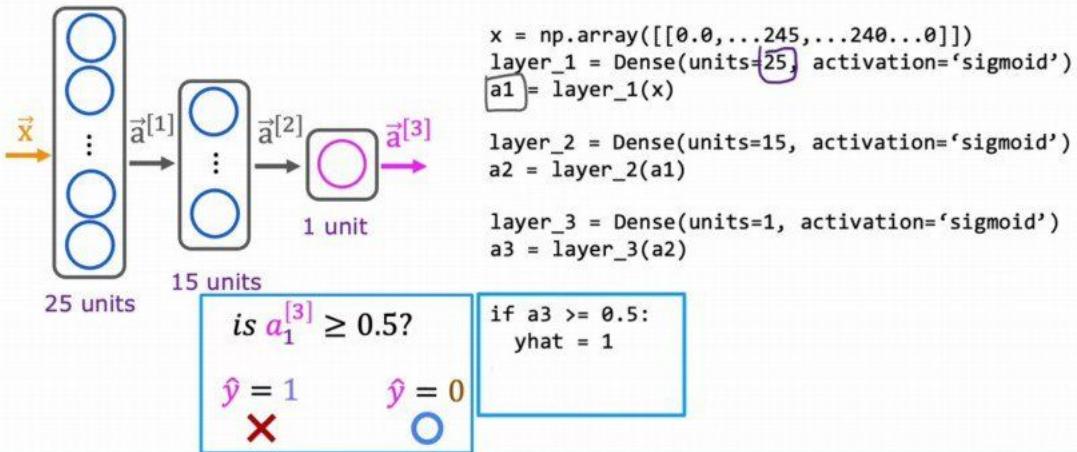
Calculation of last layer vector



## Tensorflow Implementation

- building a layer of NN in tensorflow below for handwritten digit 0/1 classification

# Model for digit classification



## Data in Tensorflow

- Inconsistency between data in Numpy and Tensorflow
- Numpy we used 1D vector, `np.array([ 200, 18 ])`
- Tensorflow uses matrices, `np.array([ [ 200, 18 ] ])`
- Representing in Matrix instead of 1D array make tensorflow run faster.
- Tensor is something like matrix

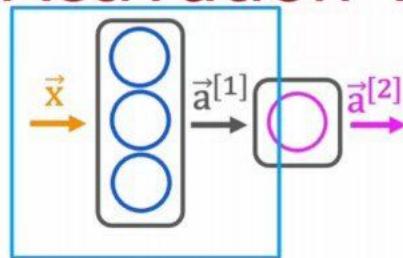
`x = np.array([[200, 17]])` → [200 17]       $1 \times 2$

`x = np.array([[200], [17]])`      → [200]  
    ↓  
    [17]       $2 \times 1$

`x = np.array([200,17])`

## Activation of Vector

# Activation vector



```
x = np.array([[200.0, 17.0]])
layer_1 = Dense(units=3, activation='sigmoid')
a1 = layer_1(x)
→      → [0.2, 0.7, 0.3]    1 x 3 matrix
→      → tf.Tensor([[0.2 0.7 0.3]], shape=(1, 3), dtype=float32)
→      → a1.numpy()
→      → array([[0.2, 0.7, 0.3]], dtype=float32)
```

## Building a Neural Network in Tensorflow

### Digit Classification Model

```
import tensorflow as tf
layer_1 = Dense( units=25, activation="sigmoid" )
layer_2 = Dense( units=15, activation="sigmoid" )
layer_3 = Dense( units=1, activation="sigmoid" )

model = Sequential ( [ layer_1, layer_2, layer_3 ] )

x = np.array( [ [ 0..., 245, ..., 17 ],
                [ 0..., 200, ..., 184 ] ] )
y = np.array( [ 1, 0 ] )

model.compile(.....)
model.fit( x, y )
model.predict( new_x )
```

## Forward Prop in Single Layer (Major Parts)

```
import tensorflow as tf
from tensorflow.keras.layers import Dense
from tensorflow.keras.models import Sequential

model = Sequential(
```

```

        [
            tf.keras.Input(shape=(2,)),
            Dense(3, activation='sigmoid', name = 'layer1'),
            Dense(1, activation='sigmoid', name = 'layer2')
        ]
    )

model.compile(
    loss = tf.keras.losses.BinaryCrossentropy(),
    optimizer = tf.keras.optimizers.Adam(learning_rate=0.01),
)

model.fit(
    Xt,Yt,
    epochs=10,
)

```

## General Implementation of Forward Prop in Single Layer

- Uppercase for Matrix
- Lowercase for Vectors and Scalars

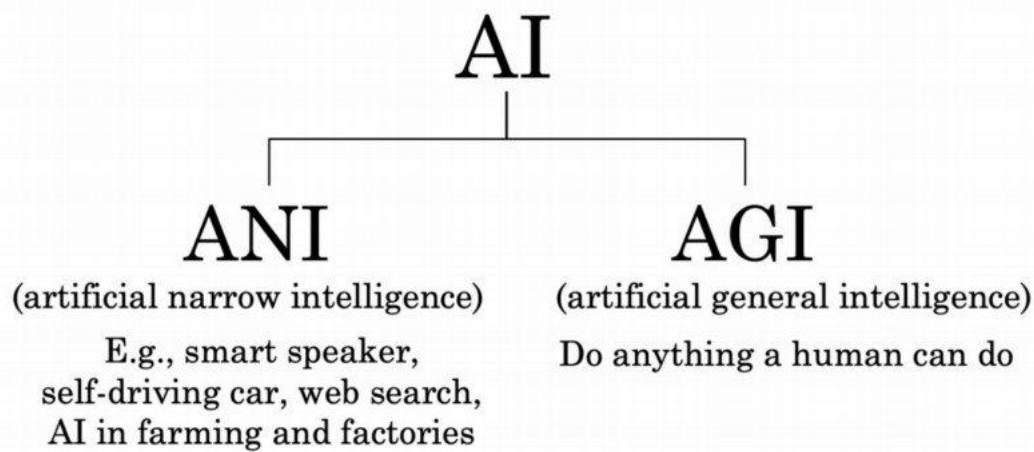
```

def my_dense(a_in, W, b, g):
    """
    Computes dense layer
    Args:
        a_in (ndarray (n, )) : Data, 1 example
        W    (ndarray (n,j)) : Weight matrix, n features per unit, j units
        b    (ndarray (j, )) : bias vector, j units
        g    activation function (e.g. sigmoid, relu..)
    Returns
        a_out (ndarray (j,)) : j units|
    """
    units = W.shape[1]
    a_out = np.zeros(units)
    for j in range(units):
        w = W[:,j]
        z = np.dot(w, a_in) + b[j]
        a_out[j] = g(z)
    return(a_out)

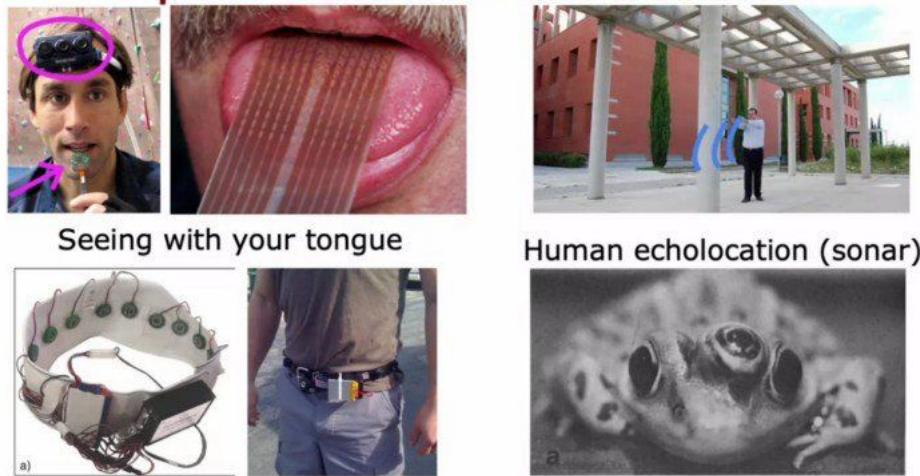
def my_sequential(x, W1, b1, W2, b2):
    a1 = my_dense(x, W1, b1, sigmoid)
    a2 = my_dense(a1, W2, b2, sigmoid)
    return(a2)

```

# Artificial General Intelligence - AGI



## Sensor representations in the brain



## Vectorization

- Matrix Multiplication in Neural Networks using Parallel Computer Hardware
- Matrix Multiplication Replaces For Loops in speed comparison

## For loops vs. vectorization

```

x = np.array([200, 17])
W = np.array([[1, -3, 5],
              [-2, 4, -6]])
b = np.array([-1, 1, 2])

def dense(a_in,W,b):
    units = W.shape[1]
    a_out = np.zeros(units)
    for j in range(units):
        w = W[:,j]
        z = np.dot(w, a_in) + b[j]
        a_out[j] = g(z)
    return a_out

```

vectorized →

```

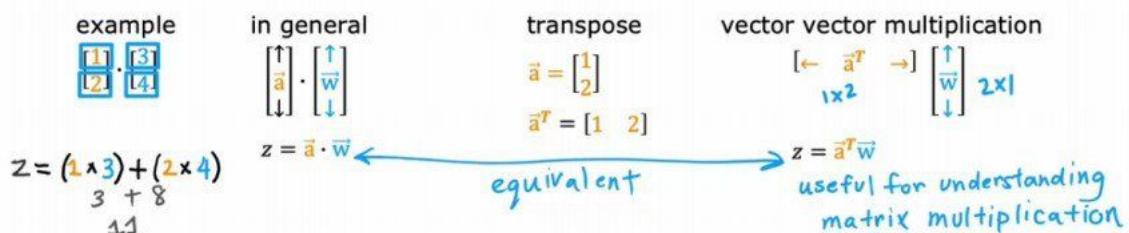
X = np.array([[200, 17]]) 2Darray
W = np.array([[1, -3, 5], same
              [-2, 4, -6]])
B = np.array([-1, 1, 2]) 1x3 2Darray
all 2Darrays
def dense(A_in,W,B):
    Z = np.matmul(A_in,W) + B
    A_out = g(Z) matrix multiplication
    return A_out
[[1,0,1]]

```

[1,0,1]

## Dot Product to Matrix Multiplication using Transpose

- $a \cdot w = a_1 * w_1 + a_2 * w_2 + \dots + a_n * w_n$
- $a \cdot w = a^T * w$



## Matrix Multiplication in Neural Networks

- Dot Product of vectors that have same length
- Matrix Multiplication is valid if col of matrix 1 = row of matrix 2
- Output will have row of matrix 1 and col of matrix 2

In code for numpy array and vectors

# Matrix multiplication in NumPy

$$A = \begin{bmatrix} 1 & -1 & 0.1 \\ 2 & -2 & 0.2 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \quad W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix} \quad Z = A^T W = \begin{bmatrix} 11 & 17 & 23 & 9 \\ -11 & -17 & -23 & -9 \\ 1.1 & 1.7 & 2.3 & 0.9 \end{bmatrix}$$

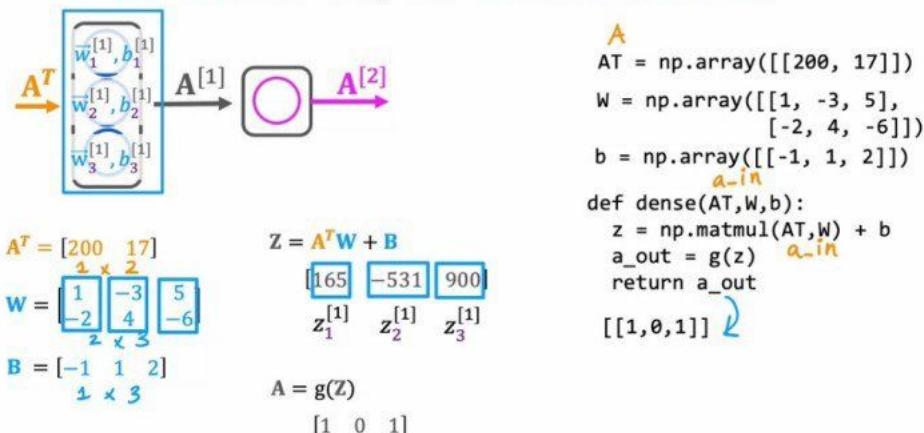
`A=np.array([[1,-1,0.1], [2,-2,0.2]])`      `W=np.array([[3,5,7,9], [4,6,8,0]])`      `Z = np.matmul(AT,W)`      *or*      `Z = AT @ W`

`AT=np.array([[1,2], [-1,-2], [0.1,0.2]])`      *result*      `[[11,17,23,9], [-11,-17,-23,-9], [1.1,1.7,2.3,0.9]]`

`AT=A.T`      *transpose*

## Dense Layer Vectorized

### Dense layer vectorized



## Tensorflow Training

- epoch is how many steps, the learning algorithm like gradient descent should run.
- BinaryCrossentropy is log loss function for logistic regression in classification problem
- if we have regression problem, use MeanSquaredError()
- loss = tf.keras.losses.MeanSquaredError()
- Derivatives of gradient descent has been calculated using backpropagation handled by model.fit

```
import tensorflow as tf
from tensorflow.keras.layers import Dense
from tensorflow.keras.models import Sequential
```

```

model = Sequential(
    [
        tf.keras.Input(shape=(2,)),
        Dense(3, activation='sigmoid', name = 'layer1'),
        Dense(1, activation='sigmoid', name = 'layer2')
    ]
)

model.compile(
    loss = tf.keras.losses.BinaryCrossentropy(),
    optimizer = tf.keras.optimizers.Adam(learning_rate=0.01),
)

model.fit(
    Xt,Yt,
    epochs=10,
)

```

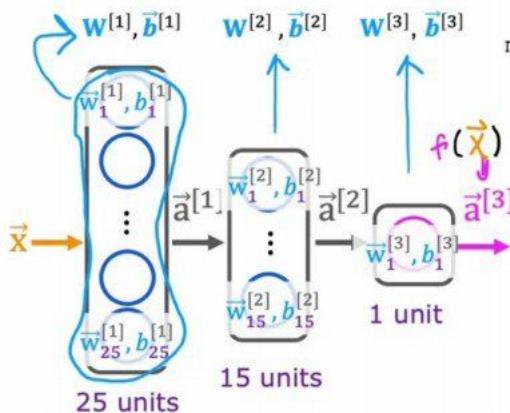
## Model Training Steps TensorFlow

<small>①</small> specify how to compute output given input $x$ and parameters $w, b$ (define model)	<small>logistic regression</small> $z = np.dot(w, x) + b$ $f_x = 1 / (1 + np.exp(-z))$	<small>neural network</small> $model = Sequential([$ $\quad Dense(...)$ $\quad Dense(...)$ $\quad Dense(...) \quad ])$
<small>②</small> specify loss and cost $L(f_{\bar{w},b}(\vec{x}), y)$ 1 example $J(\bar{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\bar{w},b}(\vec{x}^{(i)}), y^{(i)})$	<small>logistic loss</small> $loss = -y * np.log(f_x)$ $- (1-y) * np.log(1-f_x)$ $w = w - alpha * dj_dw$ $b = b - alpha * dj_db$	<small>binary cross entropy</small> $model.compile($ $\quad loss=BinaryCrossentropy())$ $model.fit(X, y, epochs=100)$
<small>③</small> Train on data to minimize $J(\bar{w}, b)$		

# 1. Create the model

define the model

$$f(\vec{x}) = ?$$



```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense

model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid'),
])
```

# 2. Loss and cost functions

handwritten digit classification problem

binary classification

$L(f(\vec{x}), y) = -y \log(f(\vec{x})) - (1-y) \log(1-f(\vec{x}))$

compare prediction vs. target

logistic loss  
also known as binary cross entropy

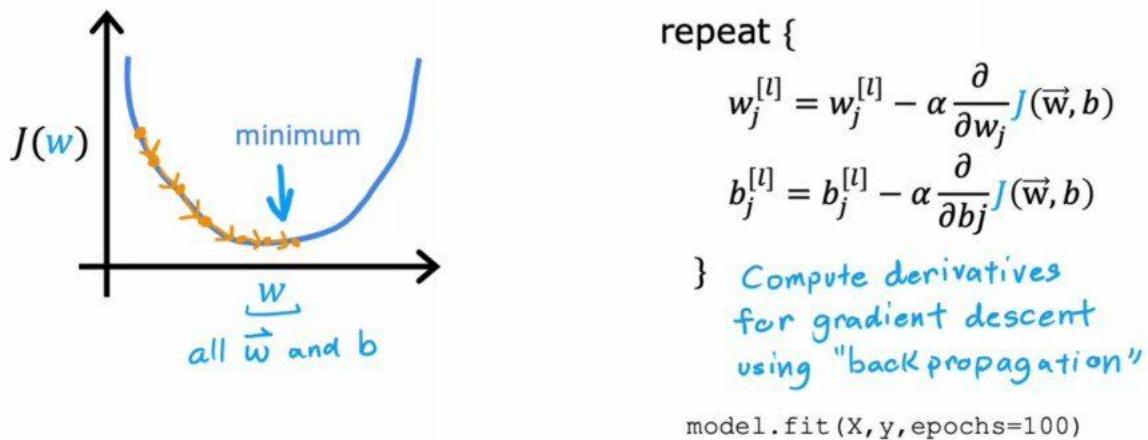
$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)})$$

$\mathbf{w}^{[1]}, \mathbf{w}^{[2]}, \mathbf{w}^{[3]}$      $\mathbf{b}^{[1]}, \mathbf{b}^{[2]}, \mathbf{b}^{[3]}$      $f_{\mathbf{W}, \mathbf{B}}(\vec{x})$

model.compile(loss= BinaryCrossentropy())  
regression  
(predicting numbers and not categories)  
model.compile(loss= MeanSquaredError())

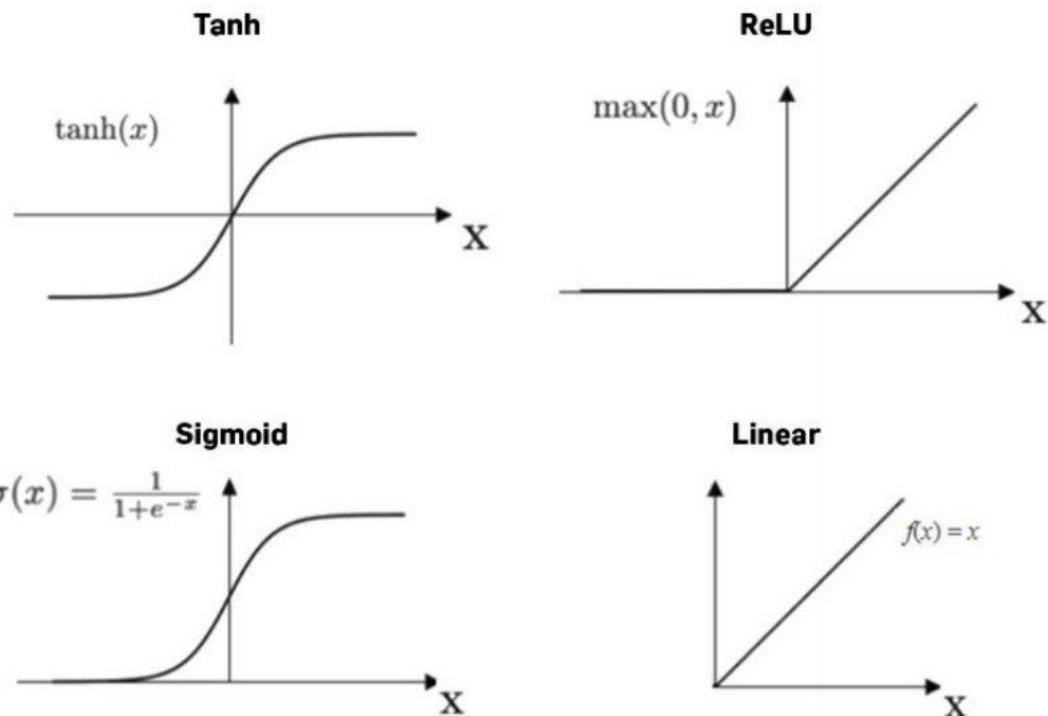
from tensorflow.keras.losses import  
BinaryCrossentropy    Keras  
from tensorflow.keras.losses import  
MeanSquaredError

### 3. Gradient descent



### Activation Function Alternatives

- awareness can not be binary, It better to use like like non negative number from zero to large number
- So best activation here is ReLU, Rectifier Linear Unit
- $g(z) = \max(0, z)$
- ReLU is common choice of training neural networks and now used more than Sigmoid
- ReLU is faster
- ReLU is flat in one side, and Sigmoid flat on both sides. So gradient descent perform faster in ReLU



## Choose Activation Function

### For Output Layer

- Binary Classification - Sigmoid
- Regression - Linear
- Regression Positive Only - ReLU

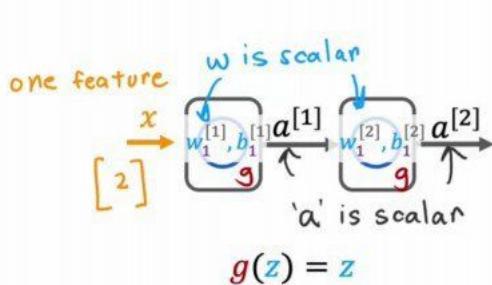
### For Hidden Layer

- ReLU as standard activation

## Why do we need activation functions?

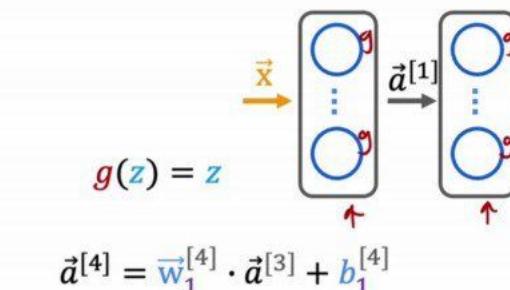
- If we use linear activation function across all the neurons the neural network is no different from Linear Regression.
- If we use linear in hidden and sigmoid in output it is similar to logistic regression.
- **Don't use linear in hidden layers. Use ReLU**
- We won't be able to fit anything complex

## Linear Example



$$\begin{aligned}
 a^{[1]} &= \underbrace{w_1^{[1]} x}_{\text{w is scalar}} + b_1^{[1]} \\
 a^{[2]} &= \underbrace{w_1^{[2]} a^{[1]}}_{\text{'a' is scalar}} + b_1^{[2]} \\
 &= w_1^{[2]} (w_1^{[1]} x + b_1^{[1]}) + b_1^{[2]} \\
 \vec{a}^{[2]} &= (\underbrace{\vec{w}_1^{[2]} \vec{w}_1^{[1]}}_{\omega}) x + \underbrace{\vec{w}_1^{[2]} b_1^{[1]} + b_1^{[2]}}_{b} \\
 \vec{a}^{[2]} &= w x + b \\
 f(x) &= w x + b \quad \text{linear regression}
 \end{aligned}$$

## Example



all linear (including output)  
 ↳ equivalent to linear regression

$$\vec{a}^{[4]} = \frac{1}{1+e^{-(\vec{w}_1^{[4]}\cdot\vec{a}^{[3]}+b_1^{[4]})}}$$

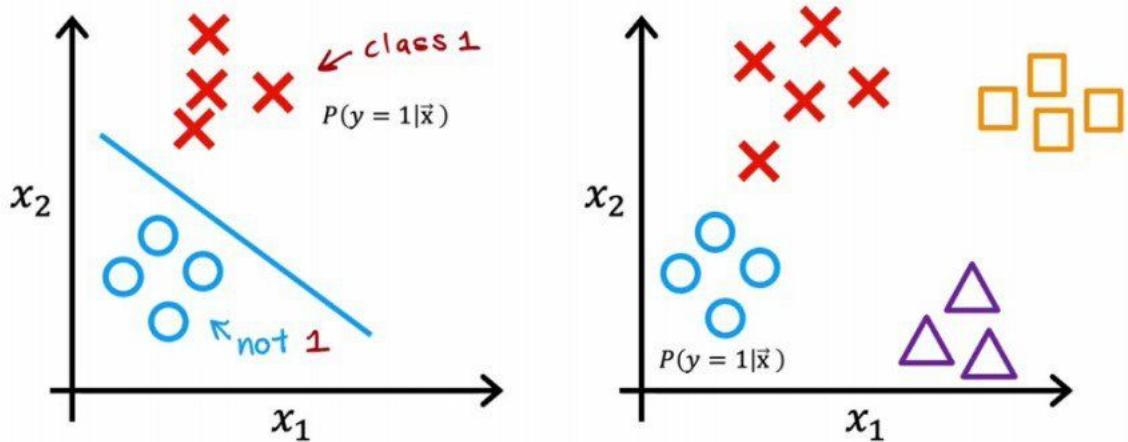
output activation is sigmoid  
 (hidden layers still linear)  
 ↳ equivalent to logistic regression

Don't use linear activations in hidden layers (use ReLU)

## Multiclass Classification

- classification with more than 2 class is called multiclass classification
- Identifying a number is multiclass since we want to classify 10 number classes (MNIST)
- for n = 2 softmax is basically logistic regression

# Multiclass classification example



Logistic regression  
(2 possible output values)

$$z = \vec{w} \cdot \vec{x} + b$$

$$\text{X } a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x}) \quad 0.11$$

$$\text{O } a_2 = 1 - a_1 = P(y=0|\vec{x}) \quad 0.29$$

Softmax regression  
(N possible outputs)  $y=1, 2, 3, \dots, N$

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j = 1, \dots, N$$

parameters  $w_1, w_2, \dots, w_N$   
 $b_1, b_2, \dots, b_N$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j|\vec{x})$$

$$\text{note: } a_1 + a_2 + \dots + a_N = 1$$

Softmax regression (4 possible outputs)  $y=1, 2, 3, 4$

$$\text{X } z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1|\vec{x}) \quad 0.30$$

$$\text{O } z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=2|\vec{x}) \quad 0.20$$

$$\text{O } z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=3|\vec{x}) \quad 0.15$$

$$\Delta z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=4|\vec{x}) \quad 0.35$$

## Cost Function of Softmax



## Softmax Function

$$F(X_i) = \frac{\text{Exp}(X_i) \quad i = 0, 1, 2, \dots k}{\sum_{j=0}^k \text{Exp}(X_j)}$$



## Sigmoid Function

$$F(X_i) = \frac{1}{1 + \text{Exp}(-X_i)}$$

- Cross entropy loss - if loss value is 1 smaller the loss.
- Cross entropy loss - if loss value is 0 larger the loss.
- **Sigmoid - BinaryCrossEntropy**
- **Softmax - SparseCategoricalCrossEntropy**

## Cost

### Logistic regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1 + e^{-z}} = P(y = 1 | \vec{x})$$

$$a_2 = 1 - a_1 = P(y = 0 | \vec{x})$$

$$\text{loss} = -y \underbrace{\log a_1}_{\text{if } y=1} - (1-y) \underbrace{\log(1-a_1)}_{\text{if } y=0}$$

$$J(\vec{w}, b) = \text{average loss}$$

### Softmax regression

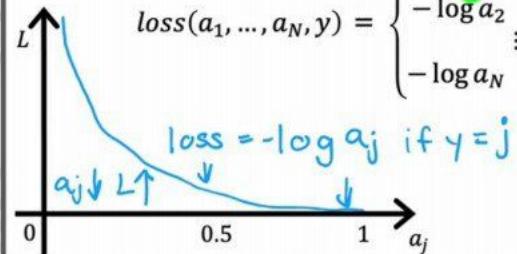
$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y = 1 | \vec{x})$$

$$\vdots$$

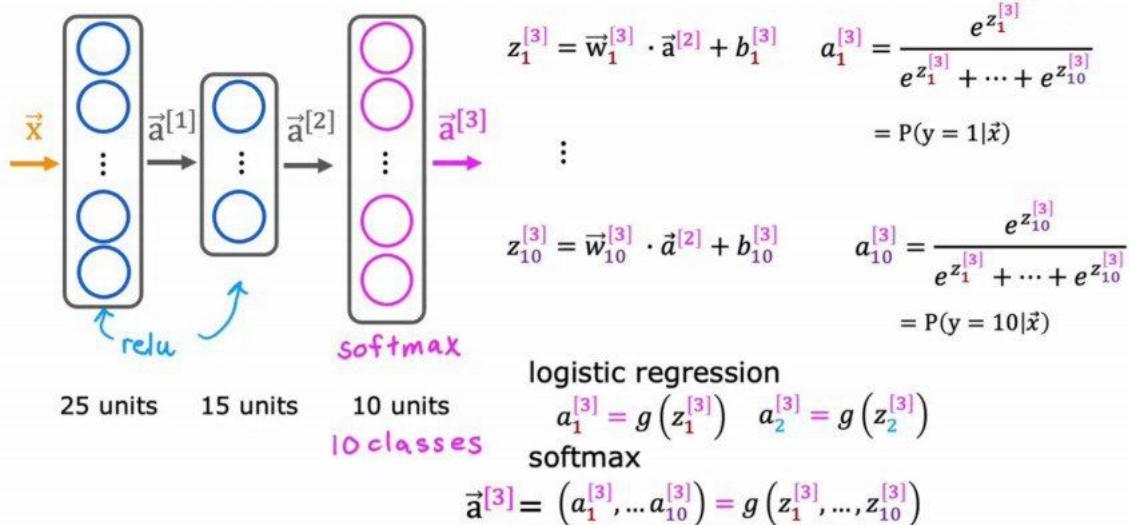
$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y = N | \vec{x})$$

#### Crossentropy loss

$$\text{loss}(a_1, \dots, a_N, y) = \begin{cases} -\log a_1 & \text{if } y = 1 \\ -\log a_2 & \text{if } y = 2 \\ \vdots \\ -\log a_N & \text{if } y = N \end{cases}$$



## Neural Network with Softmax output



## MNIST with softmax

① specify the model

$$f_{\vec{w}, b}(\vec{x}) = ?$$

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
from tensorflow.keras.losses import
    SparseCategoricalCrossentropy,
model.compile(loss= SparseCategoricalCrossentropy())
model.fit(X, Y, epochs=100)
Note: better (recommended) version later.
Don't use the version shown here!
```

② specify loss and cost

$$L(f_{\vec{w}, b}(\vec{x}), y)$$

③ Train on data to minimize  $J(\vec{w}, b)$

Don't use the above way of code for implementation. We can use alternative efficient method.

## Improved Implementation of Softmax

```
x1 = 2.0 / 10000
0.000200000000000

x2 = (1 + 1/10000) - (1 - 1/10000)
0.00019999999978
```

```
# due to memory constraints in computer rounding happens.
# inorder to avoid/reduce rounding up in softmax we can use other implementation
```

## Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

Logistic regression:

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

Original loss

$$\text{loss} = -y \log(a) - (1-y) \log(1-a)$$

model = Sequential([  
 Dense(units=25, activation='relu'),  
 Dense(units=15, activation='relu'),  
 Dense(units=10, activation='sigmoid')  
])  
model.compile(loss=BinaryCrossEntropy())

More accurate loss (in code)

$$\text{loss} = -y \log\left(\frac{1}{1 + e^{-z}}\right) - (1-y) \log\left(1 - \frac{1}{1 + e^{-z}}\right)$$

logit: z

## More numerically accurate implementation of softmax

Softmax regression

$$(a_1, \dots, a_{10}) = g(z_1, \dots, z_{10})$$

$$\text{Loss} = L(\vec{a}, y) = \begin{cases} -\log a_1 & \text{if } y = 1 \\ \vdots \\ -\log a_{10} & \text{if } y = 10 \end{cases}$$

model = Sequential([  
 Dense(units=25, activation='relu'),  
 Dense(units=15, activation='relu'),  
 Dense(units=10, activation='softmax')  
])  
model.compile(loss=SparseCategoricalCrossEntropy())

More Accurate

$$L(\vec{a}, y) = \begin{cases} -\log \frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 1 \\ \vdots \\ -\log \frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 10 \end{cases}$$

model.compile(loss=SparseCategoricalCrossEntropy(from\_logits=True))

## MNIST (more numerically accurate)

```
model    import tensorflow as tf
         from tensorflow.keras import Sequential
         from tensorflow.keras.layers import Dense
         model = Sequential([
             Dense(units=25, activation='relu'),
             Dense(units=15, activation='relu'),
             Dense(units=10, activation='linear') ])
loss     from tensorflow.keras.losses import
         SparseCategoricalCrossentropy
         model.compile(..., loss=SparseCategoricalCrossentropy(from_logits=True) )
fit      model.fit(X, Y, epochs=100)
predict  logits = model(X) ← not  $a_1 \dots a_{10}$ 
         is  $z_1 \dots z$ 
         f_x = tf.nn.softmax(logits)
```

## logistic regression (more numerically accurate)

```
model    model = Sequential([
             Dense(units=25, activation='sigmoid'),
             Dense(units=15, activation='sigmoid'),
             Dense(units=1, activation='linear')
         ])
         from tensorflow.keras.losses import
             BinaryCrossentropy
         model.compile(..., BinaryCrossentropy(from_logits=True) )
         model.fit(X, Y, epochs=100)
fit      logit = model(X)  $\rightarrow z$ 
predict  f_x = tf.nn.sigmoid(logit)
```

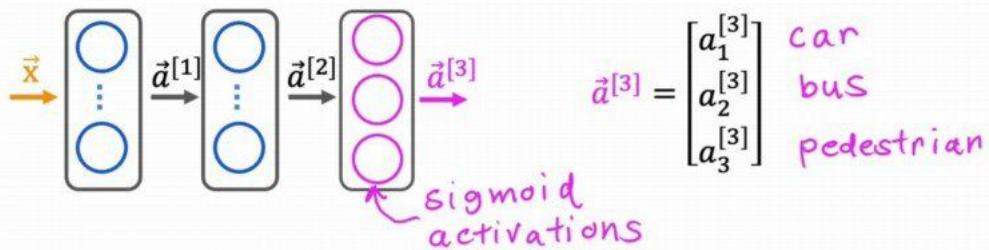
## Multi Label Classification

- In single output there are multiple class, like returning a vector of possibility
- For self driving car it have multiple labels like is there a car? is there a bus? is there a man?
- It output result as [ 1, 0, 1 ]
- We cannot make different NN for each we want it in single NN

## Multi-label Classification



Alternatively, train one neural network with three outputs



## Adam - Adaptive Moment Estimation

- Adam Optimizer is a alternative for Gradient Descent
- It is faster than GD
- learning rate is not a constant one but multiple learning rate is used
- different Alpha should be tried like so small and large even after using 0.001

## MNIST Adam model

```
model = Sequential([
    tf.keras.layers.Dense(units=25, activation='sigmoid'),
    tf.keras.layers.Dense(units=15, activation='sigmoid'),
    tf.keras.layers.Dense(units=10, activation='linear')
])

compile
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
               loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))

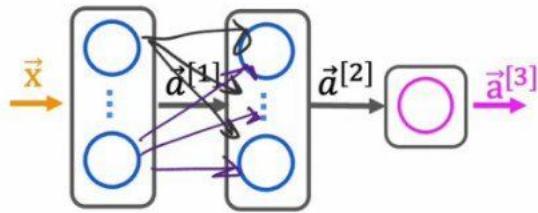
fit
model.fit(X, Y, epochs=100)
```

$\alpha = 10^{-3} = 0.001$

## Convolutional Neural Network - CNN

- It is a layer alternative for **Dense Layer**

## Dense Layer

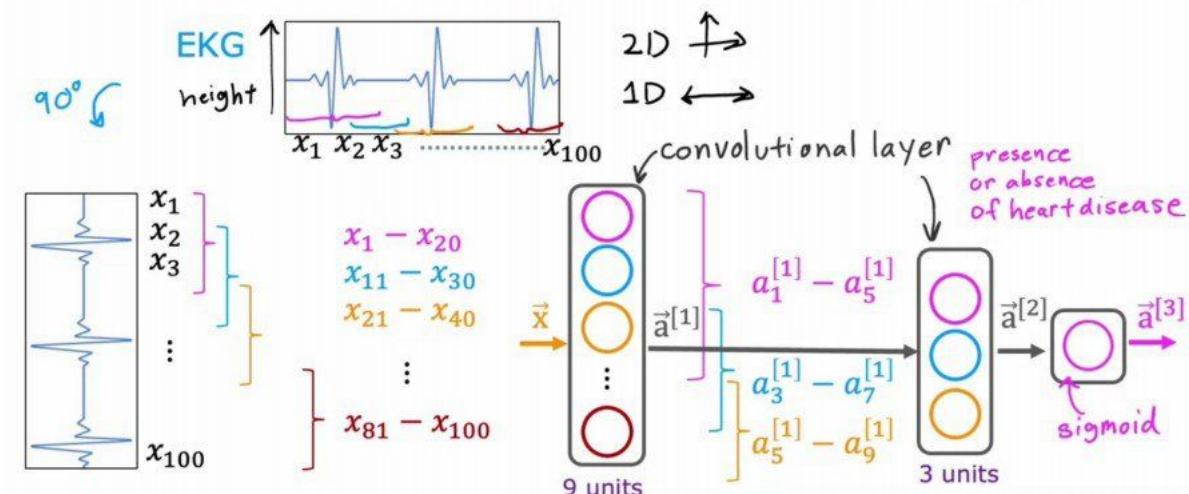


Each neuron output is a function of  
all the activation outputs of the previous layer.

$$\vec{a}_1^{[2]} = g(\vec{w}_1^{[2]} \cdot \vec{a}^{[1]} + b_1^{[2]})$$

- Hidden layer in CNN can use a set or part for creating a NN (MNIST)
- each neuron only look at part of the previous layer's input
- need less training data
- So due to this it can be fast
- we can avoid overfitting

## Convolutional Neural Network



## Debugging a learning algorithm

- Get more training data
- try small set of features
- try large set of features
- try adding polynomial features
- try decreasing and increasing learning rate and lambda (regularization parameter)
- we can use different diagnostic systems

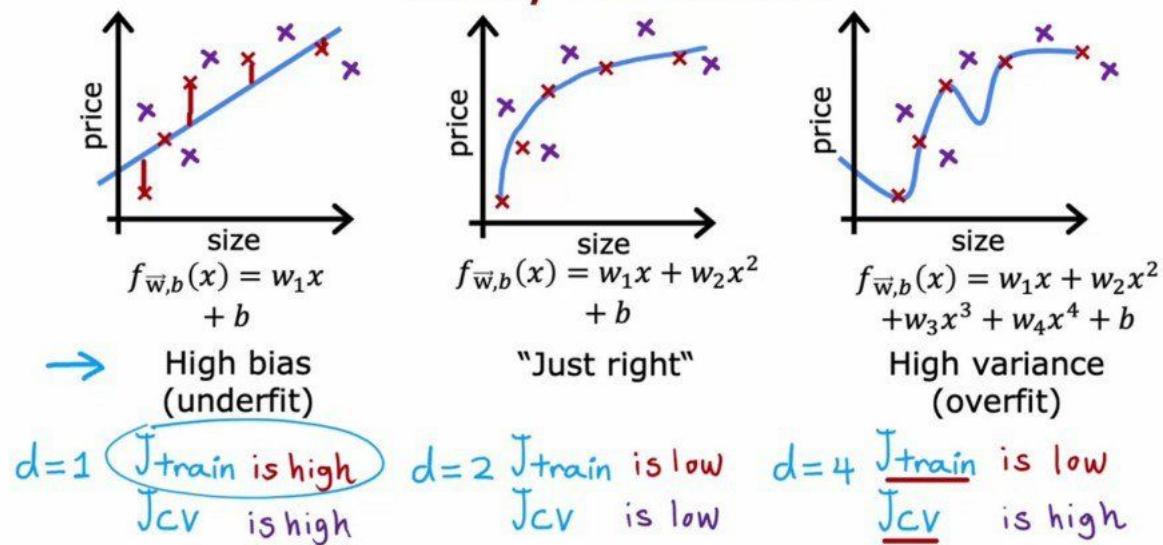
## Evaluating your Model

- Split the data training and testing.
- we can find training error and testing error
- cost function should be minimized
- MSE for Linear Regression
- Log Loss for Logistic Regression - Classification
- **For classification problem we can find what percentage of training and test set model which has predicted false, which is better than log loss**
- we can use polynomial regression and test the data for various degree of polynomial
- But to improve testing we can split data into training, cross validation set and test set
- Even for NN we can use this type of model selection
- Then we can find the cost function for each set

## Bias and Variance - Model Complexity or High Degree Polynomial

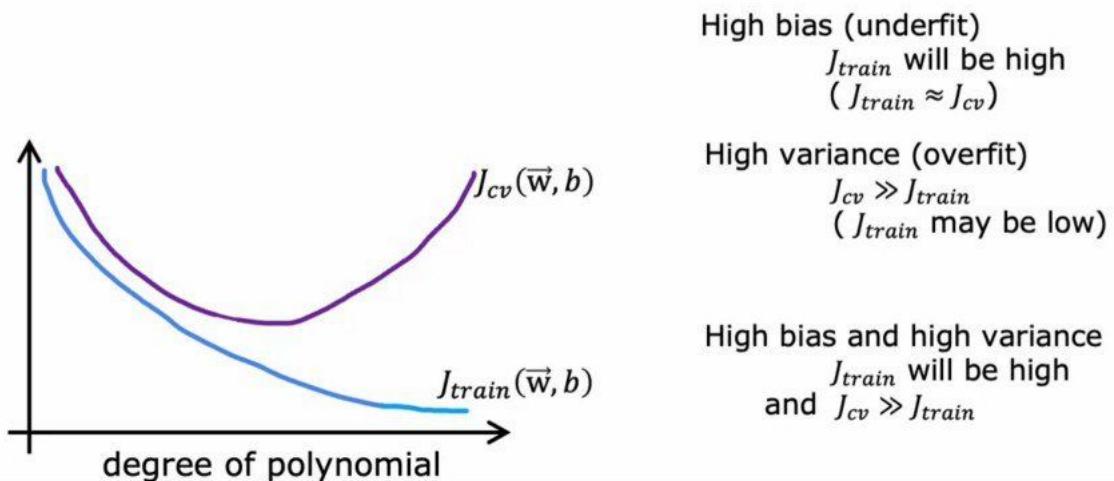
- We need low bias and low variance
- Bias - Difference between actual and predicted value
- Variance - Perform well on training set and worse on testing set
- degree of polynomial should not be small or large. It should be intermediate

## Bias/variance



## Diagnosing bias and variance

How do you tell if your algorithm has a bias or variance problem?

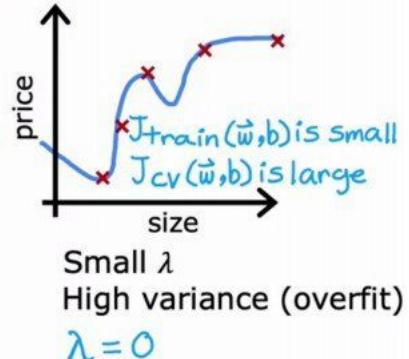
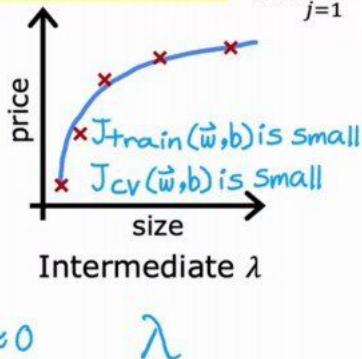
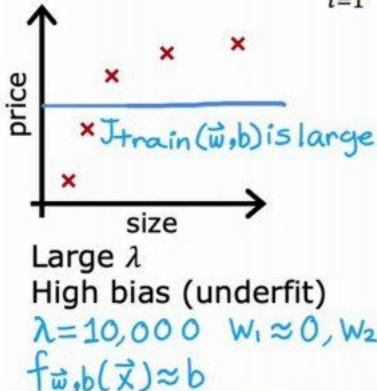


## Regularization - Bias and Variance

## Linear regression with regularization

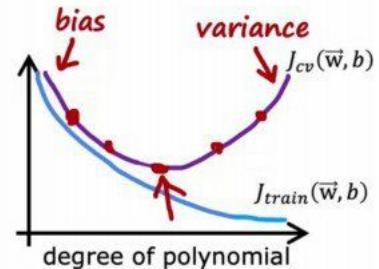
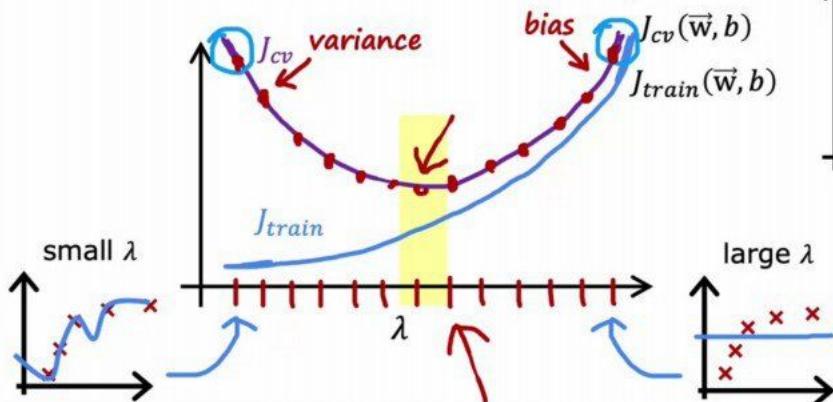
Model:  $f_{\vec{w}, b}(x) = \underline{w_1}x + \underline{w_2}x^2 + \underline{w_3}x^3 + \underline{w_4}x^4 + b$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$



### Bias and variance as a function of regularization parameter $\lambda$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$



- Underfit - low Poly degree, high lambda (regularization parameter)
- Overfit - High Poly degree, Low lambda (regularization parameter)
- Intermediate lambda and Poly degree is best

## Baseline

- First find out the baseline performance human can achieve or default baseline

## Bias/variance examples

Baseline performance	: 10.6%	0.2%	10.6%	4.4%	10.6%	4.4%
Training error ( $J_{train}$ )	: 10.8%	15.0%	15.0%	15.0%	19.7%	4.7%
Cross validation error ( $J_{cv}$ )	: 14.8%	15.5%	0.5%	19.7%	4.7%	

high variance      high bias      high bias  
                        high variance

## Learning Curve

### Bias/variance examples

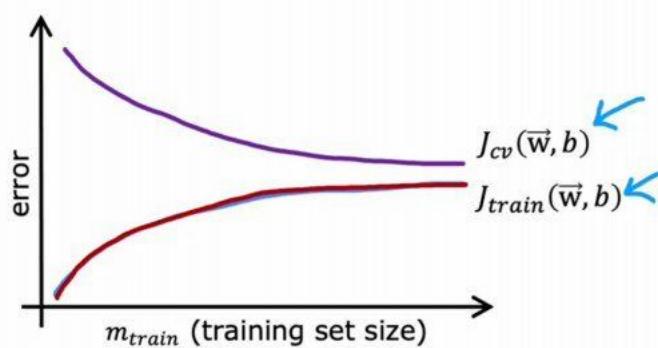
Baseline performance	: 10.6%	0.2%	10.6%	4.4%	10.6%	4.4%
Training error ( $J_{train}$ )	: 10.8%	15.0%	15.0%	15.0%	19.7%	4.7%
Cross validation error ( $J_{cv}$ )	: 14.8%	15.5%	0.5%	19.7%	4.7%	

high variance      high bias      high bias  
                        high variance

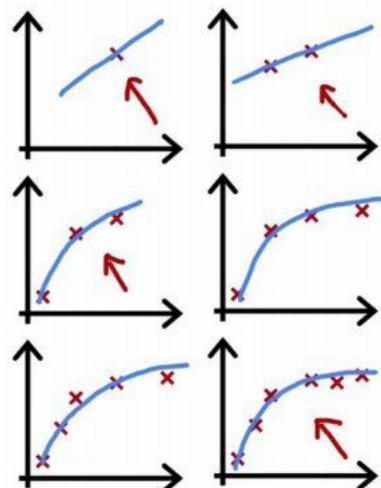
## Learning curves

$J_{train}$  = training error

$J_{cv}$  = cross validation error



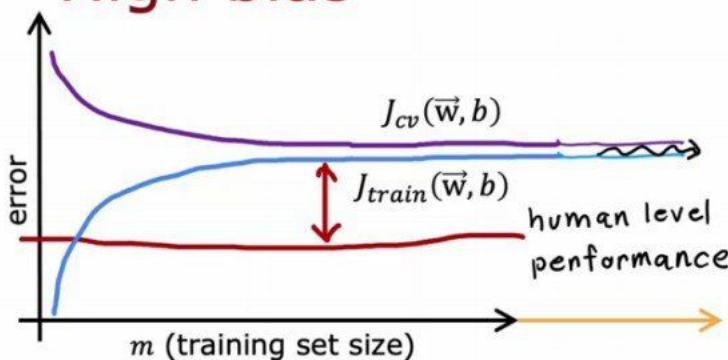
$$f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + b$$



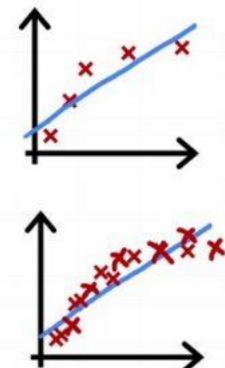
If we have high bias increasing training data is not going to help. It wont decrease the error

If we have high variance increasing training data is going to help. It will decrease the error

## High bias

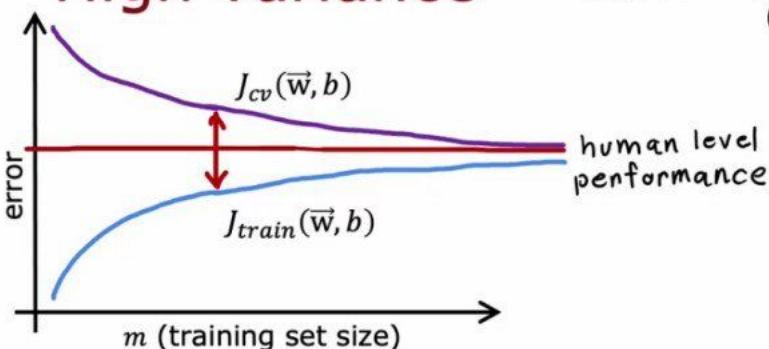


$$f_{\vec{w}, b}(x) = w_1 x + b$$

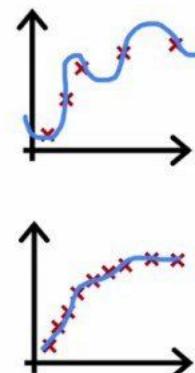


if a learning algorithm suffers from high bias, getting more training data will not (by itself) help much.

## High variance



$$f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b \quad (\text{with small } \lambda)$$



If a learning algorithm suffers from high variance, getting more training data is likely to help.

## Deciding what to try next revisited

### BIAS

- Bias is the difference between the average prediction of our model and the correct value.

- Model with high bias pays very little attention to the training data and oversimplifies the model.
- It always leads to high error on training and test data.

## VARIANCE

- Variance is the variability of model prediction for a given data point or a value which tells us spread of our data.
- Model with high variance pays a lot of attention to training data and does not generalize on test data.
- As a result, such models perform very well on training data but has high error rates on test data.

### **Overfitting - High Variance and Low Bias**

### **Underfitting - High/Low Variance and High Bias**

## **High Variance**

- Get more training data - High Variance
- try increasing lambda - Overfit High Variance
- try small set of features - Overfit High Variance

## **High Bias**

- try large set of features - Underfit High Bias
- try adding polynomial features - Underfit High Bias
- try decreasing lambda - Underfit High Bias

## **Bias and Variance Neural Networks**

- Simple Model - High Bias
- Complex Model - High Variance
- We need a model between them, that is we need to find a trade off between them.

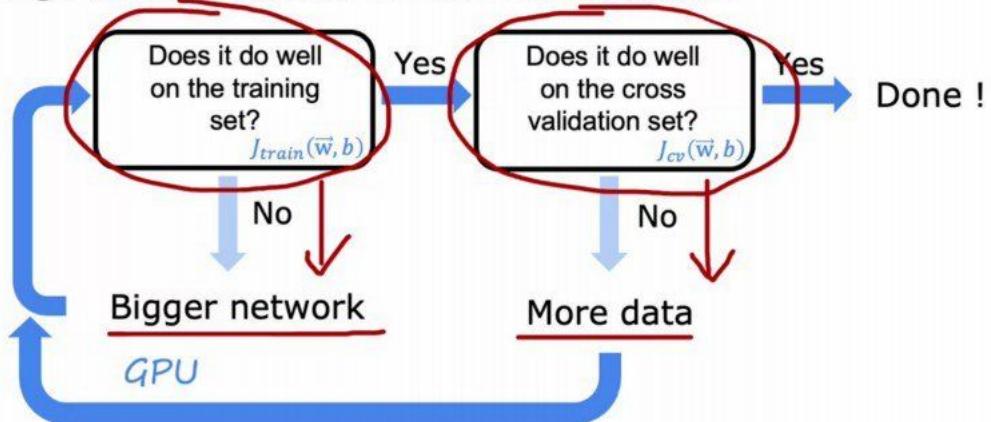
## **Neural Networks**

- Large NN are having low bias machines

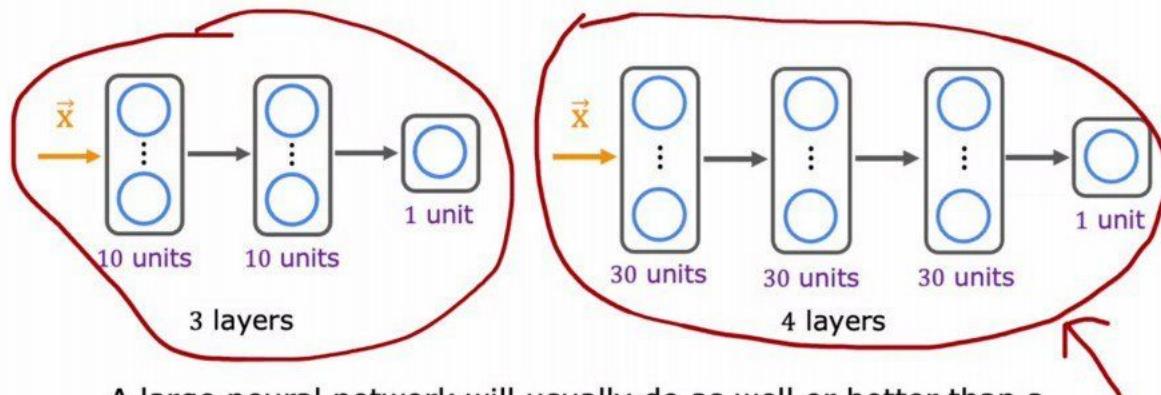
- GPU are there to speed up NN

## Neural networks and bias variance

Large neural networks are low bias machines



## Neural networks and regularization



- With best Regularization large NN will do better than smaller one.
- But large NN take time to run
- Implementing NN with regularization having lambda 0.01 (L2 regularization Ridge)

## Neural network regularization

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{\text{all weights } \mathbf{W}} (w^2)$$

Unregularized MNIST model

```
layer_1 = Dense(units=25, activation="relu")
layer_2 = Dense(units=15, activation="relu")
layer_3 = Dense(units=1, activation="sigmoid")
model = Sequential([layer_1, layer_2, layer_3])
```

Regularized MNIST model

```
layer_1 = Dense(units=25, activation="relu", kernel_regularizer=L2(0.01))
layer_2 = Dense(units=15, activation="relu", kernel_regularizer=L2(0.01))
layer_3 = Dense(units=1, activation="sigmoid", kernel_regularizer=L2(0.01))
model = Sequential([layer_1, layer_2, layer_3])
```

b

$\lambda$

## Iterative loop of ML development

- Choose architecture - model and data
- train model
- diagnostic - bias, variance, error analysis

## Building a spam classifier

Supervised learning:  $\vec{x}$  = features of email  
 $y$  = spam (1) or not spam (0)

Features: list the top 10,000 words to compute  $x_1, x_2, \dots, x_{10,000}$

$$\vec{x} = \begin{bmatrix} 0 \\ 1 \\ \text{zzz} \\ 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \quad \begin{array}{l} a \\ andrew \\ buy \\ deal \\ discount \\ \vdots \end{array}$$

From: cheapsales@buystufffromme.com  
To: Andrew Ng  
Subject: Buy now!  
  
Deal of the week! Buy now!  
Rolex w4tchs - \$100  
Medicine (any kind) - £50  
Also low cost M0rgages available.

# Building a spam classifier

How to try to reduce your spam classifier's error?

- Collect more data. E.g., "Honeypot" project.
- Develop sophisticated features based on email routing (from email header).
- Define sophisticated features from email body. E.g., should "discounting" and "discount" be treated as the same word.
- Design algorithms to detect misspellings. E.g., w4tches, med1cine, m0rtgage.



## Error Analysis

- Manually (human) going through a misclassified small dataset and analyzing the data
- We get an idea about where to focus

Error analysis involves the iterative observation, isolating, and diagnosing erroneous Machine learning (ML) predictions.

In error analysis, ML engineers must then deal with the challenges of conducting thorough performance evaluation and testing for ML models to improve model ability and performance.

Error Analysis works by

- Identification - identify data with high error rates
- Diagnosis - enables debugging and exploring the datasets further for deeper analysis
- Model debugging

This deepcheck's model error analysis check helps identify errors and diagnose their distribution across certain features and values so that you can resolve them.

```
from deepchecks.tabular.datasets.classification import adult
from deepchecks.tabular.checks import ModelErrorAnalysis

train_ds, test_ds = adult.load_data(data_format='Dataset', as_train_test=True)
model = adult.load_fitted_model()

# We create the check with a slightly lower r squared threshold to ensure that
```

```

# the check can run on the example dataset.

check = ModelErrorAnalysis(min_error_model_score=0.3)
result = check.run(train_ds, test_ds, model)
result

# If you want to only have a look at model performance at pre-defined
# segments, you can use the segment performance check.

from deepchecks.tabular.checks import SegmentPerformance
SegmentPerformance(feature_1='workclass',
                     feature_2='hours-per-week').run(validation_ds, model)

```

## Error analysis

$m_{cv} = \frac{500}{5000}$  examples in cross validation set.

Algorithm misclassifies  $\frac{100}{1000}$  of them.

Manually examine  $\frac{100}{1000}$  examples and categorize them based on common traits.

- Pharma: 21 → more data features
- Deliberate misspellings (w4tches, med1cine): 3
- Unusual email routing: 7
- Steal passwords (phishing): 18 → more data features
- Spam message in embedded image: 5

## Adding Data

- Include data where error analysis pointed

## Data Augmentation

Kind of feature engineering where we make more data from existing features

- Image Recognition
  - rotated image
  - enlarged image
  - distorted image
  - compressed image.
- Speech Recognition

- Original audio
- noisy background - crowd
- noisy background - car
- audio on a bad cellphone connection

For Image Text Recognition we can make our own data by taking screenshot of different font at different color grade

- **Synthetic Data Creation** is also now a inevitable part of majority projects.
- Focus on Data not the code

## What is Transfer Learning ?

**Transfer learning** make use of the knowledge gained while solving one problem and applying it to a different but related problem (same type of input like, image model for image and audio model for audio).

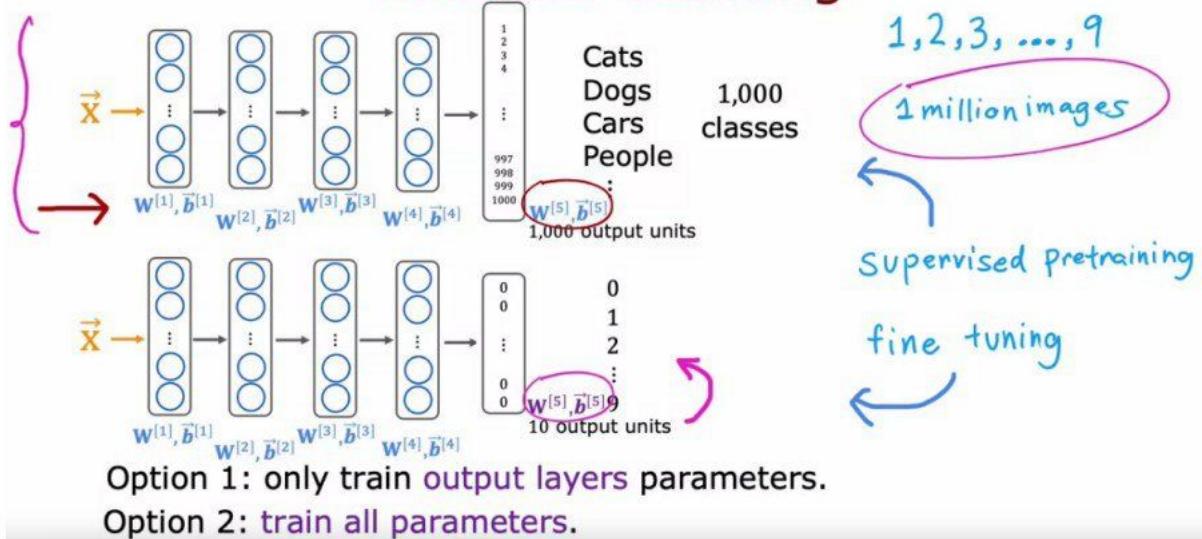
For example, knowledge gained while learning to recognize cars can be used to some extent to recognize trucks.

- We want to recognize hand written digits from 0 to 9
- We change the last output layer we wanted from the large NN all ready pre build.

## Pre Training

When we train the network on a **large dataset(for example: ImageNet)** , we train all the parameters of the neural network and therefore the model is learned. It may take hours on your GPU.

## Transfer learning



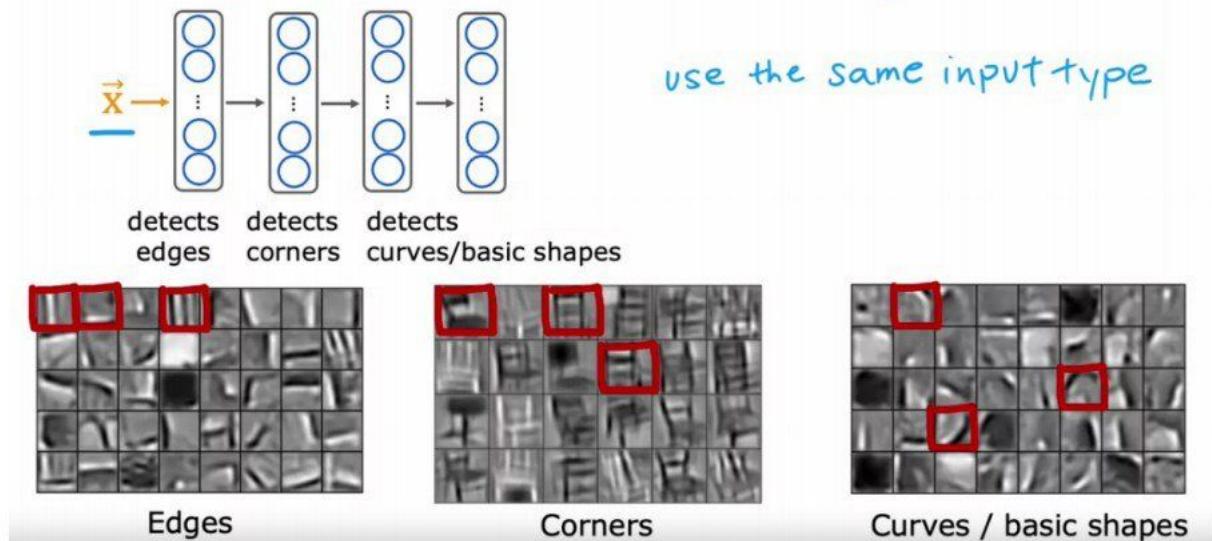
## Fine Tuning

We can give the new dataset to fine tune the pre-trained CNN. Consider that the new dataset is almost similar to the original dataset used for pre-training. Since the new dataset is similar, the same weights can be used for extracting the features from the new dataset.

1. If the new dataset is very small, it's better to train only the final layers of the network to avoid overfitting, keeping all other layers fixed. So remove the final layers of the pre-trained network. Add new layers **Retrain only the new layers**.
2. **If the new dataset is very much large, retrain the whole network** with initial weights from the pretrained model.

**How to fine tune if the new dataset is very different from the original dataset ?**

## Why does transfer learning work?



The earlier features of a ConvNet contain more **generic features** (e.g. **edge detectors** or **color blob detectors**), but later layers of the ConvNet becomes progressively more specific to the details of the **classes contained in the original dataset**.

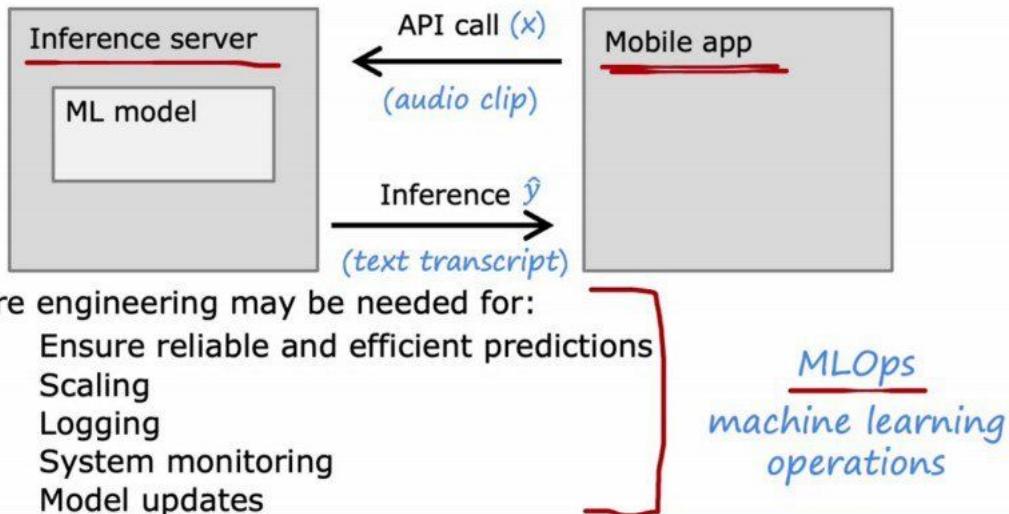
The earlier layers can help to extract the features of the new data. So it will be good if you fix the earlier layers and retrain the rest of the layers, if you got only small amount of data.

If you have large amount of data, you can retrain the whole network with weights initialized from the pre-trained network.

## Full cycle of a machine learning project

- Define Project
- Define and Collect data
- Train/Error Analysis/Iterative Improvement
- Deploy Monitor and Maintain

# Deployment



## Fairness, Bias, and Ethics

### Bias

- Hiring tool to discriminate women.
- Facial Recognition system matching dark skinned individuals to criminal mugshot.
- Biased bank loan approval
- Toxic effect of reinforcing negative stereotypes
- Using DeepFakes to create national issues/political purpose
- Spreading toxic speech for user engagement
- Using ML to build harmful products, commit fraud etc

### Guidelines

- Have a diverse team.
- Carry standard guidelines for your industry.
- Audit system against possible harm prior to deployment.
- Develop mitigation plan, monitor possible harm. (If self driving car get involved in accident)

Mitigation Plan - Reduces loss of life and property by minimizing the impact of disasters.

# Error Metrics

- For classifying a rare disease, Accuracy is not best evaluation metrics
- Precision Recall and F1Score helps to measure the classification accuracy

# Confusion Matrix in Machine Learning

Confusion Matrix helps us to display the performance of a model or how a model has made its prediction in Machine Learning.

Confusion Matrix helps us to visualize the point where our model gets confused in discriminating two classes. It can be understood well through a  $2 \times 2$  matrix where the row represents the actual truth labels, and the column represents the predicted labels.

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

# Accuracy

Simplest metrics of all, Accuracy. Accuracy is the ratio of the total number of correct predictions and the total number of predictions.

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

# Precision

Precision is the ratio between the True Positives and all the Positives. For our problem statement, that would be the measure of patients that we correctly identify having a heart/rare

disease out of all the patients actually having it.

Eg : Suppose I predicted 10 people in a class have heart disease. Out of those how many actually I predicted right.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

## Recall

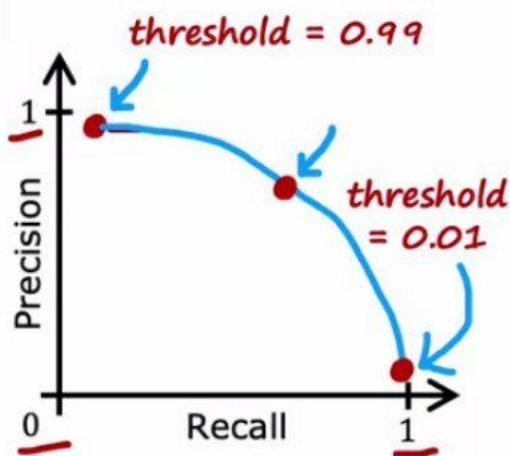
The recall is the measure of our model correctly identifying True Positives. Thus, for all the patients who actually have heart disease, recall tells us how many we correctly identified as having a heart disease.

Eg : Out of all people in a class having heart disease how many I got right prediction.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

## Trading off Precision and Recall

- If we want to predict 1 (rare disease) only if we are really confident, then
  - High Precision and Low Recall
  - We set high threshold, predict 1 if  $f(x) \geq 0.99$
- If we want to predict 1 (rare disease) when in small doubt, then
  - Low Precision and High Recall
  - We set small threshold, predict 1 if  $f(x) \geq 0.01$



## F1 Score

For some other models, like classifying whether a bank customer is a loan defaulter or not, it is desirable to have a high precision since the bank wouldn't want to lose customers who were denied a loan based on the model's prediction that they would be defaulters.

There are also a lot of situations where both precision and recall are equally important. For example, for our model, if the doctor informs us that the patients who were incorrectly classified as suffering from heart disease are equally important since they could be indicative of some other ailment, then we would aim for not only a high recall but a high precision as well.

In such cases, we use something called F1-score is used . F1-score is the Harmonic mean of the Precision and Recall

**Precision, Recall and F1 Score should be close to one, if it is close to zero then model is not working well. (General case)**

$$F1\ Score = 2 * \frac{Precision * Recall}{Precision + Recall}$$

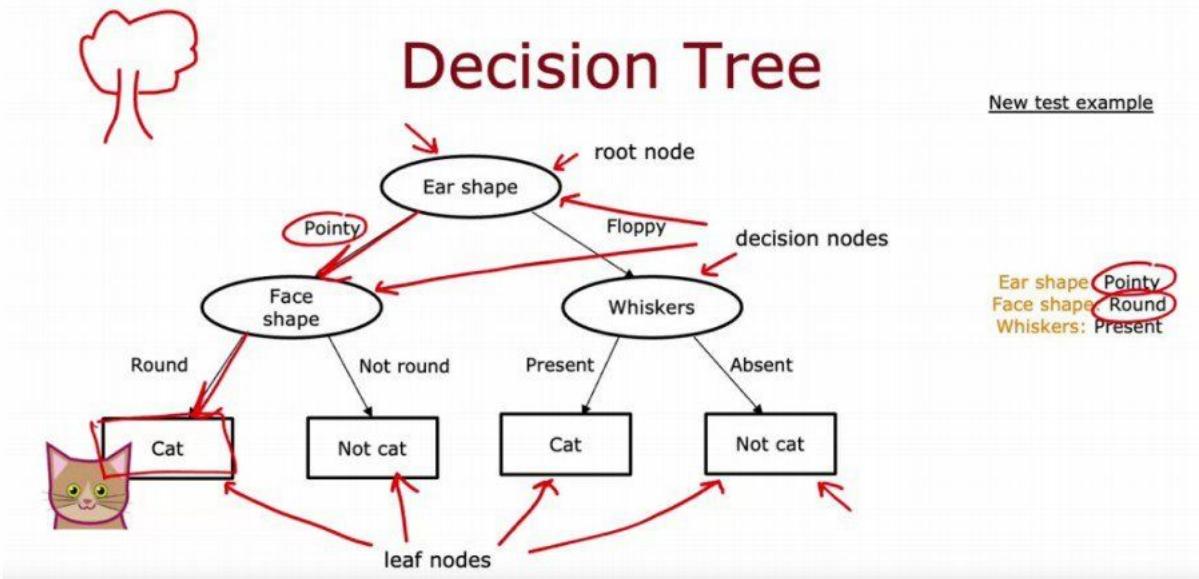
## Decision Tree Model

- If features are categorical along with a target, for a classification problem decision tree is a good model
- Starting of DT is called **Root Node**
- Nodes at bottom is **Leaf Node**
- Nodes in between them is **Decision Node**
- The purpose of DT is to find best tree from the possible set of tree

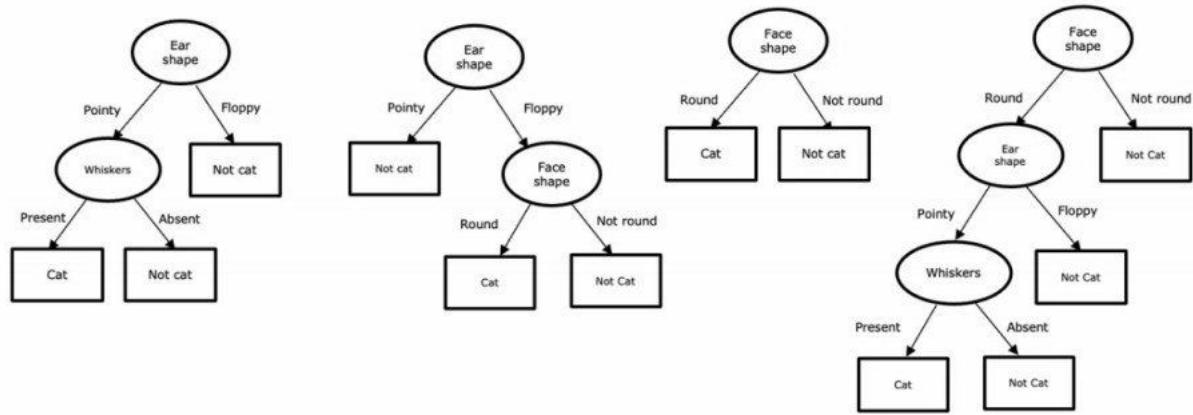
# Cat classification example

	Ear shape ( $x_1$ )	Face shape ( $x_2$ )	Whiskers ( $x_3$ )	Cat
	Pointy ↙	Round ↙	Present ↙	1
	Floppy ↙	Not round ↙	Present	1
	Floppy	Round	Absent ↙	0
	Pointy	Not round	Present	0
	Pointy	Round	Present	1
	Pointy	Round	Absent	1
	Floppy	Not round	Absent	0
	Pointy	Round	Absent	1
	Floppy	Round	Absent	0
	Floppy	Round	Absent	0

Categorical (discrete values)      X      Y



# Decision Tree



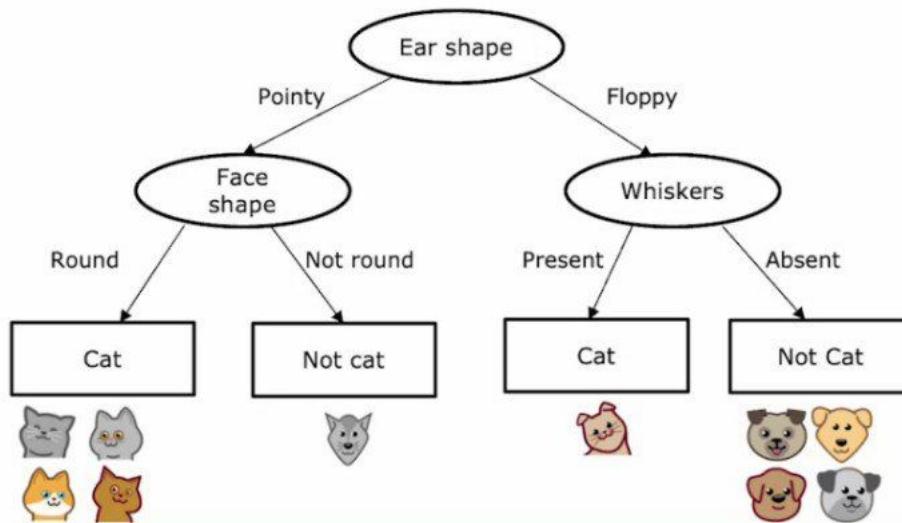
## Decision Tree Learning Process

- If we reach at a node full of specific class, we stop there and set the node as leaf node.
- we want to achieve purity, that is class full of same type (Not completely possible all time).
- In order to maximize the purity, we need to decide where we need to split on at each node.

## When to stop splitting?

- When node is 100% one class
- Reaching maximum depth of tree
- When depth increases - Overfitting
- When depth decreases - Underfitting
- when impurity score improvements are below a threshold
- Number of examples in a node is below a threshold

# Decision Tree Learning

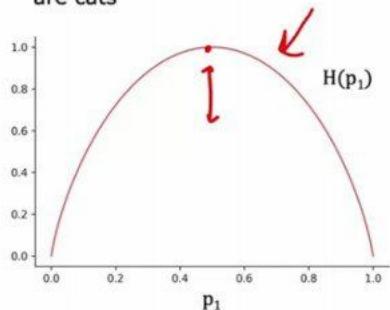


## Measuring Purity

- Entropy is measure of Impurity/randomness in data
- It starts from 0 to 1 then goes back to 0
- $P_1$  is the fraction of positive examples that are same class
- $P_0$  is opposite class
- Entropy equation is as follows;

## Entropy as a measure of impurity

$p_1$  = fraction of examples that are cats



$$p_0 = 1 - p_1$$

$$H(p_1) = -p_1 \log_2(p_1) - p_0 \log_2(p_0)$$

$$= -p_1 \log_2(p_1) - (1 - p_1) \log_2(1 - p_1)$$



Note: " $0 \log(0) = 0$ "

# Choosing a split: Information Gain

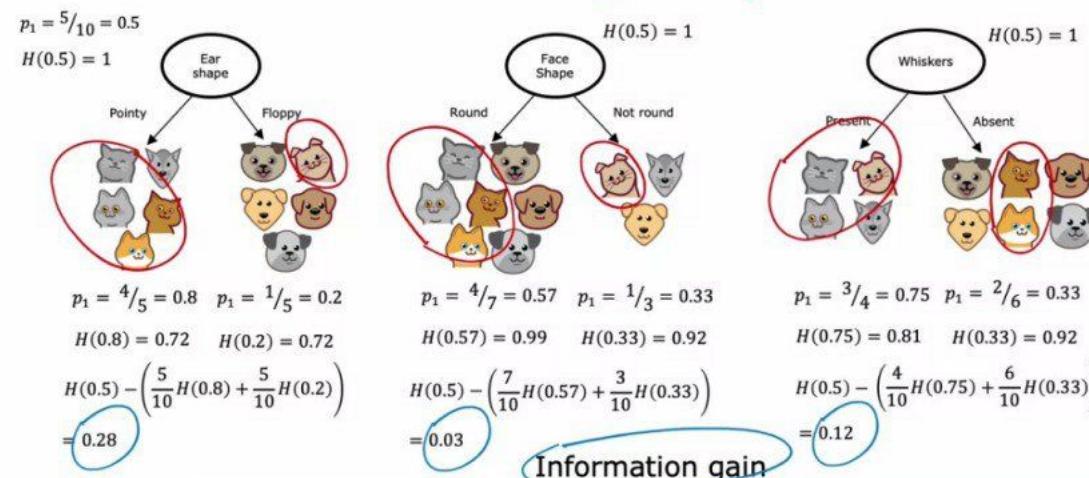
## Reduction of Entropy is Information Gain

Some Initial Steps

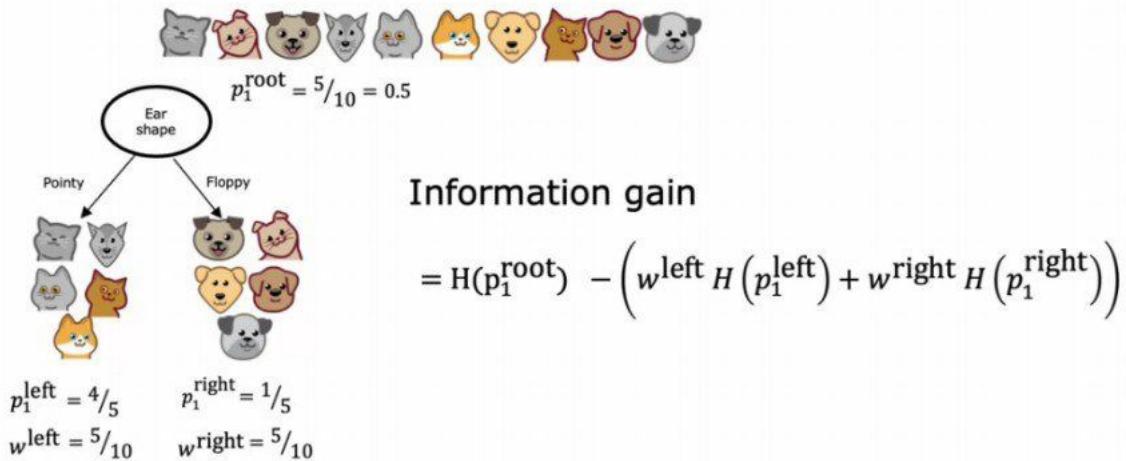
- We calculate the Entropy in each decision node - Low Entropy is Good, Shows Less Impurity
- Also number of elements in sub branch is also important - High Entropy in Lot of Samples is Worse
- We need to combine the Entropy of two side somehow
- First calculate product of ratio of elements in node with Entropy. Take sum of both of that value from each side.
- Subtract from root node P1

By this we get reduction in entropy known as Information Gain. Then pick the largest Information gain for best output in decision tree.

## Choosing a split



# Information Gain



## Putting it together

- Start with all examples at root node
- Calculate Information Gain for all features and pick one with large IG
- Split data into left and right branch
- keep repeating until,
  - Node is 100% one class
  - When splitting result in tree exceeding depth. Because it result in Overfitting
  - If information gain is less than a threshold
  - if number of examples in a node is less than a threshold

Decision Tree uses Recursive Algorithm

## Using one-hot encoding of categorical features

- One hot encoding is a process of converting categorical data variables to features with values 1 or 0.
- One hot encoding is a crucial part of feature engineering for machine learning.
- In case of binary class feature we can convert in just one column with either 0 or 1
- Can be used in Logistic/Linear Regression or NN

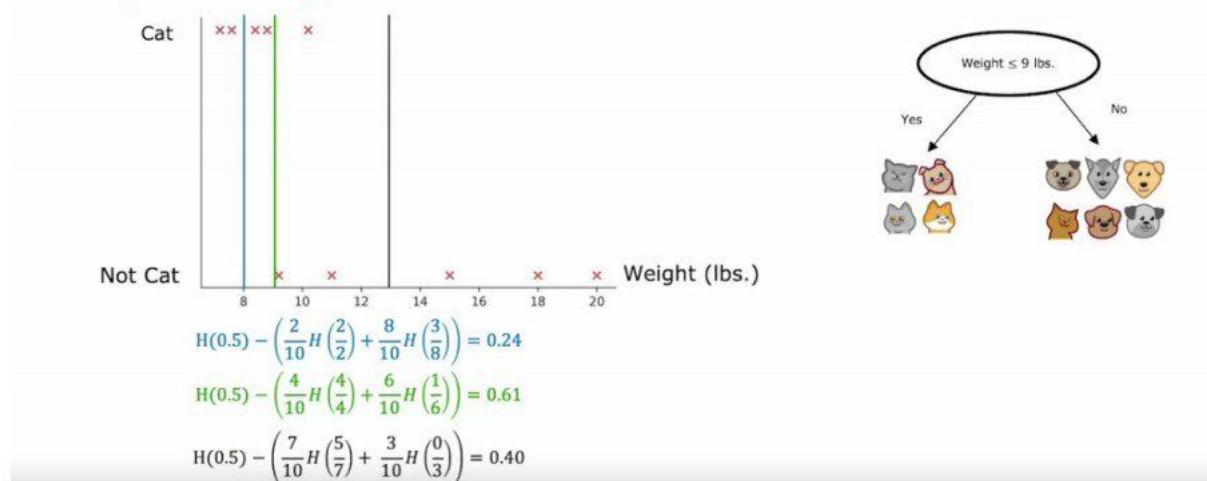
# One hot encoding and neural networks

Pointy ears	Floppy ears	Round ears	Face shape	Whiskers	Cat
1	0	0	Round 1	Present 1	1
0	0	1	Not-round 0	Present 1	1
0	0	1	Round 1	Absent 0	0
1	0	0	Not-round 0	Present 1	0
0	0	1	Round 1	Present 1	1
1	0	0	Round 1	Absent 0	1
0	1	0	Not-round 0	Absent 0	1
0	0	1	Round 1	Absent 0	1
0	1	0	Round 1	Absent 0	1
0	1	0	Round 1	Absent 0	1

## Decision Tree for Continuous valued features

- Below example weight is continuous
- We can consider weight  $\leq x$  lbs, to split into 2 classes
- The value of  $x$  can be changed and calculate Information Gain for each case.
- General case is to use 3 cases but can be changed.

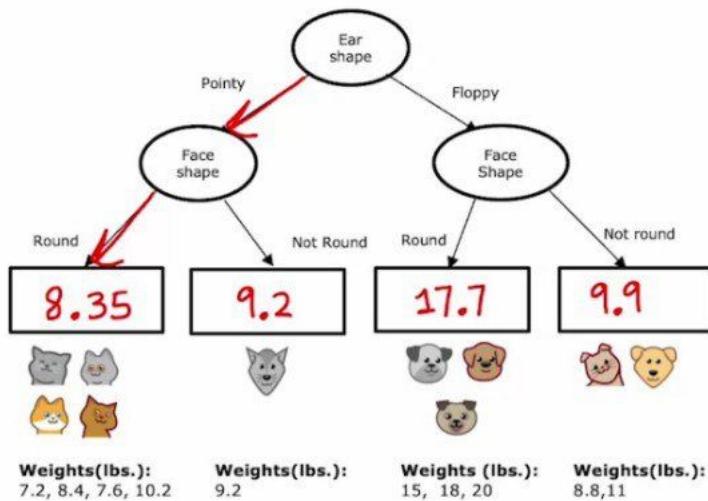
## Splitting on a continuous variable



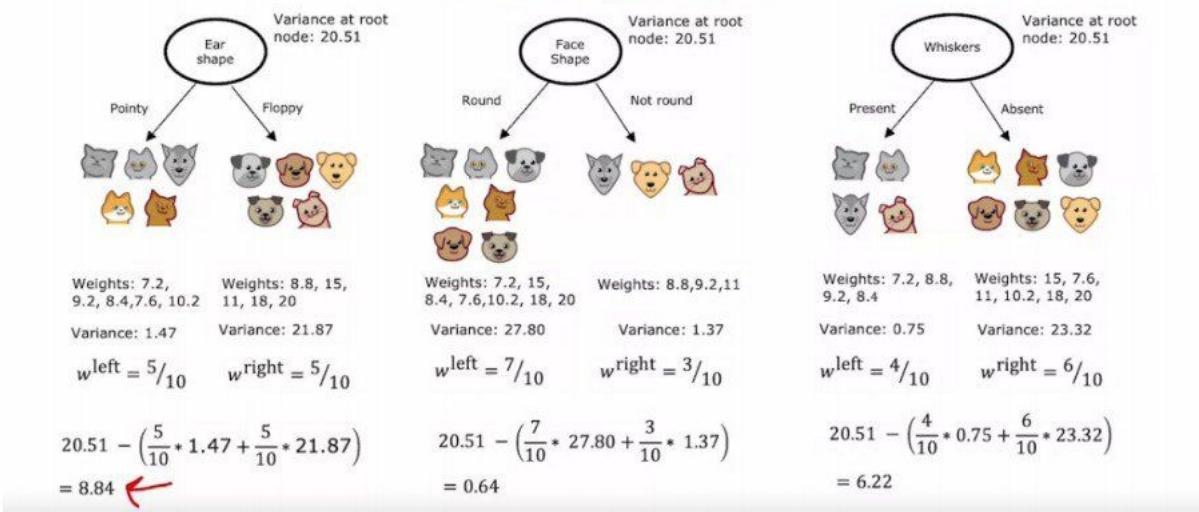
## Regression Trees

- It is same as Decision problem. At the end when we reach leaf node we calculate the average and predict
- But instead of Information Gain we calculate **Reduction in Variance**
- Calculate the variance of the leaf node then multiply by the ratio of number of elements, finally subtract it from variance of root node
- We choose split which give **Largest Reduction in Variance**

## Regression with Decision Trees



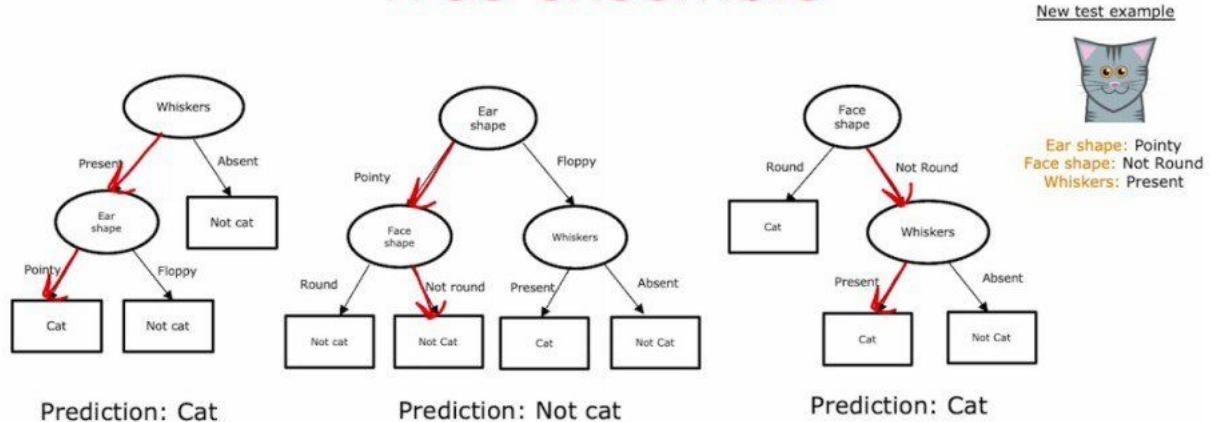
## Choosing a split



## Using Multiple Decision Trees

- To avoid the sensitivity and make model robust we can use Multiple Decision Tree called **Tree ensembles (collection of multiple tree)**
- We take majority of the Ensemble prediction

## Tree ensemble



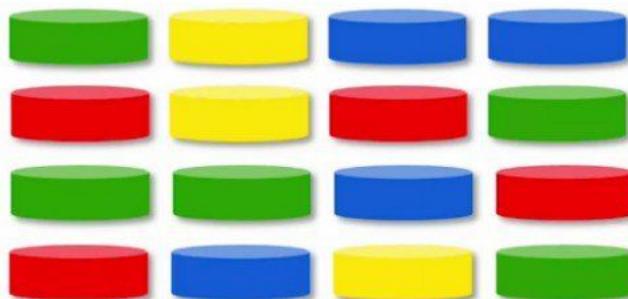
## Sampling with replacement

- Help to construct new different but similar training set.

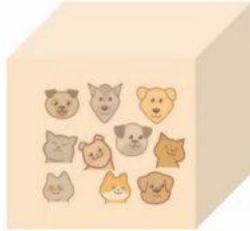
## Sampling with replacement

Tokens

Sampling with replacement:



# Sampling with replacement



Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	1
Floppy	Not round	Absent	0
Pointy	Round	Absent	1
Pointy	Not round	Present	0
Floppy	Not round	Absent	0
Pointy	Round	Absent	1
Pointy	Round	Present	1
Floppy	Not round	Present	1
Floppy	Round	Absent	0
Pointy	Round	Absent	1

## Random Forest Algorithm

- Pick a train data
- make duplicate of train data using sampling by replacement technique, Increasing this number of synthetic samples is ok. It do improve performance. But after a limit it is just a waste of GPU.
- Since we create new Decision Tree using this technique, it is also called **Bagged Decision Tree**
- Sampling helps Random Forest to understand small changes.
- Using different Decision Tree and averaging it improve robustness of Random Forest

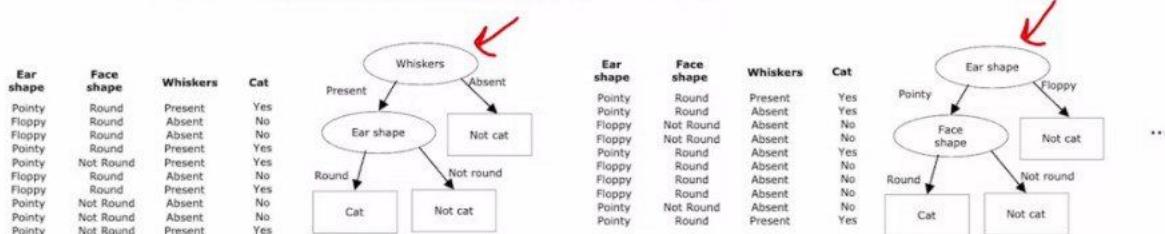
## Generating a tree sample

Given training set of size  $m$

For  $b = 1$  to  $B$

Use sampling with replacement to create a new training set of size  $m$

Train a decision tree on the new dataset



Bagged decision tree

Suppose that in our example we had three features available rather than picking from all end features, we will instead pick a random subset of K less than N features. And allow the algorithm to choose only from that subset of K features. So in other words, you would pick K features as the allowed features and then out of those K features choose the one with the highest information gain as the choice of feature to use the split. When N is large, say n is Dozens or 10's or even hundreds. A typical choice for the value of K would be to choose it to be square root of N.

## XGBoost

- Same like Random Forest, but in sampling while choosing next sample, we give more preference for picking the wrongly predicted one in the First made Decision Tree.

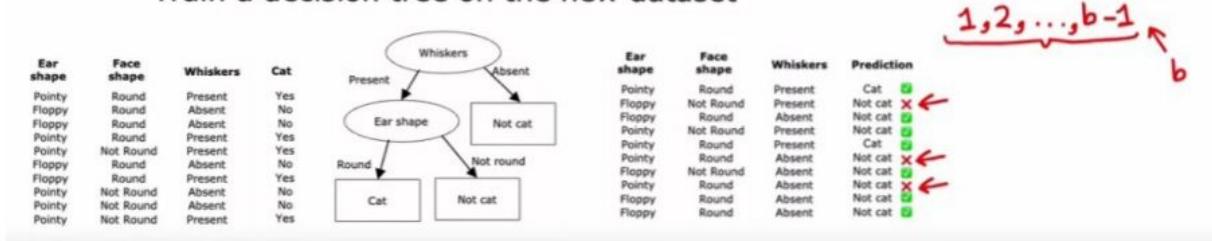
## Boosted trees intuition

Given training set of size  $m$

For  $b = 1$  to  $B$ :

Use sampling with replacement to create a new training set of size  $m$   
But instead of picking from all examples with equal  $(1/m)$  probability, make it more likely to pick misclassified examples from previously trained trees

Train a decision tree on the new dataset



## XGBoost (eXtreme Gradient Boosting)

- Open source implementation of boosted trees
- Fast efficient implementation
- Good choice of default splitting criteria and criteria for when to stop splitting
- Built in regularization to prevent overfitting
- Highly competitive algorithm for machine learning competitions (eg: Kaggle competitions)

```

# Classification
from xgboost import XGBClassifier
model = XGBClassifier()
model.fit(x_train, y_train)
y_pred = model.predict(x_test)

# Regression
from xgboost import XGBRegressor
model = XGBRegressor()
model.fit(x_train, y_train)
y_pred = model.predict(x_test)

```

## When to use Decision Trees

Decision Tree and Tree Ensembles	Neural Networks
Works well on tabular data	Works well on Tabular ( Structured and Unstructured data)
Not recommended for Images, audio and text	Recommended for Image, audio, and text
fast	slower than DT
Small Decision tree may be human interpretable	works with transfer learning
We can train one decision tree at a time.	when building a system of multiple models working together, multiple NN can be stringed together easily. We can train them all together using gradient descent.

## Unsupervised Learning, Recommenders, Reinforcement Learning - Course 3

### What is Clustering?

- Clustering is used to group unlabeled data
- There are various algorithms to perform the Clustering.