

Data Science Training – Day 4

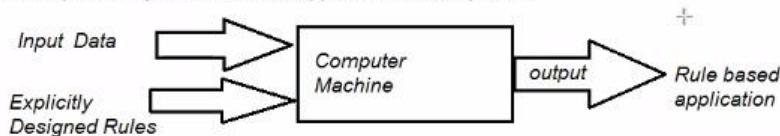
Topic: Neural Network / Deep Learning

- ANN
- CNN
- RNN

Use case:

Explicit coding the logic

*When WE Get a Use case
1) Try out the possibility of Rule Based Application Development*

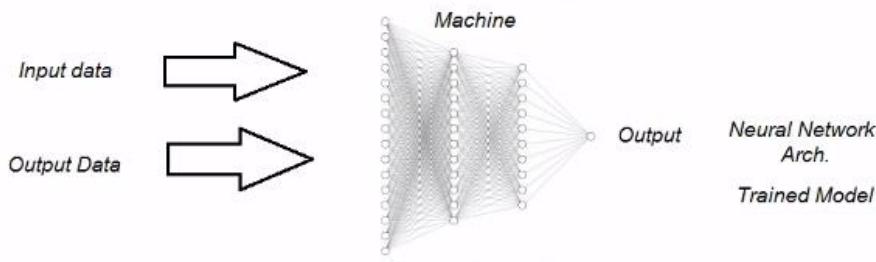


Develop the solution with explicit defined the rule-based application

<i>Independent variable</i>	<i>Dependent variable</i>
<i>Input Data</i>	<i>Output Data</i>
<i>YrsExp</i>	<i>Salary</i>
<i>Investment, Tax, Revenue</i>	<i>Profit</i>
<i>Age, Salary, City</i>	<i>Response</i>

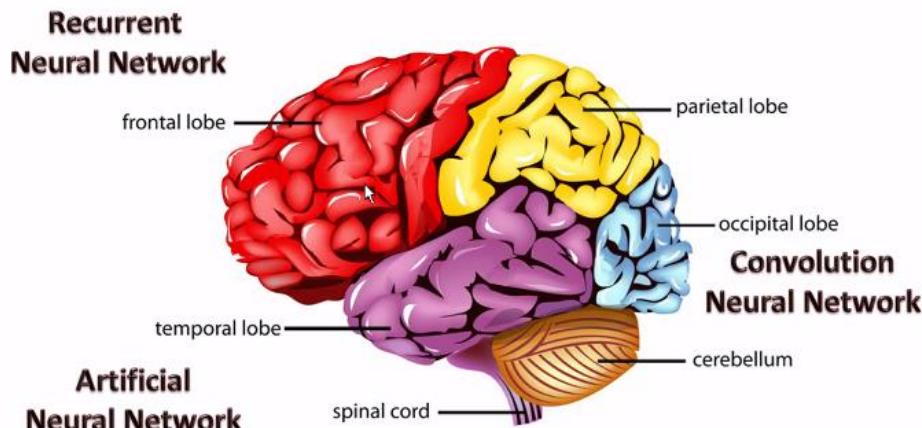
*3) We Go For Deep Learning
If We have Large Dataset
Complex Unstructured Data to deal with
High Speed Solution*

Accepts Input Data & Output Data, Feature Extraction is implicitly done by Deep Learning



Human Brain

Cerebrum



CNN - human interpretation by seeing pictures.

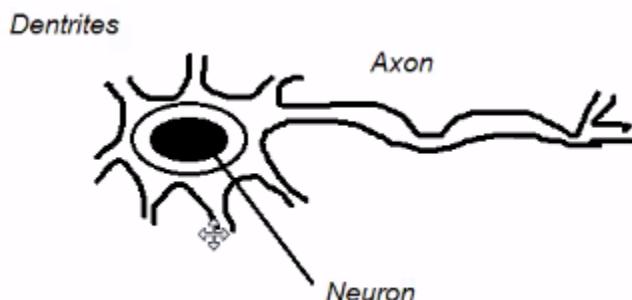
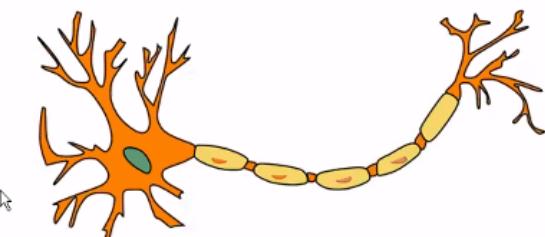
ANN -

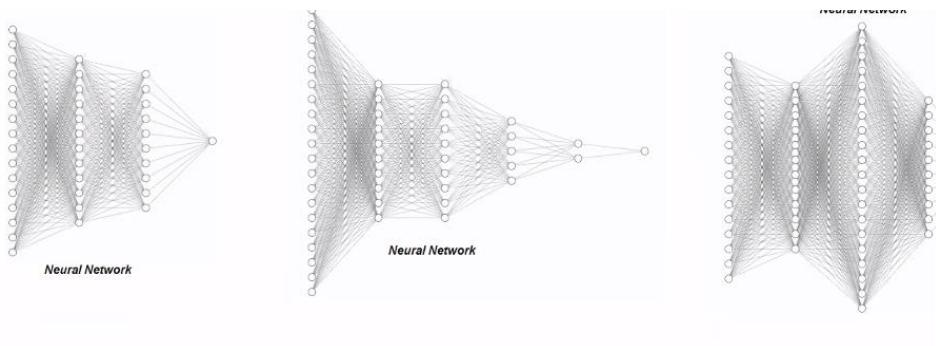
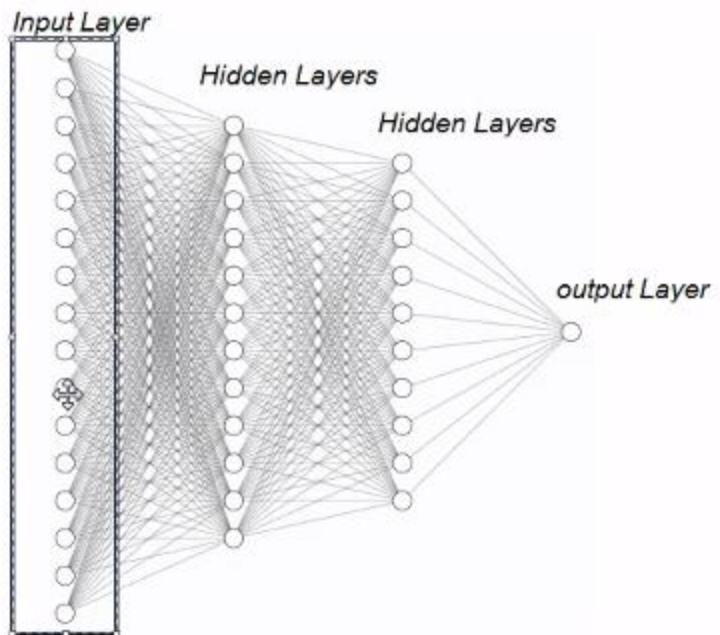
RNN – What did you have for breakfast? As time goes, you start forgetting it. On Thursday, what did you had for breakfast?

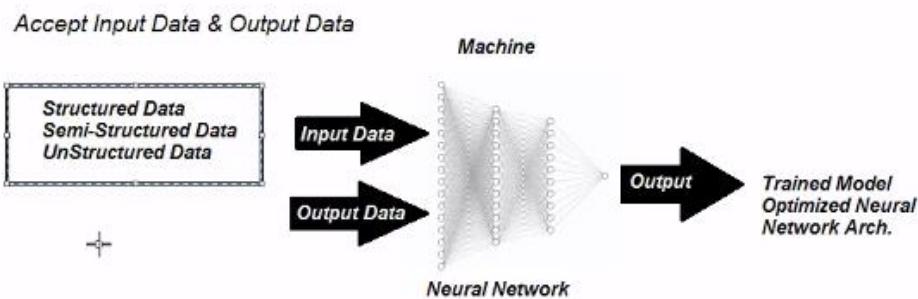
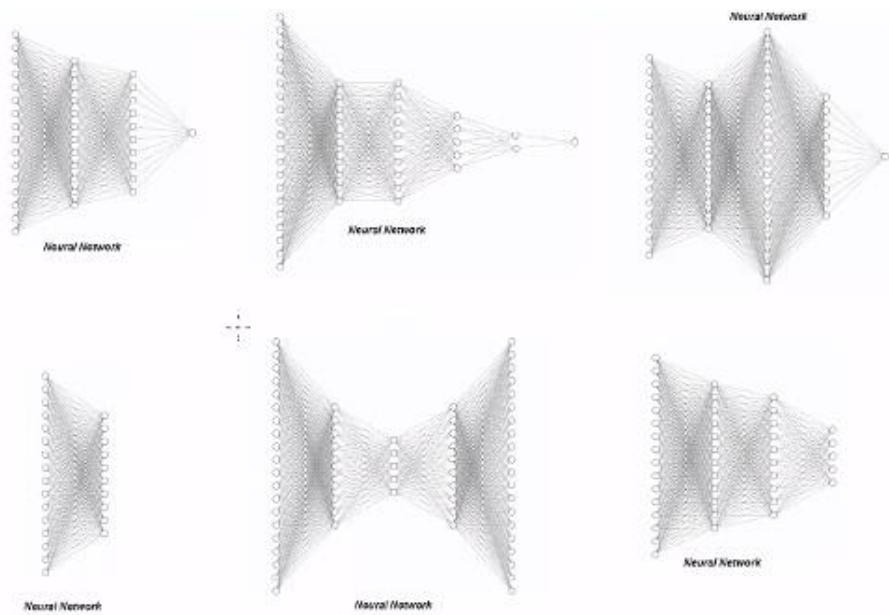
Process data in sequence, as time goes in sequence, your brain starts forgetting things.

Deep Learning

- In deep learning, these layered representations are learned via models called "neural networks", structured in literal layers stacked one after the other.
- The term "neural network" is a reference to neurobiology.

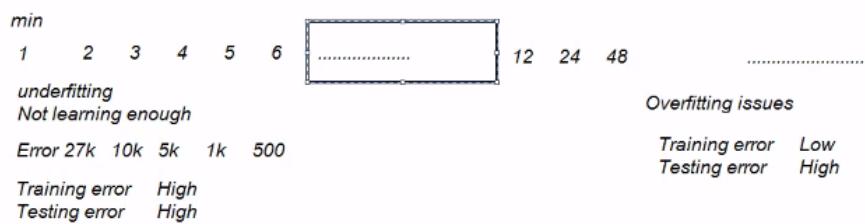
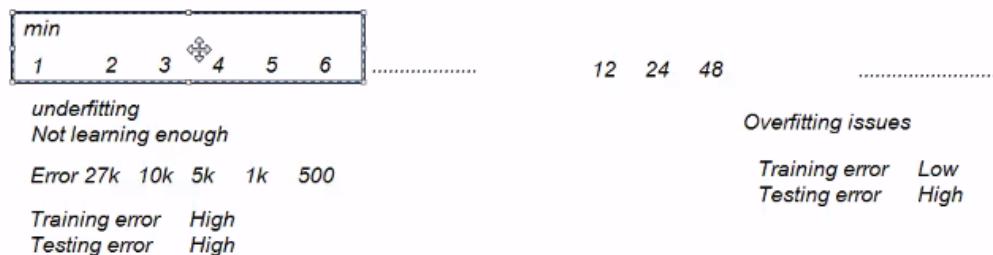
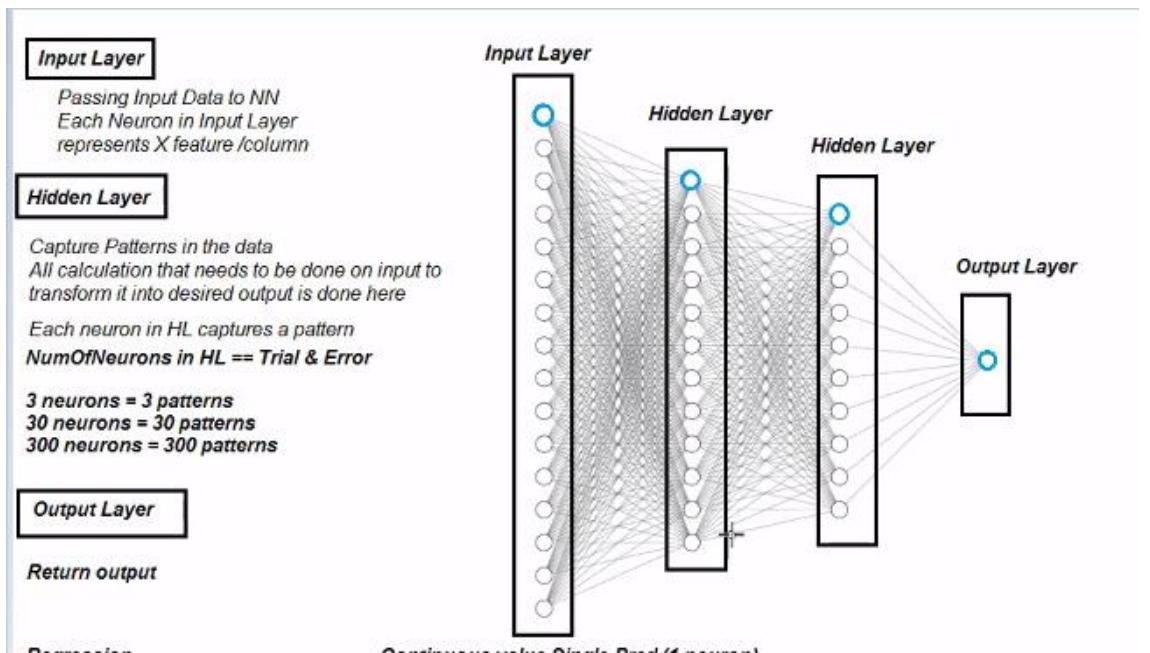




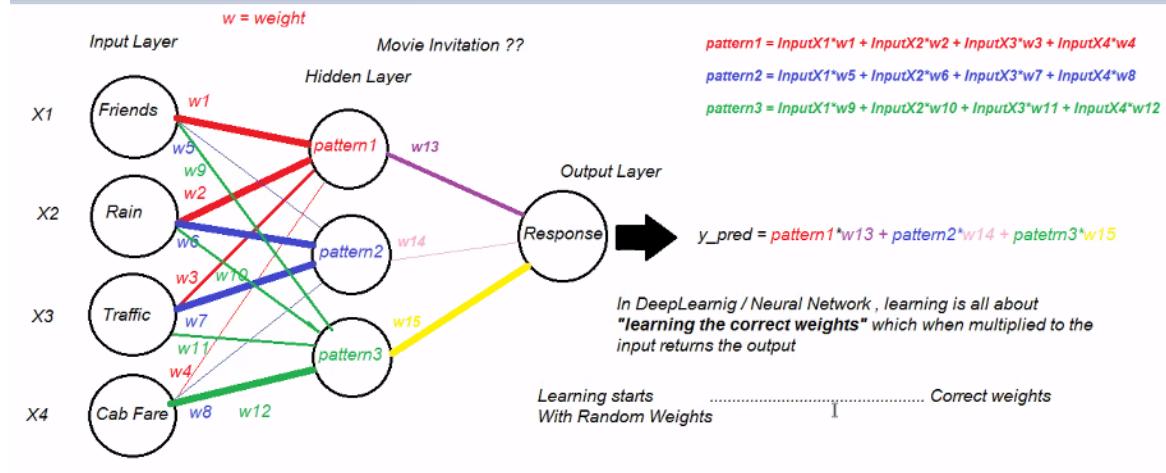


Layered structured

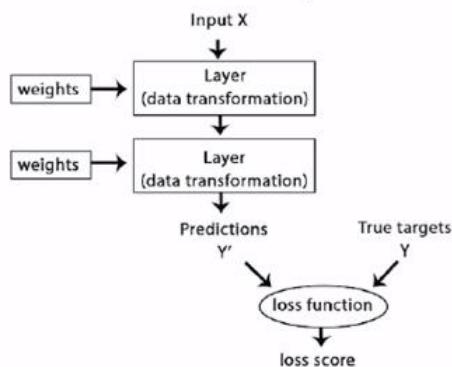
Represented as neuron



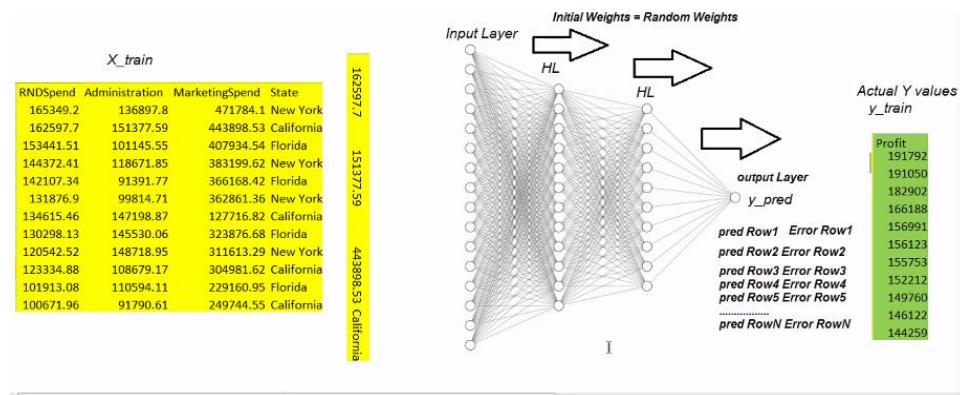
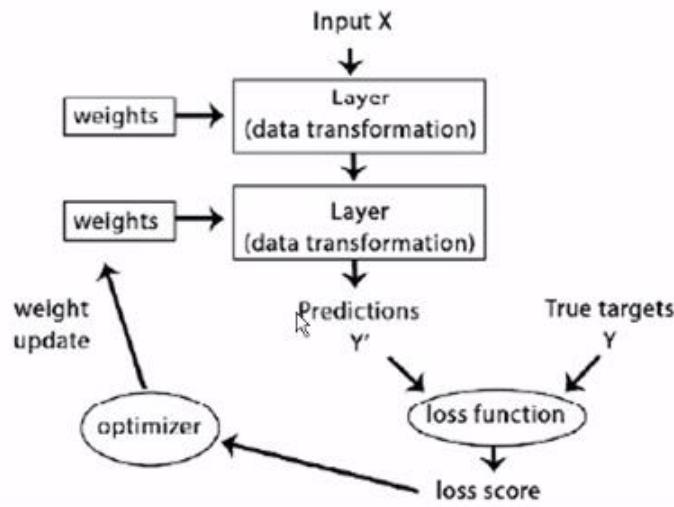
Regression	Continuous value Single Pred (1 neuron)
Binary Classification	Discrete Binary value 0/1 Pred (1 neuron)
Multi-Class Classification	Discrete Value (more than 1 neuron)

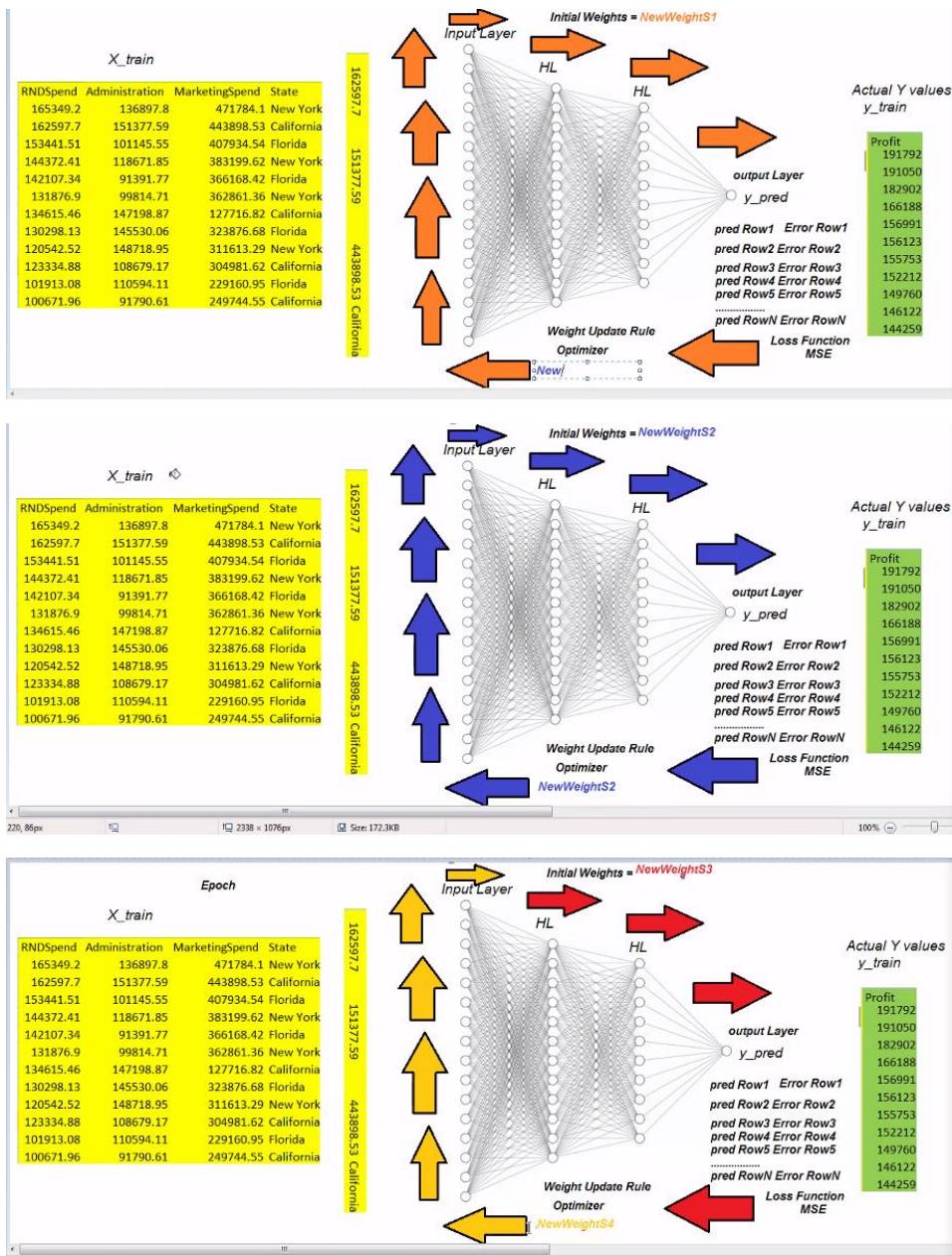


- To control the output of a neural network, you need to be able to measure how far this output is from what you expected.
- This is the job of the "**loss function**" of the network, also called "**objective function**".
- The loss function takes the predictions of the network and the true target , and computes a distance score, capturing how well the network has done on this specific example.



- The fundamental trick in deep learning is to use the loss score as a feedback signal to adjust the value of the weights by a little bit, in a direction that would lower the loss score for the current example.
- This adjustment is the job of the "optimizer", which implements what is called the "**backpropagation**" algorithm, the central algorithm in deep learning.





UnderFitting

Not Learning Enough

Good Fit

Learned well to generalize

OverFitting

Too much of Learning

TrainingAcc	TestingAcc	
60%	60%	Low Acc Underfitting
90%	90%	Good Fit
90%	95%	Good Fit
90%	75%	OverFitting

`numOfEpochs` = To be decided by developer (Trial & Error)

Step1: Initial Weights = Random Weight Allocation (All rows in training set is used against the initial weight)
 Epoch1 $y_{pred} = InputX1 * weight1 + InputX2 * weight2 + InputX3 * weight3 + \dots + InputXn * weightn$
 Error = lossFunction(y_{pred}, y_{actual})
 WeightUpdate = optimizer (weight update Rule)
 NewWeights1

Step2: Initial Weights = NewWeightS1
 Epoch2 $y_{pred} = InputX1 * weight1 + InputX2 * weight2 + InputX3 * weight3 + \dots + InputXn * weightn$
 Error = lossFunction(y_{pred}, y_{actual})
 WeightUpdate = optimizer (weight update Rule)
 NewWeightS2

```

Step3: Initial Weights = NewWeightS2
Epoch3   y_pred = InputX1*weight1 + InputX2*weight2 + InputX3*weight3 .....+InputXn*weightn
          Error = lossFunction(y_pred,y_actual)
          WeightUpdate = optimizer ( weight update Rule)
          NewWeightS3

```

StepN: Initial Weights = NewWeightSN-1
 EpochN $y_{pred} = InputX1 * weight1 + InputX2 * weight2 + InputX3 * weight3 + \dots + InputXn * weightn$
 Error = lossFunction(y_{pred}, y_{actual})
 WeightUpdate = optimizer (weight update Rule)
 NewWeightSN = Final Weights giving max accuracy and lowest error / loss

Formula at the heart of Each Neuron [Linear Spread of Data]

```
ty_pred = ActivationFunction(bias + InputX1*weight1 + InputX2*weight2 + InputX3*weight3.....+InputXn*weightn)
```

Activation Function

Linear Activation Functions (Linear Spread of data)

Non-Linear Activation Functions (Linear Spread of data & Non-Linear Spread of data)

Linear Regression - Linear Spread of Data

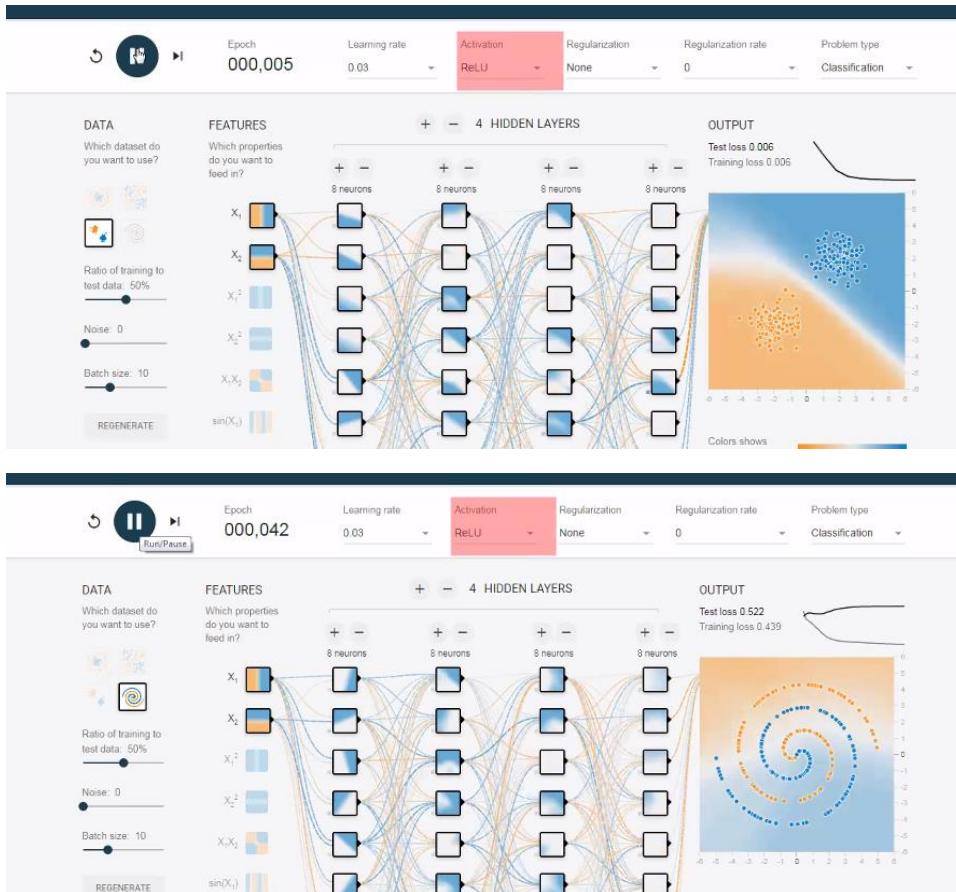
$$y_{pred} = b_0 + InputX_1 \cdot b_1 + InputX_2 \cdot b_2 + InputX_3 \cdot b_3 + \dots + InputX_n \cdot b_n$$

Bias = intercept

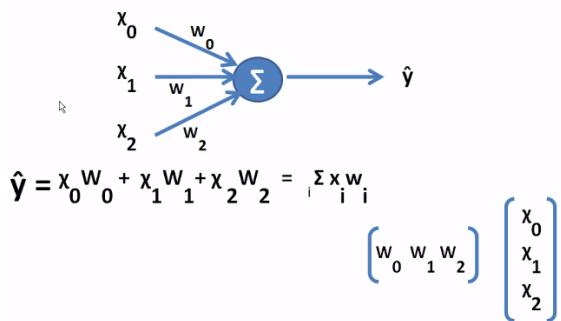
Activation function

Tensor flow playground - To see how neural network work

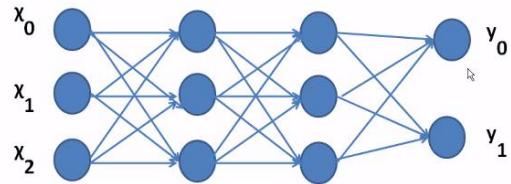
To see how activation function work



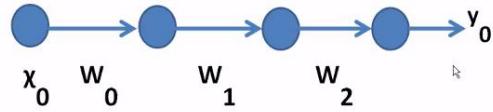
Neural Network with Single Neuron



Muti-layer Perceptron



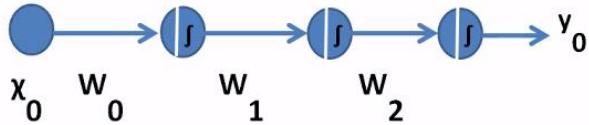
Linear Function



$$y_0 = x_0 w_0 w_1 w_2 = x_0 w_C$$

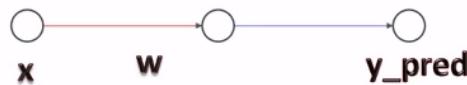
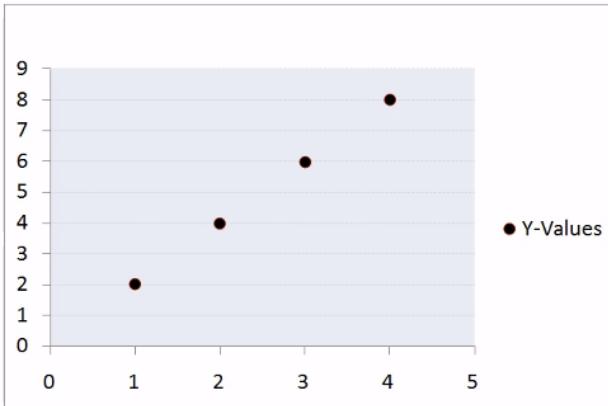


Non-Linear Function



$$y_0 = \sigma(\sigma(\sigma(x_0 w_0) w_1) w_2)$$

Simple Example



$$E = (\hat{y} - y)^2 = (xw - y)^2 \quad \text{sum of squared Error}$$

$$\frac{\partial E}{\partial w} = 2x(xw - y) \quad \text{Derivative}$$

$$w - \alpha \frac{\partial E}{\partial w} = w - \alpha 2x(xw - y) \quad \text{Update Rule}$$

- Random weight initialization $w=0.5$
- learning rating $\alpha = 0.1$
- $w_{new} = w - 0.1 * 2x(xw - y)$
- (x,y)
- $(2,4) w=0.5 \leftarrow 0.5 - 0.1*2*2(2*0.5-4) = 1.7$
- $(1,2) w = 1.7 \leftarrow 1.7 - 0.1*2*1(1*1.7-2) = 1.76$
- $(3,6) w = 1.76 \leftarrow 1.76 - 0.1*2*3(3*1.76-6) = 2.192$
- $w \sim 2$

Weight update Rule:

Weight Update Rule = oldWeight - LR * 2 * x * (x*w-y)								
	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
				0	0	0	0	0
				0	0	0	0	0
				0	0	0	0	0
				0	0	0	0	0
								0
	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
				0	0	0	0	0
				0	0	0	0	0
				0	0	0	0	0
				0	0	0	0	0

SGD:

Epoch1	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
	2	4	0.1	0.5	1.7	3.4	4	-0.6 individual Row Error
	1	2	0.1	1.7	1.76	1.76	2	-0.24 individual Row Error
	3	6	0.1	1.76	2.192	6.576	6	0.576 individual Row Error
	4	8	0.1	2.192	1.5776	6.3104	8	-1.6896 individual Row Error
								-0.4884 Epoch Error
Epoch2	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
	2	4	0.1	1.5776	1.91552	3.83104	4	-0.16896 individual Row Error
	1	2	0.1	1.91552	1.932416	1.932416	2	-0.067584 individual Row Error
	3	6	0.1	1.932416	2.0540672	6.162202	6	0.1622016 individual Row Error
	4	8	0.1	2.054067	1.88105216	7.524209	8	-0.47579136 individual Row Error
								-0.13753344 Epoch Error
Epoch3	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
	2	4	0.1	1.881052	1.976210432	3.952421	4	-0.047579136 individual Row Error
	1	2	0.1	1.97621	1.980968346	1.980968	2	-0.019031654 individual Row Error
	3	6	0.1	1.980968	2.015225324	6.045676	6	0.045675971 individual Row Error
	4	8	0.1	2.015225	1.966504288	7.866017	8	-0.133982847 individual Row Error
								-0.038729417 Epoch Error
Epoch4	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
	2	4	0.1	1.966504	1.993300858	3.986602	4	-0.013398285 individual Row Error
	1	2	0.1	1.993301	1.994640686	1.994641	2	-0.005359314 individual Row Error
	3	6	0.1	1.994641	2.004287451	6.012862	6	0.012862353 individual Row Error
	4	8	0.1	2.004287	1.990567608	7.96227	8	-0.03772957 individual Row Error
								-0.010906204 Epoch Error

BGD:

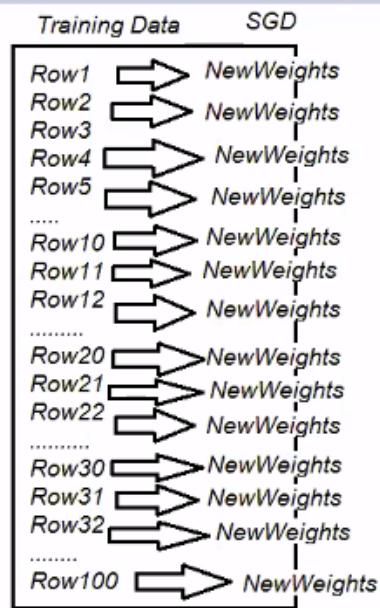
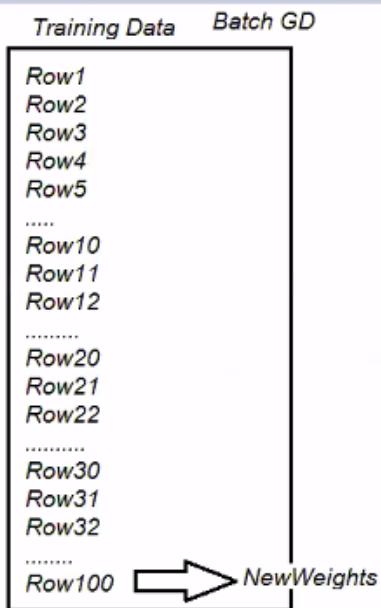
Epoch1	x	y	LR	weights	NewWeights	y_pred	y_test	Error / Loss
	2	4	0.1	0.5	1.7	3.4	4	-0.6
	1	2	0.1	0.5	0.8	0.8	2	-1.2
	3	6	0.1	0.5	3.2	9.6	6	3.6
	4	8	0.1	0.5	5.3	21.2	8	13.2
				2.75	NewWeights			3.75

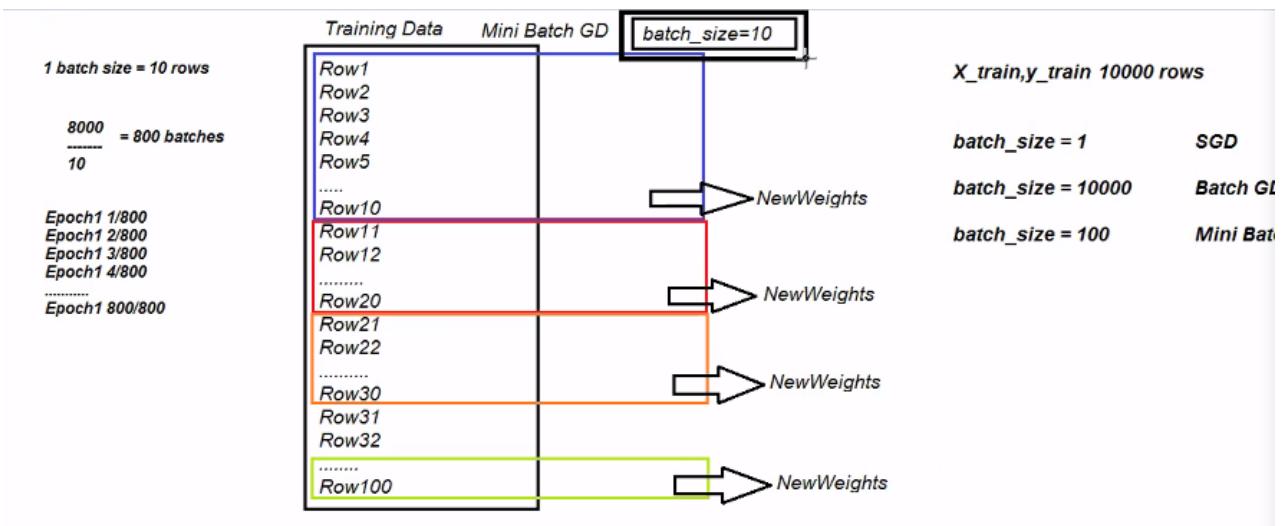
Epoch2	x	y	LR	weights	NewWeights	y_pred	y_test	Error	
	2	4	0.1	2.75	2.15	4.3	4	0.3	individual Row Error
	1	2	0.1	2.75	2.6	2.6	2	0.6	individual Row Error
	3	6	0.1	2.75	1.4	4.2	6	-1.8	individual Row Error
	4	8	0.1	2.75	0.35	1.4	8	-6.6	individual Row Error
					1.625			-1.875	Epoch Error

Epoch3	x	y	LR	weights	NewWeights	y_pred	y_test	Error	
	2	4	0.1	1.625	1.925	3.85	4	-0.15	individual Row Error
	1	2	0.1	1.625	1.7	1.7	2	-0.3	individual Row Error
	3	6	0.1	1.625	2.3	6.9	6	0.9	individual Row Error
	4	8	0.1	1.625	2.825	11.3	8	3.3	individual Row Error
					2.1875			0.9375	Epoch Error

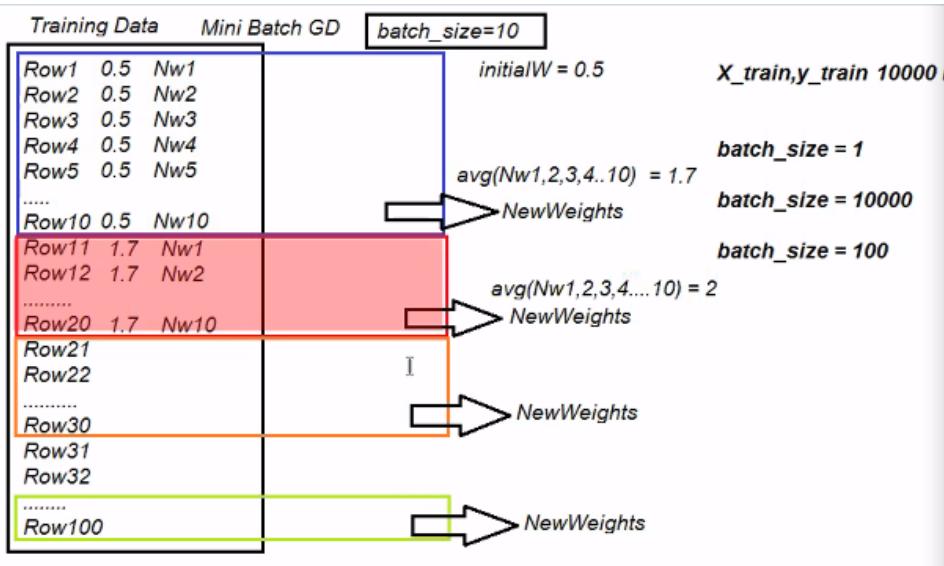
Epoch4	x	y	LR	weights	NewWeights	y_pred	y_test	Error	
	2	4	0.1		1.6	3.2	4	-0.8	
	1	2	0.1		0.4	0.4	2	-1.6	
	3	6	0.1		3.6	10.8	6	4.8	
	4	8	0.1		6.4	25.6	8	17.6	
					3			5	
									Epoch Error

Epoch5	x	y	LR	weights	NewWeights	y_pred	y_test	Error	
	2	4	0.1	1.90625	1.98125	3.9625	4	-0.0375	individual Row Error
	1	2	0.1	1.90625	1.925	1.925	2	-0.075	individual Row Error
	3	6	0.1	1.90625	2.075	6.225	6	0.225	individual Row Error
	4	8	0.1	1.90625	2.20625	8.825	8	0.825	individual Row Error
					2.046875			0.234375	Epoch Error



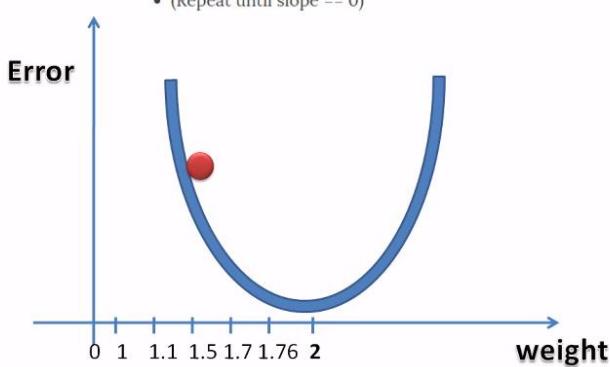


Mini batch is good approach



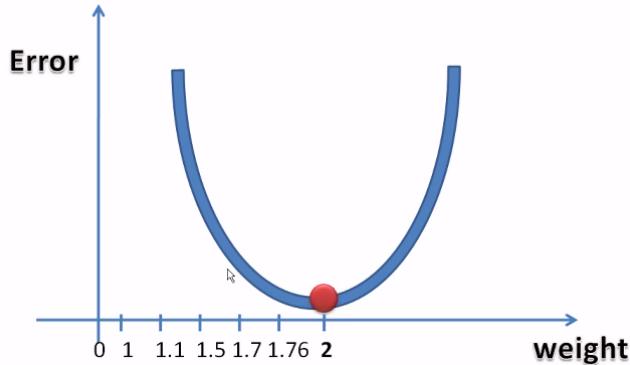
Oversimplified Gradient Descent:

- Calculate slope at current position
- If slope is negative, move right
- If slope is positive, move left
- (Repeat until slope == 0)



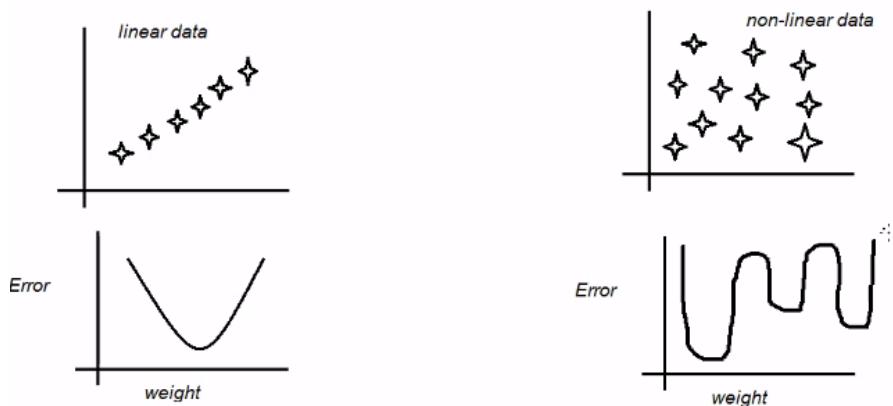
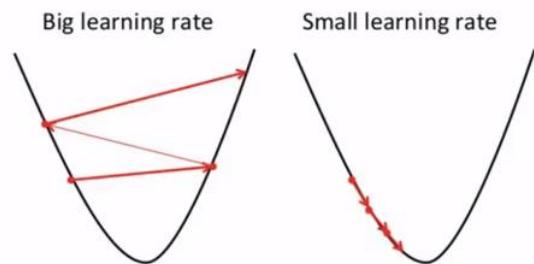
Oversimplified Gradient Descent:

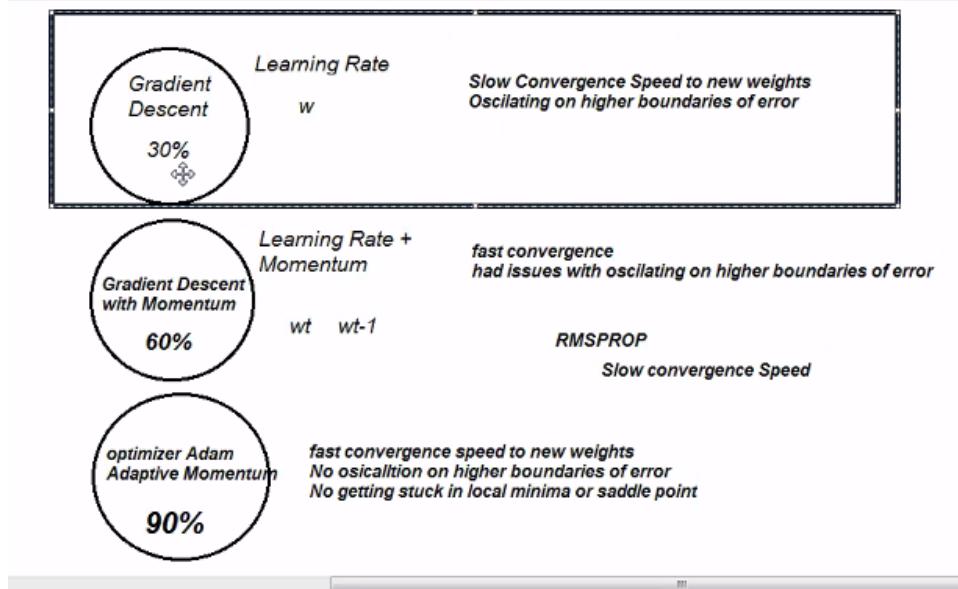
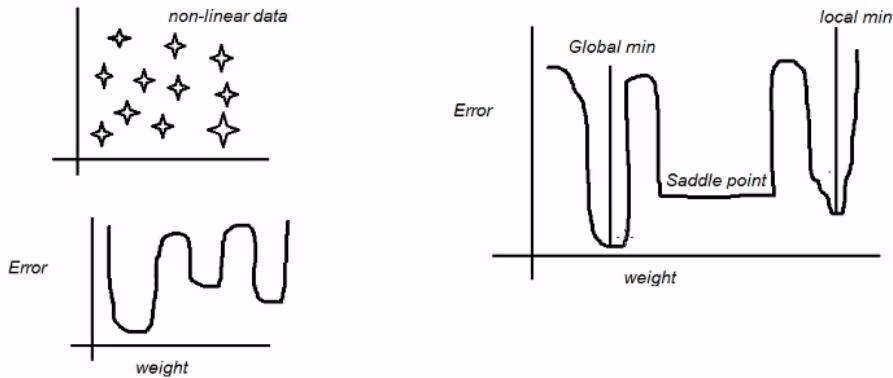
- Calculate slope at current position
- If slope is negative, move right
- If slope is positive, move left
- (Repeat until slope == 0)



Learning rate

1. if it is too small, then the model will take some time to learn.
2. if it is too large, model will converge as our pointer will shoot and we'll not be able to get to minima.

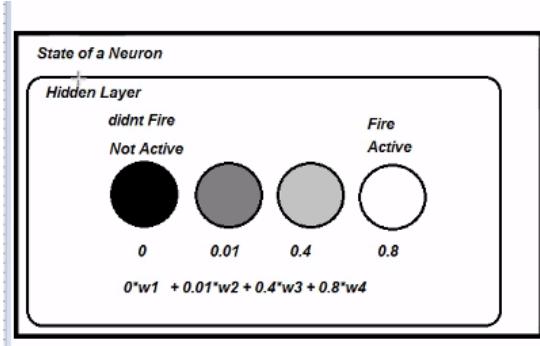




Adaptive Momentum

- Adam method is used to accelerate the gradient descent algorithm by taking into consideration the exponentially weighted average of the gradients.
- **new weight \leftarrow (old weight) - (learning rate)(gradient)**
- **new weight \leftarrow (old weight) - (learning rate)(gradient) + past gradient**
- **(accumulator) \leftarrow (old accumulator)(momentum) + gradient**
- **momentum \rightarrow weighted average of past gradients**
- **new weight \leftarrow (old weight) - (learning rate)(accumulator)**

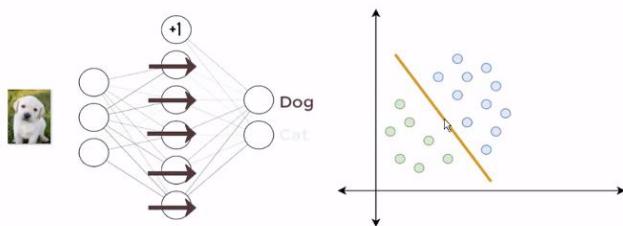
Not fire means not contribute



Activation Functions:

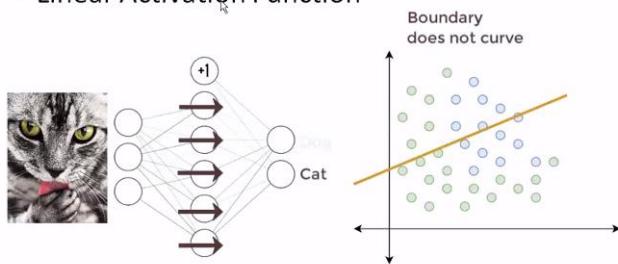
Linearly Separable Data

- Linear Activation Function

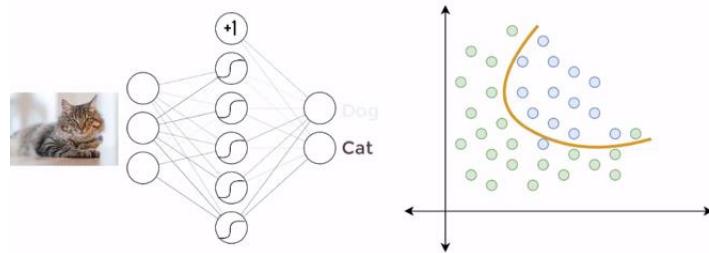


Non-Linearly Separable Data

- Linear Activation Function



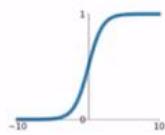
Non Linear Activation function



Activation Functions

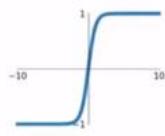
Sigmoid

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



tanh

$$\tanh(x)$$



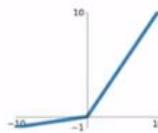
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



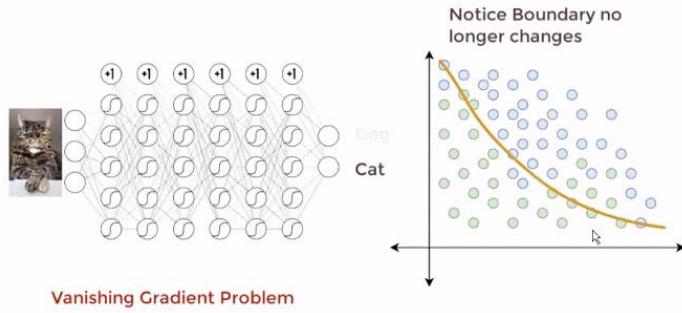
Sigmoid – squeezing in nature – Range – 0 to 1 - Put in problem in case of complex situation

Tanh –Range: -1 to 1 - squeezing in nature - Put in problem in case of complex situation

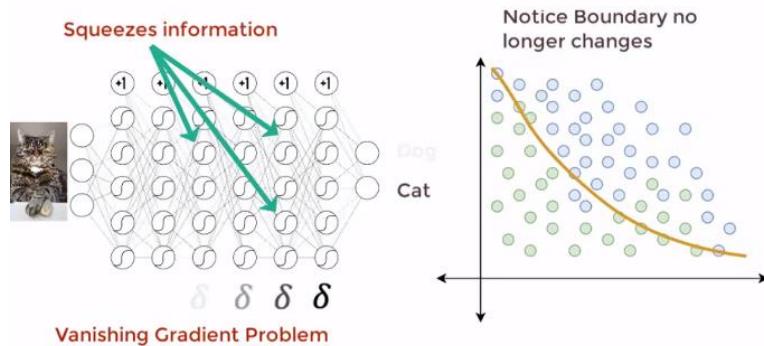
Vanishing gradient problem due to squeezing nature

ReLU - (0, x)

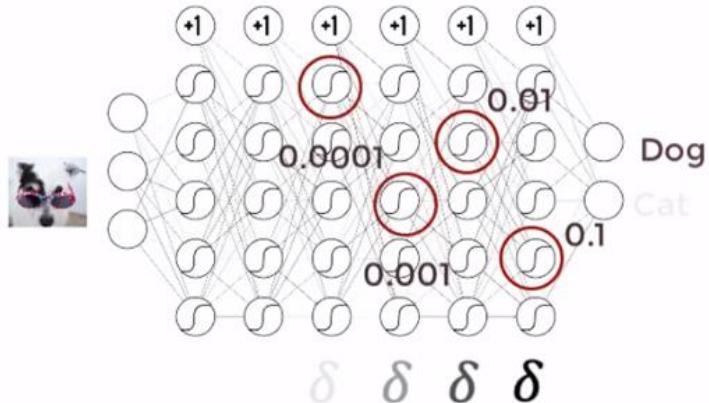
Even if the model goes through many examples it doesn't learn much



During backpropagation the gradients become smaller and smaller until they eventually vanish no gradients means no learning

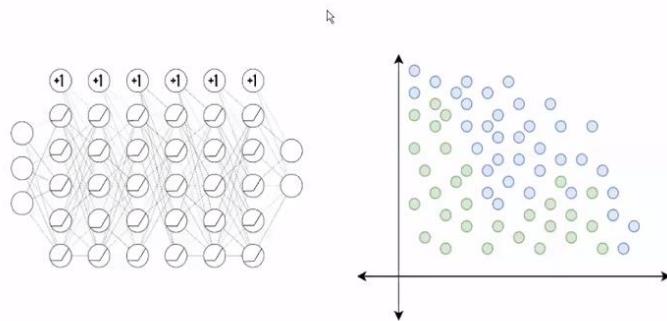


Vanishing Gradient

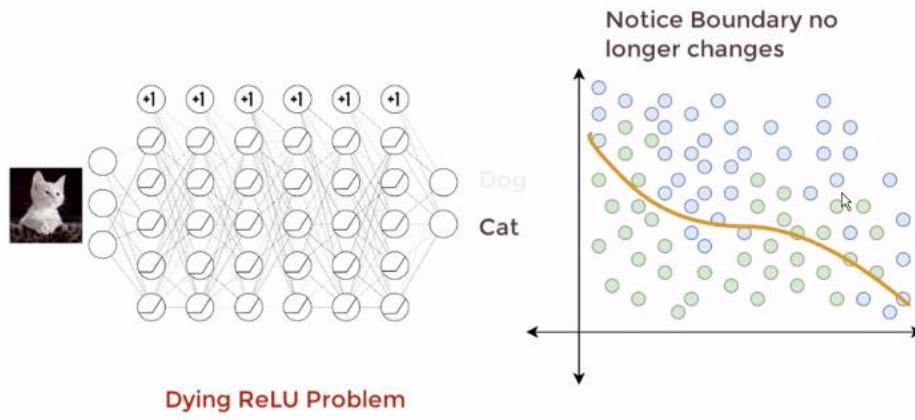


Vanishing Gradient Problem

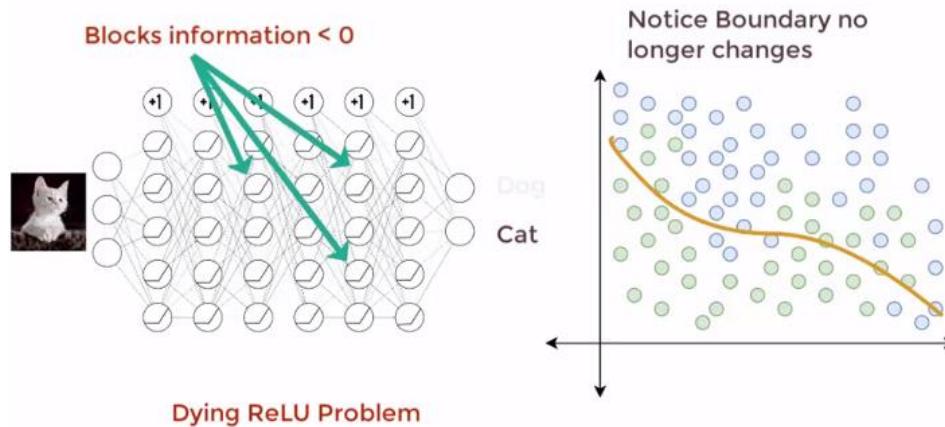
Remedy to vanishing gradient is to use activation function that doesnot squeeze values, like ReLu



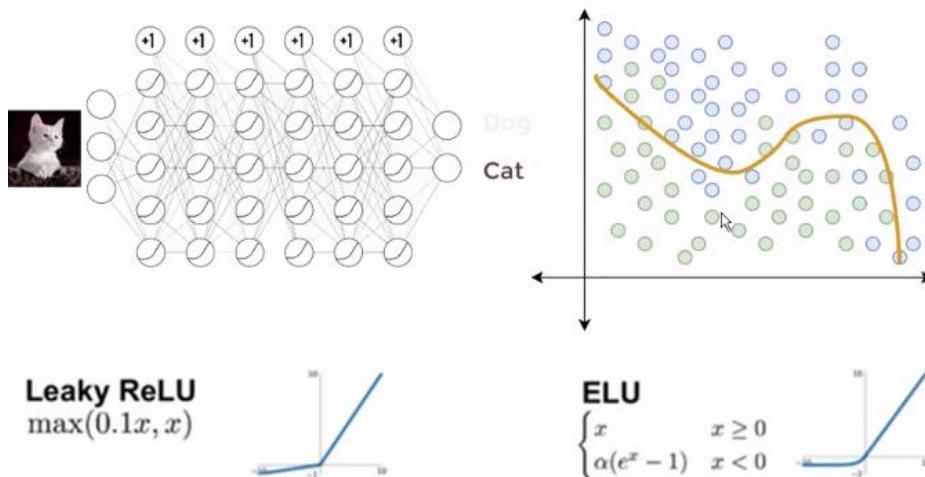
During backpropogation the gradients become smaller and smaller and smaller until they eventually vanish no gradients means no learning. Same problem like vanishing gradient



Dying ReLU causes as it blocks information less than 0



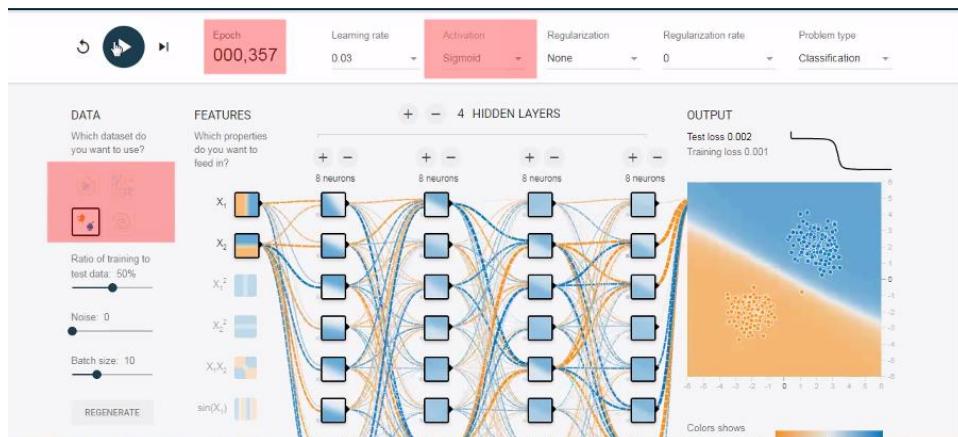
Problem is solved because of Leaky Relu or ELu



Hyper Parameter Tuning (values to be found with trial and error)

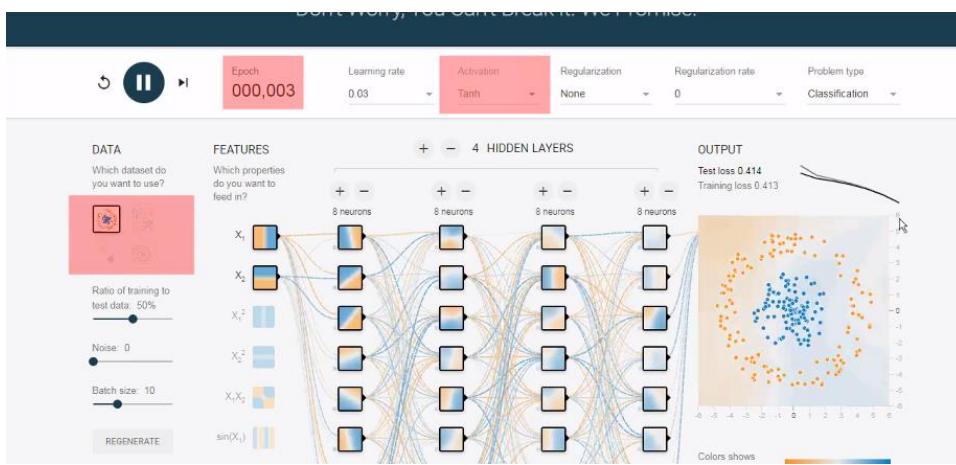
- 1) Num of Neurons in Hidden Layer
- 2) Num of Hidden Layers
- 3) Num of Epochs
- 4) Activation Functions
- 5) batch_size

Explain activation function in tensorflow playground:



Took too much time to find the boundary

Tanh – advantage – bigger boundaries

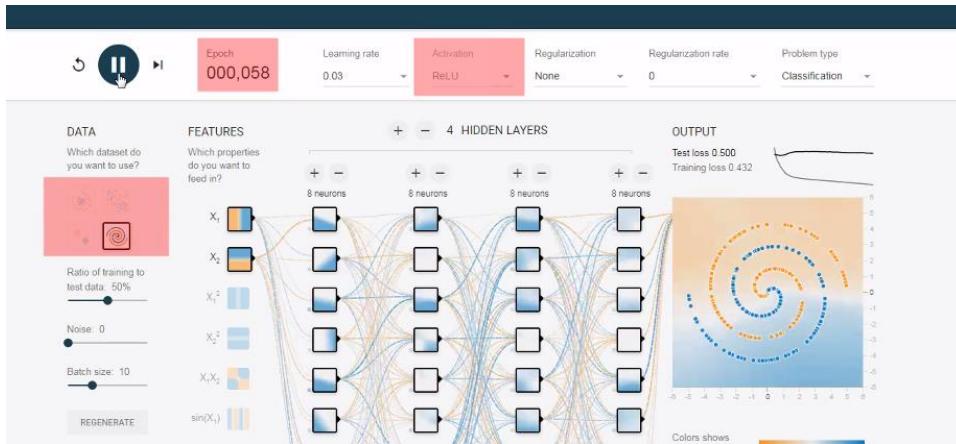


Problem of tanh:



ReLU: Identify the boundary in few numbers of epoch





Guidelines to choose the activation function:

Hyper Parameter Tuning (values to be found with trial and error)

- 1) Num of Neurons in Hidden Layer
- 2) Num of Hidden Layers
- 3) Num of Epochs
- 4) Activation Functions
- 5) batch_size

	InputLayer	HL	HL	OutputLayer	LossFunction
Regression		ReLU	ReLU	ReLU	mean_squared_error
Binary-Classification		ReLU	ReLU	Sigmoid	binary_classification
Multi-Class-Classification		ReLU	ReLU	Softmax	categorical_crossentropy

Implementation:

Step 1: Start anaconda prompt

Step 2:

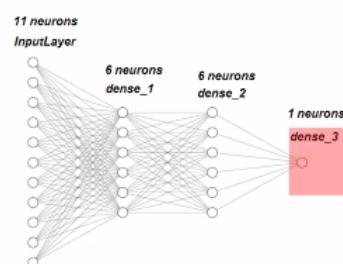
```
python -m pip install --upgrade pip
```

```
pip install keras tensorflow
```

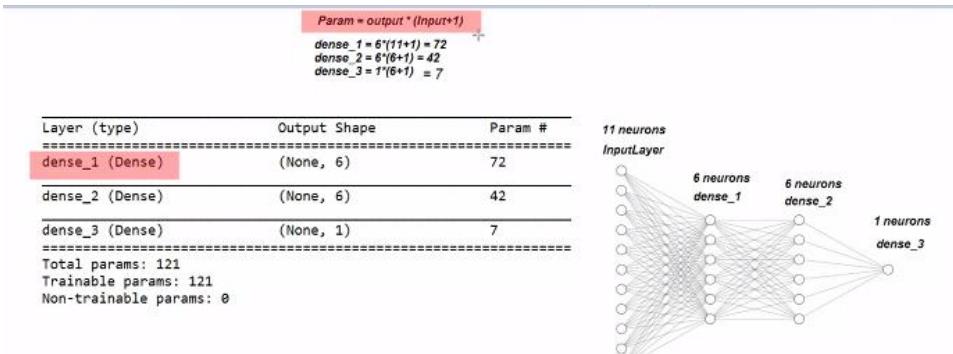
ANN:

Summary:

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 6)	72
dense_2 (Dense)	(None, 6)	42
dense_3 (Dense)	(None, 1)	7
Total params: 121		
Trainable params: 121		
Non-trainable params: 0		



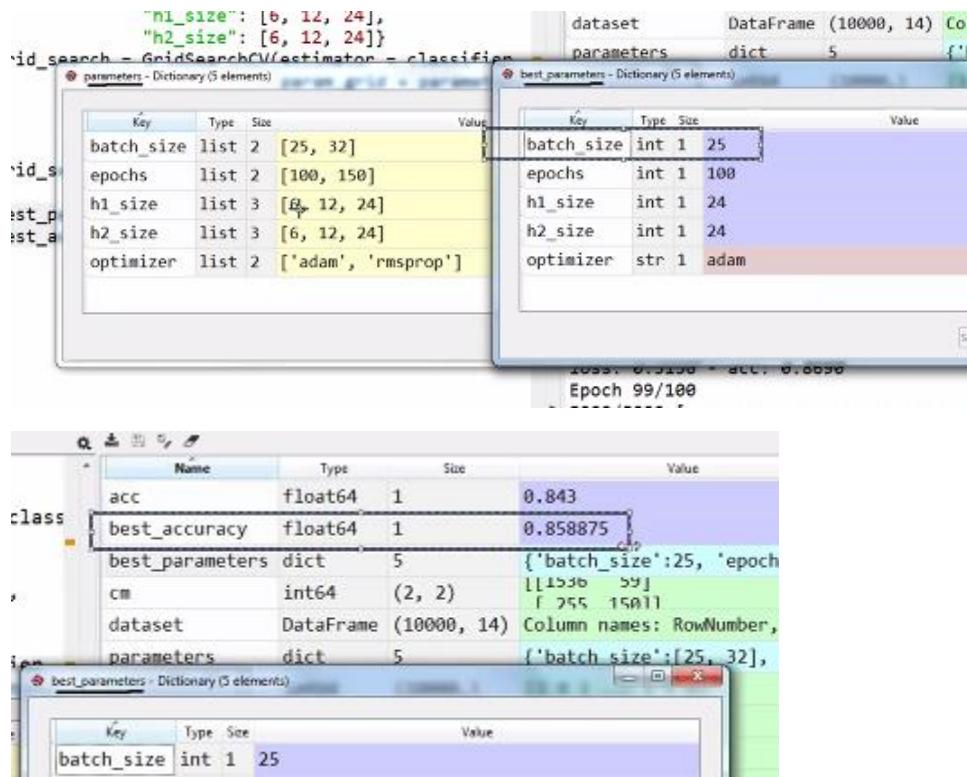
How to get the param count?



How many times I need to run?

Accuracy not stable – For that use GridSearchCV(sklearn Obj) - Trial and error with k fold cross validation

Neural network Keras object -----typecast -----sklearn obj



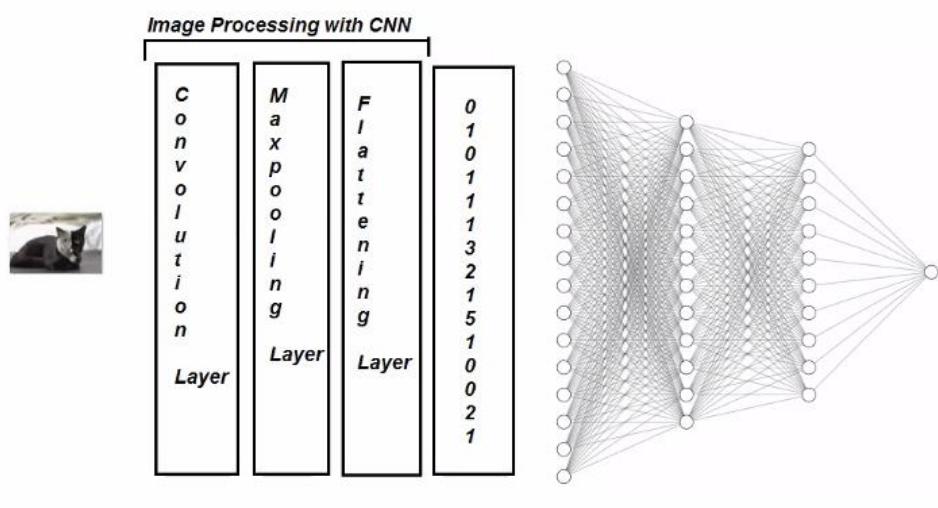
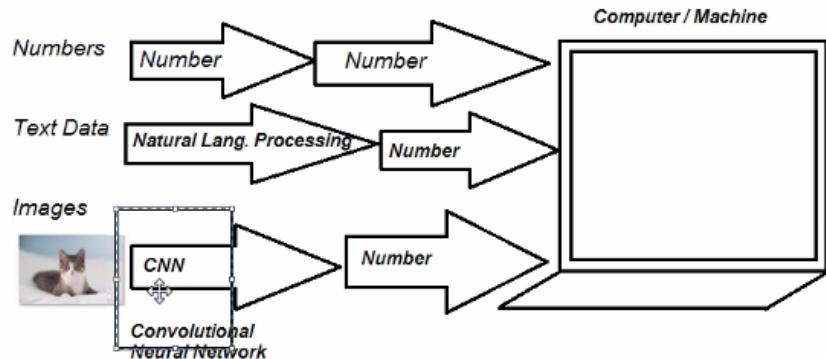
SGD:

Batch size 1 -> Slow learning -> Overfitting

BGD:

Batch size 8000 -> Fast learning -> Accuracy didn't increase – stays at 79% -> Too much of generalization

CNN:



X_{train} Input Data	y_{train} Output Data	X_{test}	y_{pred}
	Dog		????
	Cat		????
	Dog		????
	Cat		????
	Cat		

Convolution Neural Network

- Facebook uses neural nets for their automatic tagging algorithms
- Google for their photo search
- Amazon for their product recommendations
- Pinterest for their home feed personalization
- Instagram for their search infrastructure.



How Humans identify images??

- For humans, this task of recognition is one of the first skills we learn from the moment we are born and is one that comes naturally and effortlessly as adults.
- Most of the time we are able to immediately characterize the scene and give each object a label, all without even consciously noticing.

**What's
that?**

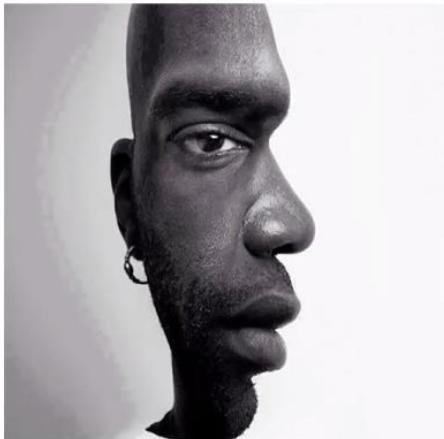


A Train

**What's
this?**



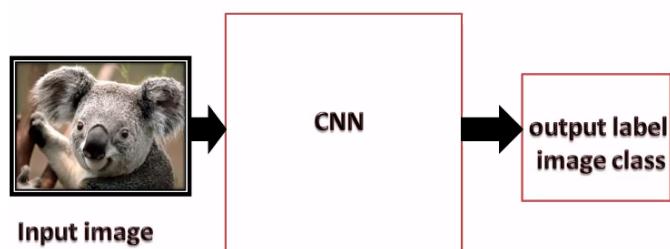
illusions



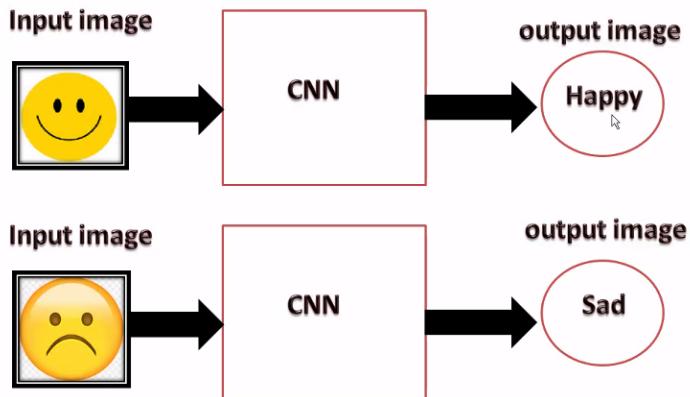
illusions



Convolution Neural Network



Convolution Neural Network



How Does Computer See an image??

When a computer sees an image ,it will see an array of pixel values. Depending on the resolution and size of the image, it will see a $32 \times 32 \times 3$ array of numbers (The 3 refers to RGB values).



What we see

What computers see

How Does Computer See an image??



What we see

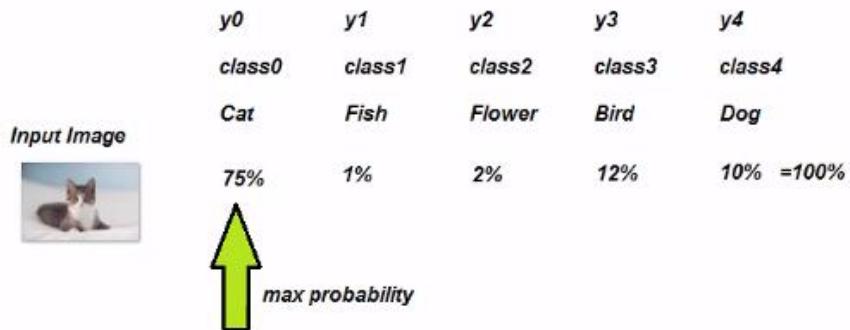
What computers see

```
06 02 22 97 98 13 09 00 97 94 05 07 01 92 12 30 77 91 08  
44 49 31 73 05 79 14 09 03 73 07 47 53 04 02 03 49 13 54 45  
81 49 31 73 05 79 14 09 03 73 07 47 53 04 02 03 49 13 54 45  
02 31 14 71 51 87 40 09 01 92 36 54 22 40 01 23 44 33 13 80  
24 47 32 49 99 03 43 02 44 79 33 53 78 04 02 22 83 37 17 12 80  
32 34 02 49 99 03 43 02 44 79 33 53 78 04 02 22 83 37 17 12 80  
47 24 20 69 02 42 12 20 05 03 99 39 42 01 02 01 92 44 49 94 21  
24 34 20 69 02 42 12 20 05 03 99 39 42 01 02 01 92 44 49 94 21  
32 34 20 69 05 76 04 21 05 03 14 06 03 02 03 01 31 33 99  
78 17 50 29 22 79 31 07 13 04 03 02 04 02 14 09 53 56 92  
24 34 20 69 05 76 04 21 05 03 14 06 03 02 03 01 31 33 99  
96 54 02 49 99 03 43 02 44 79 33 53 78 04 02 22 83 37 17 12 80  
19 80 05 69 05 99 47 04 21 73 92 13 04 02 21 77 04 09 59 56 40  
02 32 34 49 99 03 43 02 44 79 33 53 78 04 02 22 83 37 17 12 80  
88 56 49 07 37 42 20 72 03 44 33 07 47 46 50 52 32 42 93 53 49  
04 49 31 73 05 79 14 09 03 73 07 47 53 04 02 03 49 13 54 45  
20 49 34 43 72 30 23 88 34 42 39 49 32 05 03 65 74 04 34 34  
20 73 30 29 76 31 99 01 21 32 49 71 48 84 83 14 23 57 26 54  
01 70 54 75 03 51 54 09 14 32 33 42 43 48 32 02 69 19 47 49
```

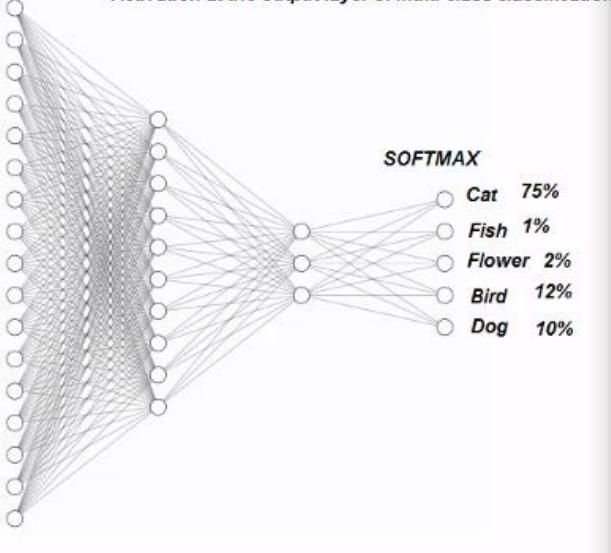
These numbers, when we perform image classification, are the only inputs available to the computer.

The idea is that you give the computer this array of numbers and it will output numbers that describe the probability of the image being a certain class (80% for flower, 15% for sky etc).

Multi-Class Classification



Activation at the output layer of multi-class classification



X_train



	<i>Cat</i>	<i>Fish</i>	<i>Flower</i>	<i>Bird</i>	<i>Dog</i>
	<i>y_train</i>				
	y_0	y_1	y_2	y_3	y_4
	1	0	0	0	0
	0	0	0	0	1
	1	0	0	0	0
	0	0	0	0	1



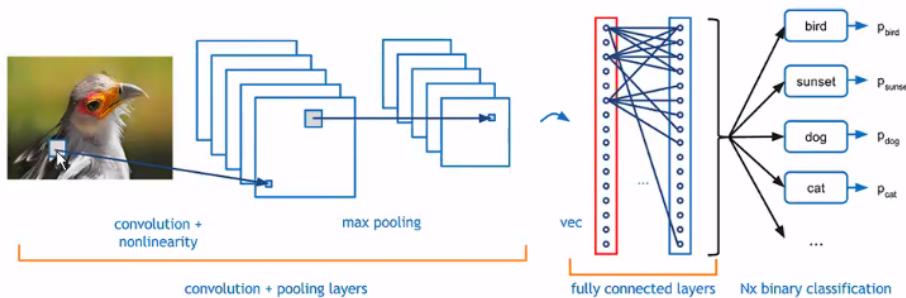


What We Want the Computer to Do

- we want the computer to do is to be able to differentiate between all the images it's given and figure out the **unique features** that make a dog a dog , that make a cat a cat or that make a flower a flower .

Convolution Neural Network

CNNs take the image, pass it through a series of convolutional, nonlinear, pooling (downsampling), and fully connected layers, and get an output.



How computer sees an image

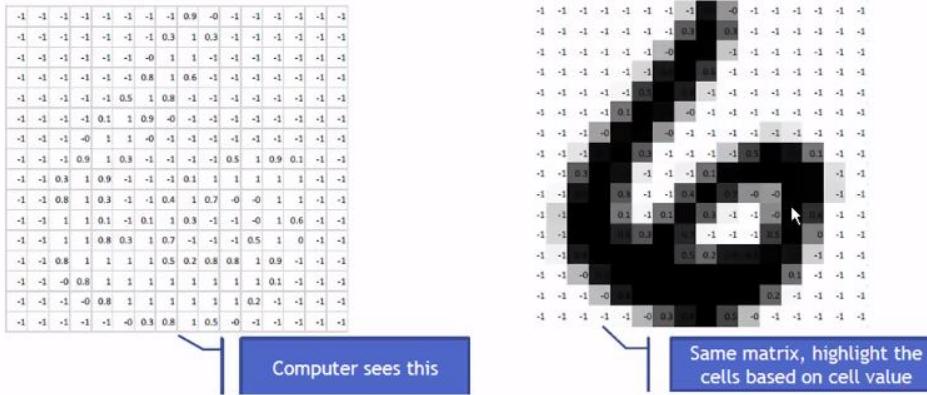


Humans see this

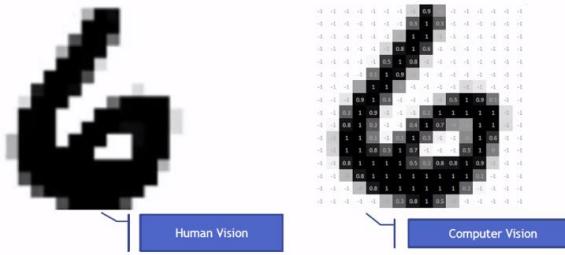
-1	-1	-1	-1	-1	-1	-1	0.9	-0	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	0.3	1	0.3	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-0	1	1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	0.8	1	0.6	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	0.5	1	0.8	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	0.3	1	0.9	-0	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-0	1	1	-0	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	0.9	1	0.3	-1	-1	-1	0.5	1	0.9	0.1	-1	-1
-1	-1	0.3	1	0.9	-1	-1	0.1	1	1	1	1	1	-1	-1
-1	-1	0.8	1	0.3	-1	-1	0.4	1	0.7	-0	-0	1	1	-1
-1	-1	1	1	0.1	-1	0.1	1	0.3	-1	-1	-0	1	0.6	-1
-1	-1	1	1	0.8	0.3	1	0.7	-1	-1	0.5	1	0	-1	-1
-1	-1	0.8	1	1	1	1	0.5	0.3	0.8	1	0.9	-1	-1	-1
-1	-1	-0.8	1	1	1	1	1	1	1	1	0.1	-1	-1	-1
-1	-1	-1	-0.8	1	1	1	1	1	1	0.2	-1	-1	-1	-1
-1	-1	-1	-1	-0	0.3	0.8	1	0.5	-0	-1	-1	-1	-1	-1

Computer sees this

How computer sees an image



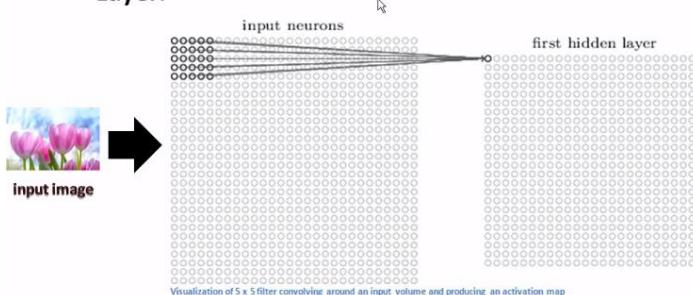
How computer sees an image



Convolutional Neural Network

• First Layer – Math Part

- The first layer in a CNN is always a **Convolutional Layer**.





Input Image



Label

Tree

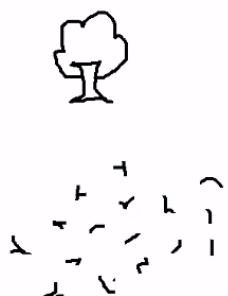


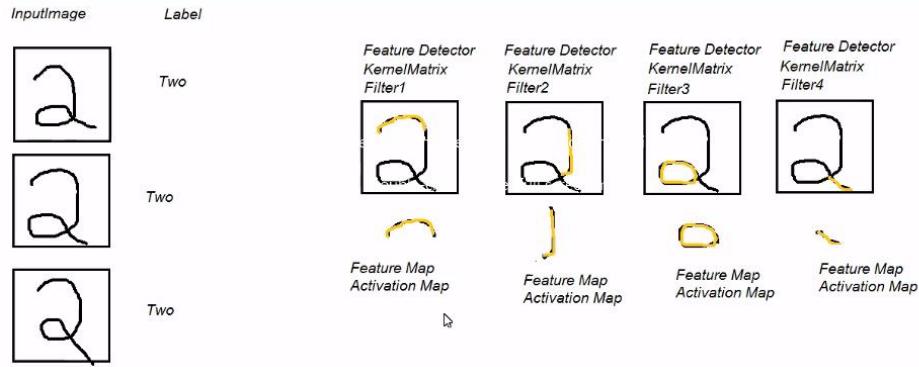
Tree



Tree

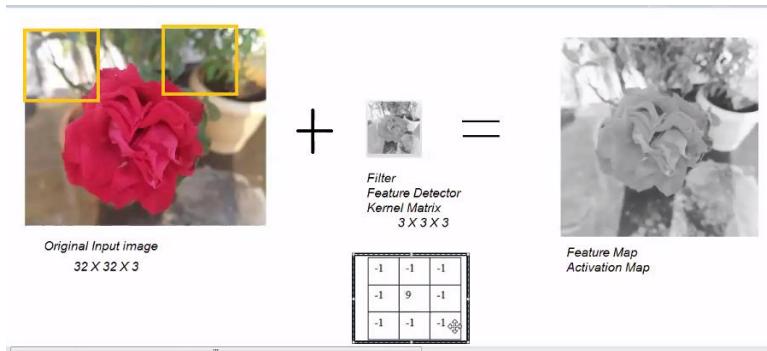
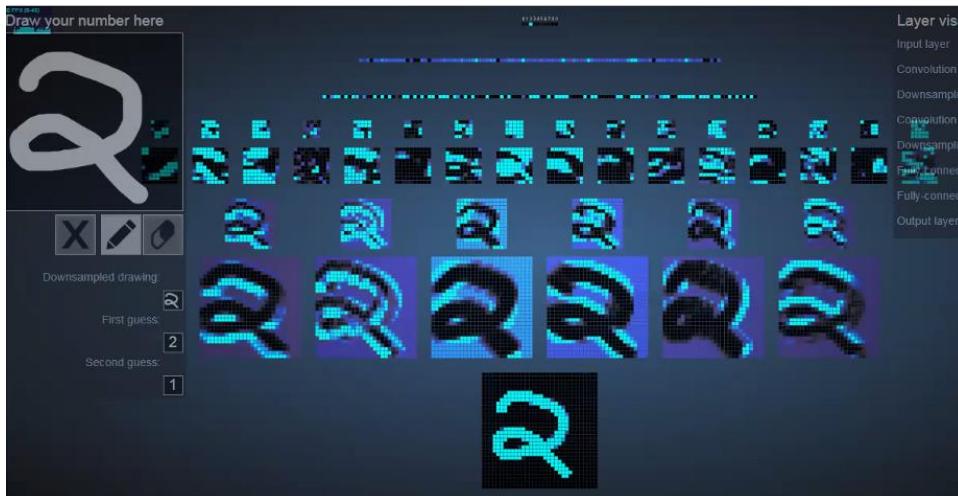
-





http://www.cs.cmu.edu/~aharley/nn_vis/cnn/2d.html

This identify the edges:



Kernel Matrix examples

1	1	1
1	1	1
1	1	1

Unweighted 3x3 smoothing kernel

0	1	0
1	4	1
0	1	0

Weighted 3x3 smoothing kernel with Gaussian blur

0	-1	0
-1	5	-1
0	-1	0

Kernel to make image sharper

-1	-1	-1
-1	9	-1
-1	-1	-1

Intensified sharper image

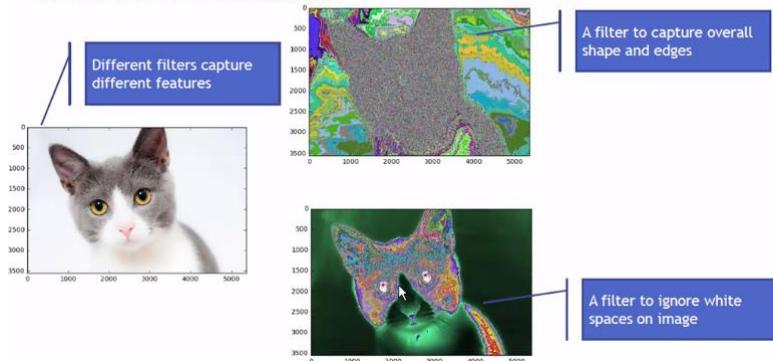


Gaussian Blur



Sharpened image

Convolved features



The depth in convolution layer

- Every filter gives us a resultant matrix (activation map)

0	0	0	0	0	0	0	0	0	0
0	1	1	1	0	0	1	1	0	
0	1	1	0	1	0	1	1	0	
0	1	0	0	1	1	1	0	0	
0	0	0	0	0	0	1	1	0	
0	1	1	0	0	0	0	1	0	
0	1	0	0	1	1	1	0	0	
0	1	1	1	0	0	1	1	0	
0	0	0	0	0	0	0	0	0	0

1	1	1
1	1	1
1	1	1

4	5	4	2	3	4	4
5	6	5	4	6	6	5
3	3	3	3	6	6	5
3	3	2	2	4	5	4
3	3	2	2	4	5	4
5	6	4	3	4	5	4
3	4	3	3	4	4	3

The depth in convolution layer

- If we apply 10 different filters then we will get 10 resultant matrices. The depth is 10

0	0	0	0	0	0	0	0	0	0
0	1	1	1	0	0	1	1	0	
0	1	1	0	1	0	1	1	0	
0	1	0	0	1	1	0	0		
0	0	0	0	0	0	1	1	0	
0	1	1	0	0	0	0	1	0	
0	1	0	0	1	1	1	0	0	
0	1	1	1	0	0	1	1	0	
0	0	0	0	0	0	0	0	0	0

The diagram illustrates the application of multiple convolutional filters to an input feature map. On the left, an input feature map is shown as a 5x10 grid of values. Multiple colored convolutional filters (red, blue, green, orange) are applied across the input. Each filter produces a corresponding output activation map, which is a 3x7 grid of values. These output maps are stacked vertically to form a final feature map with a depth of 10.

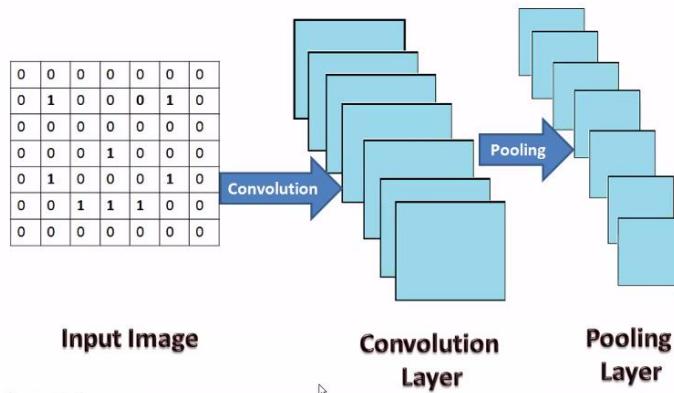
Max Pooling

0	1	0	0	0
0	1	1	1	0
1	0	1	2	1
1	4	2	1	0
0	0	1	2	1

The diagram shows the process of Max Pooling. On the left, a "Feature Map" is shown as a 5x5 grid of values. A red box highlights a 2x2 window in the second row, second column. A large blue arrow labeled "Pooling" points to the right, leading to a "Pool Feature Map" on the far right. This pool feature map is a 3x3 grid where each cell contains the maximum value from its corresponding 2x2 pooling window in the input feature map. The resulting pool feature map values are 1, 4, 0 in the first row; 1, 2, 2 in the second row; and 0, 2, 1 in the third row.

Max pooling is usually done with 2x2 windows and stride 2, so as to downsample the feature maps

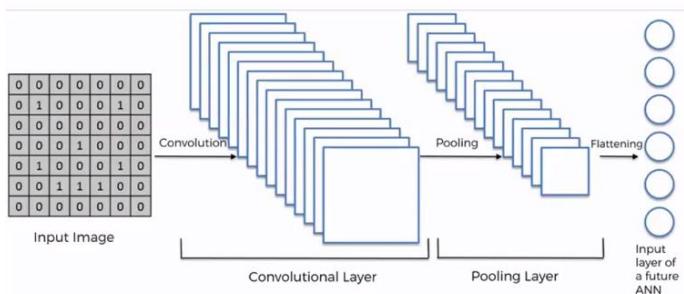
Max Pooling



Flattening



Flattening



Full Connection

