

GLOBAL PROGRAM IN MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE

1

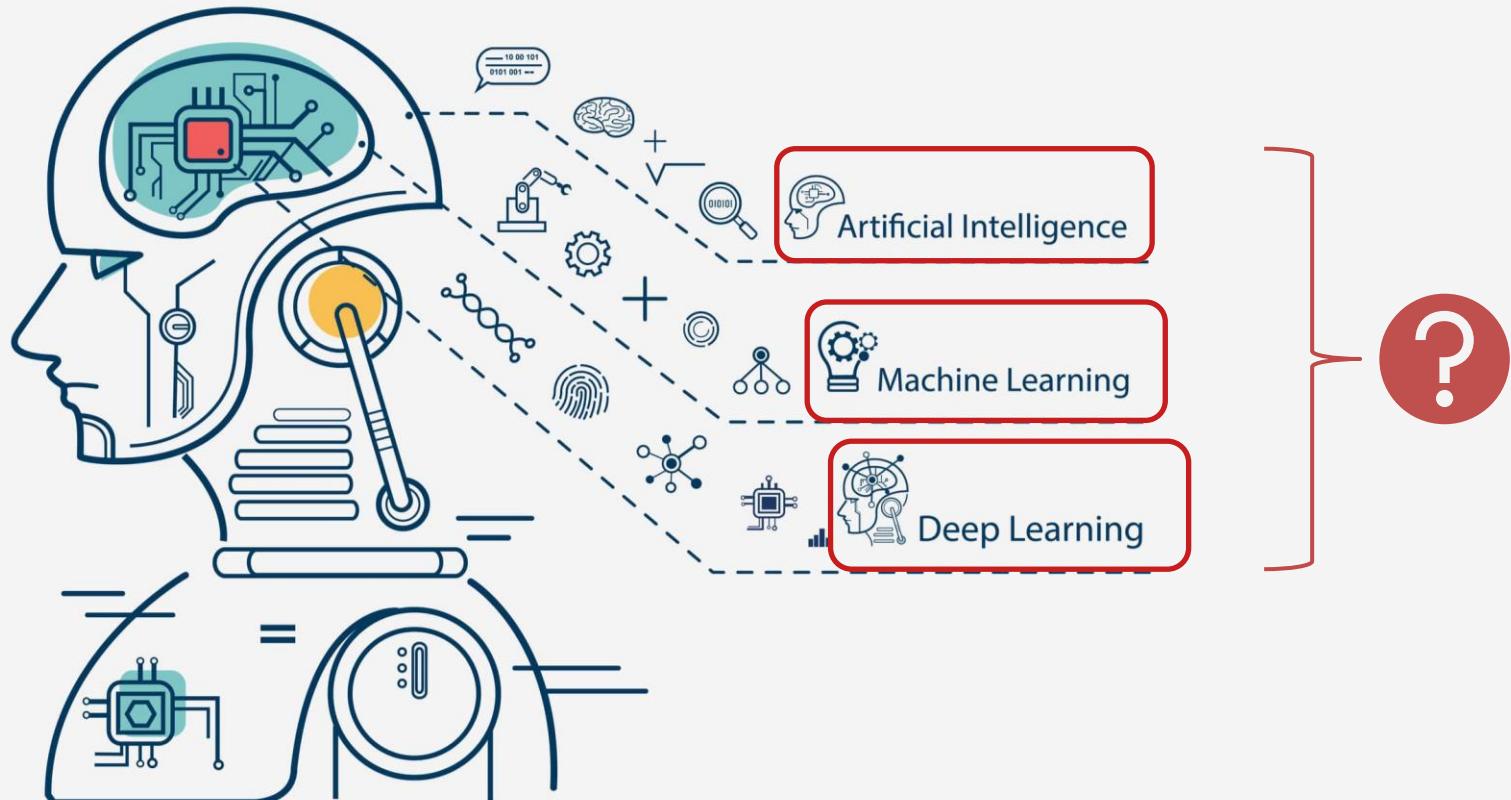
SESSION : INTRODUCTION TO NEURAL NETWORKS

Instructor : SACHIN KUMAR
Master Trainer (UpGrad)



- ❖ Introduction
 - Machine Learning & AI
 - Problems faced by traditional ML algorithms
- ❖ Perceptron – The building block
 - What is a Perceptron?
 - Activation functions
- ❖ Structure of an ANN & Feedforward in Neural Networks
 - How does an ANN look like?
 - Notations used for layers, input, outputs etc
 - Different inputs & outputs to a Neural
 - Feedforward in an Artificial Neural Network
 - Types of Neural Networks

- ❖ Backpropagation in an ANN
 - Training an ANN
 - Cost Function for Regression/Classification
 - Gradient Descent
 - Feedforward and Backprop equations
- ❖ TensorFlow playground
 - Building an ANN
 - Training an ANN
 - Changes in Parameters
- ❖ Doubt resolution



ARTIFICIAL INTELLIGENCE

Techniques that help computer mimic "cognitive" functions associated with the human mind,



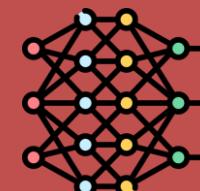
MACHINE LEARNING

Ability to learn pattern from data without being explicitly programmed



DEEP LEARNING

Extract pattern from data using Deep Neural Networks



Building Thinking Machines has been a challenge for inventors !

Easy for computers to do:

- Extensive mathematical computations
- Querying from a large corpus of text
- Combinatorial like Chess playing

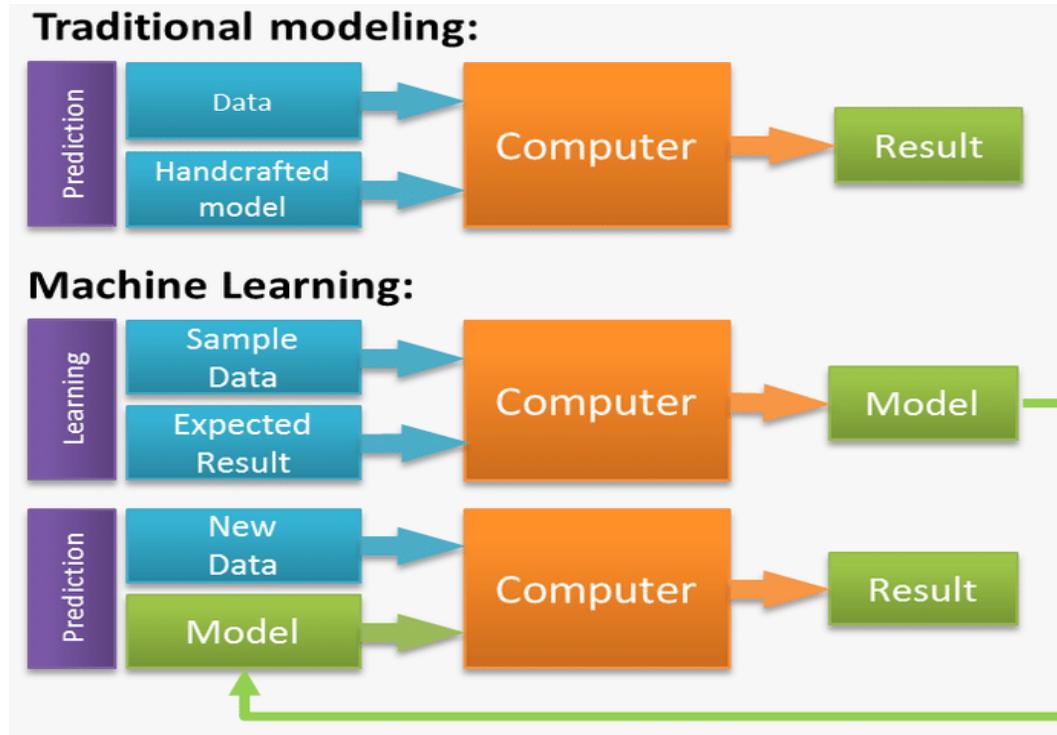
Hard for computers to do:

- Recognize images and faces
- Understand spoken language
- Recognize feelings
- Learning from few mistakes

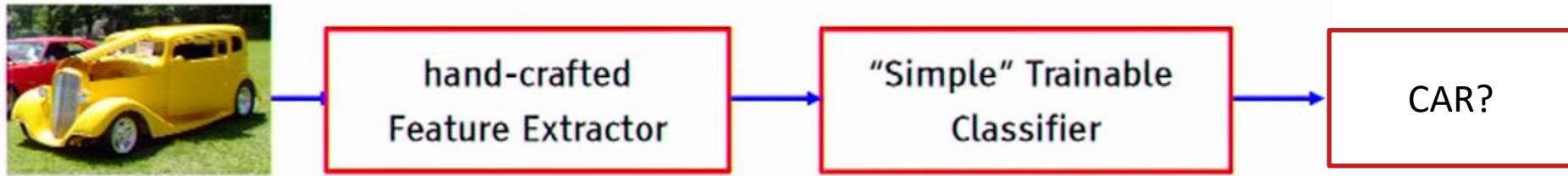


Machine Learning was a Paradigm Shift

upGrad



- ▶ Fixed/engineered features (or fixed kernel) + trainable classifier

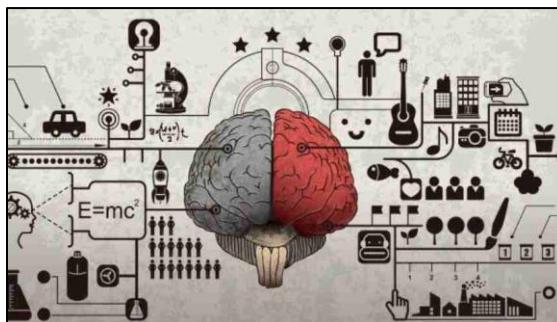


The challenge to Machine Learning proved to be solving the tasks that are easy for people to perform, but hard for people to describe formally.

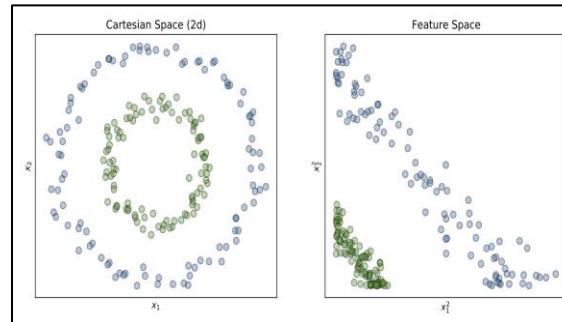
Challenges with Traditional Machine Learning



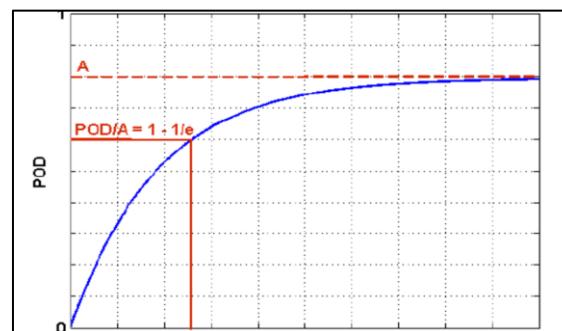
Require domain expertise to handcraft useful features



Cannot solve extremely complex problems



Results heavily depend on Feature Representation



Has an upper limit to Performance

Learning hard-to-describe problem:

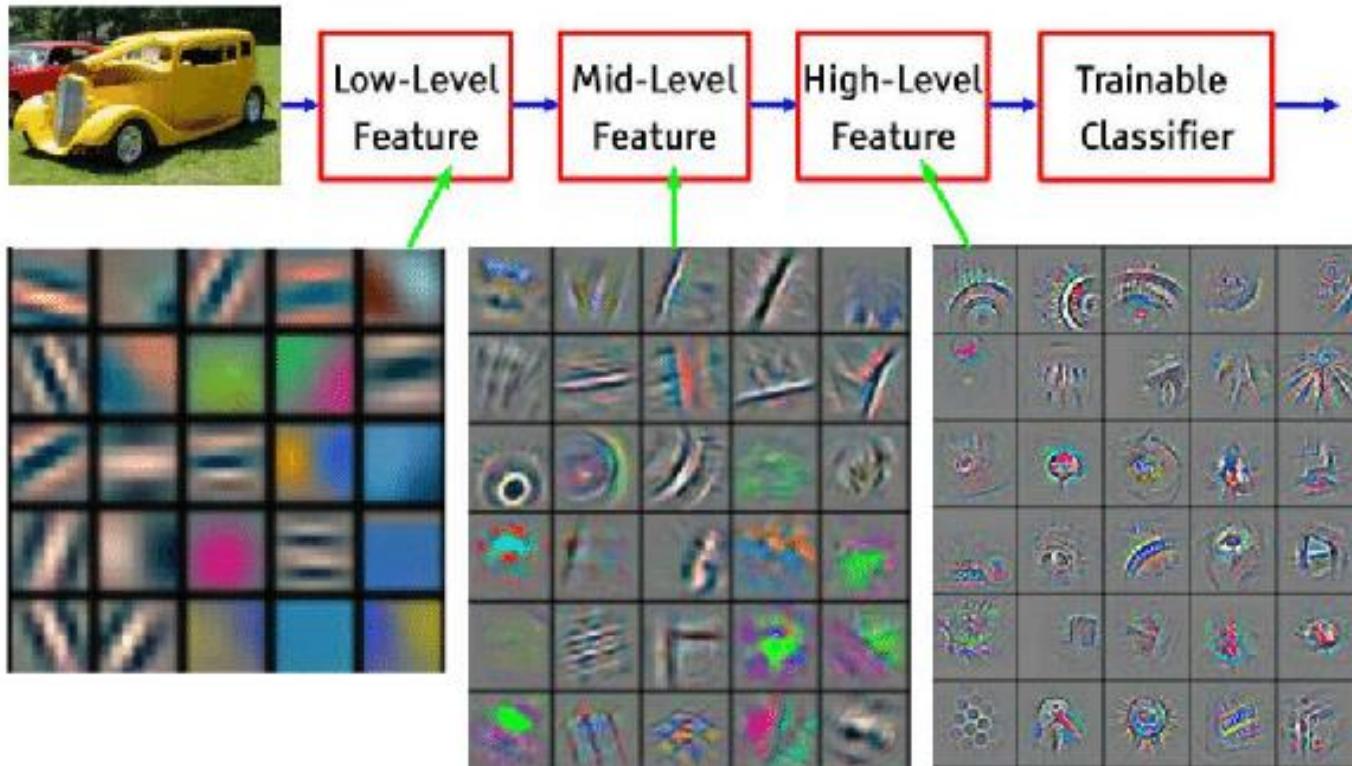
- Allow computers to learn from experience
- Understand the world in terms of a hierarchy of concepts
- Each concept defined through its relation to simpler concepts.

Advantages:

- No need for human operators to formally specify all the knowledge that the computer needs
- The hierarchy of concepts enables the computer to learn complicated concepts by building them out of simpler ones.

Deep Learning aka Hierarchical Feature Learning

upGrad





1950's

- Samuel's checker player
- Selfridge's Pandemonium



1960's

- Neural networks: Perceptron
- Pattern recognition
- Learning in the limit theory
- Minsky and Papert prove limitations of Perceptron



1970's

- Symbolic concept induction
- Winston's arch learner
- Expert systems and the knowledge acquisition bottleneck
- Quinlan's ID3
- Michalski's AQ and soybean diagnosis
- Mathematical discovery with AM



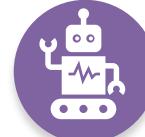
1980's

- Advanced decision tree and rule learning
- Explanation-based Learning (EBL)
- Learning and planning and problem solving
- Utility problem
- Analogy
- Cognitive Architectures
- Resurgence of neural networks (backpropagation)



1990's

- Data mining
- Adaptive software agents and web applications
- Text learning
- Reinforcement learning (RL)
- Ensembles: Bagging, Boosting, and Stacking
- Bayes Net learning



2000's

- Support vector machines & kernel methods
- Graphical models
- Statistical relational learning
- Transfer learning
- Sequence labeling
- Personalized assistants that learn
- Learning in robotics and vision



2010's

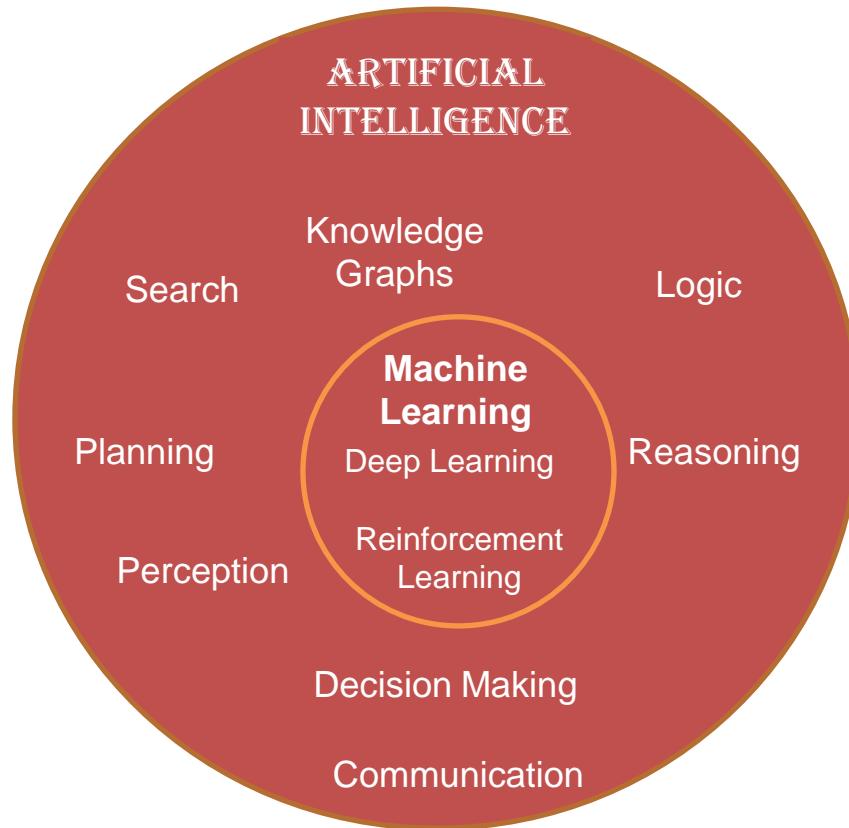
- Deep learning systems
- Learning for big data
- Bayesian methods
- Multi-task & lifelong learning
- Applications to vision, speech, social networks, learning to read



2020's

- Natural Language Understanding
- Generative Models (GANs)
- Transformers
- Large Language Models

Artificial Intelligence is a broad set of Technology



Ability of Machines to mimic "cognitive" functions that we associate with the human mind, such as "learning, seeing, interpreting and problem solving"



Examples of Artificial Intelligence Applications

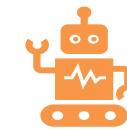
upGrad



Computer
Vision



Speech
Recognition



Personal
Assistant



Medical
Diagnosis



Chatbots



Assisted driving
Cars

Categorization of Artificial Intelligence

Narrow/Weak AI

AI created for a single or narrow set of tasks.

Computer Algorithms that perform tasks it is trained for. Cannot think or comprehend.

Used in industries due to its high efficiency and speed, over humans.

General AI

AI that can be broad and perform multiple tasks.

A type of AI that will think and function like humans

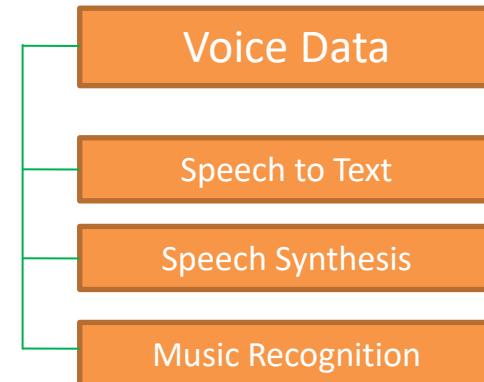
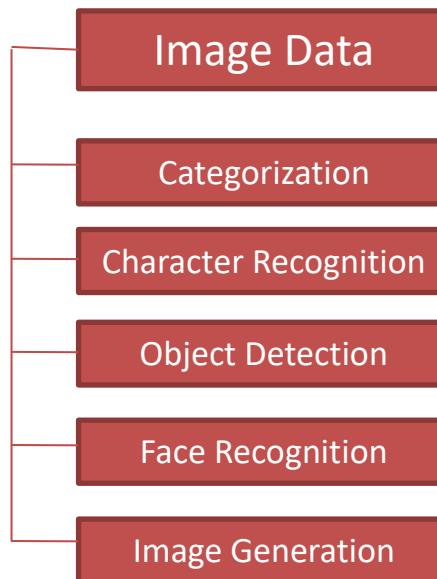
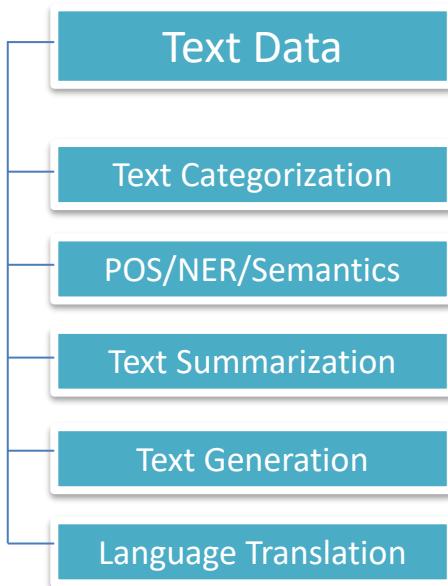
Still a far off reality, as the tools required to build it are not available yet.

Super AI

AI capable of surpassing human intelligence by manifesting cognitive skills and developing thinking skills of its own.

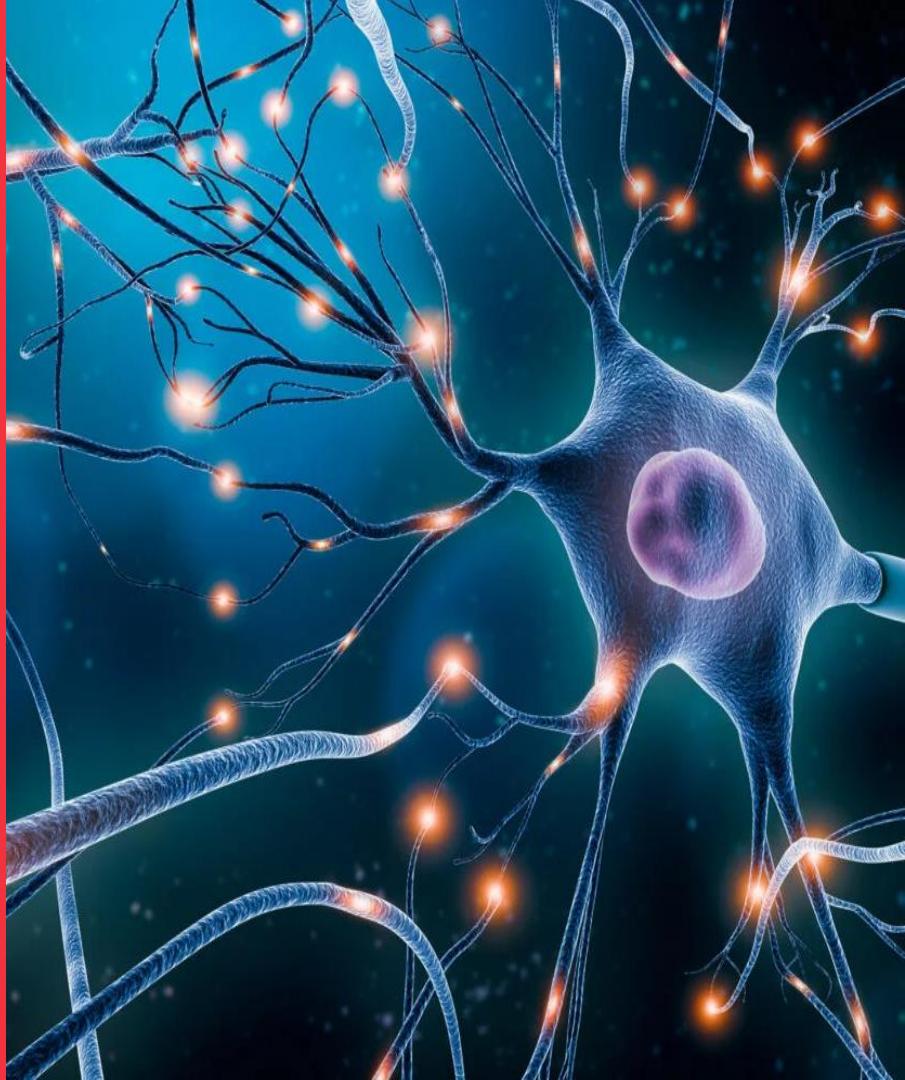
Envisioned to come from compounding growth of AI algorithms and recursive self-improvement.

Applications of Narrow AI



Perceptron and Activation Functions

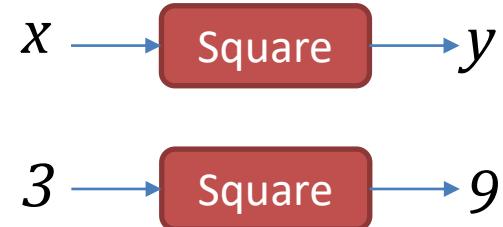
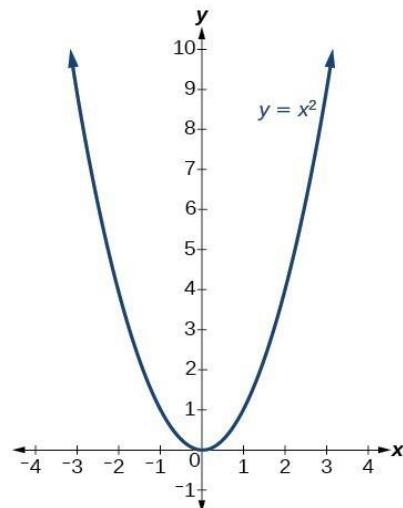
Intro to NN - Sachin Kumar



Representation of a Function

$$y = f(x)$$

$$f(x) = x^2$$

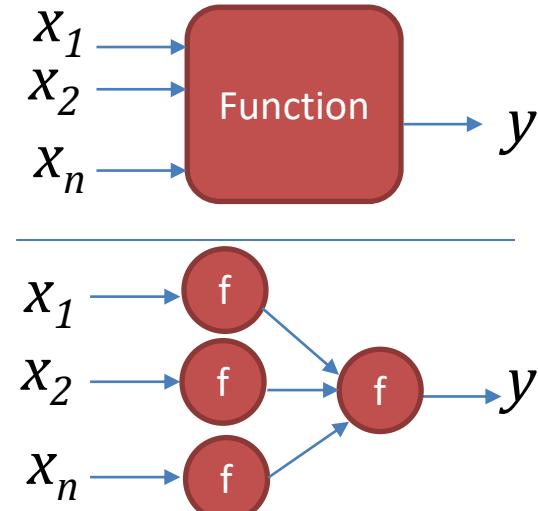
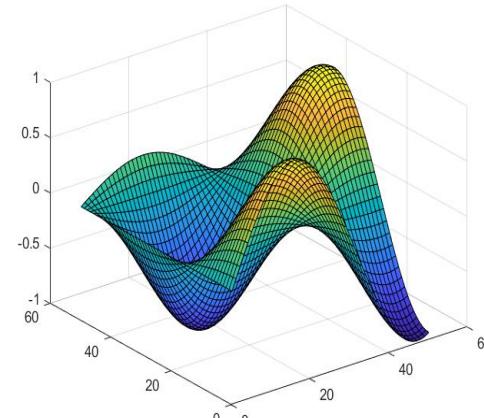


Representation of a Function – in Higher Dimension

upGrad

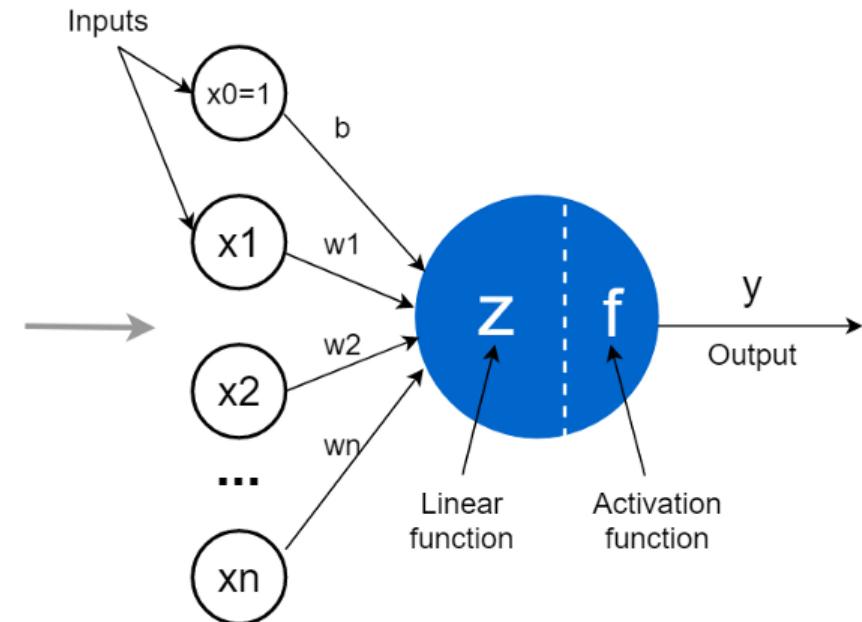
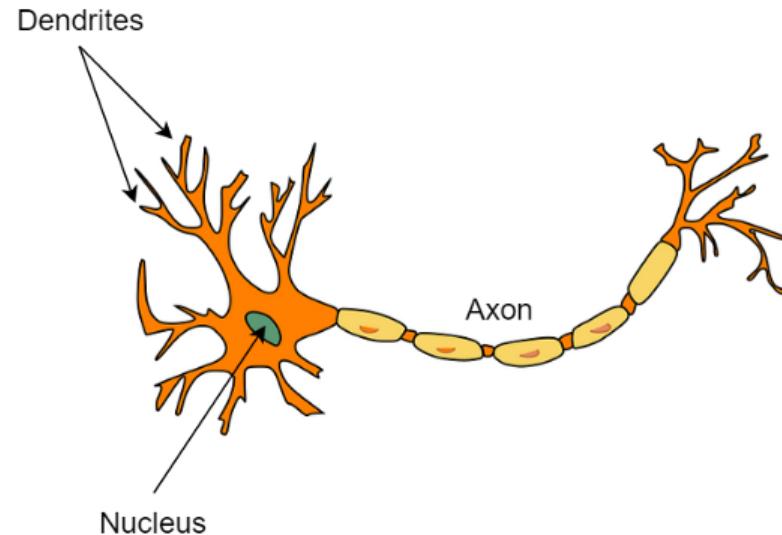
$$Y = f(W, X)$$

$$y_1 = w_1x_1 + w_2x_2 + \dots + w_nx_n$$



Artificial Neural Network - Inspired by Nature

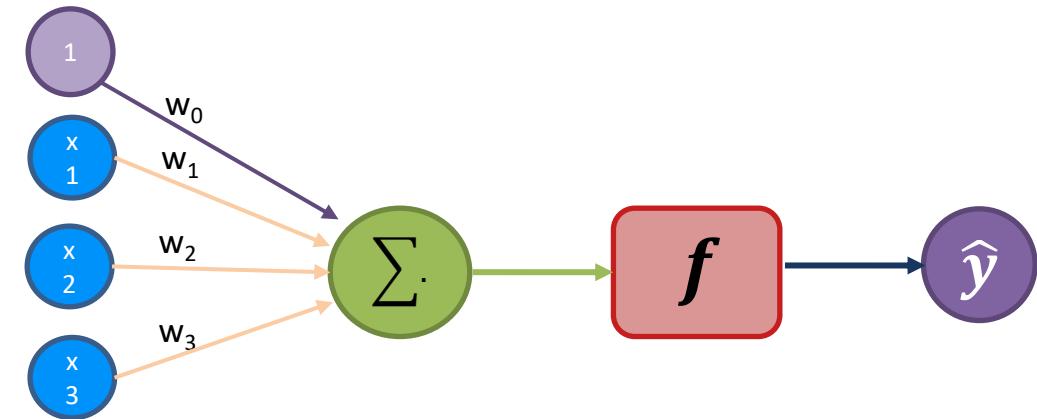
upGrad



Biological Neuron

Artificial Neuron

A Closer Look At The Perceptron



Inputs

Weights

Summation

*Activation
Function*

*Predicted
value*

*Non-linear
Function*

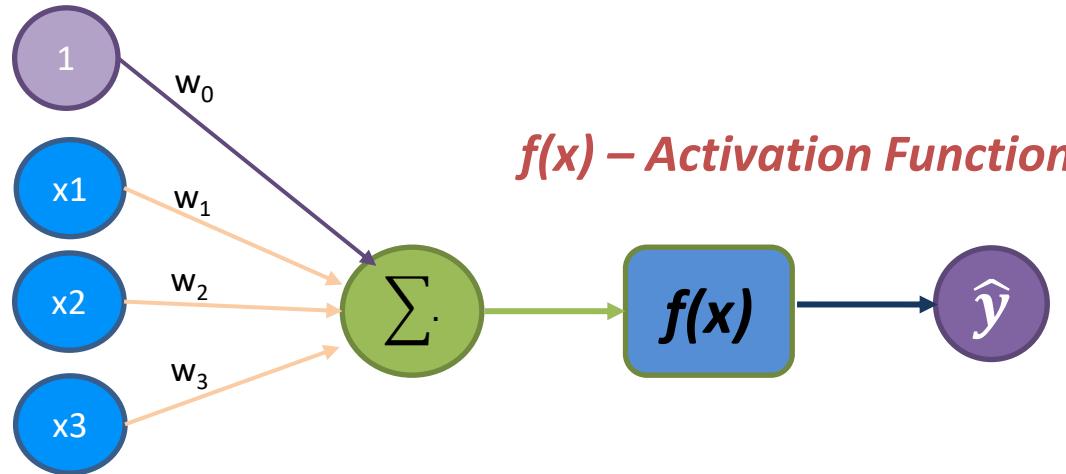
$$\hat{y} = f(\sum x_i \cdot w_i)$$

$$\hat{y} = f(w_0 + \sum x_i \cdot w_i)$$

$$\hat{y} = f(W^T \cdot X) \quad X = 1, x_1, x_2,$$

FORWARD PROPOGATION

The Perceptron Can Take Many Forms



$$\hat{y} = f(w_0 + \sum x_i \cdot w_i)$$

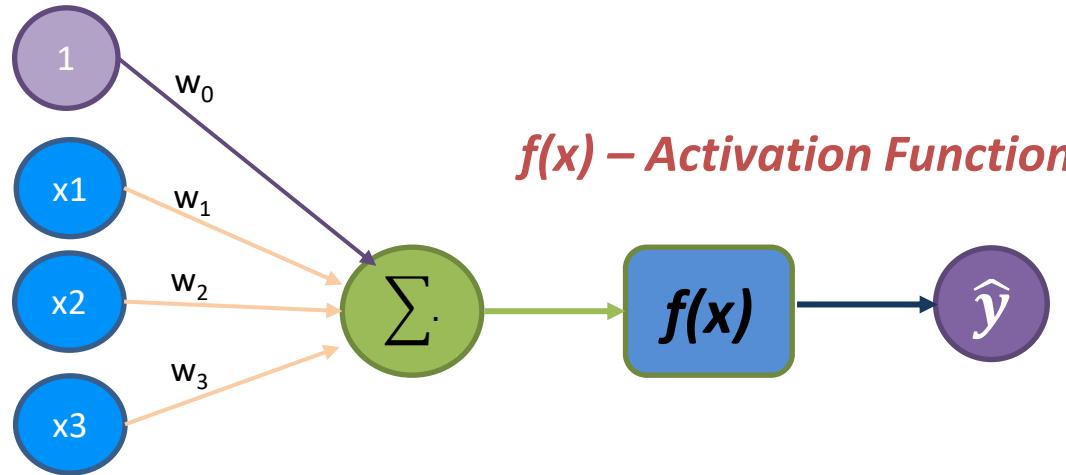
Let, $f(x) = 1$

[Linear]

$$Y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

[Linear Regression]

The Perceptron Can Take Many Forms



$$\hat{y} = f(w_0 + \sum x_i \cdot w_i)$$

Let, $f(x) = 1/(1+ e^{-x})$

[Sigmoid]

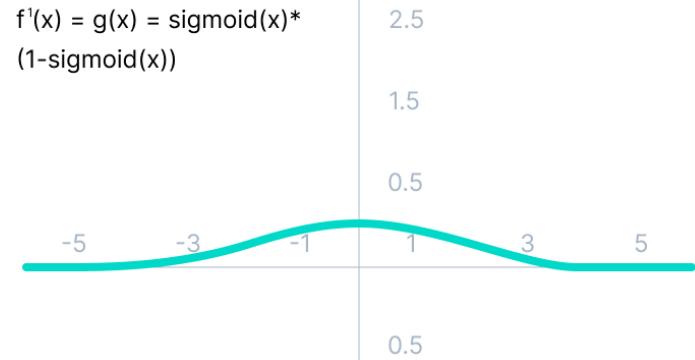
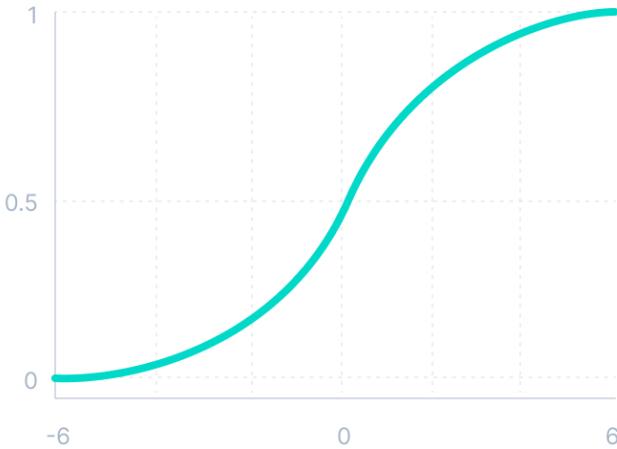
$$Y = \frac{e^{(w_0 + w_1x_1 + w_2x_2 + w_3x_3)}}{1+e^{(w_0 + w_1x_1 + w_2x_2 + w_3x_3)}}$$

[Logistic Regression]

Sigmoid Function

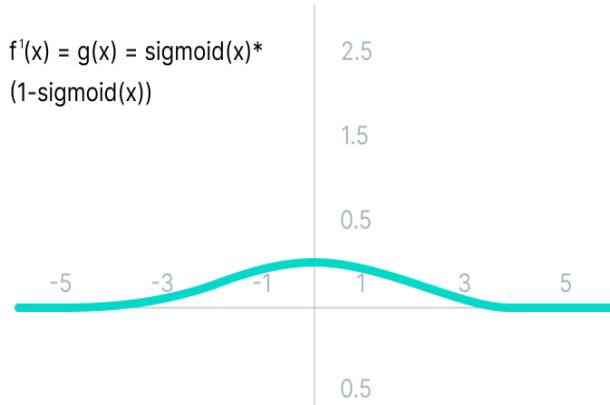
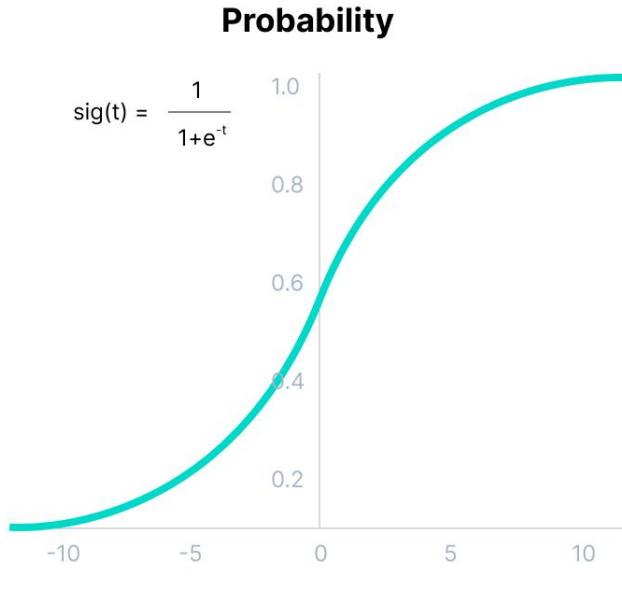
$$f(x) = \frac{1}{1 + e^{-x}}$$

Sigmoid / Logistic



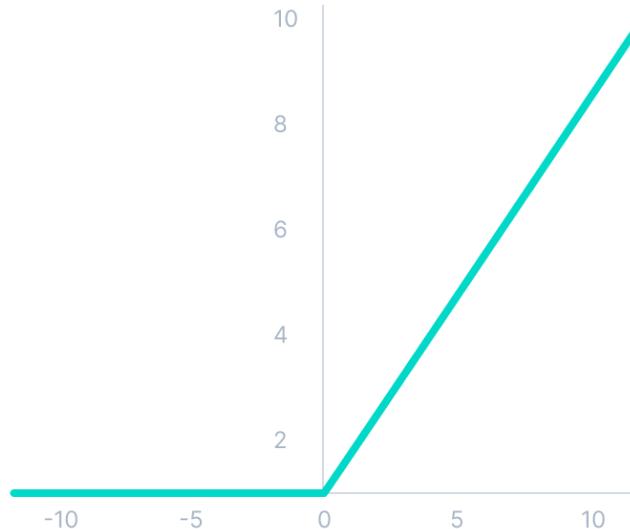
Softmax Activation Function

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

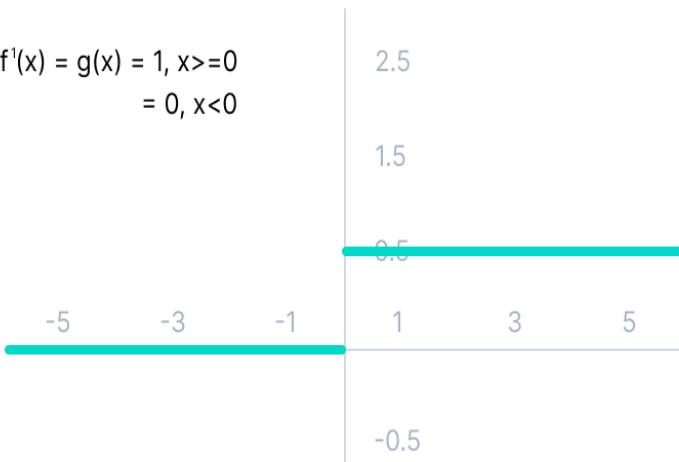


Rectified Linear Unit

$$f(x) = \max(0, x)$$

ReLU

$$\begin{aligned} f'(x) &= g(x) = 1, x \geq 0 \\ &= 0, x < 0 \end{aligned}$$



Other Activation Functions

Binary Step	Logistic	TanH	ArcTan	ReLU	PreLU	ELU	SoftPlus
$f(x)$	$f(x)$	$f(x) = \tanh(x)$	$f(x) = \tan^{-1}(x)$	$f(x) = \max(0, x)$	$f(x) = \max(x, 0)$	$f(x) = \max(\alpha x, 0)$	$f(x) = \ln(1 + e^x)$
$\begin{cases} 0 \text{ for } x < 0 \\ 1 \text{ for } x \geq 0 \end{cases}$	$\frac{1}{1 + e^{-x}}$	$\frac{2}{1 + e^{-2x}} - 1$	$\tan^{-1}(x)$	$\begin{cases} 0 \text{ for } x < 0 \\ x \text{ for } x \geq 0 \end{cases}$	$\begin{cases} \alpha x \text{ for } x < 0 \\ x \text{ for } x \geq 0 \end{cases}$	$\begin{cases} \alpha(e^{-x}-1) \text{ for } x < 0 \\ x \text{ for } x \geq 0 \end{cases}$	$\log_e(1 + e^x)$
$f'(x)$	$f'(x)$	$f'(x)$	$f'(x)$	$f'(x)$	$f'(x)$	$f'(x)$	$f'(x)$
$\begin{cases} 0 \text{ for } x \neq 0 \\ ? \text{ for } x = 0 \end{cases}$	$f(x)(1-f(x))$	$1-f(x)^2$	$\frac{1}{x^2 + 1}$	$\begin{cases} 0 \text{ for } x < 0 \\ 1 \text{ for } x \geq 0 \end{cases}$	$\begin{cases} \alpha \text{ for } x < 0 \\ 1 \text{ for } x \geq 0 \end{cases}$	$\begin{cases} f(x)+\alpha \text{ for } x < 0 \\ 1 \text{ for } x \geq 0 \end{cases}$	$\frac{1}{1 + e^{-x}}$

Artificial Neural Networks

Intro to NN - Sachin Kumar



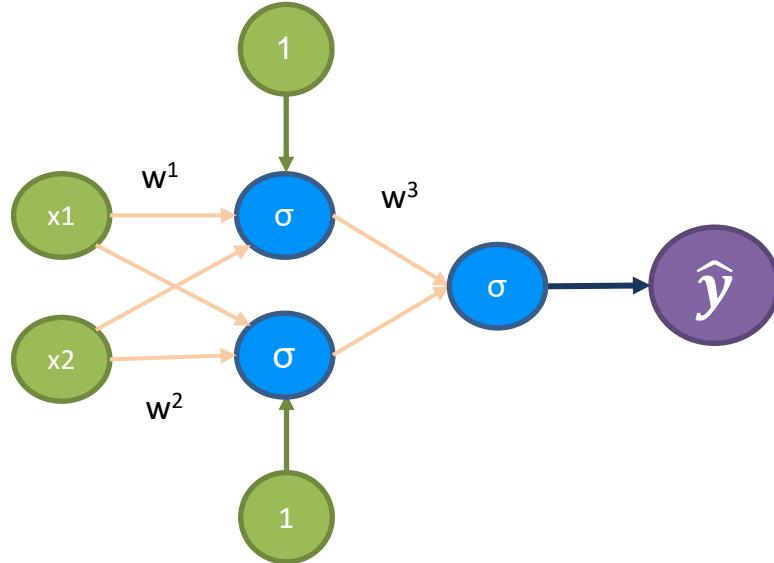
Biological Brain



Human Brain:

- 60% Fat
- 100 Billion Neurons
- 1000 trillion Synapses – connections
- Connection store memory and process information.

A Simple Multi-Layered Perceptron



$$x_1, x_2$$

$$\alpha_1 = \sigma(w_i x_i + b_1)$$

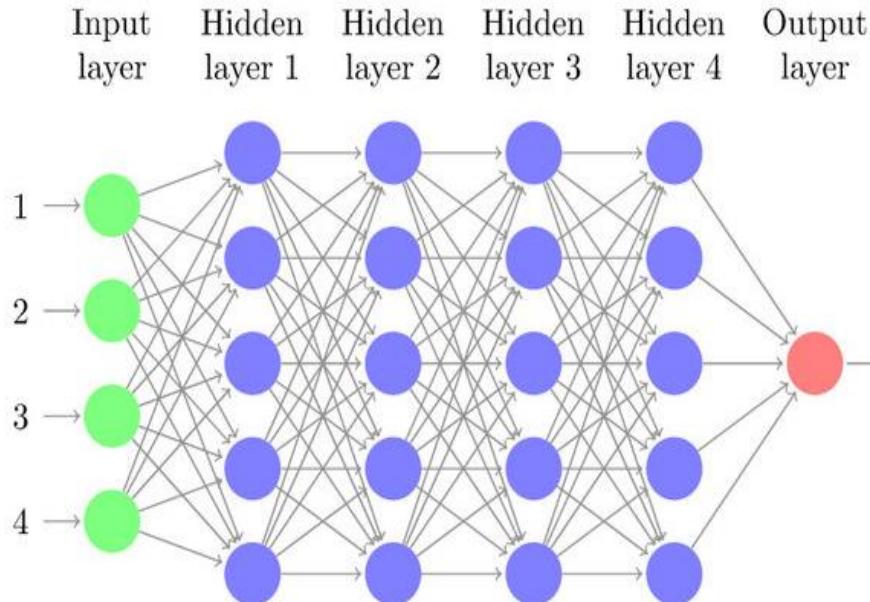
$$\alpha_2 = \sigma(w_i x_i + b_2)$$

$$\hat{y} = \sigma(w^3_1 \alpha_1 + w^3_2 \alpha_2)$$

} Input
} Intermediate
} Output

Multi-layered perceptron sum together the simple linear function, and each output a non linear function. The more layers we have in our hidden layer, the more complex non linear models we can find.

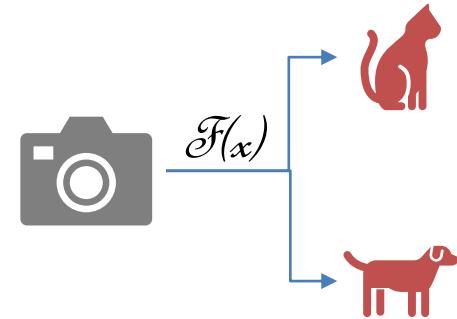
MLP – Universal Function Approximator



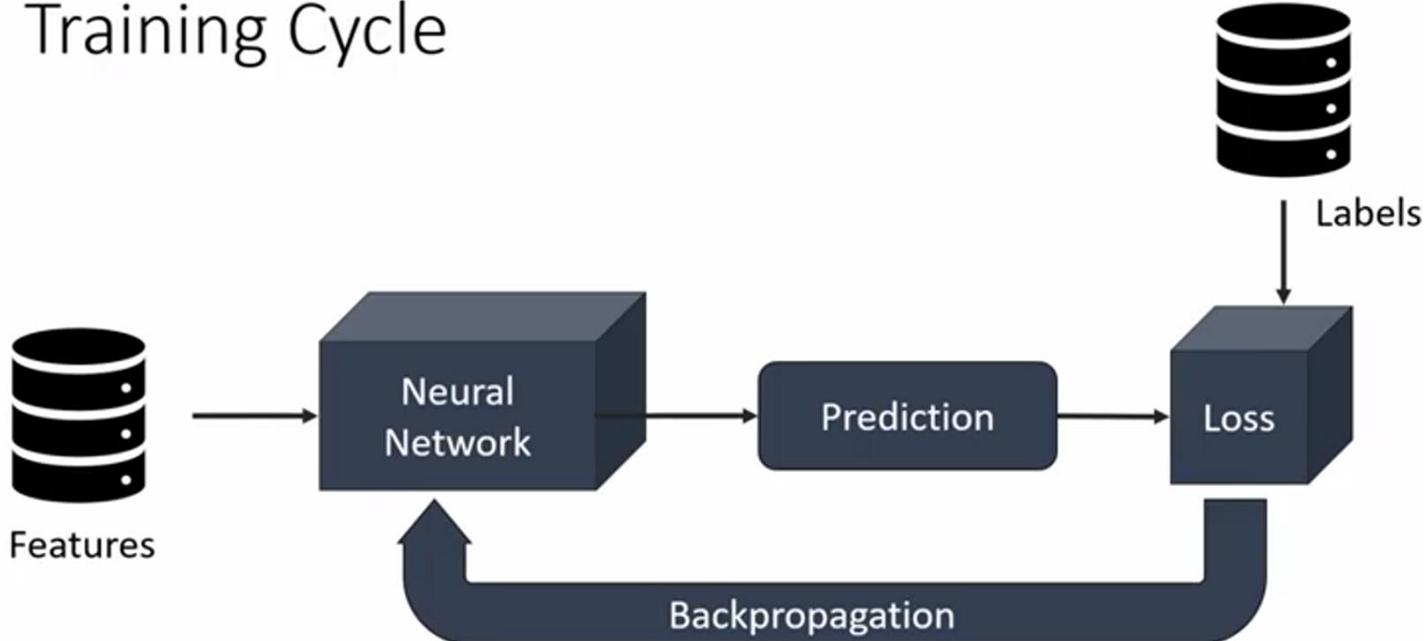
- Multi Layer Perceptron – 3 or more layers of nonlinearly-activating nodes.
- MLPs are fully connected.
 - $\text{Node } i \rightarrow w * \text{Node } j$

MLP – Universal Function Approximator

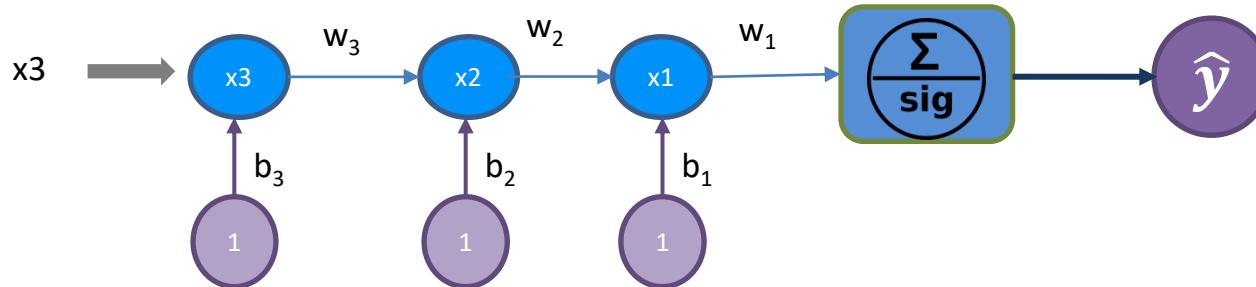
- **MLPs are universal function approximators !!**
- **Cybenko's theorem :** Given any known/unknown function, there exists a MLP that can approximate the function to an arbitrary degree of precision.
 - The network can be 1 layered with unbounded no of neurons
 - The network can be multi-layered with bounded no of neurons.
- Learning = $\delta(\text{connection weights})$ after each piece of data is processed.
- MLP learn **hierarchical feature** representations. Do not require feature engineering.



Training Cycle



Training a Simple Network



1

Define the Network

2

Define the Cost Function

3

Learn the weights/biases

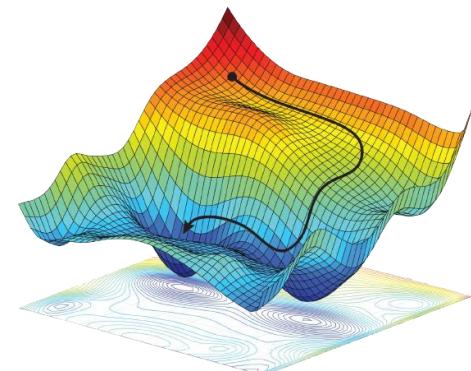
$$\hat{y} = f(w_1x_1 + b_1)$$

$$L(w): \sum (y - \hat{y})^2$$

$$x_1 = (w_2x_2 + b_2)$$

- Quantifies model performance
- # Continuous and mostly differentiable

$$x_2 = (w_3x_3 + b_3)$$



Gradient Descent

- Gradient descent is an optimization algorithm used to find the values of parameters (coefficients) of a function (f) that minimizes a cost function (cost).
- Gradient descent is best used when the parameters cannot be calculated analytically .
- It finds the parameters using an iterative method:
 - The goal is to continue to try different values for the coefficients, evaluate their cost and select new coefficients that have a slightly better (lower) cost.
 - Repeating this process enough times will lead to the coefficients that result in the minimum cost.

Gradient Descent

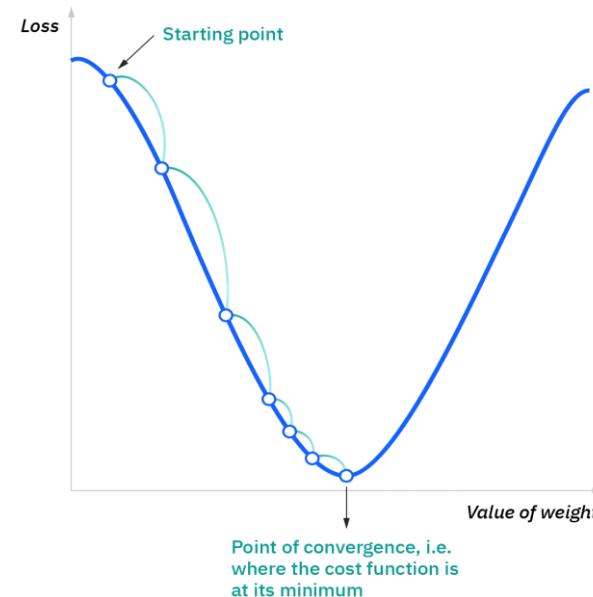
Gradient descent :

1. Estimate the error using current weights (randomly initialized)
2. Find the derivative(gradient) of the cost function.
3. Move in the direction opposite of the gradient, until convergence.

Gradient Descent Update Rule:

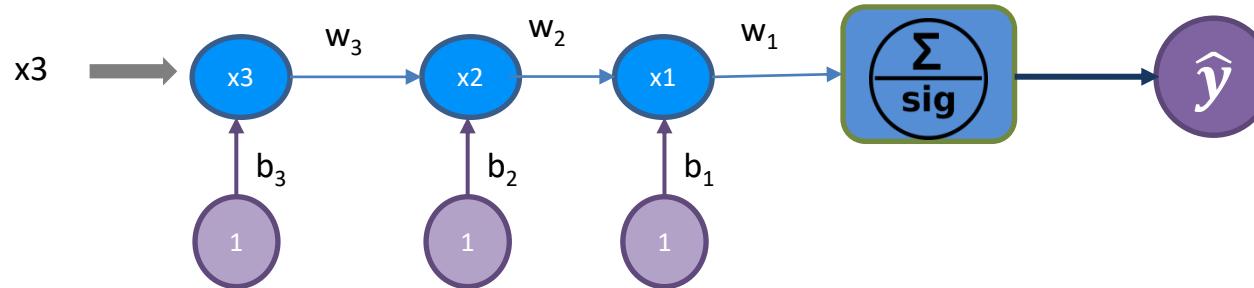
```
repeat until convergence {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$   
    (for j=1 to j=k )  
}
```

α – learning rate
 j - features



Training a Simple Network – Back Propagation

upGrad



Find the weights (w_1, w_2, w_3) such that the cost is minimized

$$L(w): \sum (y - \hat{y})^2$$

1. Start with random weights (w_1, w_2, w_3)
2. Change w_j such that

$$w_j := w_j - \alpha * \frac{\partial L}{\partial w_j}$$

BACKWARD PROPOGATION

1 $\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} * \frac{\partial y}{\partial w_1}$

2 $\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y} * \frac{\partial y}{\partial x_1} * \frac{\partial x_1}{\partial w_2}$

3 $\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial y} * \frac{\partial y}{\partial x_1} * \frac{\partial x_1}{\partial x_2} * \frac{\partial x_2}{\partial w_3}$

Cost Functions - Regression

upGrad

Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

MSE with L1 Regularization

$$\text{Cost} = \sum_{i=0}^N (y_i - \sum_{j=0}^M x_{ij} W_j)^2 + \lambda \sum_{j=0}^M |W_j|$$

MSE with L2 Regularization

$$\text{Cost} = \sum_{i=0}^N (y_i - \sum_{j=0}^M x_{ij} W_j)^2 + \lambda \sum_{j=0}^M W_j^2$$

Loss function Regularization Term

Mean Square Log Error

$$MSLE = \frac{1}{n} \sum_{i=1}^n (\log(Y_i) - \log(\hat{Y}_i))^2$$

Huber Loss Function

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

Binary Cross Entropy – Binary Classification

$$J = - \sum_{i=1}^N y_i \log(h_\theta(x_i)) + (1 - y_i) \log(1 - h_\theta(x_i))$$

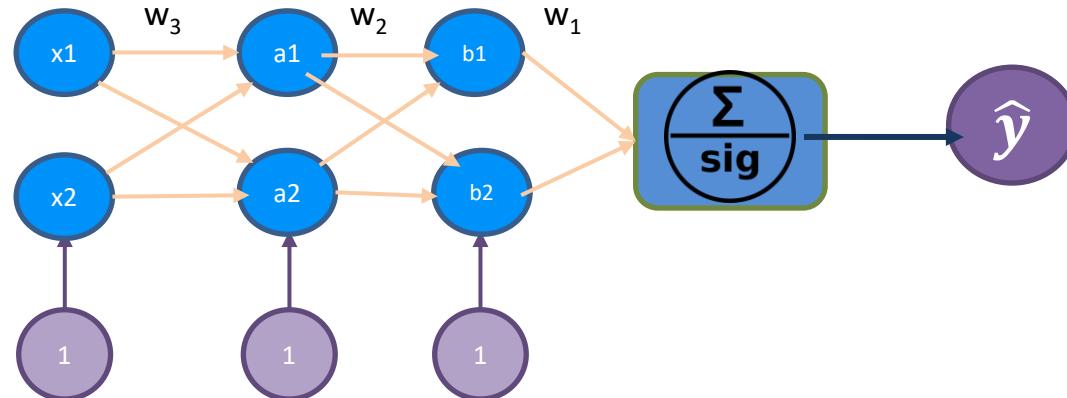
Cross Entropy – Multi Class Classification

$$\text{cross-entropy} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k t_{i,j} \log(p_{i,j})$$

Kullback Leibler Divergence Loss (KL Loss)

$$D_{KL}(P||Q) = \begin{cases} - \sum_x P(x) \cdot \log \frac{Q(x)}{P(x)} = \sum_x P(x) \cdot \log \frac{P(x)}{Q(x)}, & \text{for discrete distributions} \\ - \int P(x) \cdot \log \frac{Q(x)}{P(x)} \cdot dx = \int P(x) \cdot \log \frac{P(x)}{Q(x)} \cdot dx, & \text{for continuous distributions} \end{cases}$$

Training more Complex Networks



$$\frac{\partial L}{\partial w} = \left[\frac{\partial L}{\partial b} \right] * \left[\frac{\partial b}{\partial a} \right] * \left[\frac{\partial a}{\partial x} \right]$$

Derivative of a vector w.r.t
another vector yields a
Matrix – **Jacobian**

Linear Algebra & Matrix Calculus

Types of Gradient Descent

Batch Gradient Descent

Mini-Batch Gradient Descent

Stochastic Gradient Descent

Uses ALL training samples to calculate Error, for adjusting weights

Good for small datasets. Very time consuming with large datasets, large feature set.

Uses A random sample to calculate Error, for adjusting weights

Compromise between the two and works well in both cases. Efficient use of matrix multiplication, especially using GPUs.

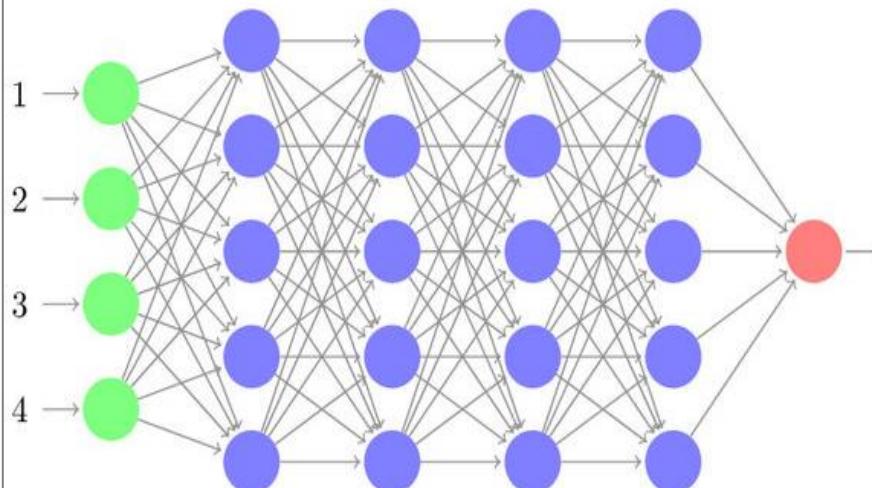
Uses ONE training Sample to calculate Error, for adjusting weights. Also called “Online” training.

Good for very big dataset with large no of features and redundant data.

Parameters:

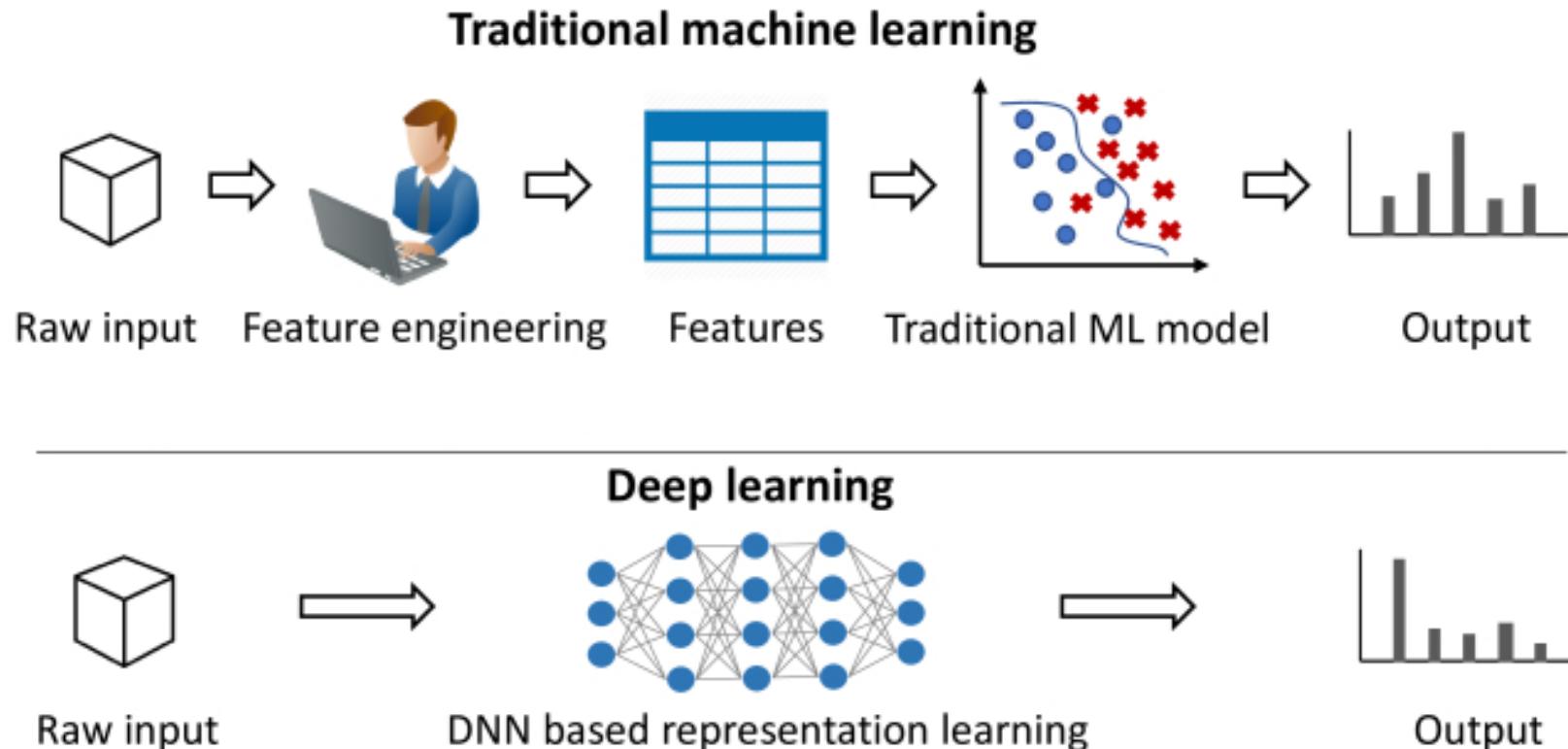
- Weights
- Biases

Input layer	Hidden layer 1	Hidden layer 2	Hidden layer 3	Hidden layer 4	Output layer
-------------	----------------	----------------	----------------	----------------	--------------

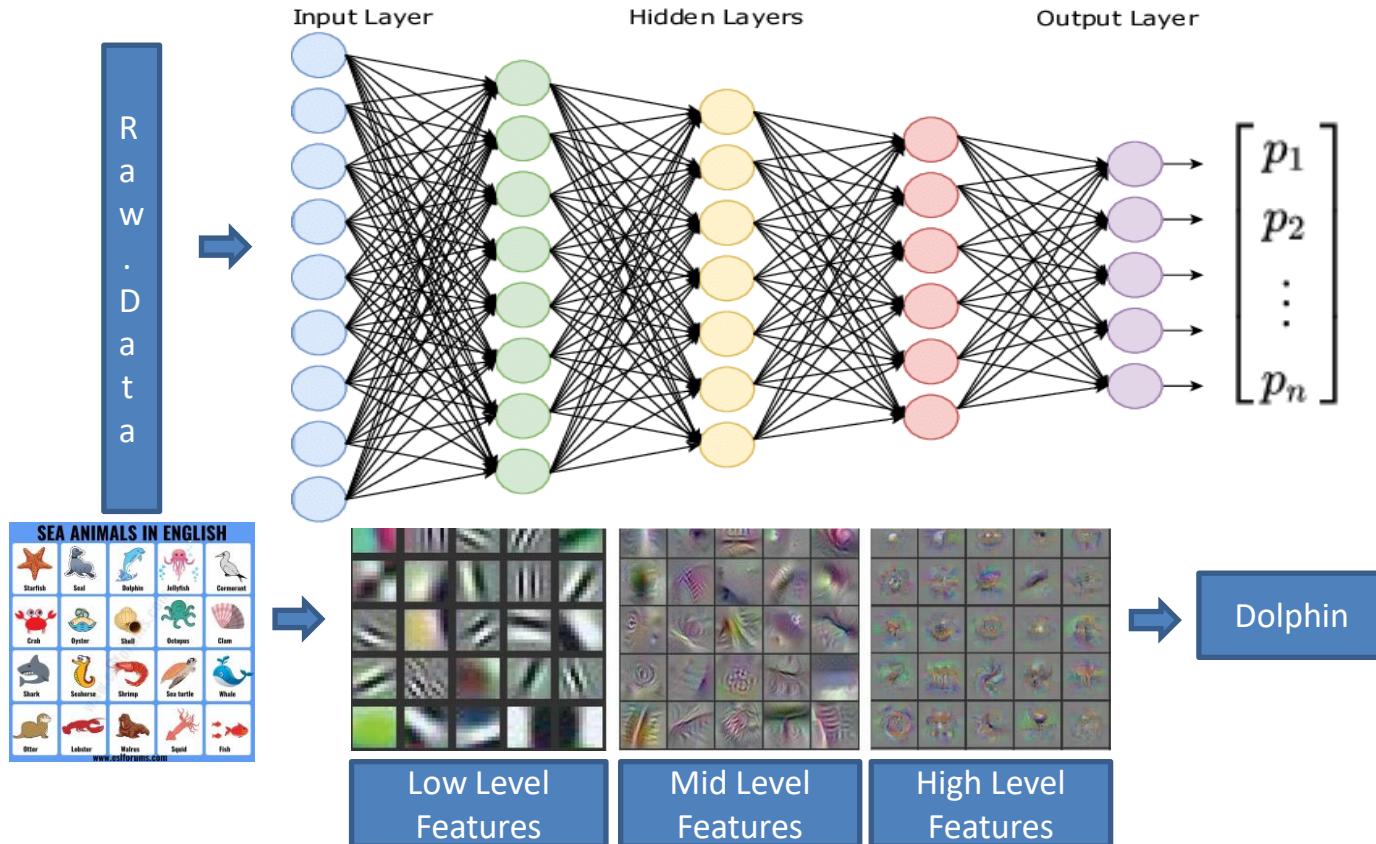


Hyper Parameters:

- Number of layers
- Number of neurons in the input, hidden and output layers
- Activation Function
- Regularization
- Learning rate
- Batch Size
- Number of epochs



CNN- Illustration of Hierarchical Feature Extraction

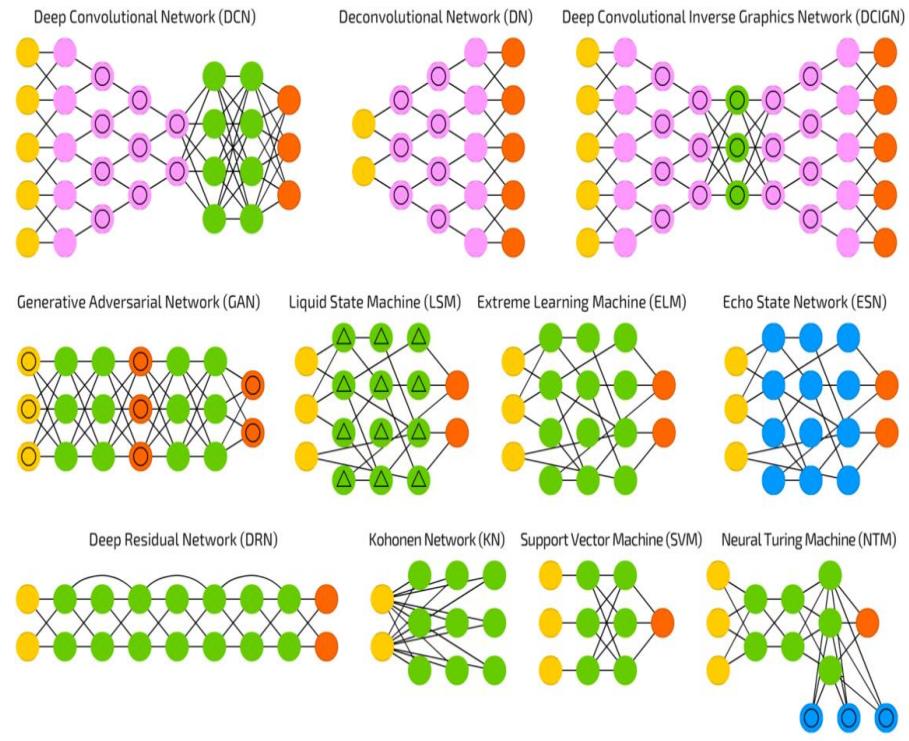
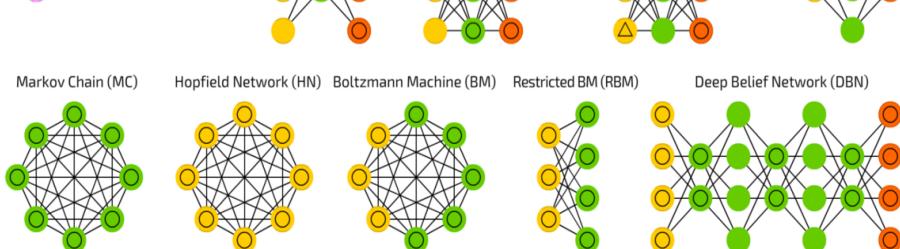


A mostly complete chart of

Neural Networks

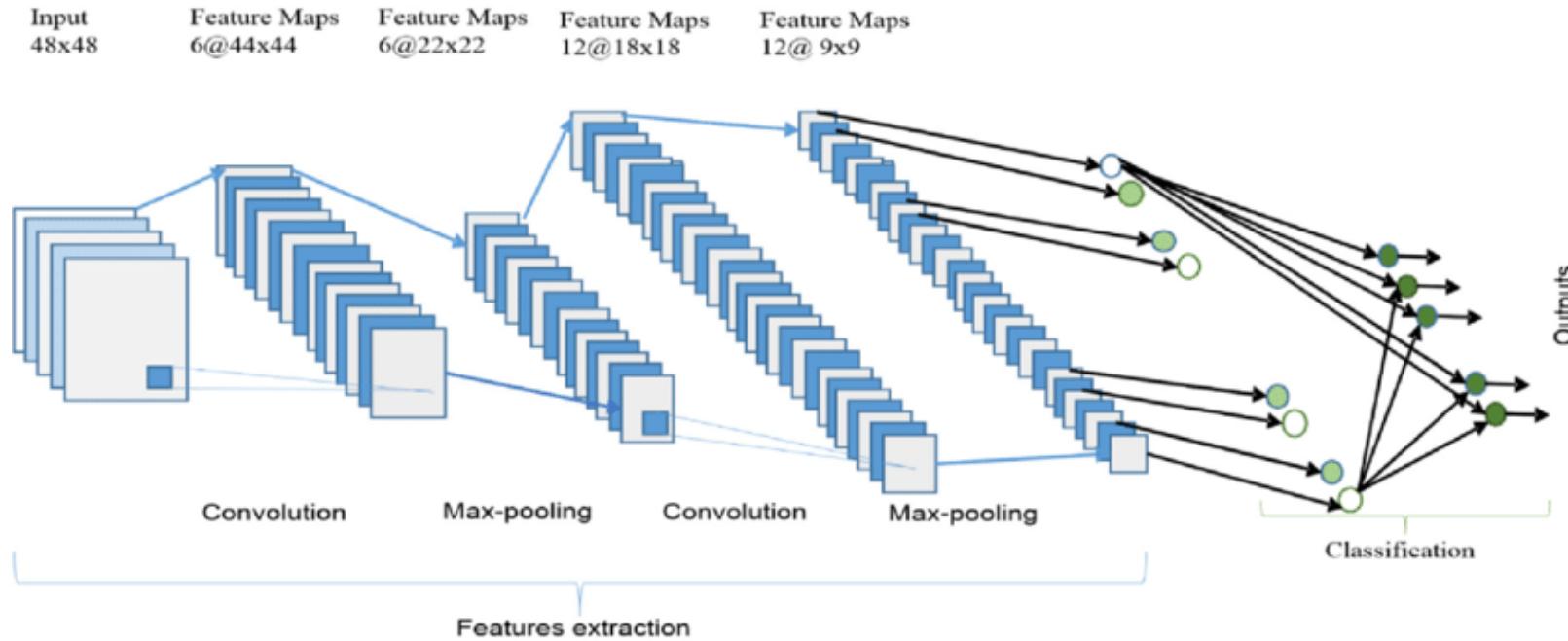
©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool



Convolutional Neural Network – Network the See

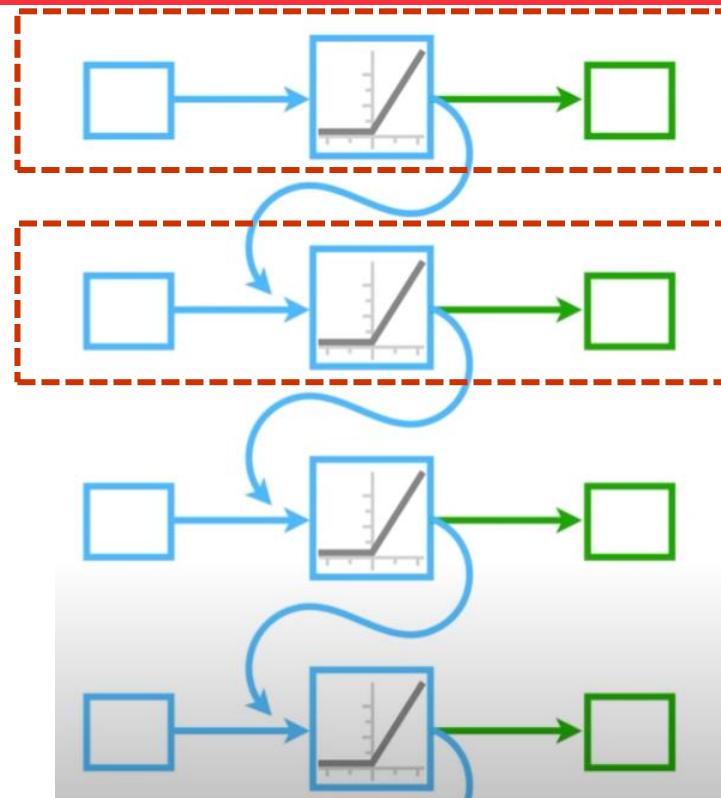
upGrad



- Sequence of **convolution** (of filters) and **pooling** of pixel values exploit spatial locality- enforce a local connectivity pattern between neurons of adjacent layers. The learned "filters" produce the strongest response to a spatially local input pattern.
- Each filter is replicated across the entire visual field – leading to translation invariance

Recurrent Neural Network – Feedback Loop

- RNN is a class of ANN with connections between nodes (cycle) - Output from some nodes to affect subsequent input to the same nodes.
- Exhibit temporal dynamic behavior.
- Process variable length sequences of inputs and generate variable sequence output.
- RNN can take a series of inputs with no predetermined limit on size and generate variable sequence output:
 1. Vector-Sequence Models — Image caption:
 2. Sequence-Vector Model — Sentiment analysis of a movie
 3. Sequence-to-Sequence Model - Language translation



Recurrent Neural Network – Applications

Speech recognition



Music generation

“The quick brown fox jumped over the lazy dog.”



Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis

AGCCCCTGTGAGGAACTAG

AG~~CCCCTGTGAGGAAC~~TAG

Machine translation

Voulez-vous chanter avec
moi?

Do you want to sing with
me?

Video activity recognition



Running

Name entity recognition

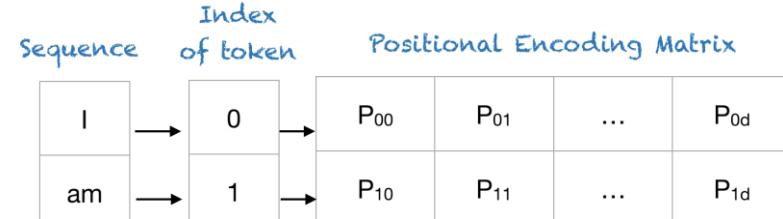
Yesterday, Harry Potter
met Hermione Granger.

Yesterday, ~~Harry~~ Potter
met ~~Hermione~~ Granger.

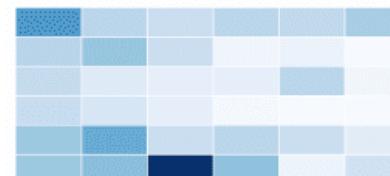
Transformer – Attention Mechanism

The transformer neural network is a novel architecture that aims to solve sequence-to-sequence tasks while handling long-range dependencies.

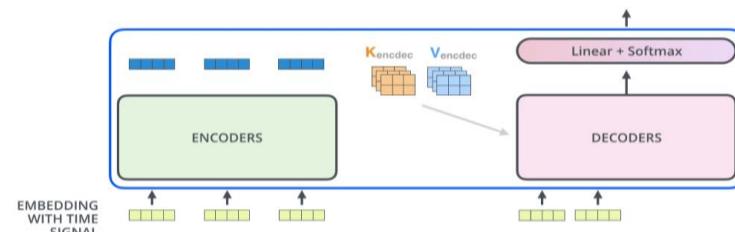
1. Positional encoding : Word embeddings are combined with their position in the sentence. This allows to input the entire sentence at once, and parallelize the process (using GPU)



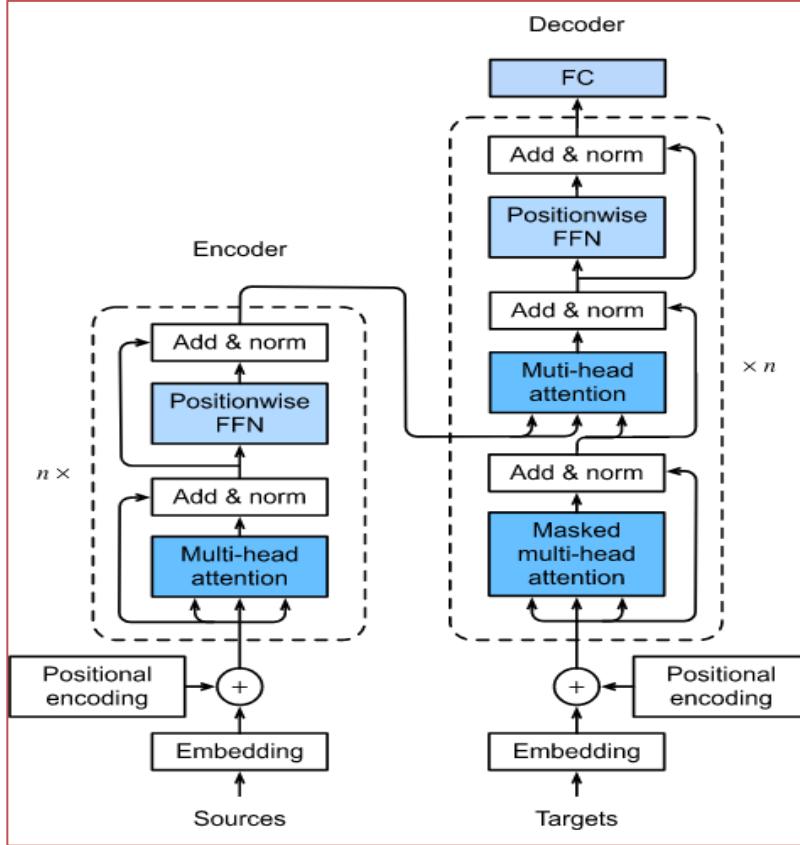
2. Self – Attention : Derive the relationship between each word with all other words in the input (using Query, Key and Value mechanism)



3. Decoding: The decoder generates the output sentence by defining the conditional probability distribution of a target sequence given a contextualized encoding sequence.

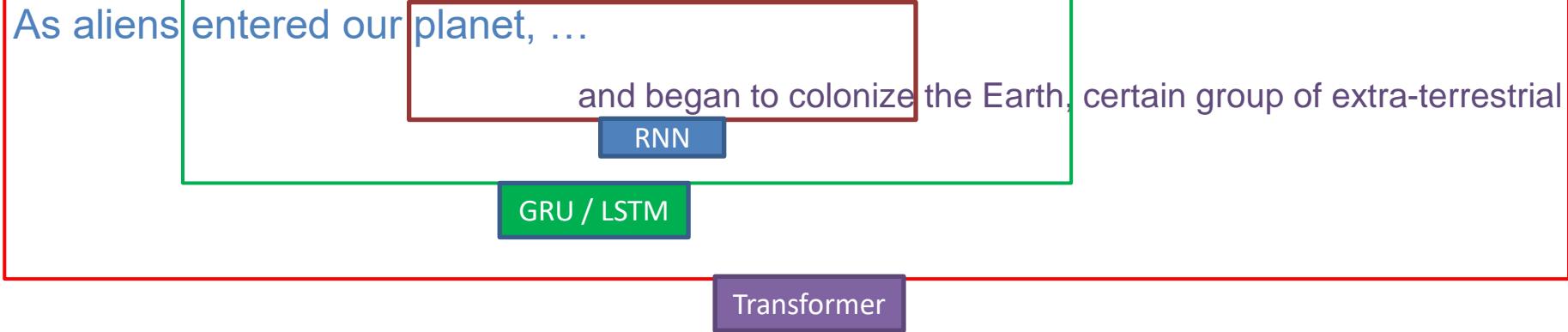


Transformer Architecture



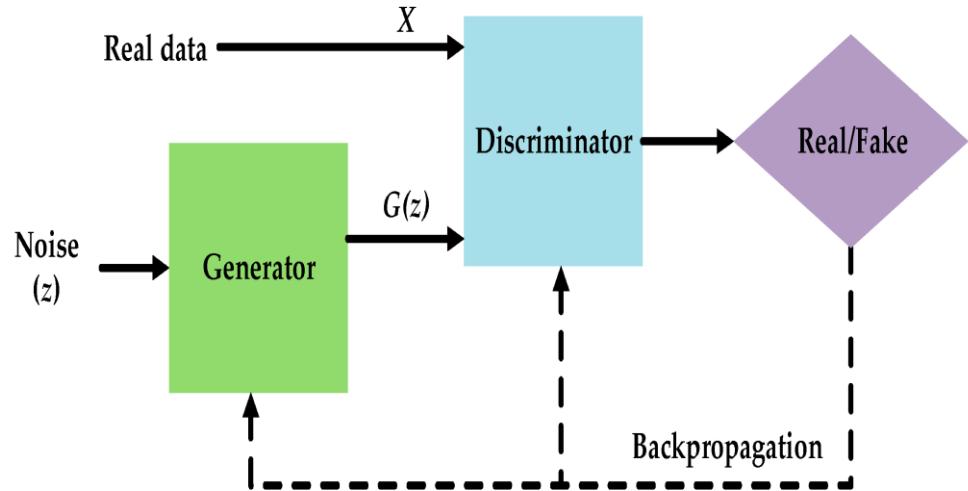
- BERT :
 - Bidirectional Encoder Representations from Transformers is a transformer-based machine learning technique for natural language processing.
 - BERT was created and published in 2018 by Google.
 - BERT is a deeply bidirectional, unsupervised language representation, pre-trained using only a plain text corpus.
- GPT-3 :
 - **Generative Pre-trained Transformer 3** is a language model that uses deep learning to produce human-like text.
 - GPT-3, which was introduced in May 2020, by OpenAI.
 - The architecture is a standard transformer network (with a few tweaks) with 175 billion parameters.

Attention – RNN vs LSTM vs Transformer

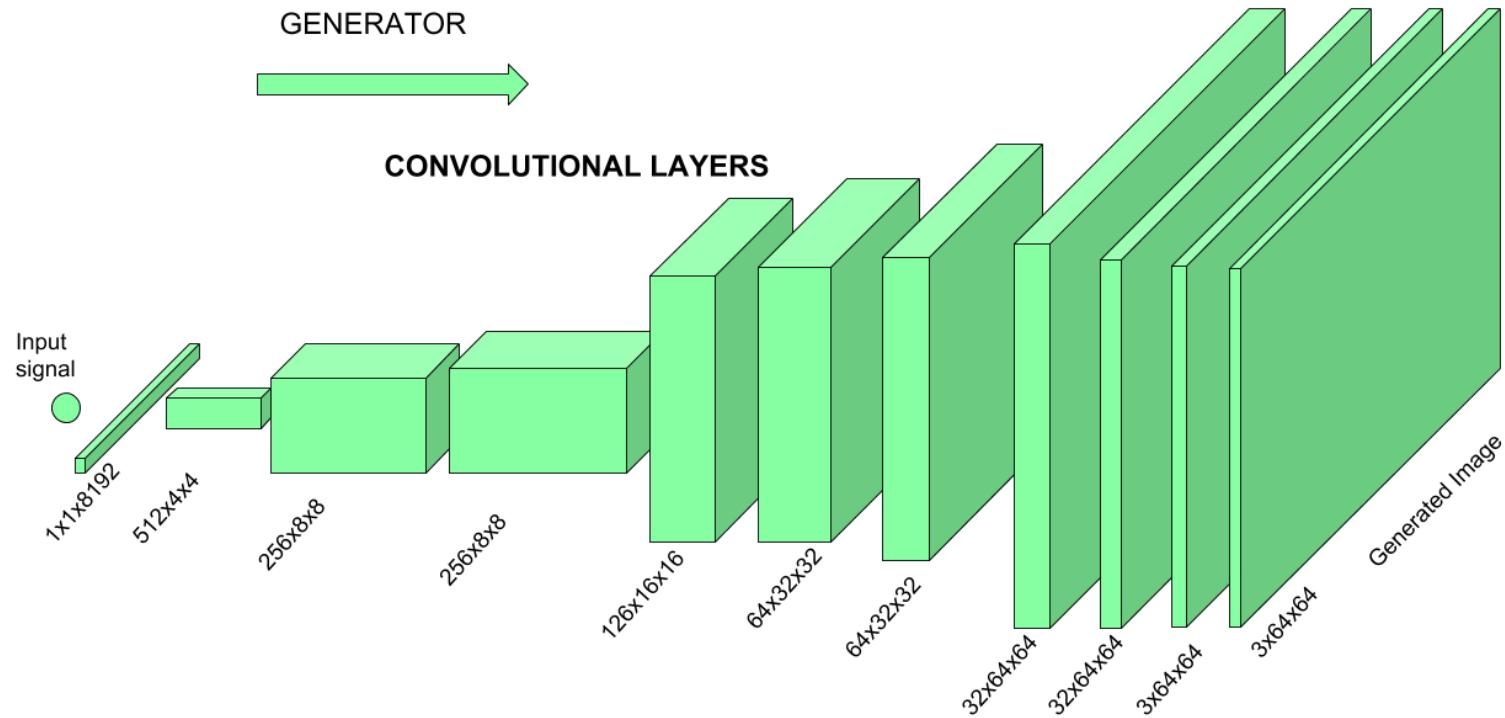


Generative Adversarial Networks

- GAN is a class of ML frameworks where two neural networks contest with each other in the form of a zero-sum game.
- Given a training set, this technique learns to generate new data with the same statistics as the training set.
- For example, a GAN trained on photographs can generate new photographs that look at least superficially authentic to human observers, having many realistic characteristics.

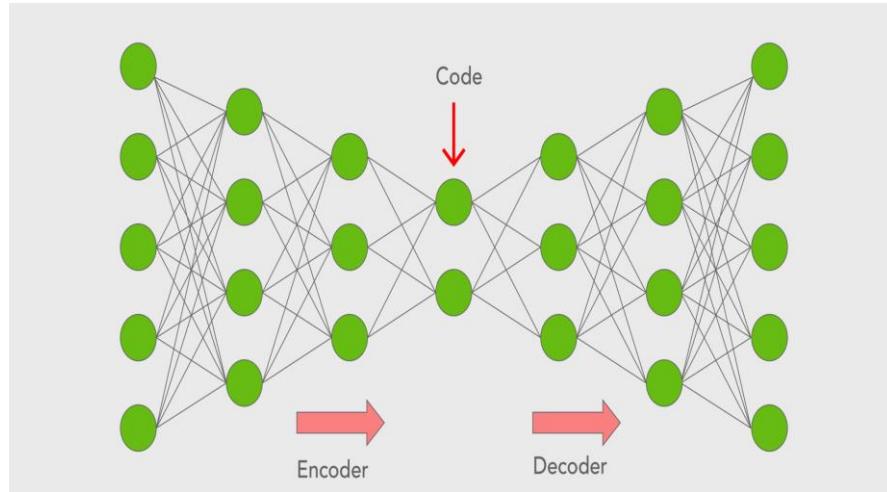


Generator of GAN



Auto Encoders

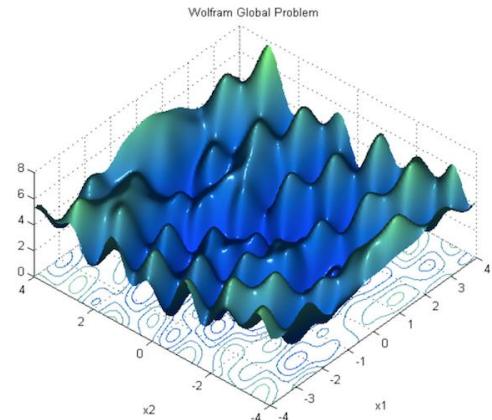
- Autoencoders operate by taking in data, compressing and encoding the data, and then reconstructing the data from the encoding representation.
- The model is trained until the loss is minimized and the data is reproduced as closely as possible.
- Through this process, an autoencoder can learn the important features of the data
- Used for dimensionality reduction, data denoising, feature extraction, image generation, sequence to sequence prediction, and recommendation systems.



Challenges of Deep Learning



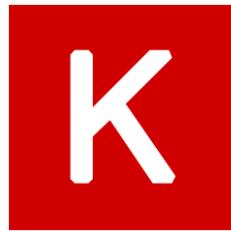
Large Dataset



Global Minima



Resource & Time



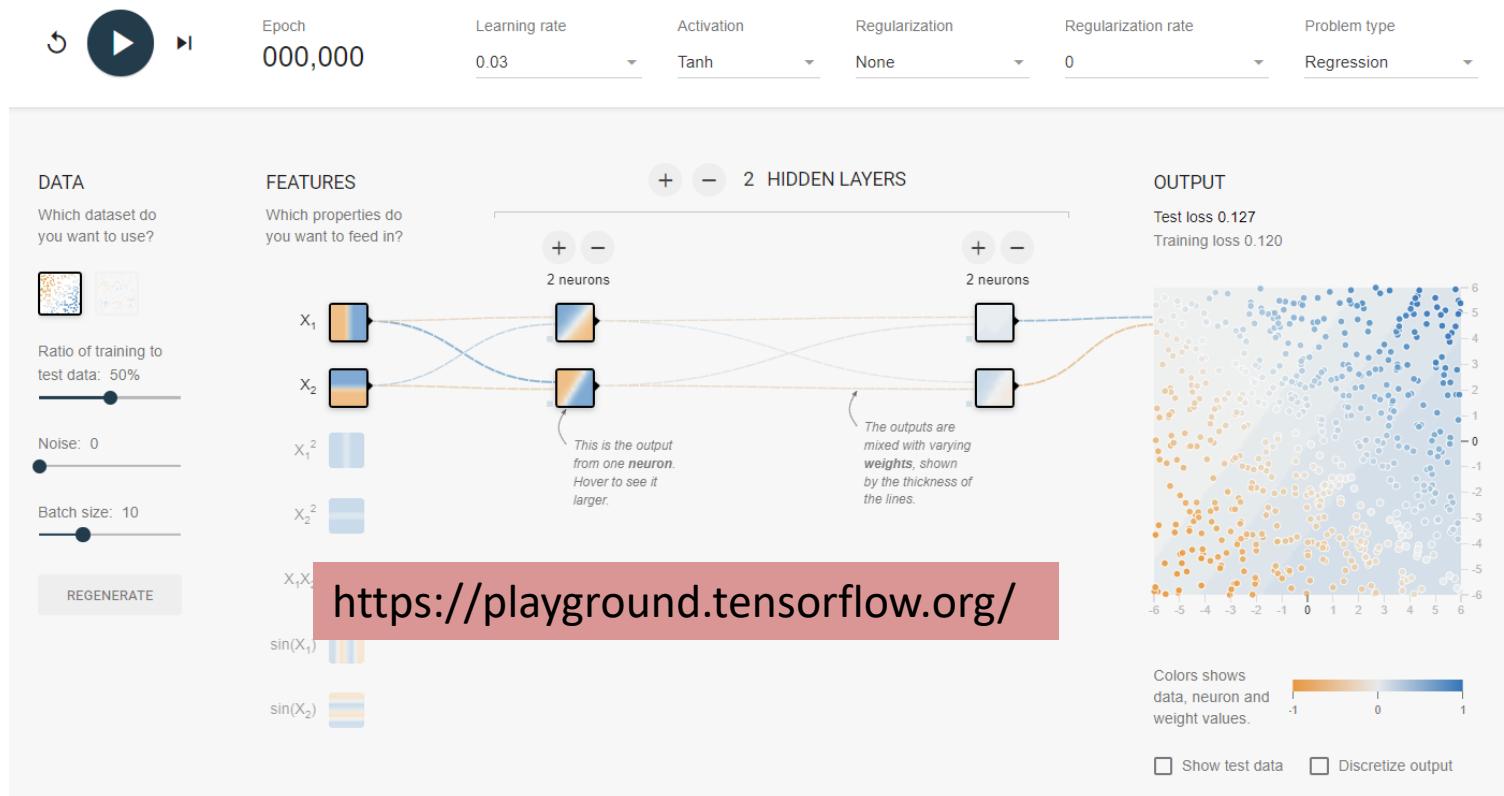
Caffe



theano

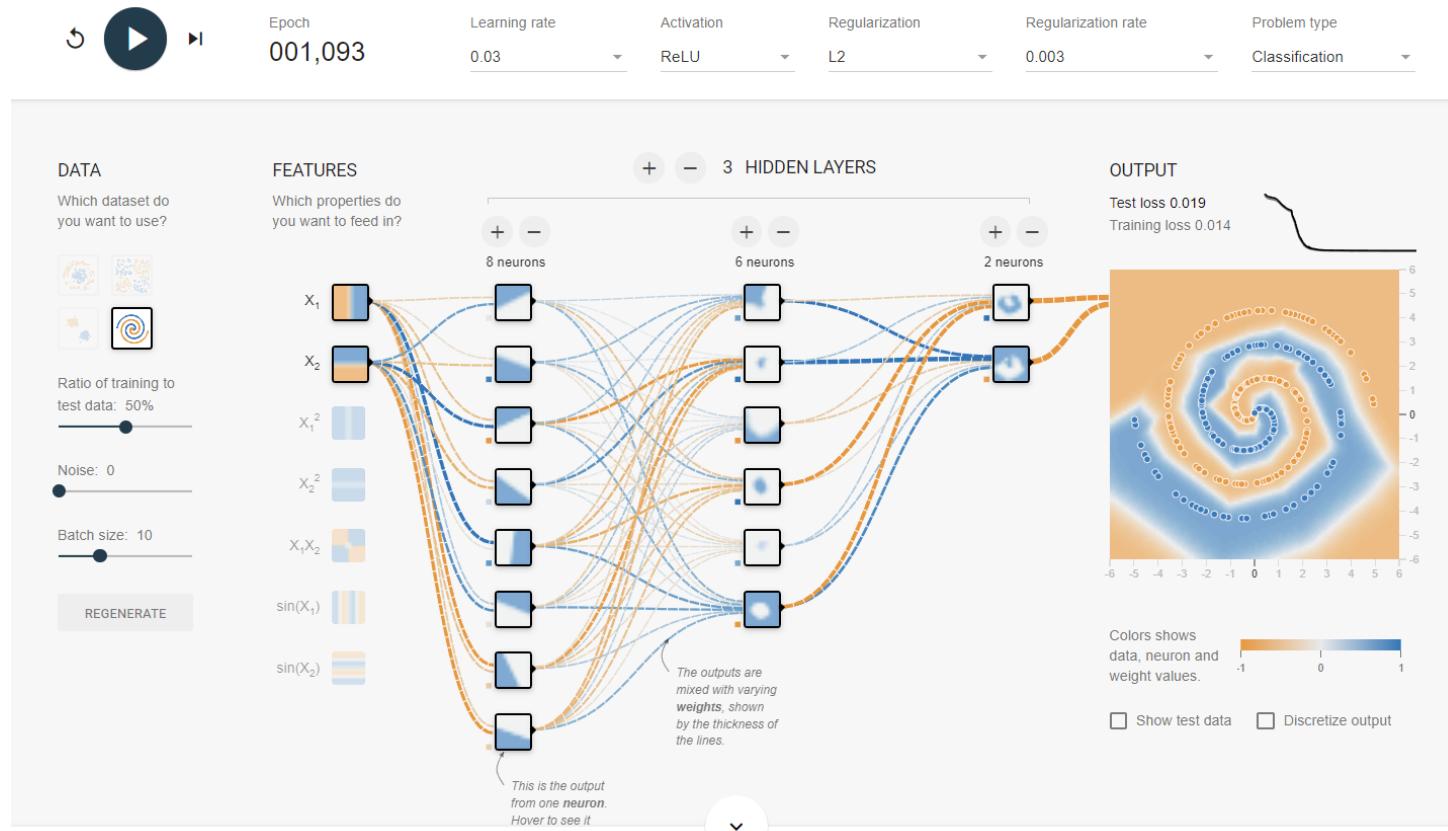
TensorFlow Playground

upGrad



Tensorflow Playground

upGrad



Discussion

upGrad





Thank You!