# AI - Past, Present & Future

## Currently Statistical Machines, NOT Cognitive Systems
## Machine (Artifical) "Intelligence" $\neq$ Human Intelligence

Jaideep Ganguly

Doctor of Science, Massachusetts Institute of Technology
Master of Science, Massachusetts Institute of Technology
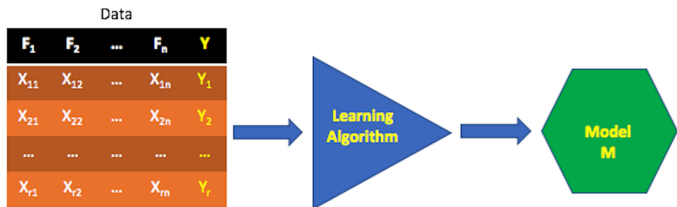Bachelor of Technology, Indian Institute of Technology, Kharagpur

March 31, 2024

# Part-1: Knowledge Based Expert Systems (1980-1995)

1. The field of AI was defined as computers performing tasks that were specifically thought of as something only humans can do.

2. In the 1980s, the expert systems were of great interest and focused on knowledge and inference mechanisms. They did a good job in their domains but were narrow in specialization and were difficult to scale.

3. Once these systems worked, they were no longer considered to be AI! For example, today the best chess players are routinely defeated by computers but chess playing is no longer really considered as AI! McCarthy referred to as the "AI effect". IBM's Watson is one such program at a level such as that of a human expert.

4. Fifty years ago Jim Slagle's (MIT) symbolic integration program (MACSYMA) was a tremendous achievement.

5. It is very hard to build a program that has "common sense" and not just narrow domains of knowledge.

# Machine Learning - Classification, Classification ...

1. In 1959, Arthur Samuel (MIT), defined ML, a subset of AI, as a *"field of study that gives computers the ability to learn without being explicitly programmed".*
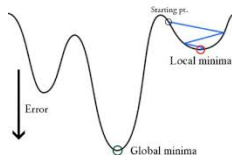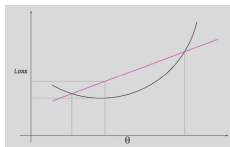


2. ML is effective for complex tasks where deterministic solution don't suffice, e.g., speech recognition, handwriting recognition, spam, fraud detection, etc. These cannot be solved manually in a large scale.

# Linear Regression - Elementary Algebra/Calculus

1. Regression is a statistical approach to find the relationship between variables between $X_i$ and $Y_i$. The model or the hypothesis is given by:

$$y_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} \cdots + \theta_n x_{in} \tag{1}$$

$$\text{Squared Error Loss (L)} = \frac{1}{2} \times \sum_{i=1}^{r} (\hat{y}_i - y_i)^2 \qquad \text{(Convex Function)} \tag{2}$$



1. Minimize by setting the partial derivative of the loss wrt $\theta_j$ to zero.

$$\frac{\partial L}{\partial \theta_0} = \sum_{i=1}^{r} (\hat{y}_i - y_i) = 0 \qquad \frac{\partial L}{\partial \theta_j} = \sum_{i=1}^{r} (\hat{y}_i - y_i)\, x_{ij} = 0 \tag{3}$$

2. Randomly assigning values to $\theta_j$.

# Logistic Regression

In many cases the value of $y_i$ need to be bounded between 0 and 1. For logistic regression, Least Squared Error will result in a *non-convex* graph with local minimums and hence is not feasible. In such cases, we use the logistic regression model as given below:

$$y_i = \frac{1}{1 + e^{-z}} \tag{4}$$

$$z_i = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \cdots + \theta_n x_n \tag{5}$$

- for z is large +ve number, $\frac{1}{e^z} = 0$; $y_i = 1$
- for $z = 0$, $y_i = 0.5$
- for z is large -ve number, $y_i = 0$

Hence, the value of $y_i$ is bounded between 0 and 1 for $z$ between $-\infty$ and $\infty$.
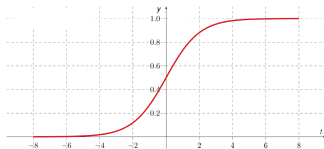


Figure: Sigmoid Curve

# Predictions & Errors



(a) Correct Fit     (b) Over fit     (c) Under fit

Predictions from the model will have differences or errors. Overfitting will occur when an excessive number of features are used than required. Underfitting occurs when an insufficient number of features are used than required.

1. Bias is the difference between average model prediction and the true target value and variance is the variation in predictions across different training data samples. Simple models with small number of features have high bias and low variance whereas complex models with large number of features have low bias and high variance.

2. Regularization is a technique used to avoid problem of overfitting. It prevents overfitting in linear models by a penalty term that penalizes large weight values.

# Information & Uncertainty - Claude Shannon (MIT)

1. **Information Content** When information is probabilistic, given $p(x)$ as probability mass function, the Information Content, $I_x$, is given by:

$$I_x = log\left(\frac{1}{p(x)}\right) = -log(p(x)) \tag{6}$$

2. **Shannon Entropy** of the random variable $X$ is defined as:

$$H(X) = \sum_x -p(x)log(p(x)) = E(I_x) \tag{7}$$

It is the *expected information content* of the measurement of $X$.

3. Cross-Entropy is defined as:

$$H(p, q) = -\sum_i p(i)log_2 q(i) \tag{8}$$

where p is the true distribution and q is the predicted distribution.
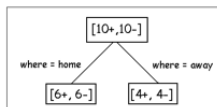
4. The *Kullback − Leibler (KL) divergence* is given by:

$$KL\ Divergence = H(p, q) - H(p)$$

## Decision Tree

1. Decision Tree is a supervised machine learning algorithm where You try to separate your data and group the samples together in the classes they belong to. You maximize the purity of the groups as much as possible each time you create a new node of the tree.

2. At each step, each branching, you want to decrease the entropy, so this quantity is computed before the cut and after the cut. If it decreases, the split is validated and we can proceed to the next step, otherwise, we must try to split with another feature or stop this branch.

3. XGBoost is a machine learning algorithm that belongs to the ensemble learning category, specifically the gradient boosting framework. It utilizes decision trees as base learners and employs regularization techniques to enhance model generalization.

## Example

Following is a table for win/loss of soccer games at home and away.



The entropy is:

$$H\left(\frac{6}{12}, \frac{6}{12}\right) = -\frac{6}{12}log_2\left(\frac{6}{12}\right) - \frac{6}{12}log_2\left(\frac{6}{12}\right) = 0.69$$

$$H\left(\frac{4}{8}, \frac{4}{8}\right) = -\frac{4}{8}log_2\left(\frac{4}{8}\right) - \frac{4}{8}log_2\left(\frac{4}{8}\right) = 0.69$$

$$H = -\frac{12}{20} \times H\left(\frac{6}{12}, \frac{6}{12}\right) - \frac{8}{20}H\left(\frac{4}{8}, \frac{4}{8}\right) = 0.69$$
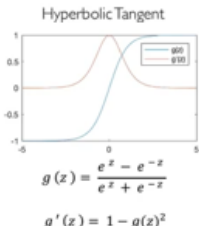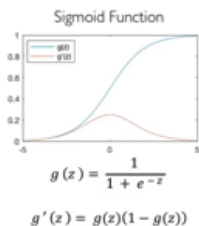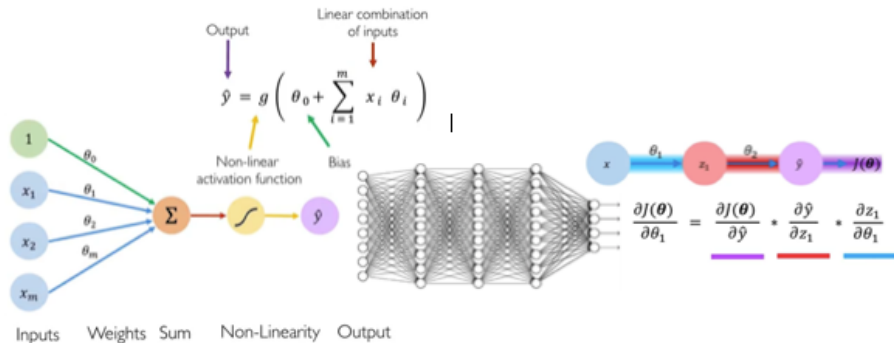


Information gain is 0.69 - 0.45 = 0.24.

# Part-2: From Perceptrons to Deep Learning (1995-2015)

1. Rosenblatt (Cornell) came up with the concept of Perceptrons, "a machine which responds like the human mind" as early as in 1957. In a critical book written in 1969, Marvin Minsky (MIT) & Seymour Papert (MIT) showed that Rosenblatt's original system was blind to simple XOR.

2. In 2006, Hinton developed Deep Learning which extends earlier important work by Yann LeCun (New York Univ). In ImageNet 2012, Hinton achieved the best accuracy in image recognition over $> 10\%$.

3. Deep Learning Success Stories - Image Recognition, Speech Comprehension, Chatbot. DNNs are suitable where the raw underlying features are not individually interpretable. This success is attributed to their ability to learn hierarchical representations, unlike traditional methods that rely upon hand-engineered features.

4. DL's innovation is to have models learn categories incrementally, attempting to nail down lower-level categories (like letters) before attempting to acquire higher-level categories (like words).
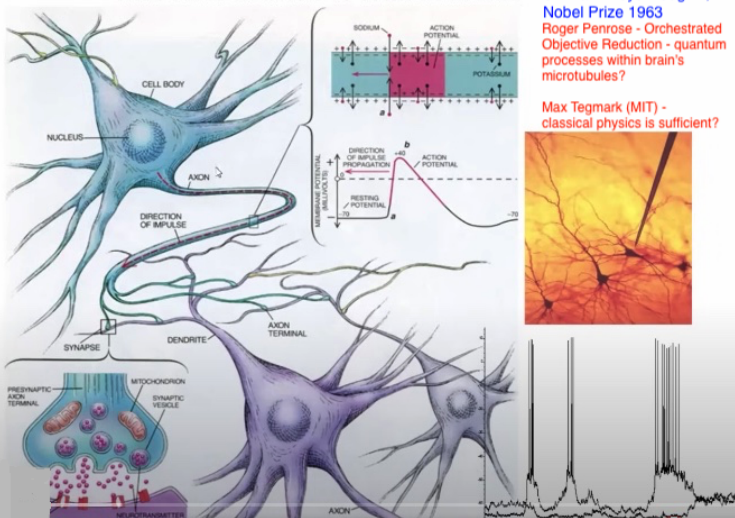
# Deep Neural Net



$$\hat{y} = g\left(\theta_0 + \sum_{i=1}^{m} x_i\,\theta_i\right)$$

Output

Linear combination of inputs

Non-linear activation function

Bias

Inputs  Weights  Sum  Non-Linearity  Output

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{\partial J(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial \theta_1}$$

**Sigmoid Function**

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

$$g'(z) = g(z)(1 - g(z))$$

**Hyperbolic Tangent**

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

$$g'(z) = 1 - g(z)^2$$

**Rectified Linear Unit (ReLU)**

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

Neurons transmit information via electrical spikes and make synapses with other neurons to form networks

Sir Alan LLoyd Hodgkin, Nobel Prize 1963

Roger Penrose - Orchestrated Objective Reduction - quantum processes within brain's microtubules?

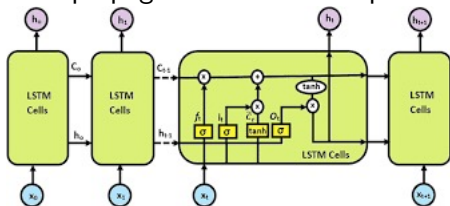Max Tegmark (MIT) - classical physics is sufficient?

# RNN

RNNs are connectionist models with the ability to selectively pass information across sequence steps while processing sequential data one element at a time.



But problems of vanishing and exploding gradients occur when backpropagating errors across many time steps.

# Long Short Term Memory - Hochreiter & Schmidhuber, RL

Can handle vanishing gradient problem faced by RNN. A separate cell state - the horizontal line running through the top of the LSTM unit allowing information to flow unchanged across many time steps. **Forget Gate**: Determines which information from the previous cell state should be forgotten. **Input Gate**: Controls the update of the cell state by selectively adding new information to it. **Output Gate**: Decides what information should be output from the cell state based on the current input and the previous cell state. Backpropagation does not require matrix multiplication



**Reinforcement Learning** - Agent learns to make decisions by interacting with the environment to achieve the goal through reward & punishment.

# **Attention Is All You Need**

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
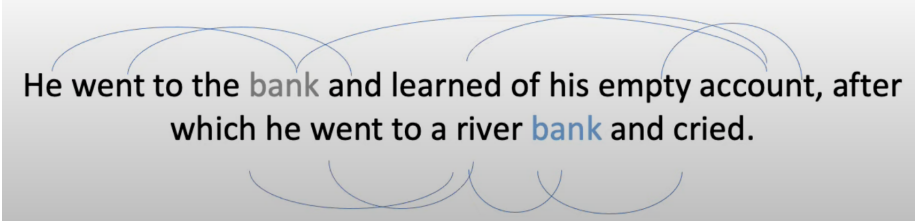Google Research
llion@google.com

**Aidan N. Gomez**[* †]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[* ‡]
illia.polosukhin@gmail.com

"Neural Machine Translation by Jointly Learning to Align and Translate"
by Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, 2014.
Introduced concept of attention mechanism and laid the foundation for
subsequent developments in NLP and DL, including the transformer
architecture introduced in "Attention Is All You Need."

# Attention



He went to the bank and learned of his empty account, after which he went to a river bank and cried.

1. Attention mechanism has an infinite reference window
2. In contrast, Recurring Neural Network (RNN) has a short reference window, Long Short Term Memory (LSTM) has a longer window. RNN does not work, LSTM has limited capability.

# Key Concepts



1. Encode, Positional encoding
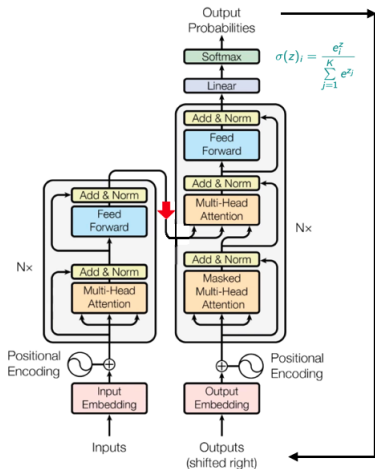$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

2. Key, Query and Value $\rightarrow$ Attention Filter

3. Multi Head Attention

4. Information preservation, normalisation

5. Decode, Masked Attention

$$\sigma(z)_i = \frac{e_i^z}{\sum\limits_{j=1}^{K} e^{z_j}}$$

Decoders are autoregressive models;
They are trained to predict the next token
after reading the preceding ones

# Query, Key & Value

input #1

| 1 | 0 | 1 | 0 |

input #2

| 0 | 2 | 0 | 2 |

input #3

| 1 | 1 | 1 | 1 |

Because every input has a dimension of 4, each set of the weights must have a shape of 4×3.
Weights are initialised randomly, it is done once before training.

Weights for **key**    Weights for **query**    Weights for **value**

[[0, 0, 1],          [[1, 0, 1],              [[0, 2, 0],
 [1, 1, 0],           [1, 0, 0],               [0, 3, 0],
 [0, 1, 0],           [0, 0, 1],               [1, 0, 3],
 [1, 1, 0]]           [0, 1, 1]]               [1, 1, 0]]

Key representation          Key representation          Key representation
for input 1:                for input 2:                for input 3:

                [0, 0, 1]                       [0, 0, 1]                      [0, 0, 1]
[1, 0, 1, 0] x [1, 1, 0] = [0, 1, 1]   [0, 2, 0, 2] x [1, 1, 0] = [4, 4, 0]   [1, 1, 1, 1] x [1, 1, 0] = [2, 3, 1]
                [0, 1, 0]                       [0, 1, 0]                      [0, 1, 0]
                [1, 1, 0]                       [1, 1, 0]                      [1, 1, 0]

Key representation:
(Vectorise)

key                         key                         key
| 0 | 1 | 1 |               | 4 | 4 | 0 |               | 2 | 3 | 1 |

                [0, 0, 1]
[1, 0, 1, 0]   [1, 1, 0]   [0, 1, 1]
[0, 2, 0, 2] x [0, 1, 0] = [4, 4, 0]
[1, 1, 1, 1]   [1, 1, 0]   [2, 3, 1]

input #1                    input #2                    input #3
| 1 | 0 | 1 | 0 |           | 0 | 2 | 0 | 2 |           | 1 | 1 | 1 | 1 |

# Multi Head Attention

# Multi Head Attention - Concatenation

$$x_0 = \frac{0.98 - 0.67}{\sqrt{0.44^2 + 0.0001}}$$

# Masking

# LLM - Summary

1. GPT does not replicate human cognitive behavior. LLMs essentially establish statistical connections among vectors representing words and extended grammatical structures with an associated probability.

2. It differs from human cognitive processes, where we use words to construct intricate and precise structures of "meanings".

3. But LLMs can do jaw-dropping things. But nobody knows exactly why. Models could seemingly fail to learn a task and then all of a sudden just get it, as if a lightbulb had switched on - "grokking".

4. Per classical statistics, the bigger a model gets, the more prone it is to overfitting. Double descent - performance of thea model initially decreases, then increases, and finally decreases again as the model complexity or dataset size increases - contradicts the traditional bias-variance tradeoff. This behavior, benign overfitting, is not fully understood. Raises basic questions on how models should be trained.

Thank You!