

# The Evolution of AI & ML

## Applications in different domains

Jaideep Ganguly, ScD (MIT)  
Head, Compass India Development Center

Distinguished speaker talk at Indian School of Business (ISB)

January 24, 2020



# About Jaideep

- ① Head of Compass India Development Center.
- ② compass.com is a startup, the India center started in December 2019.
- ③ *Compass is building the first modern real estate software platform to make search and sell experience intelligent and seamless.*
- ④ Director of Software Development, Amazon [2010 - 2019].
- ⑤ Director of Engineering, Microsoft [2006 - 2010].
- ⑥ Doctor of Science, Massachusetts Institute of Technology (MIT).
- ⑦ Master of Science, Massachusetts Institute of Technology (MIT).
- ⑧ Bachelor of Technology (Honours), Indian Institute of Technology, Kharagpur.

# Can a Computer Think?

- ① The human brain is a remarkable organ. It has enabled us understand science and advance mankind.
- ② The idea of mimicking the human brain or even improving the human cognitive functions is an alluring one and is an objective of Artificial Intelligence research.
- ③ But we are not even close in-spite of a century of research. However, it continues to have a major hold on our imagination given the potential of the rewards.
- ④ *The question of whether a computer can think is no more interesting than the question of whether a submarine can swim" - Dijkstra*
- ⑤ It is more interesting to understand the evolution of Machine Learning  
- *How did it start, here are we today and where do we go from here.*

# Symbols & Reasoning

- ① About 50,000 years ago, palaeontologists believe that some of us (1,000?), were able to deal with symbols - a major step in evolution.
- ② Noam Chomsky (MIT) thinks we were then able to create a new concept from 2 existing ideas without limiting the existing concepts.
- ③ Around 350 BC, Aristotle devised syllogistic logic, the first formal deductive reasoning system to model the way humans reason.
- ④ 2,000 years later, Bertrand Russel (Cambridge Univ) & his student Alfred Whitehead (Imperial College & Harvard) published Principia Mathematica that laid down the foundations of Mathematics.
- ⑤ John McCarthy (Stanford) championed mathematical logic in AI.
- ⑥ In 1942, Alan Turing (Cambridge Univ) showed that any form of mathematical reasoning could be processed by a machine.
- ⑦ By 1967, Marvin Minsky (MIT) declared that "within a generation, the problem of creating Artificial Intelligence would substantially be solved".
- ⑧ Clearly we are not there yet!

# Reasoning

- ① In 1961, Minsky published Steps towards Artificial Intelligence and talked about search, matching, probability, learning.
- ② In 1986, Minsky wrote the book The Society of Mind and talked about heuristic programming.
- ③ *What magical trick makes us intelligent? The trick is that there is no trick. The power of intelligence stems from our vast diversity, not from any single, perfect principle. — Marvin Minsky, The Society of Mind, p. 308*
- ④ McCarthy championed mathematical logic for artificial intelligence.
- ⑤ Newell (CMU) argued that the mind functions as a single system. He also claims the established cognitive models are vastly underdetermined by experimental data.
- ⑥ Simon (CMU) exposed the idea of bounded rationality, that rationality is limited. When individuals make decisions, by the tractability of the decision problem, the cognitive limitations of the mind, and the time available to make the decision.

# Heuristic Reasoning & Model Backed Systems

- ① In AI, model-based reasoning refers to an inference method used in expert systems based on a model of the physical world. The focus of application development is developing the model.
- ② [Terry Winograd \(MIT\)](#) built a model backed system for dialog understanding (SHRDLU).
- ③ [Patrick Winston \(MIT\)](#) built a model backed system for learning.
- ④ [Gerald Sussman \(MIT\)](#) built a model backed system for understanding blocks.
- ⑤ [Roger Schank \(Yale\)](#), in his Case-based reasoning, believed that understanding stories is the key to modeling human intelligence.
- ⑥ [David Marr \(MIT\)](#) treated vision as an information processing system. His hypothesis comprised of a computational level - what does the system, how does it do and a physical level - what neural structures and neuronal activities implement the visual system.

# Knowledge Based (KB) Expert Systems

- 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.

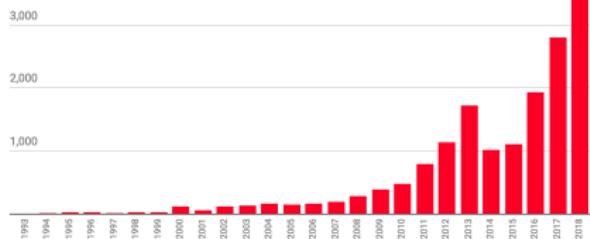
# From Perceptrons to Deep Learning

- ① [Rosenblatt \(Cornell\)](#) is credited with the concept of Perceptrons, “a machine which responds like the human mind” as early as in 1957.
- ② In a critical book written in 1969, [Marvin Minsky and Seymour Papert](#) showed that [Rosenblatt's](#) original system was blind to simple XOR.
- ③ That was incorrect and the field of “Neural Networks” disappeared!
- ④ [Geoff Hinton \(Toronto Univ\)](#) built more complex networks of neurons that allowed networks to learn more complicated functions like XOR.
- ⑤ But they learned slowly and couldn't master some of the basic things. By late 1990s, neural networks had again begun to fall out of favor.
- ⑥ In 2006, [Hinton](#) developed [Deep Learning](#) which extends earlier important work by [Yann LeCun \(New York Univ\)](#).
- ⑦ Deep learning's important 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).

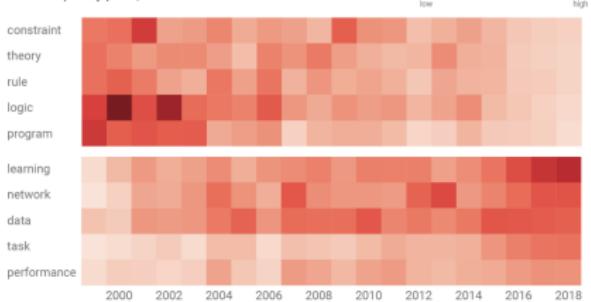
# Trends - Source MIT Technology Review

## The number of papers we downloaded from the arXiv

All of the papers available in the "artificial intelligence" section through November 18, 2018

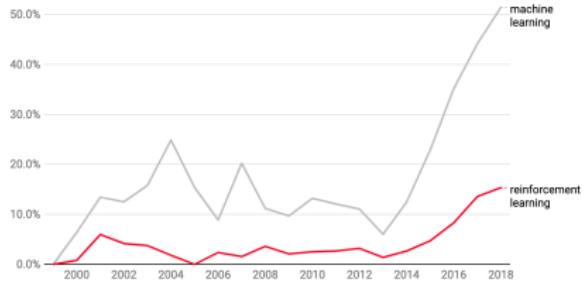


Word frequency per 1,000 words



Legend:

- neural networks
- bayesian networks
- markov methods
- evolutionary algorithms
- support vector machines

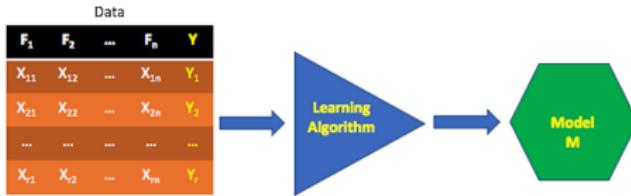


# Deep Learning (DL)

- ① DL powers Amazon's Alexa, Google's search, Facebook's news feed and Netflix's recommendation engine. Instead of manually encoding 1,000's of rules, the DL automatically extracts the rules from data.
- ② Other techniques such as Bayesian Networks, Support Vector Machines, etc., take different approaches to finding patterns in data.
- ③ In 2012, during the ImageNet competition in computer vision, [Hinton](#) achieved the best accuracy in image recognition by an astonishing margin of more than 10% using DL and sparked renewed research.
- ④ Recently, [Reinforcement Learning](#), which mimics the process of training through punishments and rewards, is in focus. In October 2015, DeepMind's AlphaGo, defeated the world champion in Go.
- ⑤ Every decade has seen focus on a different technique: Neural Networks in the late '50s and '60s, various symbolic approaches in the '70s, KB systems in the '80s, Bayesian Networks in the '90s, Support Vector machines in the '00s, and Neural Networks again in the '10s.
- ⑥ Nobody knows how to solve this problem.

# What is Machine Learning (ML)?

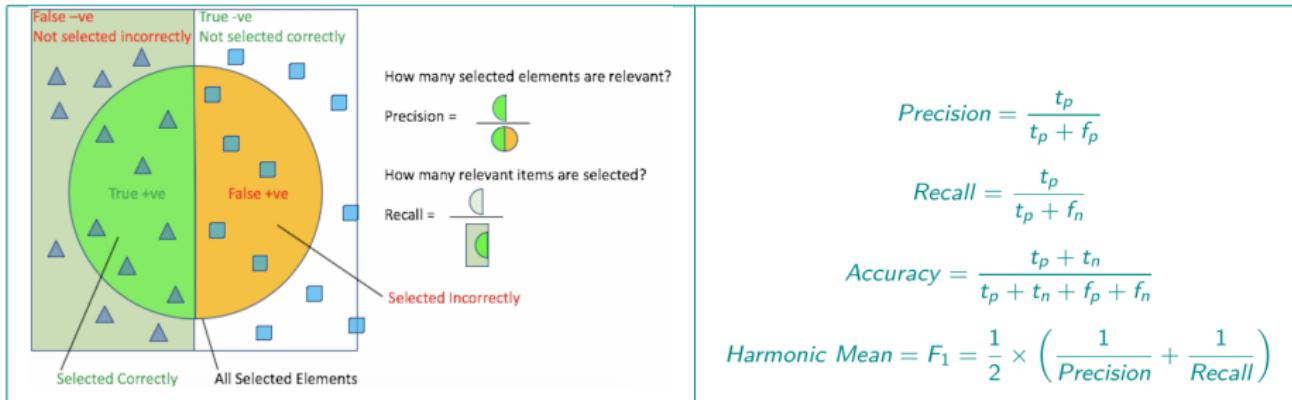
- ① In 1959, Arthur Samuel (MIT), defined ML, a subset of AI, as a “*field of study that gives computers the ability to learn without being explicitly programmed*”.



- ② ML algorithms establishes a relationship between input  $X$  to output  $Y$ . A feature is an individual measurable attribute of a phenomenon that is being observed. Given training examples  $(X_i, Y_i)$  where  $X_i$  is the feature vector and  $Y_i$  the target variable, learn a function  $F$  i.e.,  $Y_i \approx F(X_i)$  for all  $i$ . For a new sample  $X$  predict  $Y$  using  $F(X)$ .
- ③ 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.

# Types of Machine Learning & Basic Definitions

- ① ML algorithms are classified as **Supervised** or **Unsupervised**.
- ② In **supervised learning**, the machine learns from data that has already been tagged with the correct answer by an expert. This data is typically termed as the **training set**. E.g., separating spam mail.
- ③ In **unsupervised learning**, the machine groups unsorted information according to similarities without any training data. The machine automatically figures out the patterns on its own. E.g., playing GO.



# Linear Regression

- ① **Regression** is a statistical approach to find the relationship between variables between  $X_i$  and  $Y_i$ .
- ② A common function used to model training data is a **linear regression model**. The **model** or the **hypothesis** is given by:

$$y_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \dots + \theta_n x_{in} \quad (1)$$

where  $x_{ij}$  is the observed value of the feature  $x_j$ ,  $y_i$  is the predicted value of the outcome and  $\theta_i$  are constants the values of which need to be determined. For now, we limit the the values of the features  $x_{ij}$  to numbers, positive or negative. Later on, we will study techniques on how to deal with the situation where the feature value is a string.

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

This loss function is always positive, there is a minimum, is monotonically increasing from that minimum value in both positive and negative directions. Such a function is called a **convex function**.

- ③ We need to minimize some loss function over the training data.

# Gradient Descent

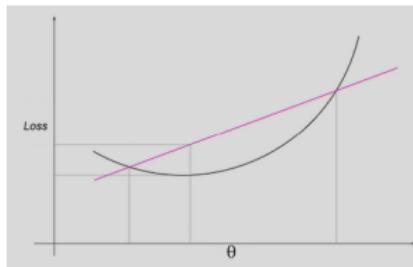


Figure: Convex Function

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

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

- ② Randomly assigning values to  $\theta_j$ . For a convex function it does not matter what the initial weights are as it will always converge.  $\hat{y}$  and  $x_{ij}$  are observed and  $y_i$  is computed. This means the slope of the loss function can be readily computed.

## Gradient Descent

Now for a given  $\theta_j$ , and keeping all other  $\theta_i$  where  $i \neq j$  constant, we change the value of  $\theta_j$  slightly and compute a new value of  $\theta_j$

$$\begin{aligned}\Delta\theta_j &= \alpha\theta_j \\ \theta_j' &= \theta_j + \Delta\theta_j\end{aligned}$$

where  $\alpha$  is small constant and is called the **learning rate**. Since the slope of  $\frac{\partial L}{\partial \theta_j}$  is already calculated from the previous equations, we can compute the new value of the loss as follows:

$$L' = L - \frac{\partial L}{\partial \theta_j}(\alpha\theta_j)$$

If  $L' < L$ , we continue incrementing  $\theta_j$  as long as the loss decreases. When the direction changes and the loss increases, we have found the  $\theta_j$  for which the loss is minimum. If  $L' > L$ , we travel in the reverse direction and follow the same process. Like this, we compute the value of each and every  $\theta_j$ .

# Stochastic versus Batch Gradient Descent

- ① **Stochastic** means a process that is linked with a random probability. In SGD, a single sample data is used to compute the  $\theta_j$ . The sample is randomly shuffled and selected for performing the iteration.
- ② In **batch gradient descent**, the gradient is calculated from the whole training set. If we have a huge dataset with millions of data points, running the batch gradient descent is expensive.
- ③ Compromise by computing the gradient for a **mini batch** at each step which performs better because of use of vectorization libraries and has smoother convergence as the gradient computed at each step uses more training examples. It is the algorithm of choice for neural nets.
- ④ In Batch Gradient Descent, the **Batch Size** refers to the total number of training examples present in a single batch and **Epoch** means that each sample in the training dataset has had an opportunity to update all the model parameters  $\theta_j$ .
- ⑤ Our objective is to find a set of  $\theta_j$  such that the loss is minimum.

# Logistic Regression

In many cases the value of  $y_i$  need to be bounded between 0 and 1. In such cases, we use the logistic regression model as given below:

$$y_i = \frac{1}{1 + e^{-z}} \quad (4)$$

$$z_i = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad (5)$$

- 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  $\infty$ .

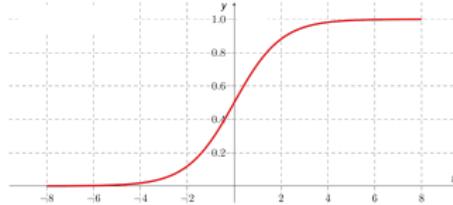
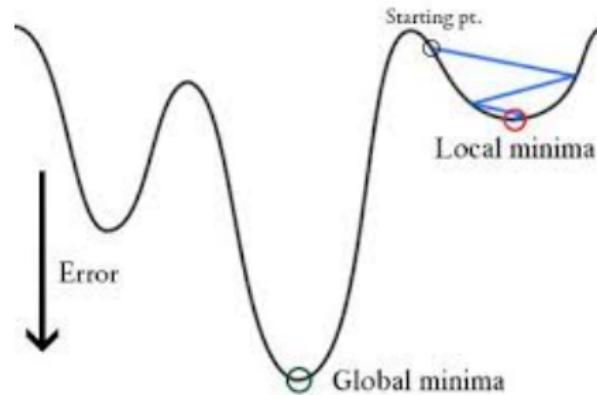


Figure: Sigmoid Curve

# Logistic Regression

For logistic regression, Least Squared Error will result in a *non-convex* graph with local minimums and hence is not feasible.



# Logistic Regression

Logistic regression is modeled as follows:

$$P(y|x, \theta) = \begin{cases} h_{\theta}(x) & \text{if } y = 1 \\ (1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \quad (6)$$

The above function can be written compactly as follows:

$$P(y|x, \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y} \quad (7)$$

Assuming that the  $r$  training examples were generated independently, we can then write down the **likelihood**  $L$  of the parameters as:

$$L(\theta) = p(\vec{y}|x, \theta) = \prod_{i=1}^r P(y^{(i)}|x^{(i)}, \theta)$$

$$L(\theta) = \prod_{i=1}^r (h_{\theta}(x^{(i)}))^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

# Log Likelihood

The log likelihood is given by:

$$l(\theta) = \log L(\theta) = \prod_{i=1}^r y^{(i)} \log(h(x^{(i)})) + (1 - y^{(i)}) \log(1 - h(x^{(i)})) \quad (8)$$

Intuitively, we want to assign more punishment when the prediction is 1 while the actual is 0 and similarly when the prediction is 0 while the actual is 1. The loss function in logistic regression is doing exactly this and hence is called **Logistic Loss**.

The sigmoid function has an interesting property. Its derivative is given by:

$$\begin{aligned} g'(x) &= \frac{d}{dx} \left( \frac{1}{1+e^{-x}} \right) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \times \left( 1 - \frac{1}{1+e^{-x}} \right) \\ &= \boxed{g(x)(1-g(x))} \end{aligned} \quad (9)$$

## Log Likelihood

Now, consider just one training example  $(x, y)$ , and take derivatives to derive the stochastic gradient ascent rule:

$$\begin{aligned}\frac{\partial(l(\theta))}{\partial\theta_j} &= y\frac{1}{g(\theta^T x)} - (1-y)\frac{1}{1-g(\theta^T x)}\frac{\partial}{\partial\theta_j}g(\theta^T x) \\&= y\frac{1}{g(\theta^T x)} - (1-y)\frac{1}{1-g(\theta^T x)}g(\theta^T x)(1-g(\theta^T x))x_j \\&= y(1-g(\theta^T x)) - (1-y)g(\theta^T x)x_j \\&= y - h_\theta(x)x_j\end{aligned}$$

With the above equation, we are now in a position to compute the gradient and follow the same steps to determine  $\theta_j$  corresponding to minimum loss as in the case of linear regression.

# Predictions & Errors

- ① Predictions from the model will have differences or errors.



Figure: Correct Fit

- ② **Overfitting** will occur when an excessive number of features are used than required.

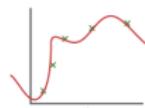
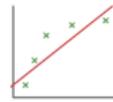


Figure: Overfit

- ③ **Underfitting** occurs when an insufficient number of features are used than required.



# Regularization

- ➊ 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.
- ➋ 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.
- ➌ A regression model that uses L1 regularization (or L1 Norm) technique is called Lasso Regression and model that uses L2 regularization (or L2 Norm) is called Ridge Regression.
- ➍ The key difference between these two is the penalty term.

# Regularization - Lasso & Ridge

- ① **Lasso Regression (Least Absolute Shrinkage and Selection Operator)** adds the absolute value of magnitude of coefficient as a penalty term to the loss function. Here the highlighted part represents L1 regularization element where  $\lambda$  is a constant.

$$L + \boxed{\lambda \sum_{j=1}^n |\theta_j|} \quad (10)$$

- ② **Ridge regression** adds the squared magnitude of coefficient as a penalty term to the loss function. Here the highlighted part represents L2 regularization element.

$$L + \boxed{\lambda \sum_{j=1}^n \theta_j^2} \quad (11)$$

- ③ Lasso shrinks the less important feature's coefficient to zero thus removing some features altogether. Hence, this works well for feature selection where we have a large number of features. However, the L1-norm does not have an analytical solution. L2-norm solutions can be computed efficiently.

# SoftMax

- ① In the training data the values corresponding to the features  $x_j$  can have numerical values of different magnitudes. For example, one feature may have values ranging between 0 ... 10, while another feature may have values ranging between 10,000 ... 100,000. Some of these values can also be negative and the components will obviously not add up to 1.
- ② The softmax function is a function that takes as input a vector of K real numbers, and normalizes it into a probability distribution consisting of K probabilities in the interval 0...1 with the components adding up to 1.

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad (12)$$

The feature values are normalized with the softmax function prior to inputting in a regression model.

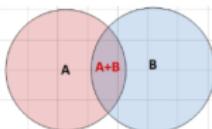
# One Hot Encoding

- 1 In the linear regression and the logistic regression models that we studied so far can only deal with numerical data. It cannot handle features that have text as values. In **one-hot encoding**, the string encoded variable is replaced with new variables of boolean type. For example, a particular feature, say color, can have "red", "green" and "blue" as possible values. This feature color, will be replaced by 3 features, red, blue and green which hold the values 1 or 0. We can then encode color, in terms of red, blue and green as follows:

red	green	blue
1	0	0
0	1	0
0	0	1

Note that if the boolean value of the feature is 1, it must be 0 for all the other features

# Conditional Probability - Naive Bayes



$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) = P(B|A) \times \frac{P(A)}{P(B)}$$

$$P(B|A) = P(A|B) \times \frac{P(B)}{P(A)}$$

(13)

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

# Naive Bayes

$$\begin{aligned}P(\text{Play} = \text{Yes}) &= 9/14 \\P(\text{Play} = \text{No}) &= 5/14 \\P(\text{Outlook} = \text{Sunny} \mid \text{Play} = \text{Yes}) &= 2/9 \\P(\text{Temperature} = \text{Cool} \mid \text{Play} = \text{Yes}) &= 3/9 \\P(\text{Humidity} = \text{High} \mid \text{Play} = \text{Yes}) &= 3/9 \\P(\text{Wind} = \text{Strong} \mid \text{Play} = \text{Yes}) &= 3/9 \\P(\text{Outlook} = \text{Sunny} \mid \text{Play} = \text{No}) &= 3/5 \\P(\text{Temperature} = \text{Cool} \mid \text{Play} = \text{No}) &= 1/5 \\P(\text{Humidity} = \text{High} \mid \text{Play} = \text{No}) &= 4/5 \\P(\text{Wind} = \text{Strong} \mid \text{Play} = \text{No}) &= 3/5\end{aligned}$$

Given,  $X = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$ , do we play?

$$\begin{aligned}P(\text{Yes} \mid X) &= P(\text{Sunny} \mid \text{Yes}) \times P(\text{Cool} \mid \text{Yes}) \times P(\text{High} \mid \text{Yes}) \\&\quad \times P(\text{Strong} \mid \text{Yes}) \times P(\text{Play} = \text{Yes}) = 0.0053\end{aligned}$$

$$\begin{aligned}P(\text{No} \mid x) &= P(\text{Sunny} \mid \text{No}) \times P(\text{Cool} \mid \text{No}) \times P(\text{High} \mid \text{No}) \\&\quad \times P(\text{Strong} \mid \text{No}) \times P(\text{Play} = \text{No}) = 0.0206\end{aligned}$$

Since  $P(\text{Yes} \mid X) < P(\text{No} \mid x)$ , we label  $X$  to be No.

# Latent Dirichlet Allocation (LDA)

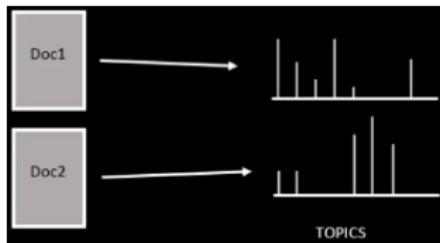
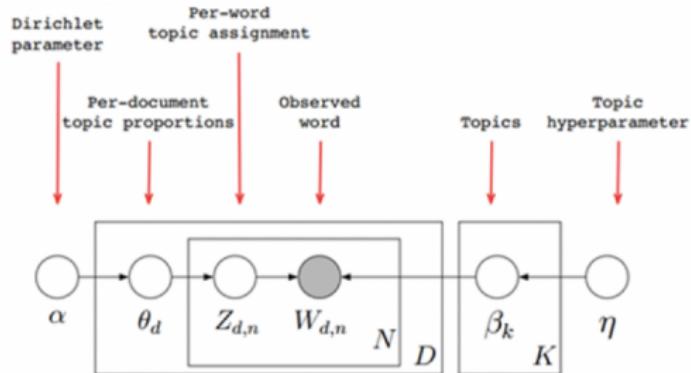


Figure: Topic Distribution

- ① LDA is an unsupervised machine learning algorithm. LDA uses a generative process that treats a document as bag of words and learns the topics in the document.
- ② The words are observable and the topics in the document are hidden and hence the term latent.
- ③ We are directly observing the words and not the topics, so the topics themselves are latent variables (along with the distributions themselves). Documents are probability distribution over latent topics and topics are probability distribution over words.

# Latent Dirichlet Allocation (LDA)



Each piece of the structure is a random variable.

Figure: Plate Notation

The plate notation is a concise way of visually representing the dependencies among the model parameters.  $\alpha$  is the parameter in the Dirichlet prior on the per document distribution. A high  $\alpha$  implies that a document is likely to contain most of the topics.  $\beta$  is the parameter in the Dirichlet prior on the per topic word distribution. A high  $\beta$  implies that a topic is likely to contain most of the words.

# LDA

- ①  $w$  represents a word
- ②  $W$  represents a document (a vector of words);  $w = w_1, w_2, \dots, w_N$
- ③  $\alpha$  is the parameter of the Dirichlet distribution;

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$$

- ④  $\theta$  is per document topic proportions
- ⑤  $Z$  is a vector of topics, where if the  $i$  th element of  $z$  is 1 then  $w$  draws from the  $i$  th topic  $D$  is the total set of documents
- ⑥ LDA assumes that new documents are generated by choosing a topic mixture for the document, (eg., 20% topic-1, 30%, topic-2 and 50% topic-3)) and generate words by first pick a topic based on the document's multinomial distribution and then pick a word based on the topics topic's multinomial distribution.

## LDA Steps

- ① The first step is to decide that there are  $k$  topics.
- ② Next, from a given document, assign each word to one of the  $k$  topics randomly.
- ③ Assume that all topic assignments except for the current word are correct. Calculate two ratios.
  - Proportion of words in the document  $d$  that are currently assigned to topic  $t$ .
  - Proportion of assignments to topic  $t$  over all documents that come from this word.
  - Assign the current word to the topic for which the probability  $p$  is maximum.

We continue repeating this process of re-assigning words until we reach a steady state. Mathematically, the LDA is expressed as follows:

$$p(\theta, z|w, \alpha, \beta) \propto p(\theta, z, w|\alpha, \beta) = p(w|z, \beta)p(z|\theta)p(\theta|\alpha) \quad (14)$$

Now, given a new document, how do we determine what topics does it consist of? We can score potential topics using the cross-entropy, i.e., KL distance.

# Information & Uncertainty

- ① In Claude Shannon's (MIT) information theory, one bit of information reduces the uncertainty by 2. Similarly, if 3 bits of information are sent, then the reduction in uncertainty by  $2^3$ , i.e., 8. This is intuitive. With 3 bits, there could be 8 possible values and so if a particular set of bits are transmitted, 8 possibilities are eliminated with 1 certainty.
- ② **Information Content** When the information is probabilistic, the self-information  $I_x$ , or Information Content of *measuring a random variable X as outcome x is defined as:*

$$I_x = \log \left( \frac{1}{p(x)} \right) = -\log(p(x)) \quad (15)$$

where  $p(x)$  is probability mass function.

- ③ **Shannon Entropy of the random variable X is defined as:**

$$H(X) = \sum_x -p(x)\log(p(x)) = E(I_x) \quad (16)$$

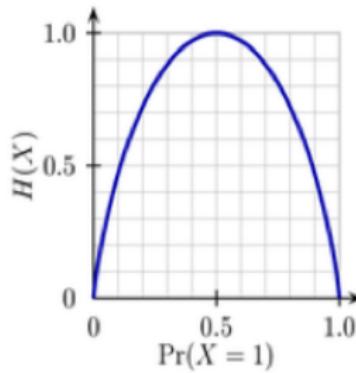
It is the *expected information content* of the measurement of  $X$ .

## Example

- ① Suppose S has 25 examples, 15 positive and 10 negatives [15+, 10-]. Then the entropy of S relative to this classification is computed as:

$$E(S) = -(15/25)\log_2(15/25) - (10/25)\log_2(10/25) = 0.67$$

In a 2 class system, The entropy is 0 if the outcome is *certain*. The entropy is maximum if we have no knowledge of the system, i.e., when any outcome is equally possible.



# Cross Entropy - Kullback–Leibler (KL) divergence

- ① Cross-Entropy is defined as:

$$H(p, q) = - \sum_i p(i) \log_2 q(i) \quad (17)$$

where  $p$  is the true distribution and  $q$  is the predicted distribution. If the predictions are perfect, then the cross-entropy is same as the entropy. If the prediction differs, then there is a divergence which is known as *Kullback – Leibler (KL) divergence*. Hence,

$$\text{Cross Entropy} = \text{Entropy} + \text{KL Divergence}$$

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

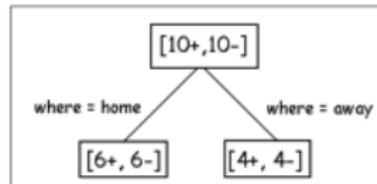
Hence, if the predicted distribution is closer to true distribution when KL divergence is low.

# Decision Tree

- ① Decision Tree is a supervised machine learning algorithm. In decision trees, the goal is to tidy the data. You try to separate your data and group the samples together in the classes they belong to. You know their label since you construct the trees from the training set. You maximize the purity of the groups as much as possible each time you create a new node of the tree (meaning you cut your set in two). Of course, at the end of the tree, you want to have a clear answer.
- 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 Before and after the decision, the sets are different and have different sizes. Still, entropy can be compared between these sets, using a weighted sum.

## Example

Consider the following table for win/loss of soccer games at home and away.



The entropy is:

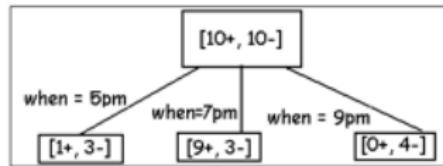
$$H\left(\frac{6}{12}, \frac{6}{12}\right) = -\frac{6}{12} \log_2\left(\frac{6}{12}\right) - \frac{6}{12} \log_2\left(\frac{6}{12}\right) = 0.69$$

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

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

## Example

Partitioning by when, we have:



$$H(5pm) = H\left(\frac{1}{4}, \frac{3}{4}\right) = -\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right) = 0.56$$

$$H(7pm) = H\left(\frac{9}{12}, \frac{3}{12}\right) = -\frac{9}{12} \log_2\left(\frac{9}{12}\right) - \frac{3}{12} \log_2\left(\frac{3}{12}\right) = 0.56$$

$$H(9pm) = H\left(\frac{0}{4}, \frac{4}{4}\right) = -\frac{0}{4} \log_2\left(\frac{0}{4}\right) - \frac{4}{4} \log_2\left(\frac{4}{4}\right) = 0.00$$

Entropy after partition is:

$$H = -\frac{4}{20} H\left(\frac{1}{4}, \frac{3}{4}\right) - \frac{12}{20} H\left(\frac{9}{12}, \frac{3}{12}\right) - \frac{4}{20} H\left(\frac{0}{4}, \frac{4}{4}\right) = 0.45$$

Hence, the Information gain is  $0.69 - 0.45 = 0.24$ .

# Decision Trees

- ① We can apply the same approach for other columns to find the most dominant factor on decision. This way we construct the decision tree that has the least entropy.
- ② Decision trees can generate understandable rules and classify without much computation. It can handle both continuous and categorical variables. Also, it provides a clear indication of which fields are most important for prediction or classification.
- ③ However, it is not suitable for prediction of continuous attributes. It performs poorly with many class and small data. It is computationally expensive to train.

# Random Forest

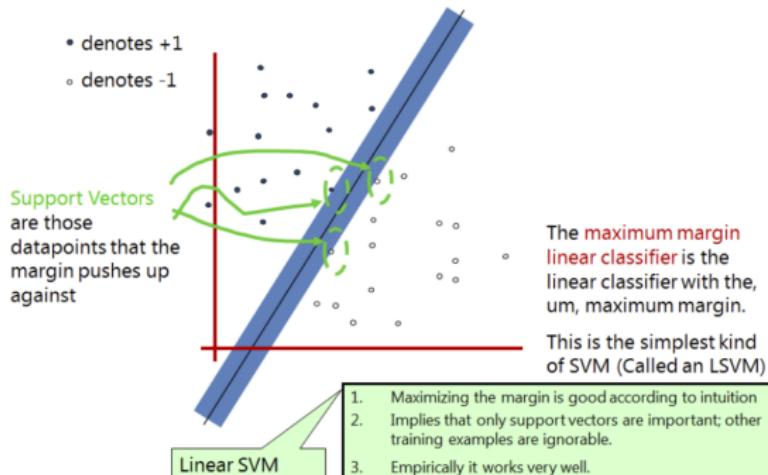
- ① Random forest builds multiple decision trees and merges them together to get a more accurate and stable prediction. The fundamental idea behind a random forest is to combine many decision trees into a single model. Individually, predictions made by decision trees (or humans) may not be accurate, but combined together, the predictions will be closer to the mark on average.
- ② Why the name 'random forest?' Well, much as people might rely on different sources to make a prediction, each decision tree in the forest considers a random subset of features when forming questions and only has access to a random set of the training data points.
- ③ This increases diversity in the forest leading to more robust overall predictions and the name 'random forest.' When it comes time to make a prediction, the random forest takes an average of all the individual decision tree estimates. (This is the case for a regression task, such as our problem where we are predicting a continuous value of temperature.)

# Clustering

- ① Clustering is an unsupervised learning technique
- ② Create the first K initial clusters from the dataset by choosing K rows of data randomly from the dataset.
- ③ Each record is assigned to the nearest cluster (the cluster which it is most similar to) using a measure of distance or similarity.
- ④ Re-calculate the arithmetic mean for each cluster.
- ⑤ Repeat step 2.
- ⑥ Stop when arithmetic mean for the clusters stabilize.

# Support Vector Machines (SVM)

- ① SVM is a supervised machine learning algorithm. SVM performs classification tasks by constructing hyperplanes in a multidimensional space that separates cases of different class labels.
- ② SVM supports both regression and classification tasks and can handle multiple continuous and categorical variables.



# SVM



Figure: Plane

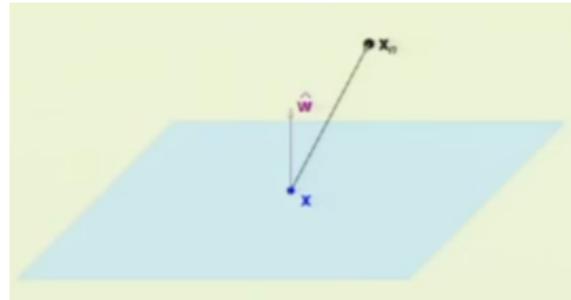


Figure: Normal to the Plane

# SVM

- ① The equation of a plane is given by:

$$w^T x + b = 0 \quad (18)$$

For two points  $x'$  and  $x''$  on the plane, we have:

$$\begin{aligned} w^T x' + b &= 0 \\ w^T x'' + b &= 0 \\ w^T(x' - x'') &= 0 \end{aligned}$$

The normal distance  $d$  of a point  $x_n$  from  $x$  is given by:

$$d = \|\hat{w}(x_n - x)\| = \frac{1}{\|w\|} (w^T x_n + b - w^T x - b) = \frac{2}{\|w\|} \quad (19)$$

# SVM

Our goal is to maximize the distance which is the same as minimizing  $\frac{1}{2}w^T w$  subject to the condition:

$$wx_i + b \geq 1 \text{ for } y_i = +1 \quad (20)$$

$$wx_i + b \leq -1 \text{ for } y_i = -1 \quad (21)$$

i.e., for all i:

$$y_i(wx_i + b) \geq 1 \quad (22)$$

We will now use the Lagrange formulation to minimize:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2}w^T w - \sum_{n=1}^N \alpha_n(y_n(w^T x_n + b) - 1) \quad (23)$$

with respect to w and b and maximize with respect to each  $\alpha_n > 0$

# SVM

$$\nabla \mathcal{L} = w - \sum_{n=1}^N \alpha_n y_n x_n = 0 \quad (24)$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{n=1}^N \alpha_n y_n = 0 \quad (25)$$

$$w = \sum_{n=1}^N \alpha_n y_n x_n \quad (26)$$

$$\sum_{n=1}^N \alpha_n y_n = 0 \quad (27)$$

$$\boxed{\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} - \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m} \quad (28)$$

maximize above eqn w.r.t  $\alpha$  subject to  $\alpha_x > 0$  for  $n = 1, 2, \dots, N$ .

# SVM

$$\max_{\alpha} \sum_{n=1}^N \alpha_n - \frac{1}{2} - \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m \quad (29)$$

$$\min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m - \sum_{n=1}^N \alpha_n \quad (30)$$

$$\min_{\alpha} \frac{1}{2} \alpha^T \begin{vmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \cdots & y_1 y_N x_1^T x_N \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 & \cdots & y_2 y_N x_2^T x_N \\ \cdots & \cdots & \cdots & \cdots \\ y_N y_1 x_N^T x_1 & y_N y_2 x_N^T x_2 & \cdots & y_N y_N x_N^T x_N \end{vmatrix} \alpha + (-1^T) \alpha \quad (31)$$

subject to  $y^T \alpha = 0; 0 < \alpha \leq \infty$ .

$$\boxed{\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha} \quad (32)$$

# SVM

subject to Karush Kuhn Tucker (KKT) condition:  $y^T \alpha = 0; \alpha \geq 0$ . We then compute w:

$$w = \sum_{n=1}^N \alpha_n y_n x_n \quad (33)$$

and solve for b:

$$y_n(wx_n + b) = 1 \quad (34)$$

# Non Linearity

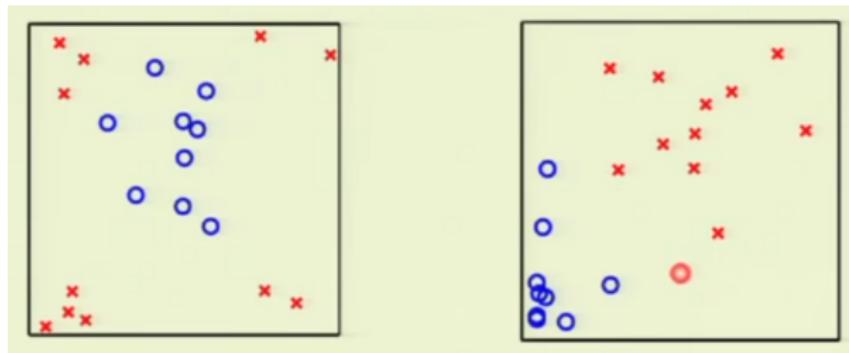


Figure: Transform to linear space

For non-linear systems, we transform to z space:

$$\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} - \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m z_n^T z_m \quad (35)$$

and we will have support vectors in the transformed space.

# Deep Neural Net

- ① Deep Learning Success Stories - Image Recognition, Speech Comprehension, Chatbot
- ② 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.
- ③ DNNs are greedily built by stacking restricted Boltzmann machines, and Convolutional Neural Networks, which exploit the local dependency of visual information, and have demonstrated record-setting results on many important applications.
- ④ Neural networks are powerful learning models that achieve state-of-the-art results in a wide range of supervised and unsupervised machine learning tasks for classification & regression.
- ⑤ Simple linear models tend to under-fit, and often under utilize computing resources with 8-GPU Hydra.

# Deep Neural Net

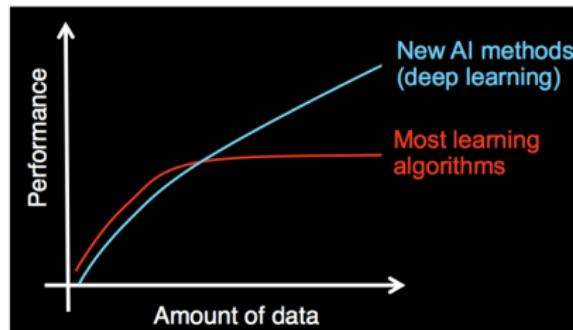
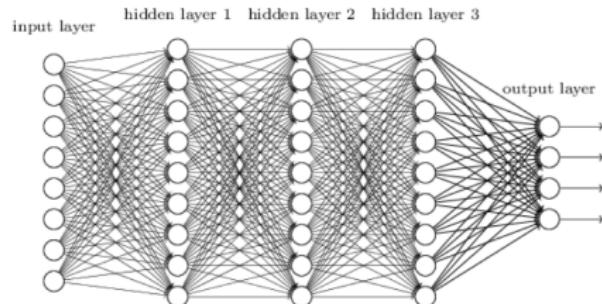
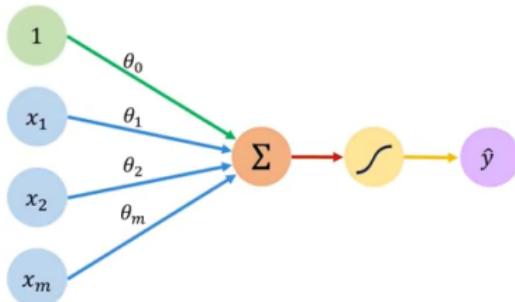


Figure: Deep Learning

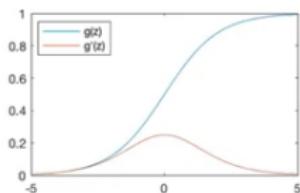
# Deep Neural Net



Linear combination of inputs  
Output  
 $\hat{y} = g \left( \theta_0 + \sum_{i=1}^m x_i \theta_i \right)$   
Non-linear activation function  
Bias

Inputs    Weights    Sum    Non-Linearity    Output

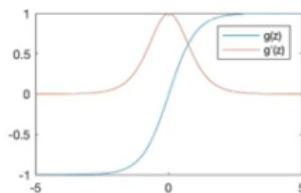
Sigmoid Function



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

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

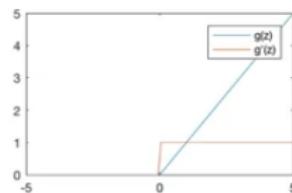
Hyperbolic Tangent



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

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

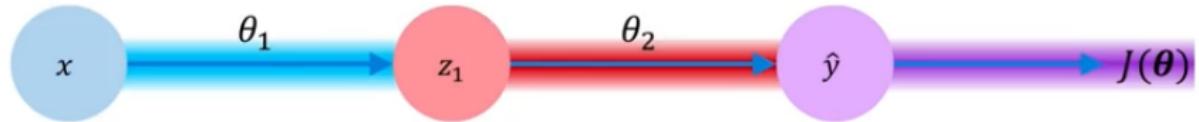
Rectified Linear Unit (ReLU)



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

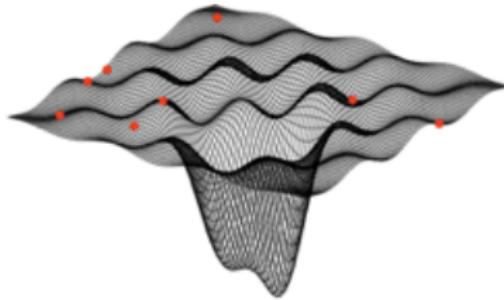
$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Back Propagation



$$\frac{\partial J(\theta)}{\partial \theta_1} = \underbrace{\frac{\partial J(\theta)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial \theta_1}}_{\text{cyan}}$$

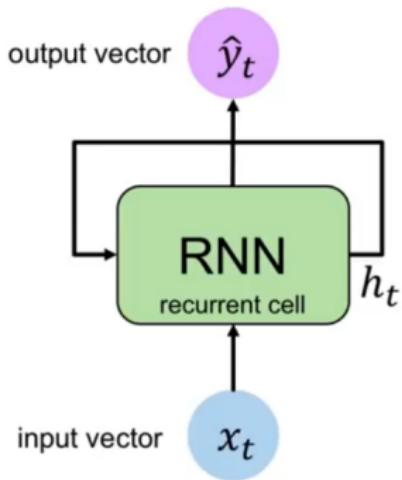
Figure: Loss Minimization



# Recurrent Neural Network (RNN)

- ① Dialogue systems, self-driving cars, robotic surgery, speech, etc. among others require an explicit model of sequentiality or time – a combination of classifiers or regressors cannot provide these functionalities.
- ② SVM, logistic regression, feedforward networks have proved very useful without explicitly modeling time. But the assumption of independence precludes modeling long range dependencies.
- ③ Frames from video, snippets of audio, and words pulled from sentences – the independence assumption fails.
- ④ RNNs are connectionist models with the ability to selectively pass information across sequence steps while processing sequential data one element at a time. They can model input and/or output consisting of sequences of elements that are not independent. In other words can handle variable length sequences, track long term dependencies, maintain information about order and share parameters across the sequence.

# RNN



Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(h_{t-1}, x_t)$$

cell state      function parameterized by W  
old state

Figure: Same function and set of parameters are used in every step

Source - MIT 6.S191- Recurrent Neural Networks

# RNN

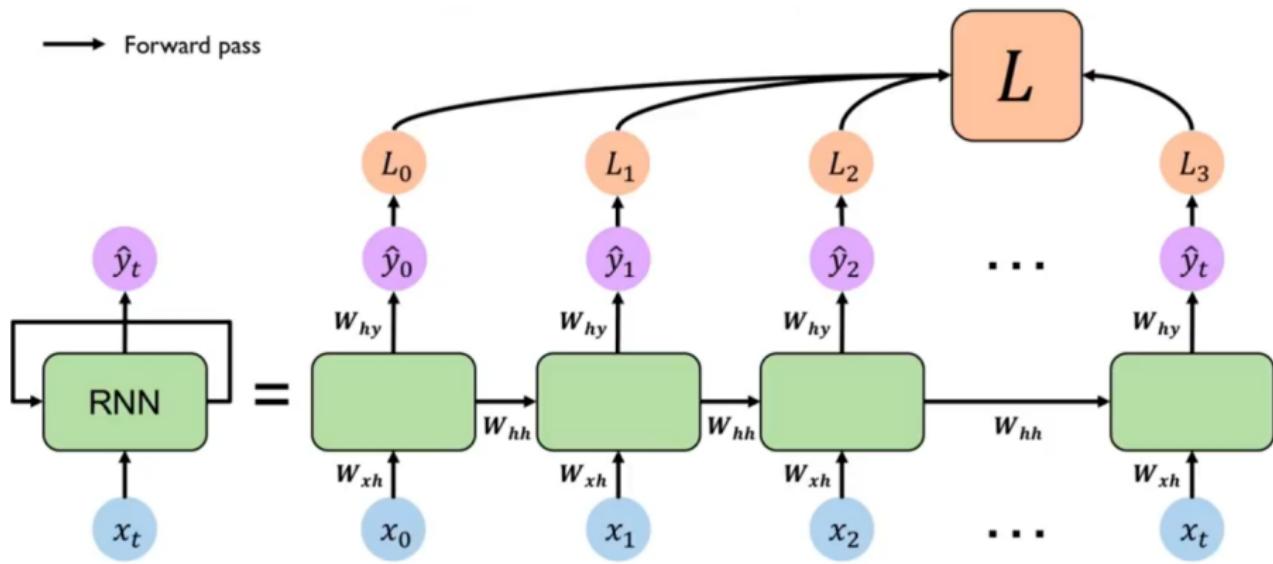


Figure: Computational Graph

# RNN

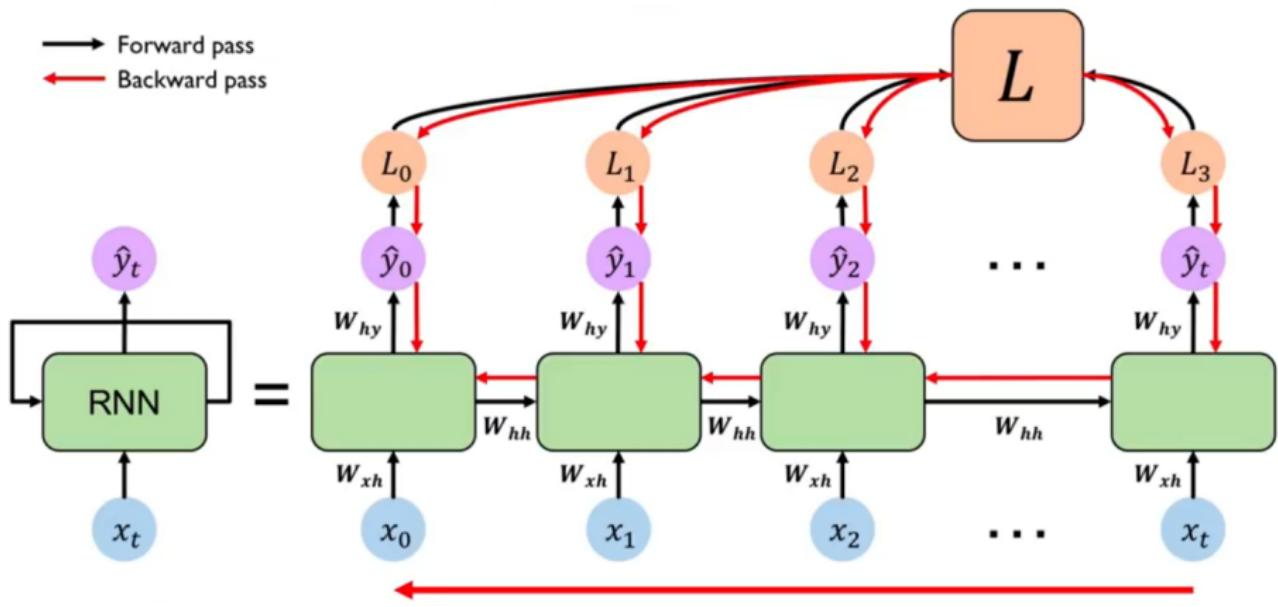


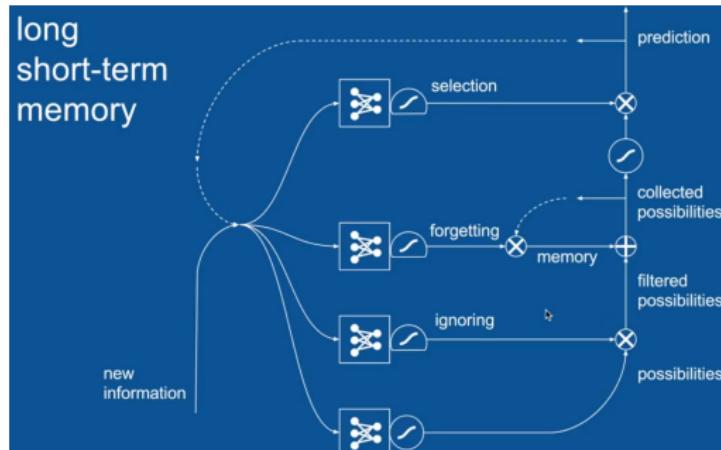
Figure: Backpropagation over time

# The Problem of Vanishing / Exploding Gradients

- ① Learning with recurrent networks is challenging due to the difficulty of learning long-range dependencies, [Bengio et al. 1994], [Hochreiter et al, 2001].
- ② Problems of vanishing and exploding gradients occur when backpropagating errors across many time steps.
- ③ Which of the two phenomena occurs depends on whether the weight of the recurrent edge  $|w_{jj}| > 1$  or  $|w_{jj}| < 1$  and on the activation function in the hidden node. For a sigmoid activation function, the vanishing gradient problem is more pressing, but with a rectified linear unit  $\max(0, x)$ , it is easier to imagine the exploding gradient.

# Long Short Term Memory (LSTM) - Hochreiter & Schmidhuber

- ① Maintain a separate cell state
- ② Use gates
  - ① Forget
  - ② Selectively update
  - ③ Output
- ③ Backpropagation does not require matrix multiplication



# Reinforcement Learning

- ①  $\hat{P}$  is sampled and is NOT predicted.
- ②  $P_i$  is predicted and is a function of the weights and is differentiable with respect to the weights.
- ③ Cross Entropy Loss has to be minimized, this will lead to determining the weights.
- ④ Policy Gradients

- ① Discounted Rewards for the  $i^{th}$  move.  $R_i$  reduces geometrically over previous steps.

$$\text{loss}_i = -R_i \log(P_i^{\text{chosen}}) \quad (36)$$

- ② Normalized Rewards across a batch of moves

$$R_i = \frac{R_i - \bar{R}}{\text{stdev}(R)} \quad (37)$$

# Generative Adversarial Network (GAN)

- ① GAN is a neural network architecture that allows neural networks to generate data and learn a probability distribution of a data set by pitting two neural networks against each other.
- ② The generator acts like a painting forger. It tries to create images that look like real. The discriminator tries to detect the fakes.
- ③ The generator keeps getting better at making fakes, while the discriminator gets better at detecting fakes. The 2 models compete, and after many iterations, the generator creates images indistinguishable from the real data set.

