



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

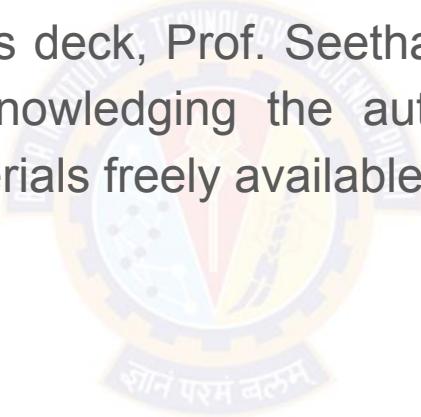
Deep Neural Network

AIML Module 1

Seetha Parameswaran

BITS Pilani

The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Course Logistics



What we Learn.... (Module Structure)

1. Fundamentals of Neural Network
2. Multilayer Perceptron
3. Deep Feedforward Neural Network
4. Improve the DNN performance by Optimization and Regularization
5. Convolutional Neural Networks
6. Sequence Models
7. Attention Mechanism
8. Neural Network search
9. Time series Modelling and Forecasting
10. Other Learning Techniques

Text book

- Dive into Deep Learning by Aston Zhang, Zack C. Lipton, Mu Li, Alex J. Smola. https://d2l.ai/chapter_introduction/index.html



Course Logistics

- Refer Canvas for the following

- Handout
- Schedule for Webinar
- Schedule of Quiz, and Assignments.
- Evaluation scheme
- Session Slide Deck
- Demo Lab Sheets
- Quiz-I, Quiz-II
- Assignment-I, Assignment-II
- Sample QPs



- Lecture Recordings

- Available on Microsoft teams

Honour Code

All submissions for graded components must be the result of your original effort. It is strictly prohibited to copy and paste verbatim from any sources, whether online or from your peers. The use of unauthorized sources or materials, as well as collusion or unauthorized collaboration to gain an unfair advantage, is also strictly prohibited. Please note that we will not distinguish between the person sharing their resources and the one receiving them for plagiarism, and the consequences will apply to both parties equally.

In cases where suspicious circumstances arise, such as identical verbatim answers or a significant overlap of unreasonable similarities in a set of submissions, will be investigated, and severe punishments will be imposed on all those found guilty of plagiarism.

In case of Queries regarding the course....

Step 1: Post in the discussion forum.

- Read through the existing post and if you find any topic similar to your concern, add on to the existing discussion.
- Avoid duplication of queries or issues.

Step 2: Email to the IC at seetha.p@pilani.bits-pilani.ac.in if the query or issue is not resolved within 1 weeks time. Turn around for response to the email is 48hrs.

- In the subject pl mention the phrase "DNN" clearly.
- Use BITS email id for correspondence.

PATIENCE is highly APPRECIATED :)

What is Deep Learning?



Definitions

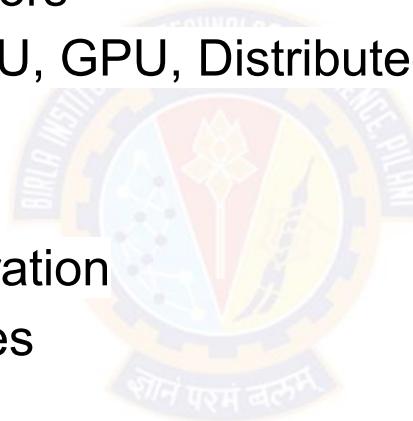
- Deep Learning is a type of **machine learning** based on **artificial neural networks** in which multiple layers of processing are used to extract progressively higher level features from data.
- Deep learning is a method in artificial intelligence (AI) that teaches computers to process data in a way that is **inspired by the human brain**.
- Deep learning is a machine learning technique that teaches computers to do what comes naturally to humans: **learn by example**.
- Deep learning is a subset of machine learning, which is essentially a neural network with **three or more layers**.
- Deep Learning gets its name from the fact that we add more "**Layers**" to learn from the data.

Why Deep Learning?



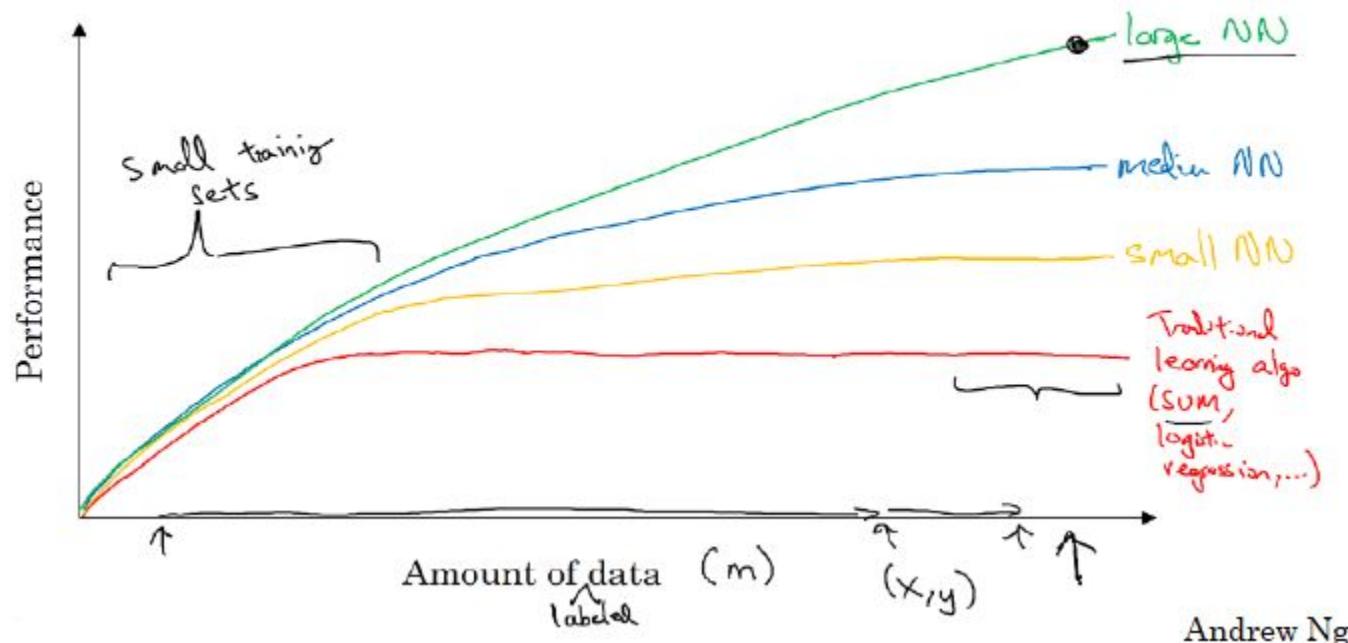
Why Deep Learning?

- Large amounts of data
- Lots and lots of unstructured data like images, text, audio, video
- Cheap, high-quality sensors
- Cheap computation - CPU, GPU, Distributed clusters
- Cheap data storage
- Learn by examples
- Automated feature generation
- Better learning capabilities
- Scalability
- Advance analytics can be applied

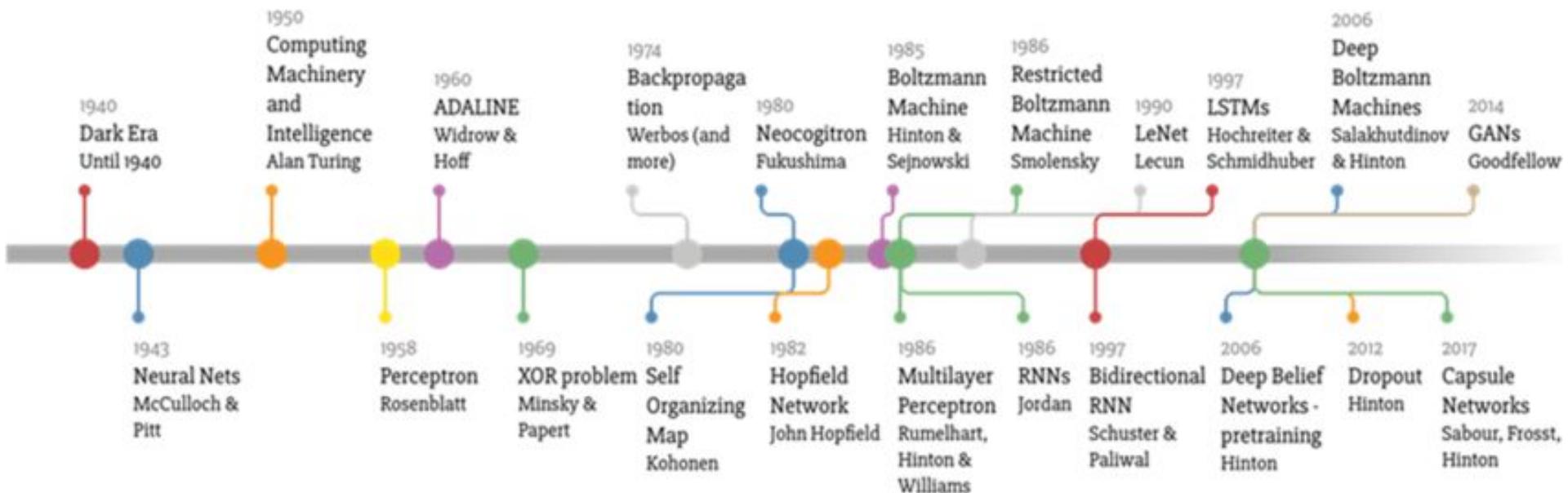


Why deep learning?

Scale drives deep learning progress



Deep Learning Timeline



Applications of Deep Learning



Breakthroughs with Neural Networks

TECHNEWSWORLD

EMERGING TECH

SEARCH

Computing Internet IT Mobile Tech Reviews Security Technology Tech Blog

Reader Service

Microsoft AI Beats Humans at Speech Recognition

By Richard Adhikari

Oct 20, 2016 11:40 AM PT

Print

Email



Image: Adobe Stock

Microsoft's Artificial Intelligence and Research Unit earlier this week reported that its speech recognition technology had surpassed the performance of human transcriptionists.



[Most Popular](#) [Newsletters](#) [News Alerts](#)

How do you feel about Black Friday and Cyber Monday?

- They're great -- I get a lot of bargains!
- The deals are too spread out -- I'd prefer just one day.
- They're a fun way to kick off the holiday season.
- I don't like the commercialization of Thanksgiving Day.
- They're crucial for the retail industry and the economy.
- The deals typically aren't that good.

[Vote to See Results](#)

E-Commerce Times

[Black Friday Shoppers Hungry for New Experiences, New Tech](#)

[Pay TV's Newest Innovation: Giving Users Control](#)

[Apple Celebrates Itself in \\$300 Coffee Table Tome](#)

[AWS Enjoys Top Perch in IaaS, PaaS Markets](#)

[US Comptroller Gears Up for Blockchain and](#)

Breakthroughs with Neural Networks

The Keyword Latest Stories Product News Topics

TRANSLATE NOV 15, 2016

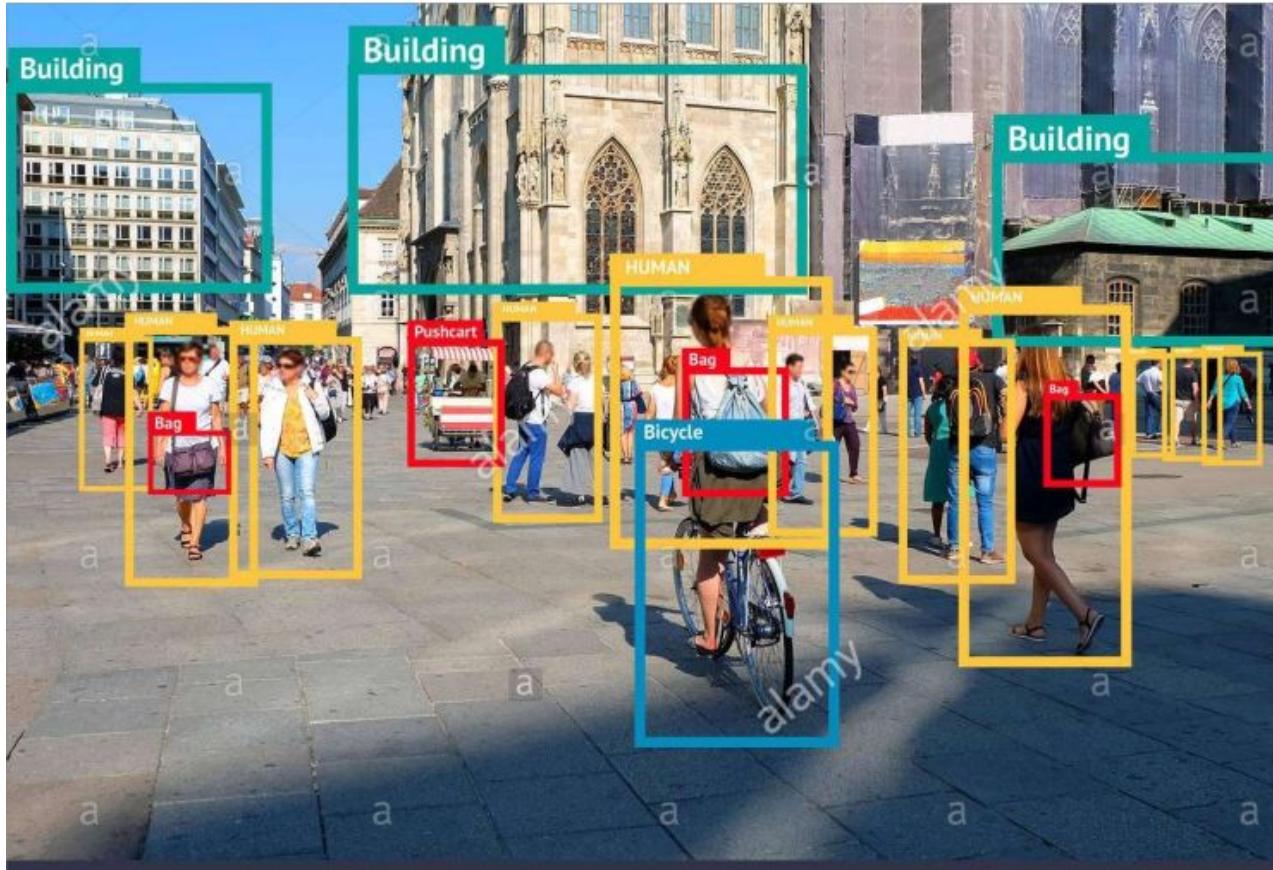
Found in translation: More accurate, fluent sentences in Google Translate

Barak Turovsky
PRODUCT LEAD, GOOGLE TRANSLATE

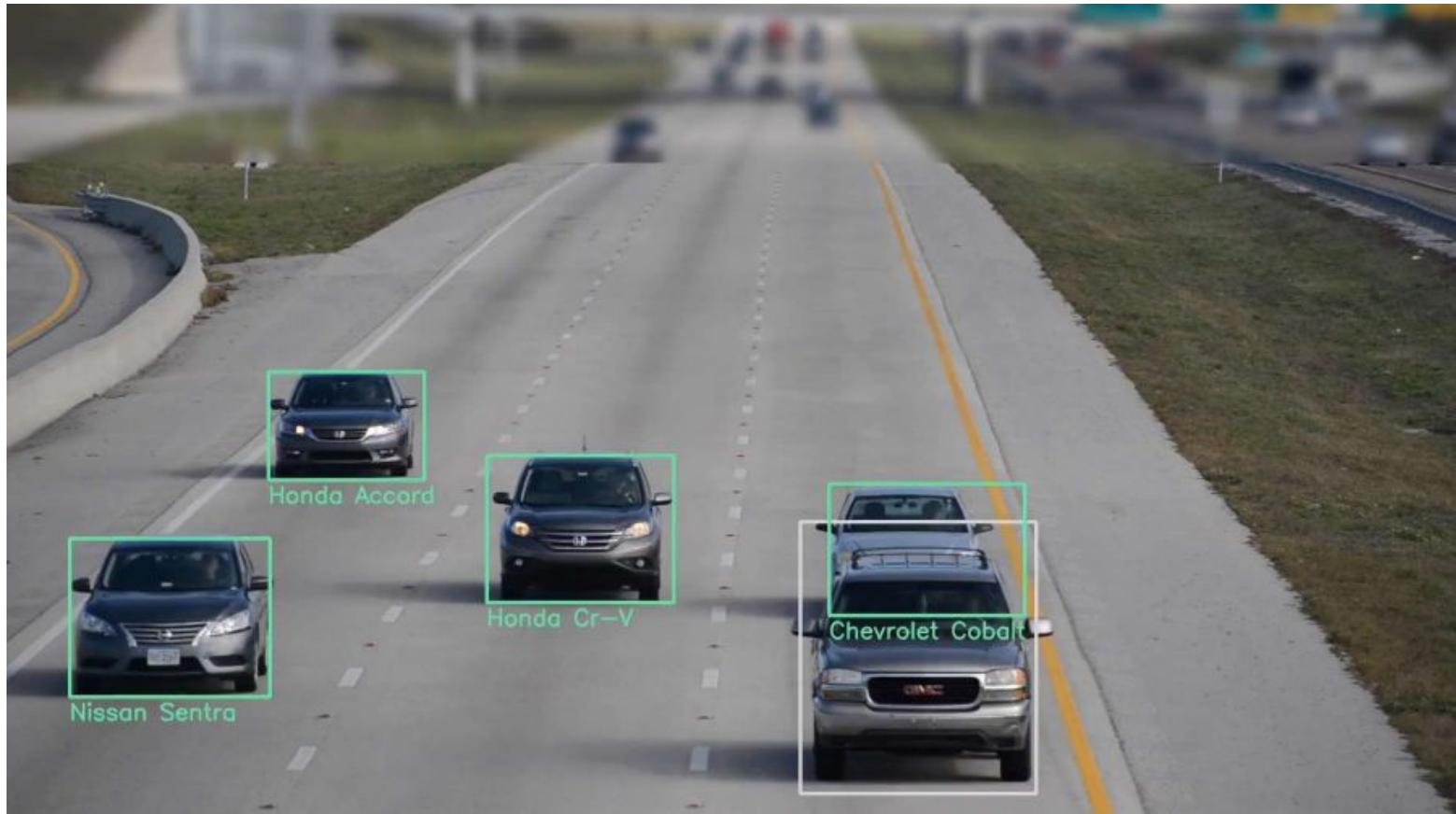
In 10 years, Google Translate has gone from supporting just a few languages to 103, connecting strangers, reaching across language barriers and even helping



Breakthroughs with Neural Networks



Breakthroughs with Neural Networks



Breakthroughs with Neural Networks

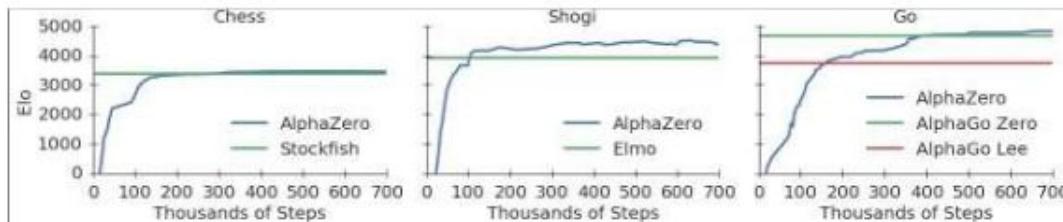


Figure 1: Training *AlphaZero* for 700,000 steps. Elo ratings were computed from evaluation games between different players when given one second per move. **a** Performance of *AlphaZero* in chess, compared to 2016 TCEC world-champion program *Stockfish*. **b** Performance of *AlphaZero* in shogi, compared to 2017 CSA world-champion program *Elmo*. **c** Performance of *AlphaZero* in Go, compared to *AlphaGo Lee* and *AlphaGo Zero* (20 block / 3 day) (29).

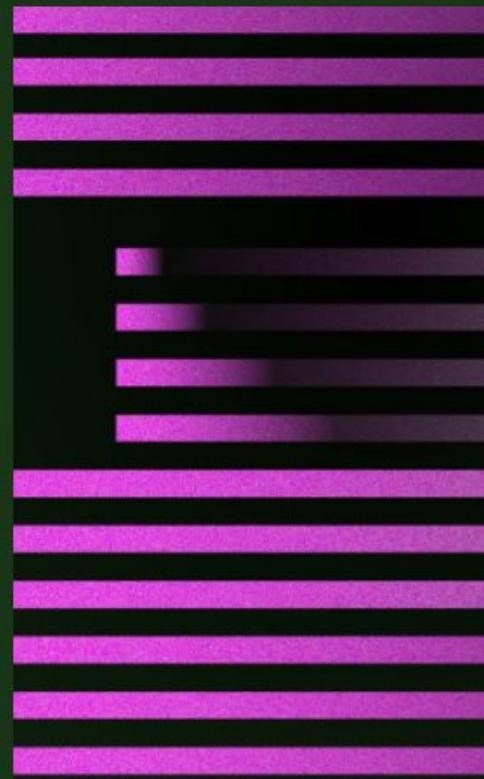
Breakthroughs with Neural Networks

Introducing ChatGPT

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

[Try ChatGPT ↗](#)

[Read about ChatGPT Plus](#)



Many more applications....

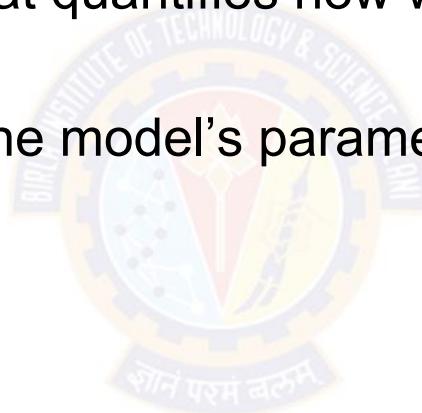
- Write a program that predicts tomorrow's weather given geographic information, satellite images, and a trailing window of past weather.
- Write a program that takes in a question, expressed in free-form text, and answers it correctly.
- Write a program that given an image can identify all the people it contains, drawing outlines around each.
- Write a program that presents users with products that they are likely to enjoy but unlikely, in the natural course of browsing, to encounter.

Key components of DL problem



Core components of DL problem

1. The **data** that we can learn from.
2. A **model** of how to transform the data.
3. An **objective function** that quantifies how well (or badly) the model is doing.
4. An **algorithm** to adjust the model's parameters to optimize the objective function.



Data

- Collection of examples.
- Data has to be converted to an useful and a suitable numerical **representation**.
- Each example (or data point, data instance, sample) typically consists of a set of attributes called **features** (or covariates), from which the model must make its predictions.
- In the supervised learning problems, the attribute to be predicted is designated as the **label** (or target).
- We need **right** data.

Data

- **Dimensionality** of data

- Each example has the same number of numerical values. This data consist of fixed-length vectors. Eg: Image
- The constant length of the vectors as the dimensionality of the data.
- Text data has varying-length data.

Model

- Model denotes the **computational machinery** for ingesting data of one type, and spitting out predictions of a possibly different type.
- Deep learning models consist of many successive transformations of the data that are chained together top to bottom, thus the name deep learning.

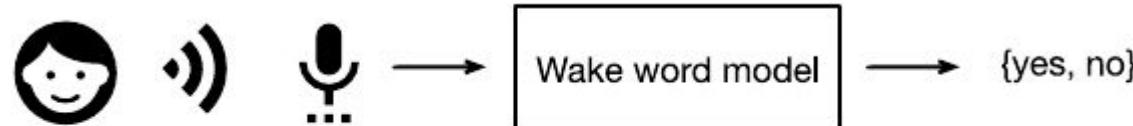


Fig. 1.1.1: Identify a wake word.

Objective Function

- **Learning means improving at some task over time.**
- A formal mathematical system of learning machines is defined using formal measures of how good (or bad) the models are. These formal measures are called as **objective functions**.
- By convention, objective functions are defined so that lower is better.
- Because lower is better, these functions are sometimes called **loss functions**.

Loss Functions

- To predict numerical values (regression), the most common loss function is squared error.
- For classification, the most common objective is to minimize error rate, i.e., the fraction of examples on which our predictions disagree with the ground truth.

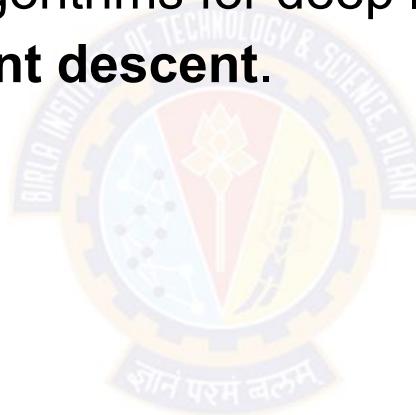


Loss Functions

- Loss function is defined **with respect to the model's parameters** and **depends upon the dataset**.
- We learn the best values of our model's parameters by **minimizing the loss** incurred on a set consisting of some number of examples collected for training. However, doing well on the training data does not guarantee that we will do well on unseen data. I.e **Model has to generalize better**.
- When a model performs well on the training set but fails to generalize to unseen data, we say that it is **overfitting**.

Optimization Algorithms

- Optimization Algorithm is an algorithm capable of **searching for the best possible parameters for minimizing the loss function.**
- Popular optimization algorithms for deep learning are based on an approach called **gradient descent**.



Kinds of Machine Learning Problems



Learning Problems

1. Supervised Learning
2. Unsupervised Learning
 - o Situations where we feed the model a giant dataset containing only the features
 - o Type and number of questions we could ask is limited only by our creativity.
 - o Eg: Clustering, Generate synthetic data (uses GAN)
 - o Self-supervised learning leverages unlabeled data to provide supervision in training, such as by predicting some withheld part of the data using other parts.
 - o Eg: Fill in the blanks
3. Reinforcement Learning
 - o Develop an agent that interacts with an environment and takes actions over a series of time steps.

Supervised Learning

- Task of predicting labels given input features.
- Each **feature–label pair** is called an example.
- Goal is to produce **a model that maps any input to a label prediction**.
- The supervision comes into play because for choosing the parameters, we (the supervisors) provide the model with a dataset consisting of labeled examples, where each example is matched with additional probability of a label given input features.

Supervised Learning

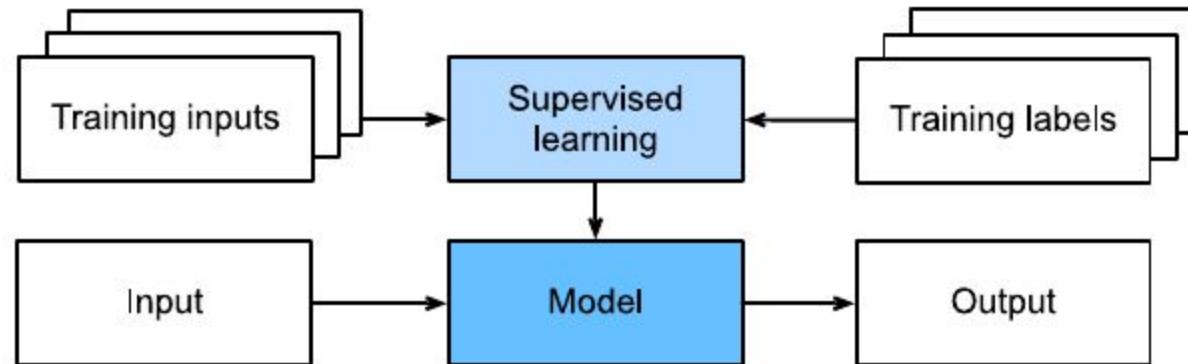
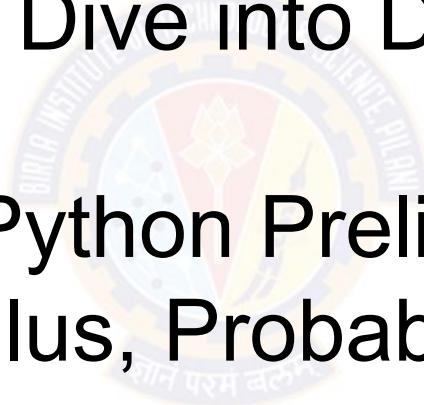


Fig. 1.3.1: Supervised learning.

Supervised Learning Tasks

- Regression
 - how much or how many question
- Classification
 - Binary classification
 - Multi-class classification
 - Multi-label classification
- Tagging
 - Multi-label classification
- Recommender systems
- Sequence Learning
 - Tagging and Parsing
 - Automatic speech recognition
 - Text to speech
 - Machine translation





Reading from TB Dive into Deep Learning

- Chapter 1
- Chapter 2 for Python Prelims, Linear Algebra, Calculus, Probability

Next Session: What is Neural Network?





BITS Pilani
Pilani | Dubai | Goa | Hyderabad

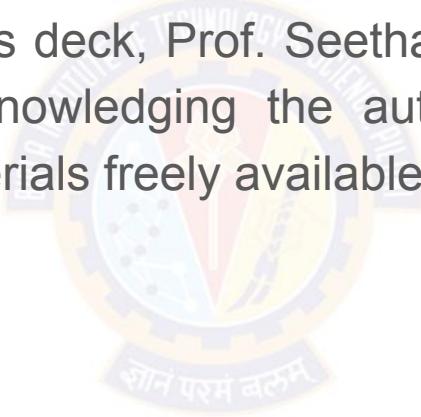
Deep Learning

AIML Module 1

Seetha Parameswaran

BITS Pilani

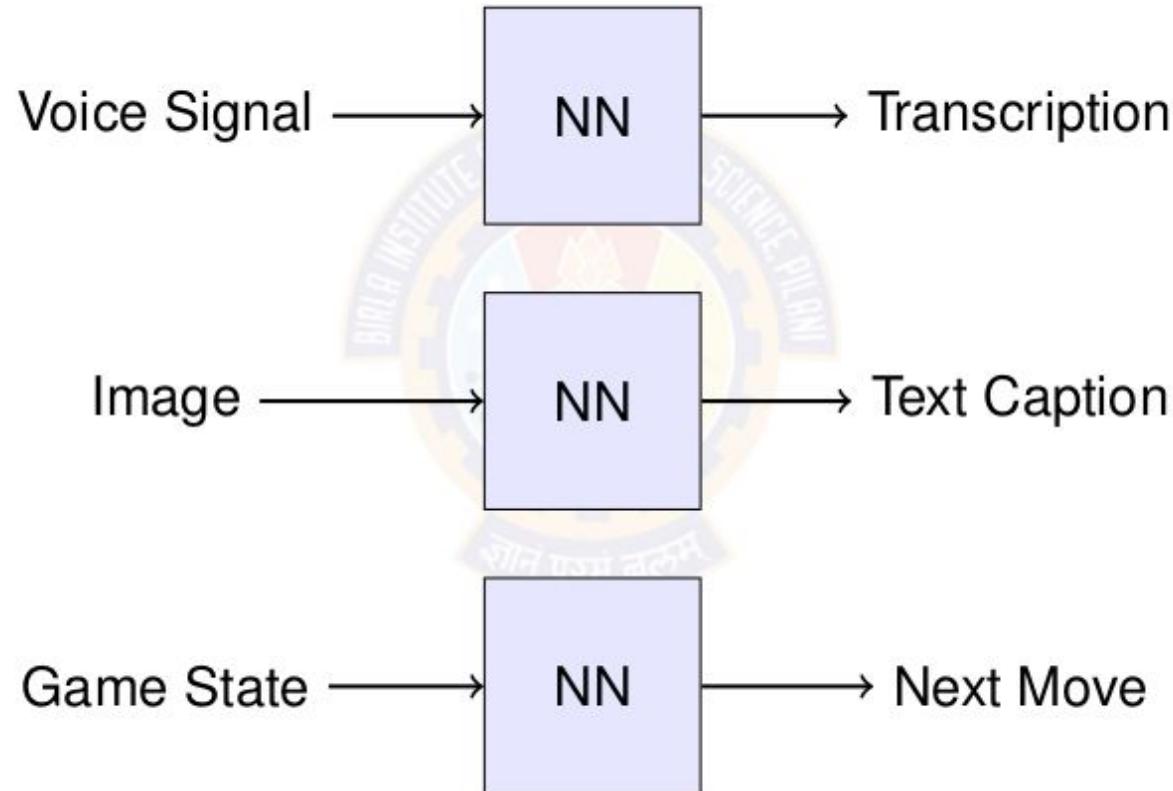
The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Artificial Neural Network



What are Neural Networks?



What are Neural Networks?

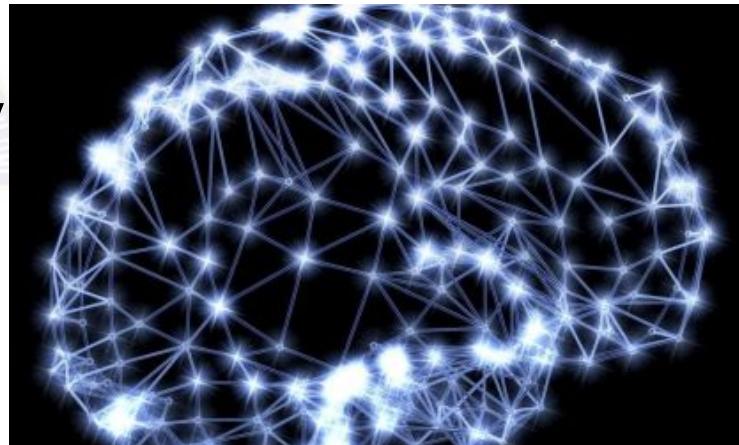
- It begins with human brain.



- Humans learn, solve problems, recognize patterns, create, think deeply about something, meditate and many many more.....
- Humans learn through **association**. [Refer to Associationism for more details.]

Observation: The Brain

- The brain is a mass of interconnected neurons.
- Number of neurons is approximately 10^{10} .
- Connections per neuron is approximately 10^4 to 10^5 .
- Neuron switching time is approximately 0.001 second.
- Scene recognition time is 1 second.
- 100 inference steps doesn't seem like enough. Lot of parallel computation.

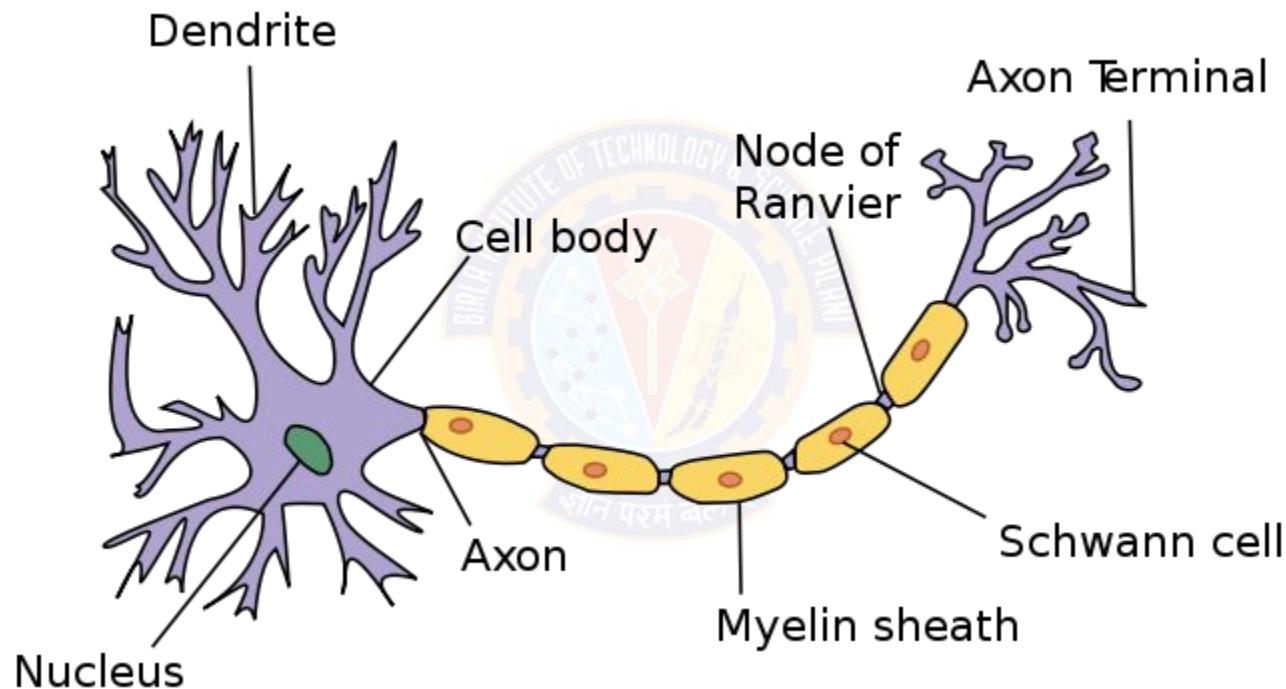


Brain: Interconnected Neurons

- Many neurons connect in to each neuron.
- Each neuron connects out to many neurons.

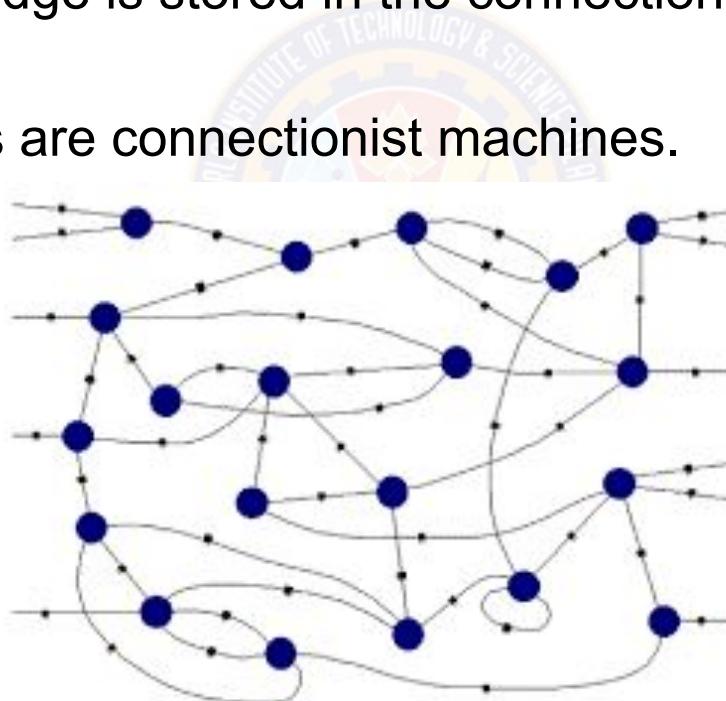


Biological Neuron



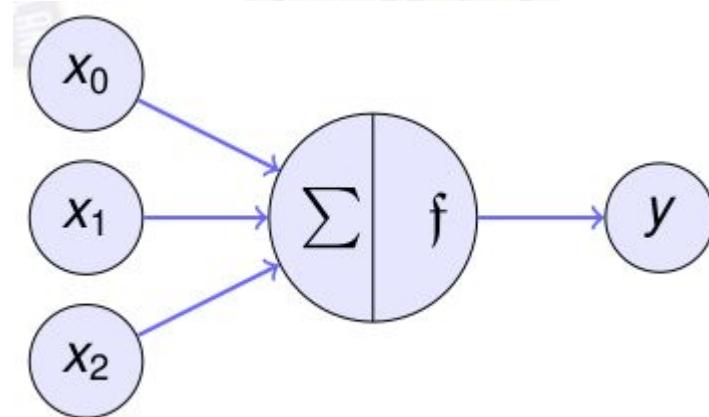
Connectionist Machines

- Network of processing elements, called artificial neural unit.
- The neurons are interconnected to form a network.
- All world knowledge is stored in the connections between these elements.
- Neural networks are connectionist machines.



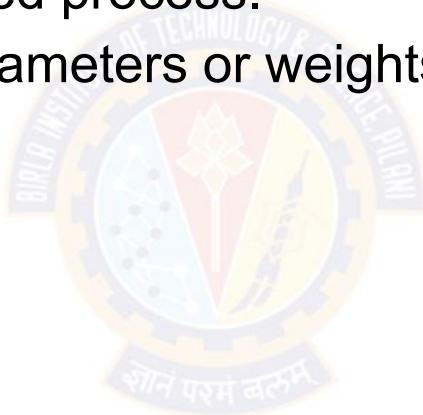
What are Artificial Neurons?

- Neuron is a processing element inspired by how the brain works.
- Similar to biological neuron, each artificial neuron will do some computation. Each neuron is interconnected to other neurons.
- Similar to brain, the interconnections between neurons store the knowledge it learns. The knowledge is stored as parameters.



Properties of Artificial Neural Nets (ANNs)

- Many neuron-like threshold switching units.
- Many weighted interconnections among units.
- Highly parallel, distributed process.
- Emphasis on tuning parameters or weights automatically.



When to consider Neural Networks?

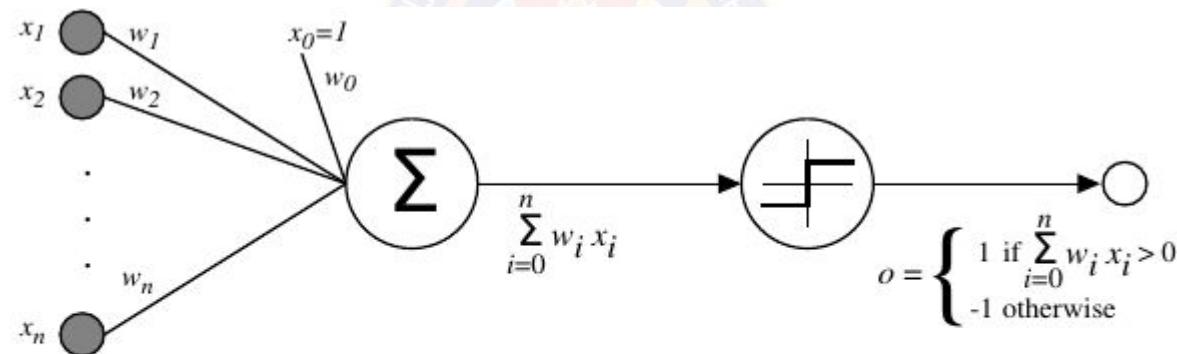
- Input is high-dimensional discrete or real-valued (e.g. raw sensor input).
- Possibly noisy data. Data has lots of errors.
- Output is discrete or real valued or a vector of values.
- Form of target function is unknown.
- Human readability, in other words, explainability, of result is unimportant.
- Examples:
 - Speech phoneme recognition
 - Image classification
 - Financial prediction

Perceptron



Perceptron

- One type of ANN system is based on a unit called a perceptron.
- A perceptron takes a vector of real-valued inputs, calculates a linear combination of these inputs, then outputs a 1 if the result is greater than some threshold and -1 otherwise.

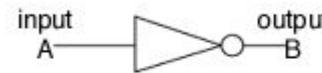


$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

NOT Logic Gate

Question:

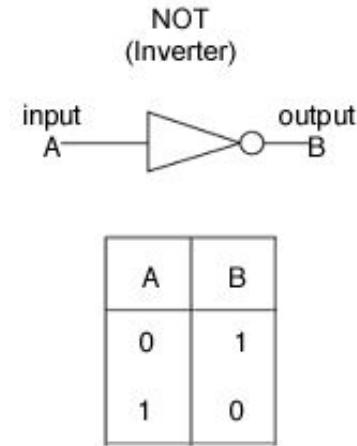
- How to represent NOT gate using a perceptron?
- What are the parameters for the NOT perceptron?
- Data is given below.



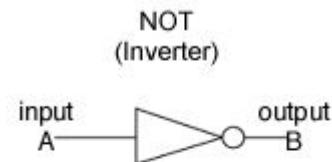
A	B
0	1
1	0

Perceptron for NOT gate

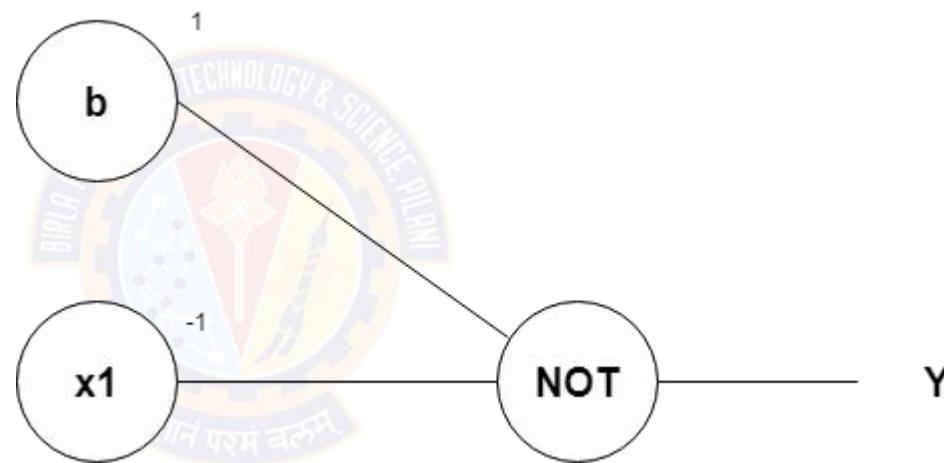
- One input say x_1 .
- One output say $o(X)$.
- Perceptron equation is $o(X) = \text{sign}(w_0 + w_1 x_1)$.
 - This should be > 0 for output to be 1.
- For each row of truth table, the equations are
 - $w_0 + 0 > 0$
 - $w_0 + w_1 < 0$
- One solution is $w_0 = 1$ and $w_1 = (-1)$. (Intuitive solution)
- This give a beautiful linear decision boundary.



Solution for NOT gate



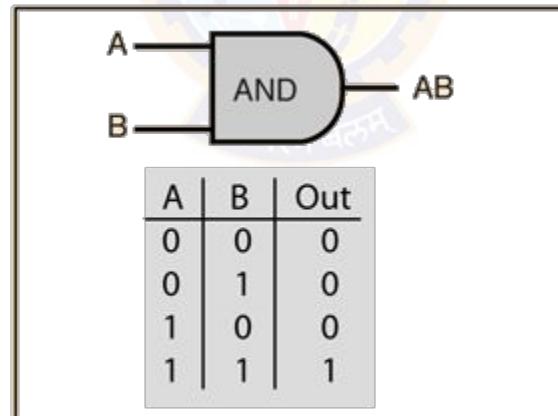
A	B
0	1
1	0



AND Logic Gate

Question:

- How to represent AND gate using a perceptron?
- What are the parameters for the AND perceptron?
- Data is given below.



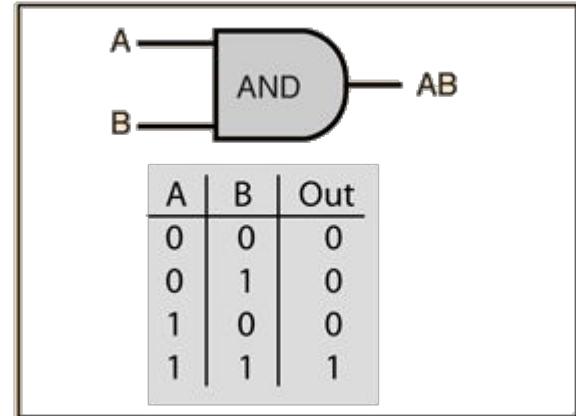
Perceptron for AND gate

- Two inputs say x_1 and x_2 .
- One output say $o(X)$.
- Perceptron equation is

$$o(X) = \text{sign}(w_0 + w_1 x_1 + w_2 x_2).$$

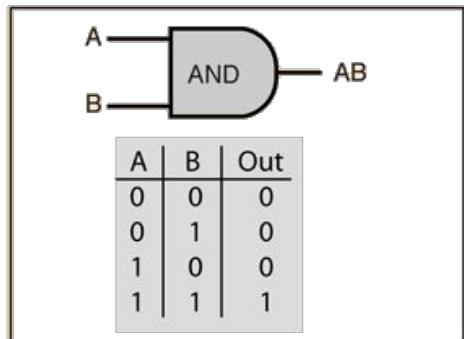
- This should be > 0 for output to be 1.

- For each row of truth table, the equations are
 - $w_0 + 0 + 0 < 0$
 - $w_0 + 0 + w_2 < 0$
 - $w_0 + w_1 + 0 < 0$
 - $w_0 + w_1 + w_2 > 0$
- One solution is $w_1 = w_2 = 1$ and $w_0 = (-1)$.
- This give a beautiful linear decision boundary.

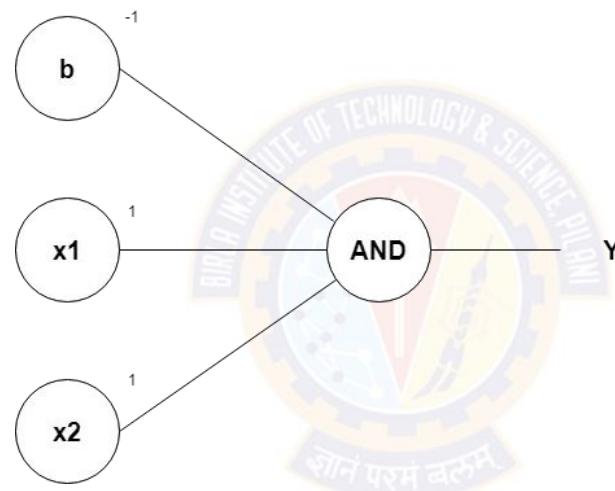


Perceptron for AND gate

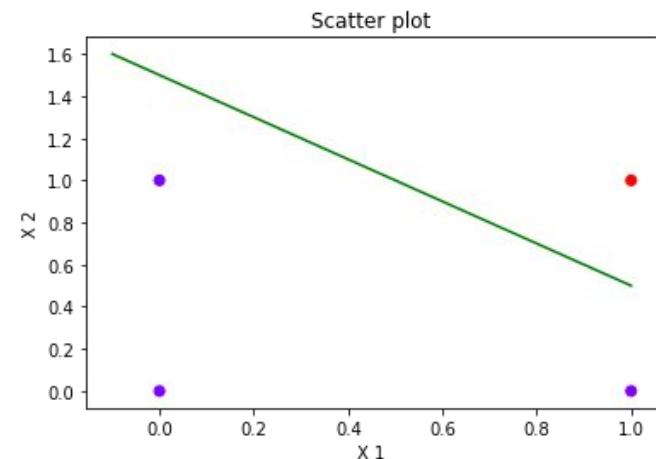
DATA



MODEL

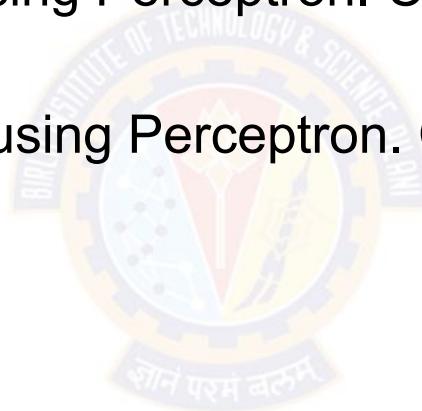


DECISION BOUNDARY



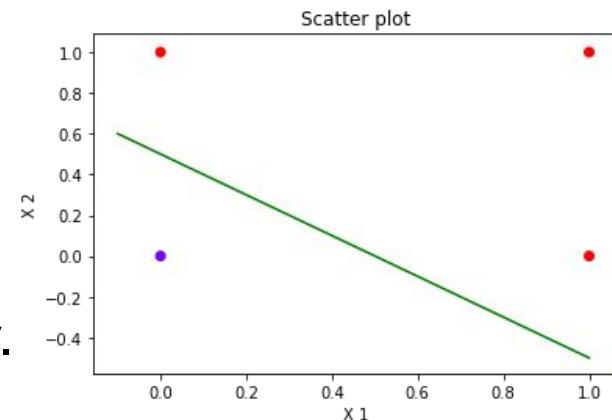
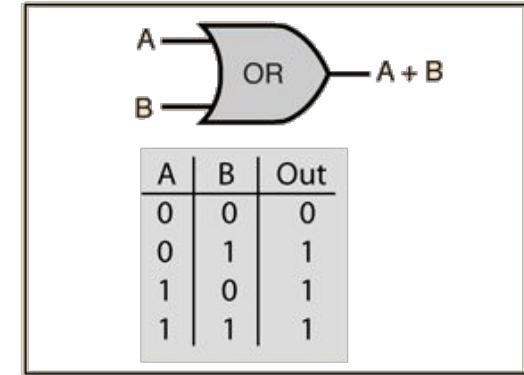
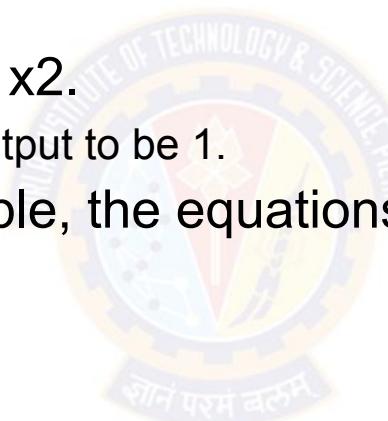
Exercise

1. Represent OR gate using Perceptron. Compute the parameters of the perceptron.
2. Represent NOR gate using Perceptron. Compute the parameters of the perceptron.
3. Represent NAND gate using Perceptron. Compute the parameters of the perceptron.



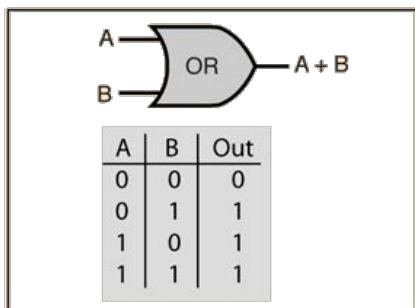
Perceptron for OR gate

- Two inputs say x_1 and x_2 .
- One output say $o(X)$.
- Perceptron equation is
- $o(X) = w_0 + w_1 x_1 + w_2 x_2$.
 - This should be > 0 for output to be 1.
- For each row of truth table, the equations are
 - $w_0 + 0 + 0 < 0$
 - $w_0 + 0 + w_2 > 0$
 - $w_0 + w_1 + 0 > 0$
 - $w_0 + w_1 + w_2 > 0$
- One solution is $w_1 = w_2 = 20$ and $w_0 = (-10)$.
- This give a beautiful linear decision boundary.

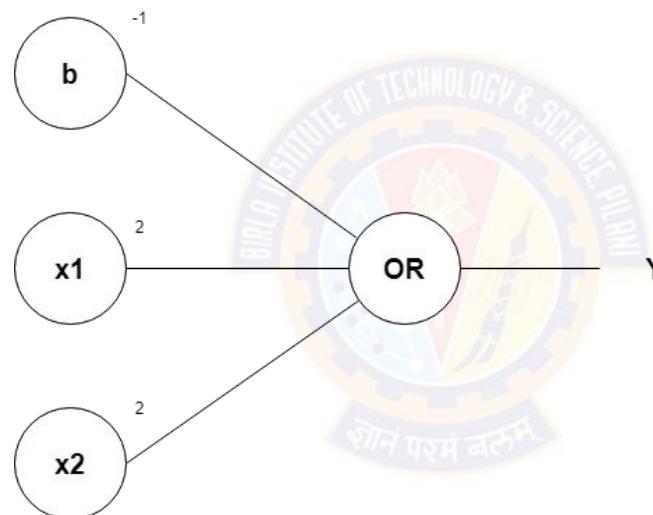


Perceptron for OR gate

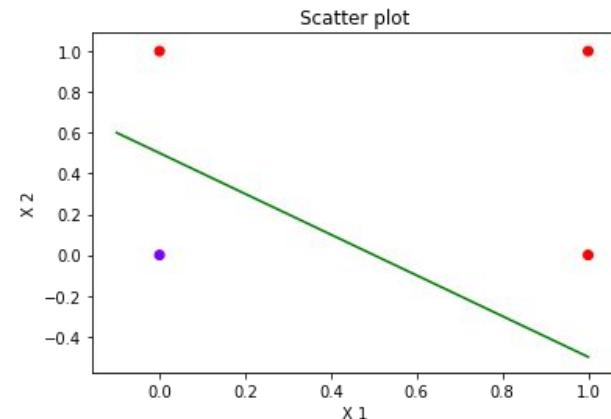
DATA



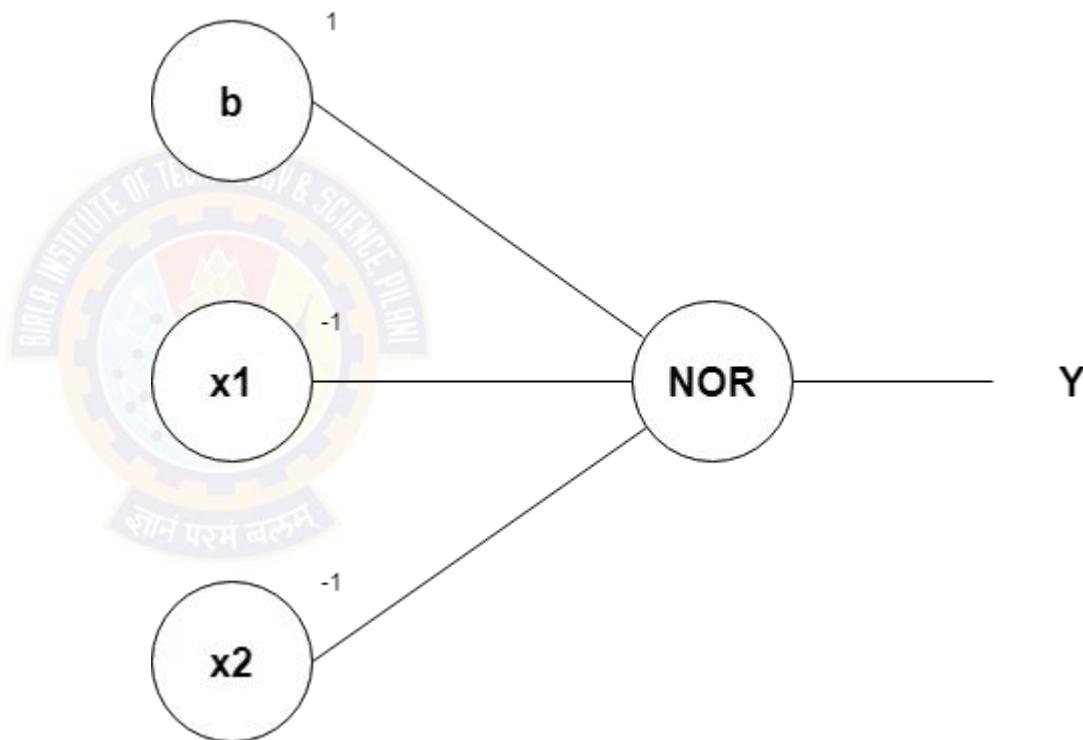
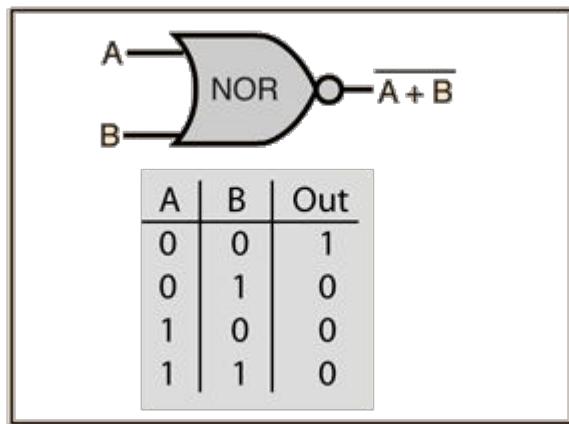
MODEL



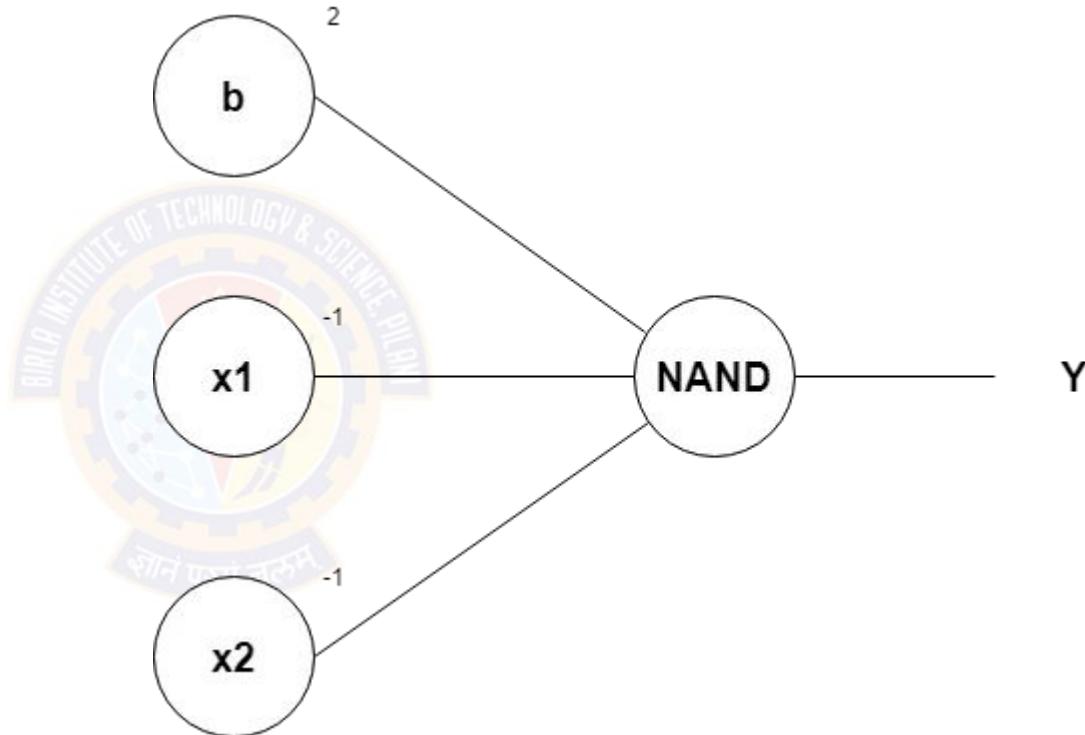
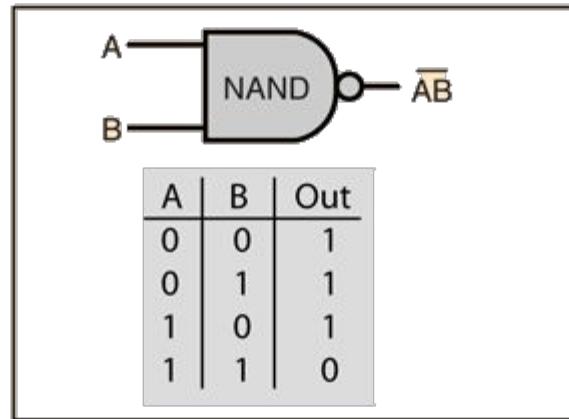
DECISION BOUNDARY



Perceptron for NOR Gate



Perceptron for NAND Gate



Perceptron Learning Algorithm



Perceptron Learning Algorithm

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

Where:

- $t = c(\vec{x})$ is target value
- o is perceptron output
- η is small constant (e.g., .1) called *learning rate*

Convergence of Perceptron Learning Algorithm

Can prove it will converge

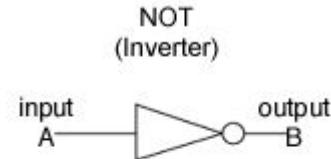
- If training data is linearly separable.
- Learning rate is sufficiently small
 - The role of the learning rate is to moderate the degree to which weights are changed at each step.
 - It is usually set to some small value (e.g., 0.1) and is sometimes made to decay as the number of weight-tuning iterations increases.

Perceptron Learning Algorithm for NOT gate

- Assume $w_0 = w_1 = 0$. Let the learning rate = $\eta = 1$.
- Equation is $o(x) = \text{sign}(w_0 + w_1 x_1)$

Epoch 1

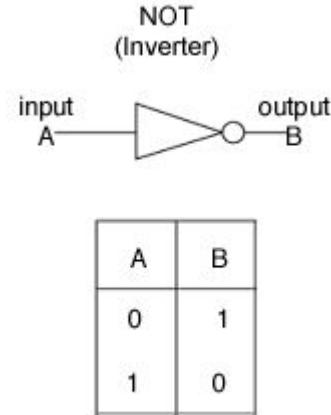
- For first example, $x_1 = 0$ and $d(x) = 1$
 - sum = $0 + 0 = 0$
 - $o(x) = \text{sign}(\text{sum}) = -1$
 - Not equal to $d(x)$. So update parameters.
 - $\Delta w_1 = \eta (d(x) - o(x)) x_1 = 1 (1 - (-1)) 0 = 0$
 - $\Delta w_0 = \eta (d(x) - o(x)) = 1 (1 - (-1)) = 2$ (for w_0 ; x_0 is assumed as 1)
 - New $w_1 = \text{old } w_1 + \Delta w_1 = 0 + 0 = 0$
 - New $w_0 = \text{old } w_0 + \Delta w_0 = 0 + 2 = 2$
- Use the above weights for next example.



A	B
0	1
1	0

Perceptron Learning Algorithm for NOT gate

- Second example uses $w_0 = 2$ and $w_1 = 0$. $\eta = 1$.
- For second example, $x_1 = 1$ and $d(x) = -1$
 - sum = $2 + 0 = 2$
 - $o(x) = \text{sign}(\text{sum}) = 1$
 - Not equal to $d(x)$
 - So update parameters.
 - $\Delta w_1 = \eta (d(x) - o(x)) x_1 = 1 (-1-1) 1 = -2$
 - $\Delta w_0 = \eta (d(x) - o(x)) = 1 (-1-1) = -2$
 - New $w_1 = \text{old } w_1 + \Delta w_1 = 0 + -2 = -2$
 - New $w_0 = \text{old } w_0 + \Delta w_0 = 2 + -2 = 0$
- Use the above weights for next epoch.

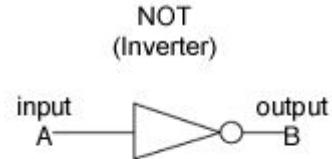


Perceptron Learning Algorithm for NOT gate

- From previous epoch $w_0 = 0$ and $w_1 = -2$. $\eta = 1$.

Epoch 2

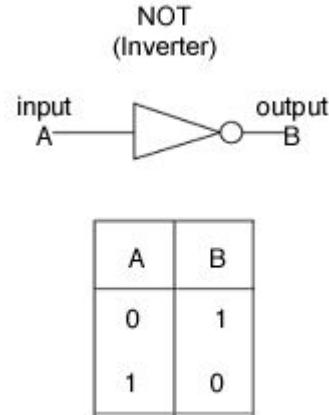
- For first example, $x_1 = 0$ and $d(x) = 1$
 - $\text{sum} = 0 + 0 = 0$
 - $o(x) = \text{sign}(\text{sum}) = -1$
 - Not equal to $d(x)$. So update parameters.
 - $\Delta w_1 = \eta (d(x) - o(x)) x_1 = 1 (1 - (-1)) 0 = 0$
 - $\Delta w_0 = \eta (d(x) - o(x)) = 1 (1 - (-1)) = 2$ (for w_0 ; x_0 is assumed as 1)
 - New $w_1 = \text{old } w_1 + \Delta w_1 = -2 + 0 = -2$
 - New $w_0 = \text{old } w_0 + \Delta w_0 = 0 + 2 = 2$



A	B
0	1
1	0

Perceptron Learning Algorithm for NOT gate

- Second example uses $w_0 = 2$ and $w_1 = -2$. $\eta = 1$.
- For second example, $d(x) = 0$
 - $\text{sum} = 2 + -2 = 0$
 - $o(x) = \text{sign}(\text{sum}) = -1$
 - Equal to $d(x)$. So NO need to update parameters.
- The Algorithm converges with $w_0 = 2$ and $w_1 = -2$.



Demo Python Code

<https://colab.research.google.com/drive/1DUVcOoUIWhl8GQKc6AWR1wi0LaMeNkgD?usp=sharing>



Student pl note:

The Python notebook is shared for anyone who has access to the link and the access is restricted to use BITS email id. So pl do not access from non-BITS email id and send requests for access.

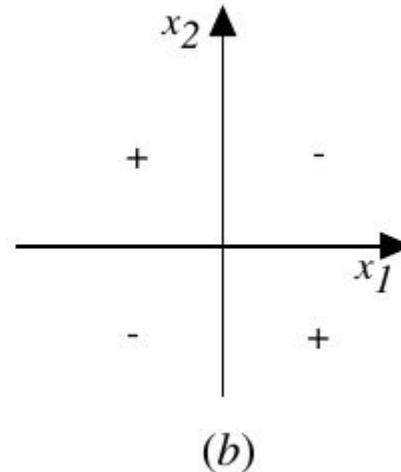
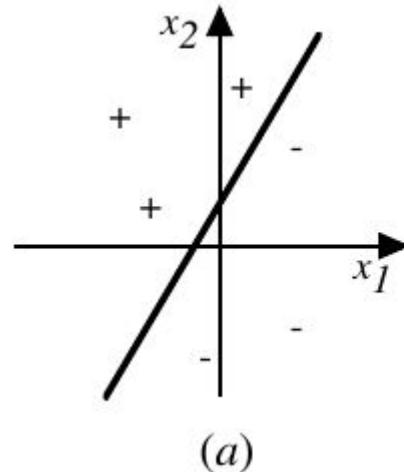
Exercise

1. Represent OR gate using Perceptron. Compute the parameters of the perceptron using perceptron learning algorithm.
2. Represent AND gate using Perceptron. Compute the parameters of the perceptron using perceptron learning algorithm.



Representational Power of Perceptrons

- A perceptron represents a hyperplane decision surface in the n-dimensional space of examples.
- The perceptron outputs a 1 for examples lying on one side of the hyperplane and outputs a -1 for examples lying on the other side.

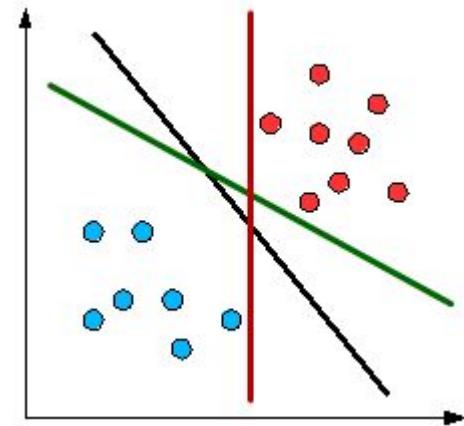


Linearly Separable Data

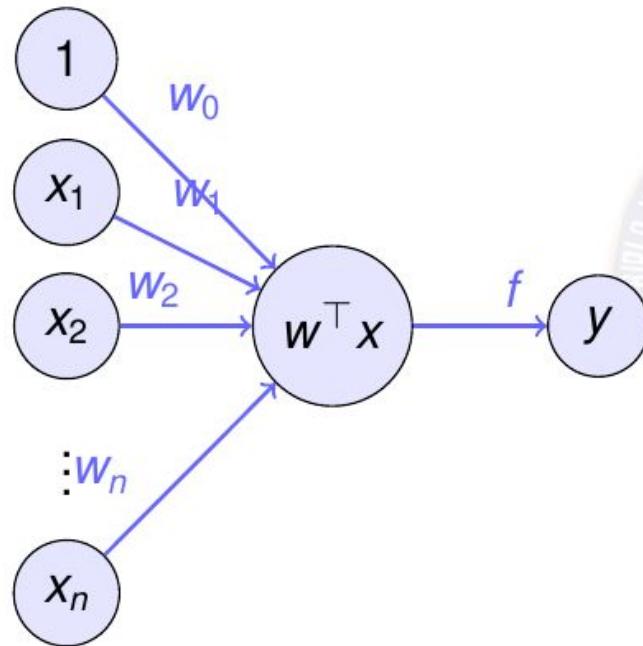


Linearly Separable Data

- Two sets of data points in a two dimensional space are said to be linearly separable when they can be completely separable by a single straight line.
- In general, **two groups of data points are separable in a n-dimensional space if they can be separated by an (n-1) dimensional hyperplane.**
- A straight line can be drawn to separate all the data examples belonging to class +1 from all the examples belonging to the class -1. Then the two-dimensional data are clearly linearly separable.
- An infinite number of straight lines can be drawn to separate the class +1 from the class -1.



Perceptron for Linearly Separable Data



$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$\begin{aligned} \sum &= w_0 + w_1 x_1 + \cdots + w_n x_n \\ &= w^\top x \end{aligned}$$

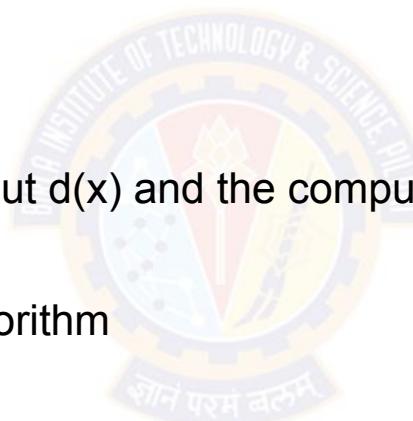
$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w^\top x > 0 \\ -1 & \text{otherwise.} \end{cases}$$

Challenge: How to learn these n parameters $w_1 \dots w_n$?

Solution: Use Perceptron learning algorithm.

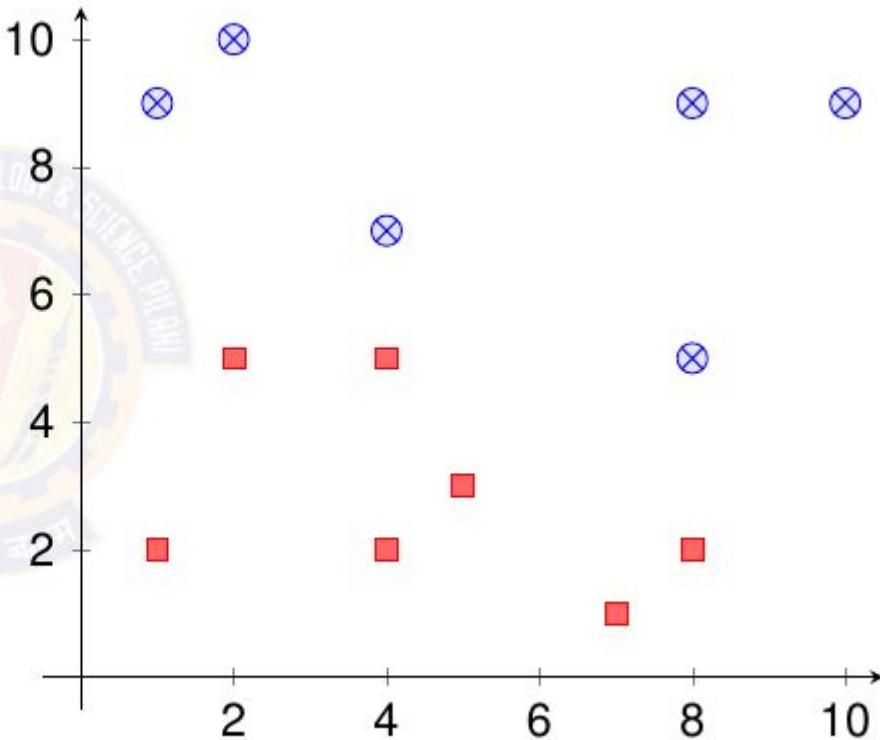
Perceptron and its Learning - Review

- Data
 - Truth tables or set of examples
- Model
 - Perceptron
- Objective Function
 - deviation of desired output $d(x)$ and the computed output $o(x)$
- Learning Algorithm
 - Perceptron Learning algorithm

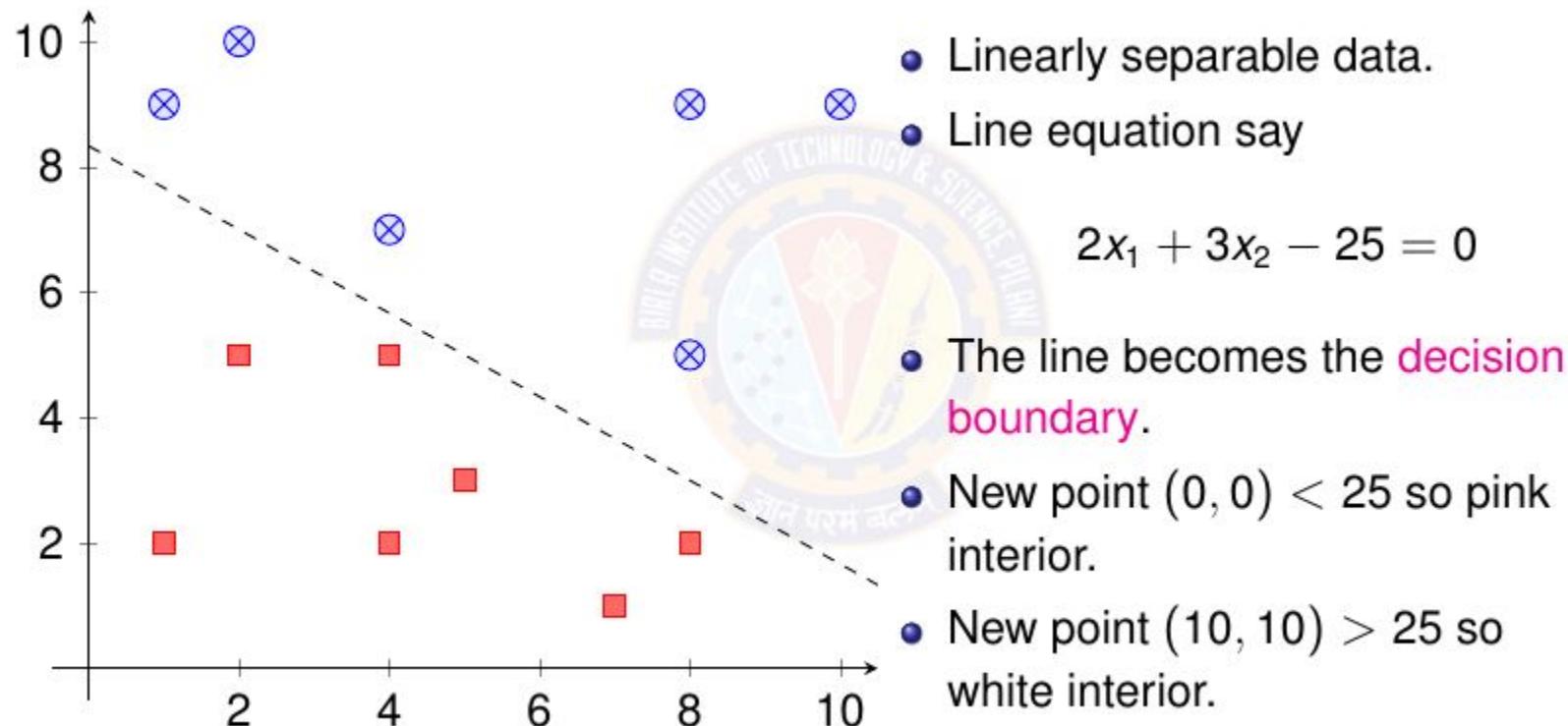


Example of Linearly Separable data

x_1	x_2	y
1	9	white
10	9	white
4	7	white
4	5	pink
5	3	pink
8	9	white
4	2	pink
2	5	pink
7	1	pink
2	10	white
8	5	white
1	2	pink
8	2	pink



Example of Linearly Separable data

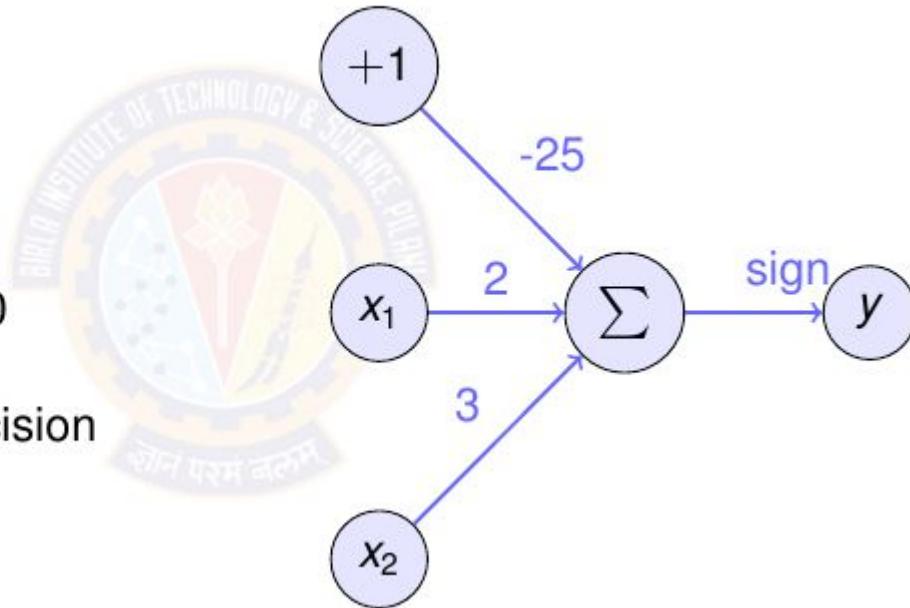


Example of Linearly Separable data

- Linearly separable data.
- Line equation say

$$2x_1 + 3x_2 - 25 = 0$$

- The line becomes the decision boundary.

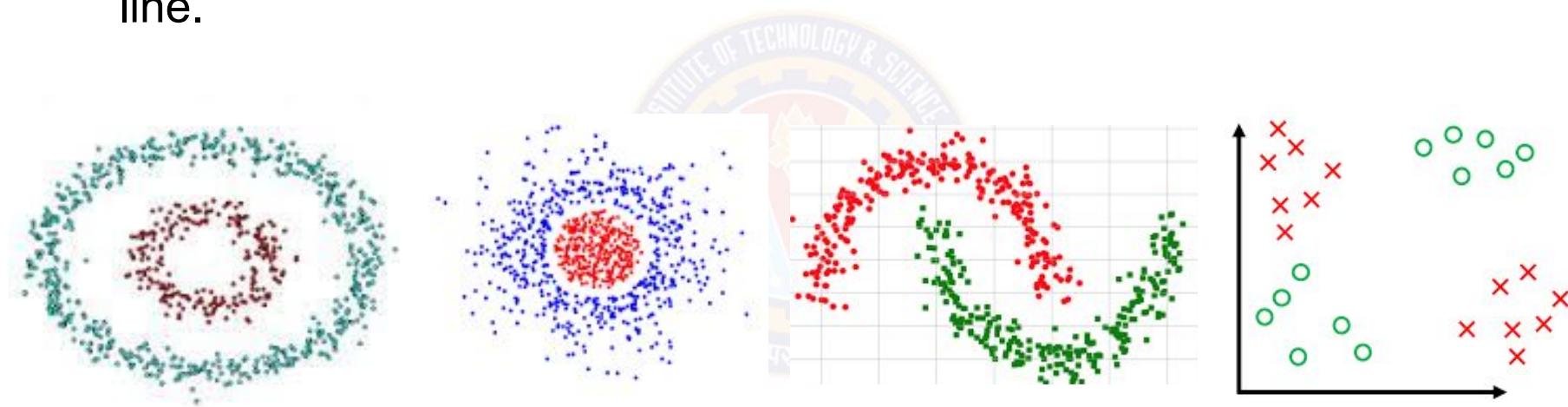


Non-Linearly Separable Data



Non-linearly Separable Data

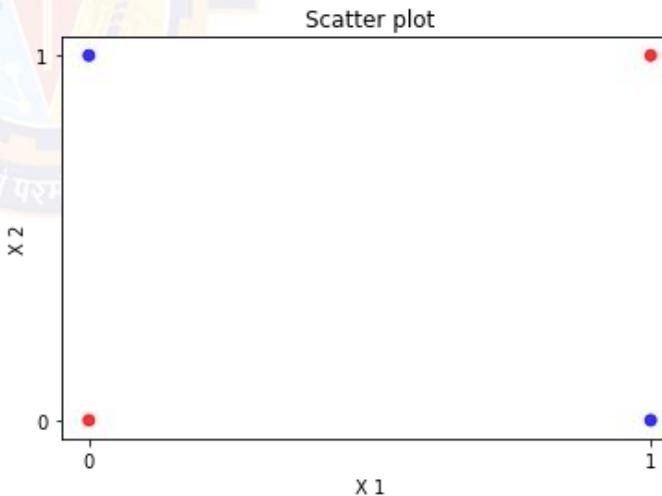
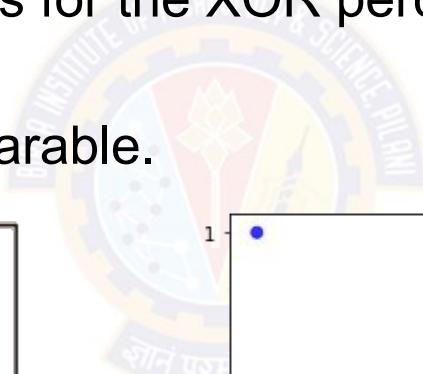
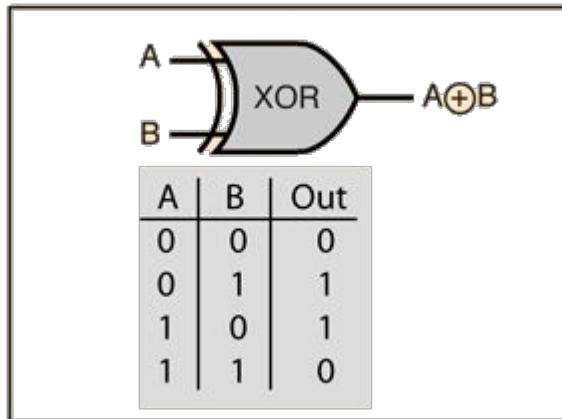
- Two groups of data points are non-linearly separable in a 2-dimensional space if they cannot be easily separated with a linear line.



Perceptron for XOR Gate

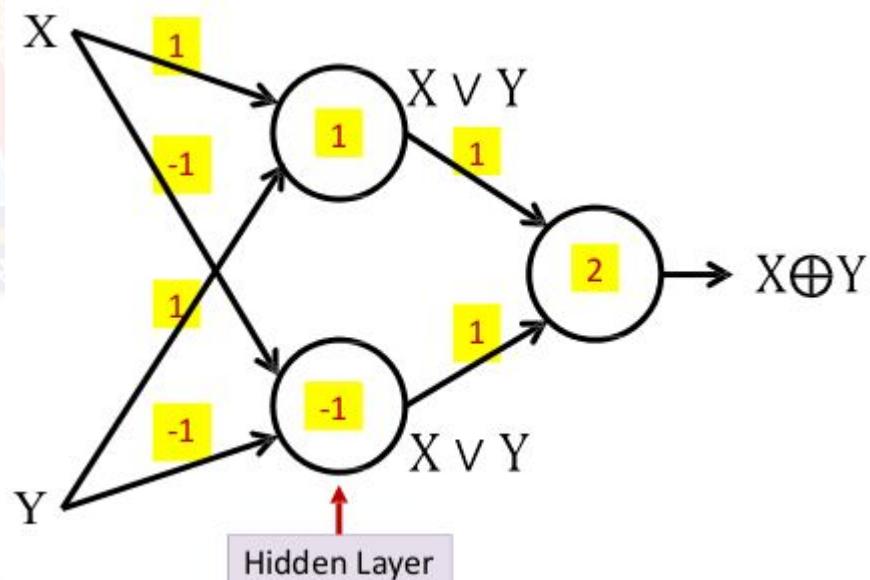
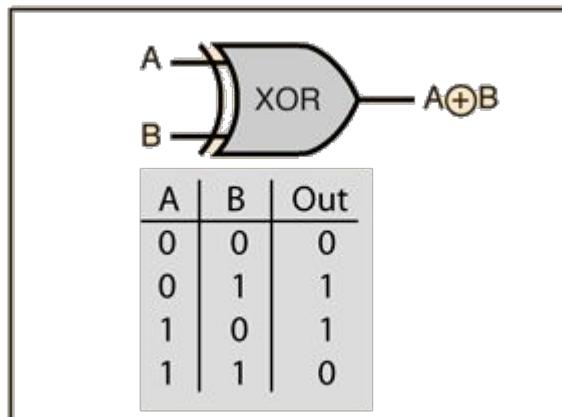
Question:

- How to represent XOR gate using a perceptron?
- What are the parameters for the XOR perceptron?
- Data is given below.
- Data is non-linearly separable.



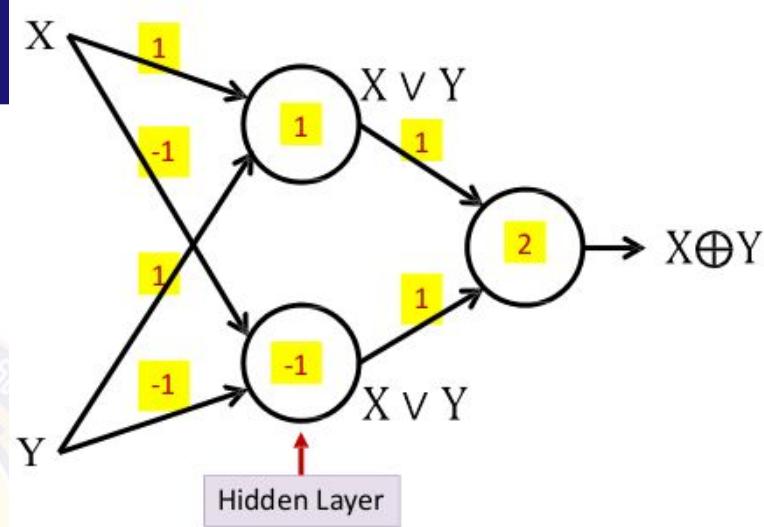
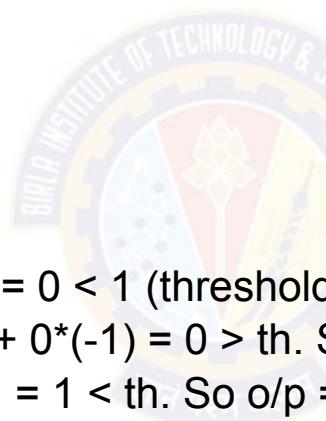
Solution for XOR

- Qn: How to represent XOR gate using a perceptron?
- Ans: Use Multilayer Perceptron (MLP)
- Introduce another layer in between the input and output.
- This in-between layer is called hidden layer.



MLP for XOR

MLP works :)



- Input (0,0)
 - First neuron = $0*1 + 0*1 = 0 < 1$ (threshold). So o/p = 0.
 - Second neuron = $0*(-1) + 0*(-1) = 0 > th.$ So o/p = 1
 - Third neuron = $0*1 + 1*1 = 1 < th.$ So o/p = 0. The desired output.
- Input (1,0)
 - First neuron = $1*1 + 0*1 = 1 = th.$ So o/p = 1
 - Second neuron = $1*(-1) + 0*(-1) = -1 = th.$ So o/p = 1
 - Third neuron = $1*1 + 1*1 = 2 < th.$ So o/p = 1. The desired output.

Solution of XOR Data

- Data
 - Truth table
- Model
 - Multi-layered Perceptron.
- Challenge
 - How to learn the parameters and threshold?
- Solution for learning
 - Use gradient descent algorithm



Gradient Descent Algorithm



Gradient Descent Algorithm

To understand, consider simpler *linear unit*, where

$$o = w_0 + w_1 x_1 + \cdots + w_n x_n$$

Let's learn w_i 's that minimize the squared error

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

where D is set of training examples.

Gradient

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}] \quad \text{i.e.,} \quad \Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Gradient Descent Algorithm

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\&= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\&= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\&= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ \frac{\partial E}{\partial w_i} &= \sum_d (t_d - o_d) (-x_{i,d})\end{aligned}$$

Gradient Descent Algorithm

Each training example is a pair of the form $\langle \vec{x}, t \rangle$, where \vec{x} is the vector of input values, and t is the target output value. η is the learning rate (e.g., .05).

Gradient-Descent(*training_examples*, η)

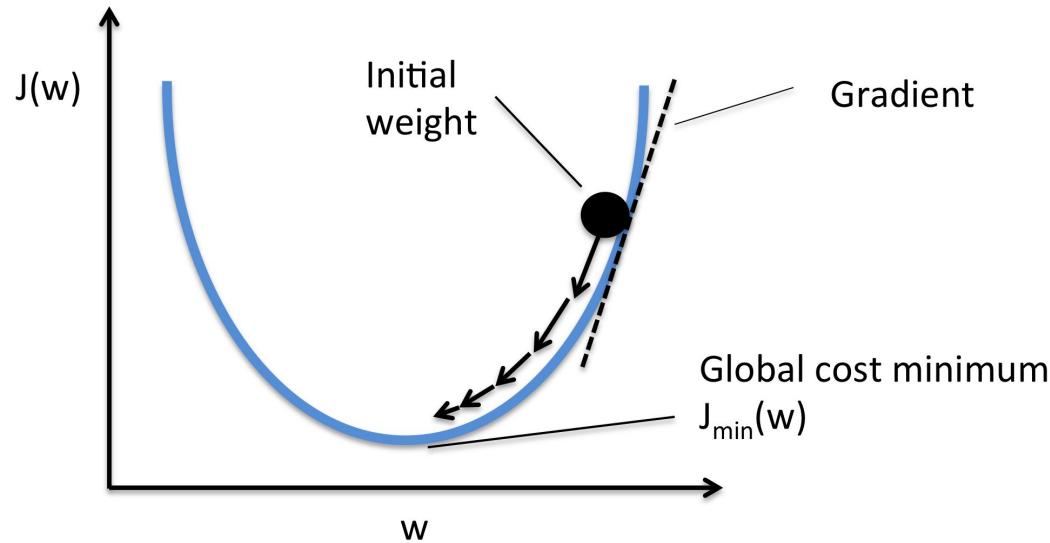
- Initialize each w_i to some small random value
- Until the termination condition is met, Do
 - ▶ Initialize each Δw_i to zero.
 - ▶ For each $\langle \vec{x}, t \rangle$ in *training_examples*, Do
 - ★ Input the instance \vec{x} to the unit and compute the output o
 - ★ For each linear unit weight w_i , Do

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i$$

- ▶ For each linear unit weight w_i , Do

$$w_i \leftarrow w_i + \Delta w_i$$

Gradient Descent Algorithm



Numerical Example

Equation is $y = (x+5)^2$

Initialize Parameters :

$$X_0 = 3$$

When will it be minimum?

$$\text{Learning rate} = 0.01$$

Use gradient descent algorithm

Assume starting point as 3 and

Learning rate as 0.01.

$$\frac{dy}{dx} = \frac{d}{dx}(x + 5)^2 = 2 * (x + 5)$$

Iteration 1 :

$$X_1 = X_0 - (\text{learning rate}) * \left(\frac{dy}{dx}\right)$$

$$X_1 = 3 - (0.01) * (2 * (3 + 5)) = 2.84$$

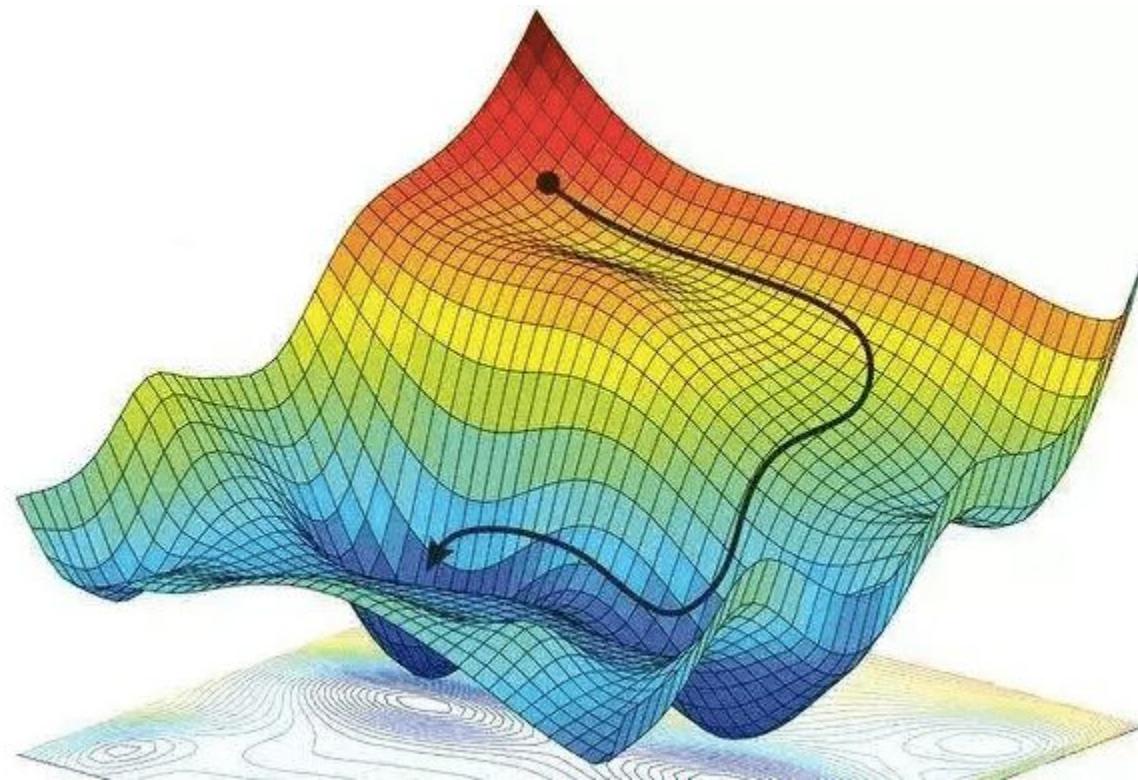
Iteration 2 :

$$X_2 = X_1 - (\text{learning rate}) * \left(\frac{dy}{dx}\right)$$

$$X_2 = 2.84 - (0.01) * (2 * (2.84 + 5)) = 2.6832$$

Gradient Descent Algorithm

Contour or hyperplane



Incremental (Stochastic) Gradient Descent

Incremental mode Gradient Descent

- Do until satisfied
 - ▶ For each training example d in D
 - ★ Compute the gradient $\nabla E_d[\vec{w}]$
 - ★ $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d[\vec{w}]$

$$E_d[\vec{w}] \equiv \frac{1}{2}(t_d - o_d)^2$$

Incremental Gradient Descent can approximate Batch Gradient Descent arbitrarily closely if η made small enough.

MultiLayer Perceptron (MLP)

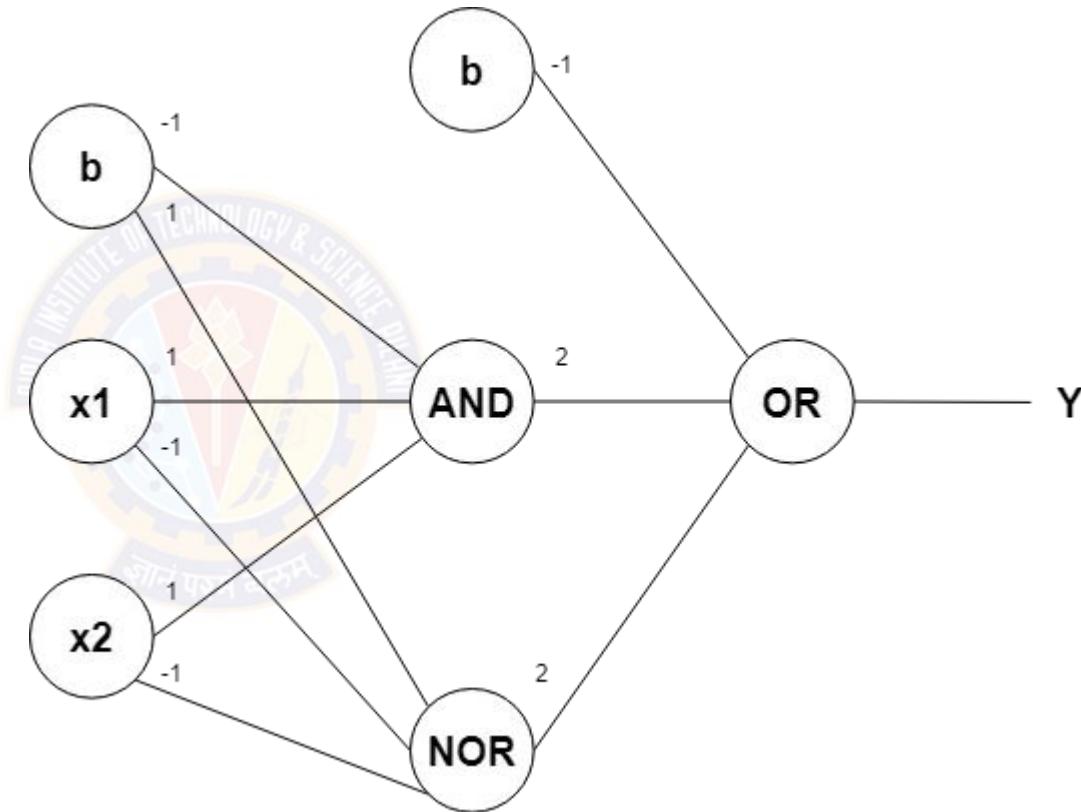
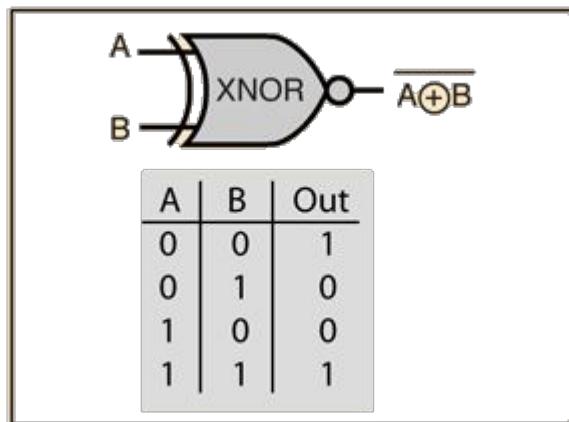


MultiLayer Perceptron (MLP)

Requires

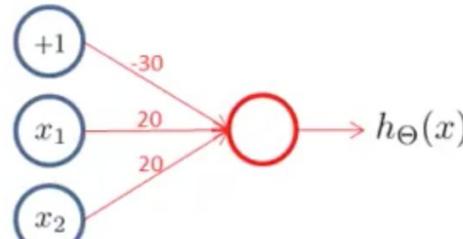
1. **Forward Pass** through each layer to compute the output.
2. Compute the deviation or error between the desired output and computed output in the forward pass (first step). This morphs into **objective function**, as we want to minimize this deviation or error.
3. The deviation has to be send back through each layer to compute the delta or change in the parameter values. This is achieved using **back propagation algorithm**.
4. **Update** the parameters.

Perceptron for XNOR Gate

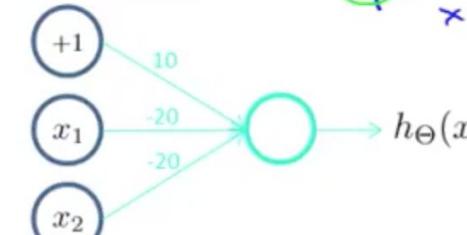


XNOR

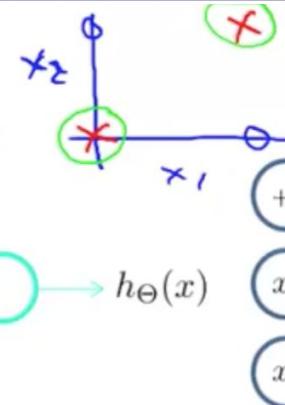
Putting it together: x_1 XNOR x_2



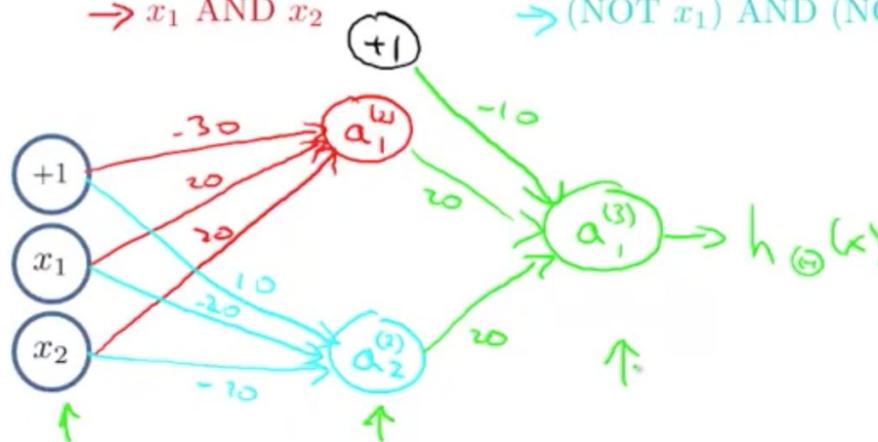
$\rightarrow x_1 \text{ AND } x_2$



$\rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$



$\rightarrow x_1 \text{ OR } x_2$



x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1 ↪
0	1	0	0	0 ↪
1	0	0	0	0 ↪
1	1	1	0	1 ↪

Ref:

Chapter 4 of Book: Machine Learning by
Tom M. Mitchell



Next Session:
How to Train MLP?
Power of MLP





BITS Pilani
Pilani | Dubai | Goa | Hyderabad

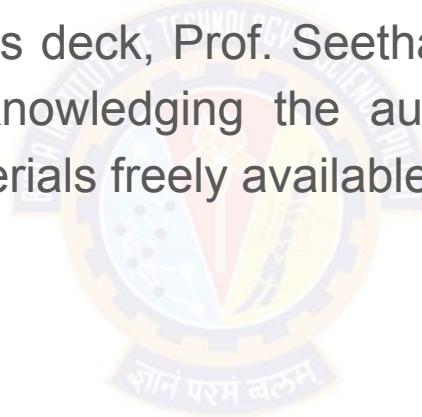
Deep Learning

DSE Module 2

Seetha Parameswaran

BITS Pilani

The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Single Perceptron for Regression Real Valued Output

Linear Regression Example

- Suppose that we wish to estimate the prices of houses (in dollars) based on their area (in square feet) and age (in years).
- The linearity assumption just says that the target (price) can be expressed as a weighted sum of the features (area and age):

$$price = w_{\text{area}} * \text{area} + w_{\text{age}} * \text{age} + b$$

- w_{area} and w_{age} are called **weights**, and b is called a **bias**.
- The weights determine the influence of each feature on our prediction and the bias just says what value the predicted price should take when all of the features take value 0.

Data

- Data

- The dataset is called a **training dataset or training set**.
- Each row is called an **example** (or data point, data instance, sample).
- The thing we are trying to predict is called a **label** (or target).
- The independent variables upon which the predictions are based are called **features** (or covariates).
 - m – number of training examples
 - i – index of i^{th} example
 - $x^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}]$ – features of i^{th} example
 - $y^{(i)}$ – label corresponding to i^{th} example

Affine transformations and Linear Models

- The equation of the form

$$\hat{y} = m_1x_1 + \dots + m_dx_d + b$$

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b \quad \mathbf{x} \in \mathbb{R}^d \quad \mathbf{w} \in \mathbb{R}^d$$

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + b \quad \mathbf{X} \in \mathbb{R}^{n \times d} \quad \hat{\mathbf{y}} \in \mathbb{R}^n$$

- is an **affine transformation** of input features, which is characterized by a linear transformation of features via weighted sum, combined with a translation via the added bias.
- Models whose output prediction is determined by the affine transformation of input features are **linear models**.
- The affine transformation is specified by the chosen weights (w) and bias (b).

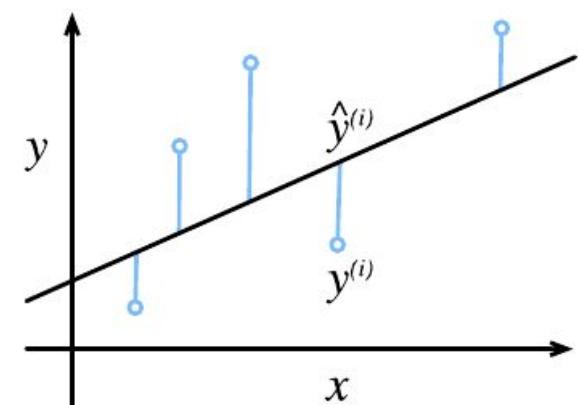
Loss Function

- Loss function is a quality measure for some given model or a measure of fitness.
- The loss function quantifies the distance between the real and predicted value of the target.
- The loss will usually be a non-negative number where smaller values are better and perfect predictions incur a loss of 0.
- The most popular loss function in regression problems is the squared error.

Squared Error Loss Function

- The most popular loss function in regression problems is the squared error.
- For each example, $l^{(i)}(\mathbf{w}, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$
- For the entire dataset of n examples
 - average (or equivalently, sum) the losses

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)})^2.$$



- When training the model, find parameters (\mathbf{w}^*, b^*) that minimize the total loss across all training examples

$$\mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\operatorname{argmin}} L(\mathbf{w}, b)$$

Minibatch Stochastic Gradient Descent (SGD)

- Apply Gradient descent algorithm on a random minibatch of examples every time we need to compute the update.
- In each iteration,
 - Step 1: randomly sample a minibatch \mathcal{B} consisting of a fixed number of training examples.
 - Step 2: compute the derivative (gradient) of the average loss on the minibatch with regard to the model parameters.
 - Step 3: multiply the gradient by a predetermined positive value η and subtract the resulting term from the current parameter values.

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \partial_{\mathbf{w}} l^{(i)}(\mathbf{w}, b) = \mathbf{w} - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \mathbf{x}^{(i)} \left(\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)} \right),$$

$$b \leftarrow b - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \partial_b l^{(i)}(\mathbf{w}, b) = b - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \left(\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)} \right).$$

Training using SGD Algorithm

- Initialize parameters (\mathbf{w}, b)
- Repeat until done
 - Compute gradient $\mathbf{g} \leftarrow \partial_{(\mathbf{w}, b)} \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} l(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w}, b)$
 - Update parameters $(\mathbf{w}, b) \leftarrow (\mathbf{w}, b) - \eta \mathbf{g}$



PS: The number of epochs and the learning rate are both hyperparameters. Setting hyperparameters requires some adjustment by trial and error.

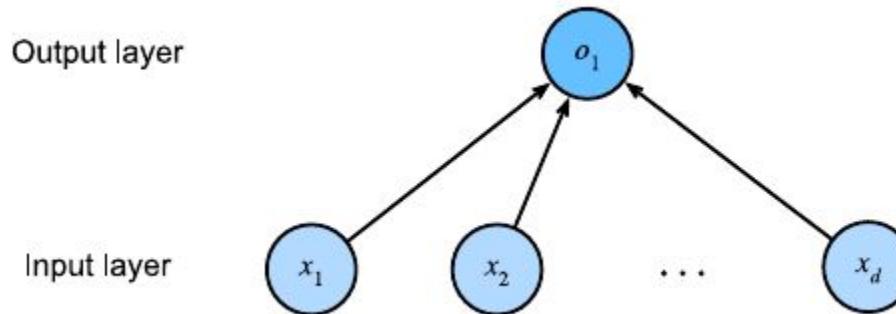
Prediction

- Estimating targets given features is commonly called prediction or inference.
- Given the learned model, values of target can be predicted, for any set of features.



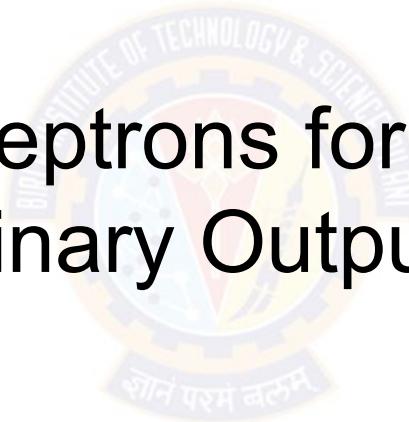
Single-Layer Neural Network

Linear regression is a single-layer neural network



- Number of inputs (or feature dimensionality) in the input layer is d .
- The inputs are x_1, \dots, x_d .
- The output is o_1 .
- Number of outputs in the output layer is 1.
- Number of layers for the neural network is 1. (conventionally we do not consider the input layer when counting layers.)
- Every input is connected to every output, This transformation is a **fully-connected layer** or **dense layer**.

Multiple Perceptrons for Classification Binary Outputs

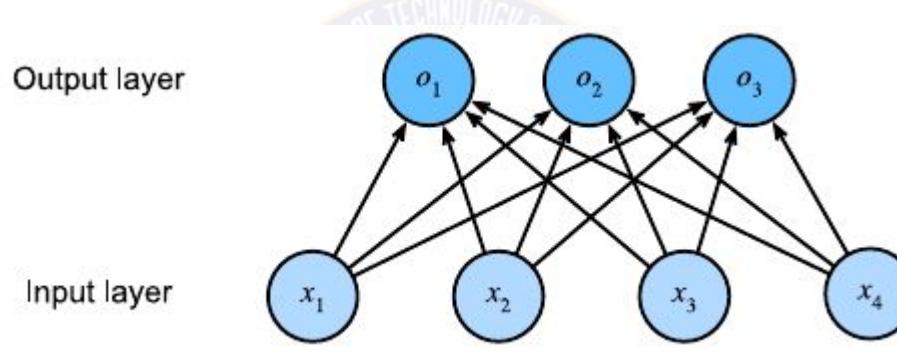


Classification Example

- Each input consists of a 2×2 grayscale image.
- Represent each pixel value with a single scalar, giving four features x_1, x_2, x_3, x_4 .
- Assume that each image belongs to one among the categories “square”, “triangle”, and “circle”.
- How to represent the labels?
 - Use label encoding. $y \in \{1, 2, 3\}$, where the integers represent {circle, square, triangle} respectively.
 - Use one-hot encoding. $y \in \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$.
 - y would be a three-dimensional vector, with $(1, 0, 0)$ corresponding to “circle”, $(0, 1, 0)$ to “square”, and $(0, 0, 1)$ to “triangle”.

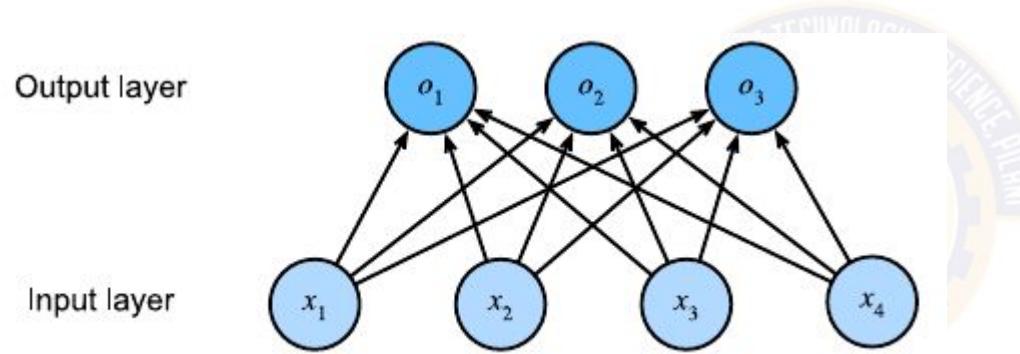
Network Architecture

- A model with multiple outputs, one per class. Each output will correspond to its own affine function.
 - 4 features and 3 possible output categories



Network Architecture

- 12 scalars to represent the weights and 3 scalars to represent the biases
- compute three logits, o_1 , o_2 , and o_3 , for each input.
- weights is a 3×4 matrix and bias is 1×4 matrix



$$o_1 = x_1w_{11} + x_2w_{12} + x_3w_{13} + x_4w_{14} + b_1,$$

$$o_2 = x_1w_{21} + x_2w_{22} + x_3w_{23} + x_4w_{24} + b_2,$$

$$o_3 = x_1w_{31} + x_2w_{32} + x_3w_{33} + x_4w_{34} + b_3.$$

$$\mathbf{O} = \mathbf{XW} + \mathbf{b},$$

$$\hat{\mathbf{Y}} = \text{softmax}(\mathbf{O}).$$

Softmax Operation

- Interpret the outputs of our model as probabilities.
 - Any output \hat{y}_j is interpreted as the probability that a given item belongs to class j . Then choose the class with the largest output value as our prediction $\text{argmax}_j \hat{y}_j$.
 - If \hat{y}_1 , \hat{y}_2 , and \hat{y}_3 are 0.1, 0.8, and 0.1, respectively, then predict category 2.
 - To interpret the outputs as probabilities, we must guarantee that, they will be nonnegative and sum up to 1.
- The softmax function transforms the outputs such that they become nonnegative and sum to 1, while requiring that the model remains differentiable.

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{o}) \quad \text{where} \quad \hat{y}_j = \frac{\exp(o_j)}{\sum_k \exp(o_k)}$$

- first exponentiate each logit (ensuring non-negativity) and then divide by their sum (ensuring that they sum to 1)
- Softmax is a nonlinear function.

Log-Likelihood Loss Function / Cross-Entropy loss

- The softmax function gives us a vector \hat{y} , which we can interpret as estimated conditional probabilities of each class given any input x ,
 - $\hat{y}_1 = P(y = \text{cat} | x)$.
- Compare the estimates with reality by checking how probable the actual classes are according to our model, given the features:

$$P(\mathbf{Y} | \mathbf{X}) = \prod_{i=1}^n P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$$

- Maximize $P(Y | X) =$ minimizing the negative log-likelihood

$$-\log P(\mathbf{Y} | \mathbf{X}) = \sum_{i=1}^n -\log P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}) = \sum_{i=1}^n l(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

$$l(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^q y_j \log \hat{y}_j.$$

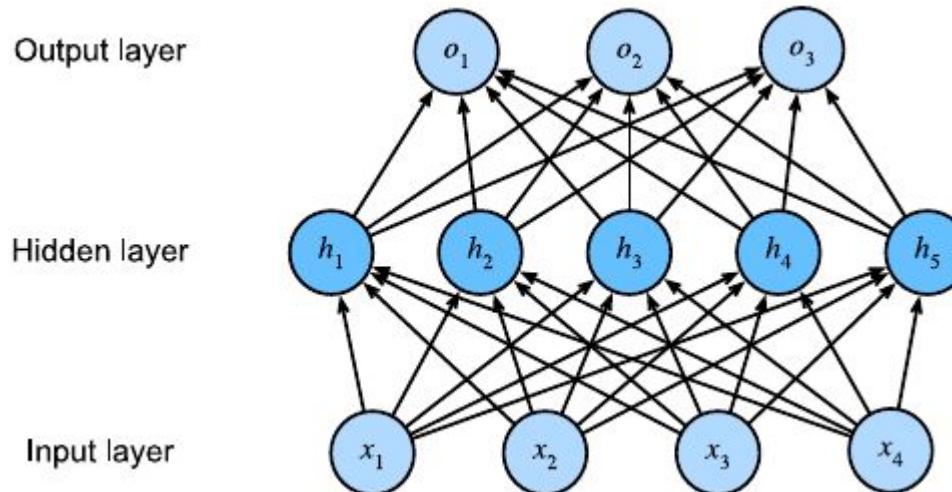
Multi Layered Perceptrons (MLP)



Multilayer Perceptron

- With deep neural networks, use the data to jointly learn both a representation via hidden layers and a linear predictor that acts upon that representation.
- Add many hidden layers by stacking many fully-connected layers on top of each other. Each layer feeds into the layer above it, until we generate outputs.
- The first ($L-1$) layers learns the representation and the final layer is the linear predictor. This architecture is commonly called a **multilayer perceptron (MLP)**.

MLP Architecture



- MLP has 4 inputs, 3 outputs, and its hidden layer contains 5 hidden units.
- Number of layers in this MLP is 2.
- The layers are both fully connected. Every input influences every neuron in the hidden layer, and each of these in turn influences every neuron in the output layer.
- The outputs of the hidden layer are called as **hidden representations** or hidden-layer variable or a hidden variable.

Nonlinearity in MLP

- Input is $\mathbf{X} \in \mathbb{R}^{n \times d}$ with n examples where each example has d features.
- Hidden Layer has h hidden units $\mathbf{H} \in \mathbb{R}^{n \times h}$
- Hidden-layer weights $\mathbf{W}^{(1)} \in \mathbb{R}^{d \times h}$ and biases $\mathbf{b}^{(1)} \in \mathbb{R}^{1 \times h}$
- Output-layer weights $\mathbf{W}^{(2)} \in \mathbb{R}^{h \times q}$ and biases $\mathbf{b}^{(2)} \in \mathbb{R}^{1 \times q}$
- Output is $\mathbf{O} \in \mathbb{R}^{n \times q}$
- A nonlinear activation function σ has to be applied to each hidden unit following the affine transformation. The outputs of activation functions are called **activations**.

$$\mathbf{H} = \sigma(\mathbf{X}\mathbf{W}^{(1)} + \mathbf{b}^{(1)})$$

$$\mathbf{O} = \text{softmax}(\mathbf{X}\mathbf{W}^{(2)} + \mathbf{b}^{(2)})$$

Activation Functions



Activation function

- Activation function of a neuron defines the output of that neuron given an input or set of inputs.
- Activation functions decide whether a neuron should be activated or not by calculating the weighted sum and adding bias with it.
- They are differentiable operators to transform input signals to outputs, while most of them add non-linearity.
- Artificial neural networks are designed as universal function approximators, they must have the ability to calculate and learn any nonlinear function.

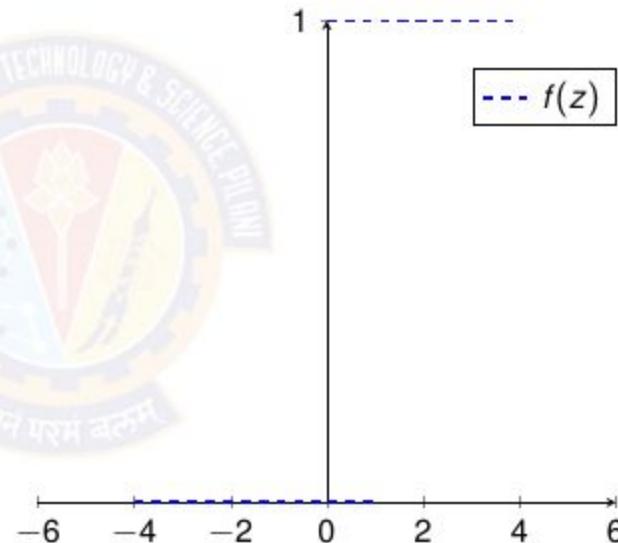
1. Step Function

Function: $f(z) = \begin{cases} 0 & \text{for } z \leq 0 \\ 1 & \text{for } z > 0 \end{cases}$

Derivative: $f'(z) = 0$

Range : $(-\infty, +\infty)$

- Linear function
- Derivative is 0. So neuron cannot learn.

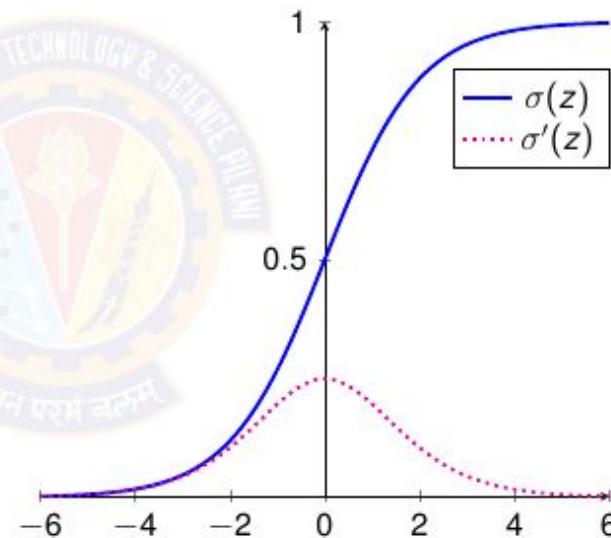
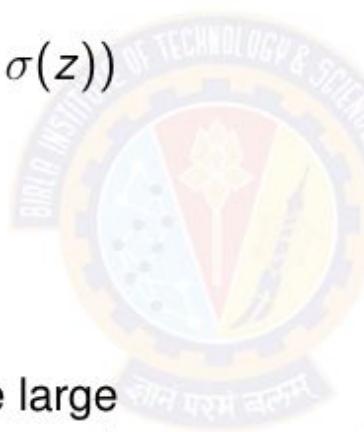


2. Sigmoid (Logistic) Activation Function

Function: $\sigma(z) = \frac{1}{1 + e^{-z}}$

Derivative: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$

Range : $(0, 1)$



- Non-Linear function
- Small changes in z will be large in $f(z)$. This means it is a good classifier.
- Leads to vanishing gradient. So the learning is minimal.

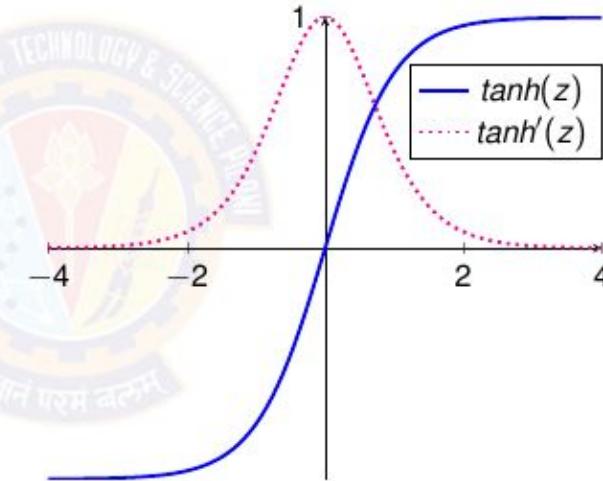
Squashing function as: it squashes any input in the range $(-\infty, \infty)$ to some value in the range $(0, 1)$

3. Tanh (hyperbolic tangent) Activation Function

$$\text{Function: } \tanh(z) = \frac{(e^z - e^{-z})}{(e^z + e^{-z})}$$

$$\text{Derivative: } \tanh'(z) = (1 - \tanh^2(z))$$

Range : $(-1, 1)$



- More efficient because it has a wider range for faster learning and grading.
- The tanh activation usually works better than sigmoid activation function for hidden units because the mean of its output is closer to zero, and so it centers the data better for the next layer.
- Issues with tanh
 - computationally expensive
 - lead to vanishing gradients

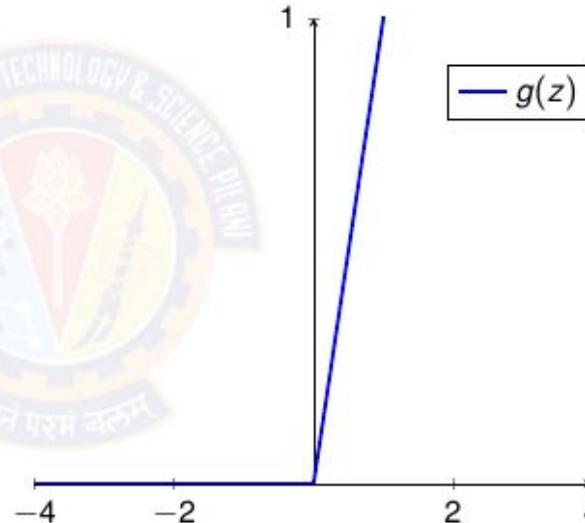
4. ReLU Activation Function

Rectified Linear Unit

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

Derivative: $g'(z) = \begin{cases} 0 & \text{for } z \leq 0 \\ 1 & \text{for } z > 0 \end{cases}$

Range : $[0, \infty]$



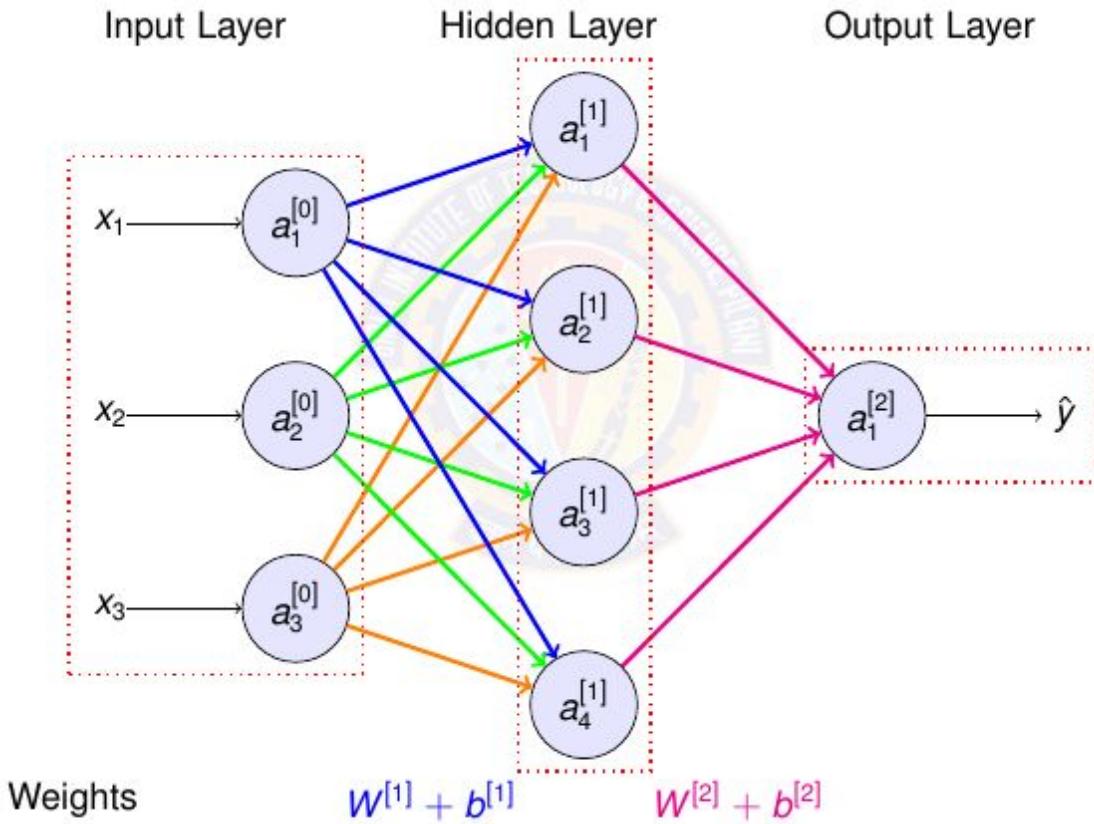
- If we combine two ReLU units, we can recover a piecewise linear approximation of the Sigmoid function.
- Some ReLU variants: Softplus (Smooth ReLU), Noisy ReLU, Leaky ReLU, Parametric ReLU and Exponential ReLU (ELU).
- Advantages
 - Fast Learning and Efficient computation
 - Fewer vanishing gradient problems
 - Sparse activation
 - Scale invariant (max operation)
- Disadvantages
 - Leads to exploding gradient.

Comparing Activation Functions

Activation Function	Sigmoid	Tanh	ReLU
Linearity	Non Linear	Non Linear	Non Linear
Activation Function	$\sigma(x) = \frac{1}{1+e^{-x}}$	$tanh(z) = \frac{(e^z - e^{-z})}{(e^z + e^{-z})}$	$g(x) = \max(0, z)$
Derivative	$\sigma(z)(1 - \sigma(z))$	$(1 - tanh^2(z))$	$1 \text{ if } z > 0 \text{ else } 0$
Symmetric Function	No	Yes	No
Range	$[0, 1]$	$[-1, 1]$	$[0, \infty]$
Vanishing Gradient	Yes	Yes	No
Exploding Gradient	No	No	Yes
Zero Centered	No	Yes	No

Training MLP

Two Layer Neural Network



Compute the Activations

For Hidden layer

$$z_1^{[1]} = w_1^{[1]} a^{[0]} + b^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]} a^{[0]} + b^{[1]}$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]} a^{[0]} + b^{[1]}$$

$$a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]} a^{[0]} + b^{[1]}$$

$$a_4^{[1]} = \sigma(z_4^{[1]})$$

In the matrix form


$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \underbrace{\begin{bmatrix} \dots w_1^{[1]\top} \dots \\ \dots w_2^{[1]\top} \dots \\ \dots w_3^{[1]\top} \dots \\ \dots w_4^{[1]\top} \dots \end{bmatrix}}_{4 \times 3} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{3 \times 1} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{4 \times 1}$$

$$z^{[1]} = \begin{bmatrix} w_1^{[1]\top} a^{[0]} + b_1^{[1]} \\ w_2^{[1]\top} a^{[0]} + b_2^{[1]} \\ w_3^{[1]\top} a^{[0]} + b_3^{[1]} \\ w_4^{[1]\top} a^{[0]} + b_4^{[1]} \end{bmatrix} \leftarrow 4 \times 1 \text{ matrix}$$

$$a^{[1]} = \sigma(z^{[1]}) \leftarrow 4 \times 1 \text{ matrix}$$

Vectoring Forward Propagation

Vectorizing for one training example

$$a^{[0]} = X$$

$$Z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(Z^{[2]})$$

$$\hat{y} = a^{[2]}$$

Vectorizing for all training examples

$$\text{Example 1 : } X^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$$

$$\text{Example 2 : } X^{(2)} \longrightarrow a^{2} = \hat{y}^{(2)}$$

$$\text{Example 3 : } X^{(2)} \longrightarrow a^{[2](3)} = \hat{y}^{(3)}$$

:

:

$$\text{Example m : } X^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$$

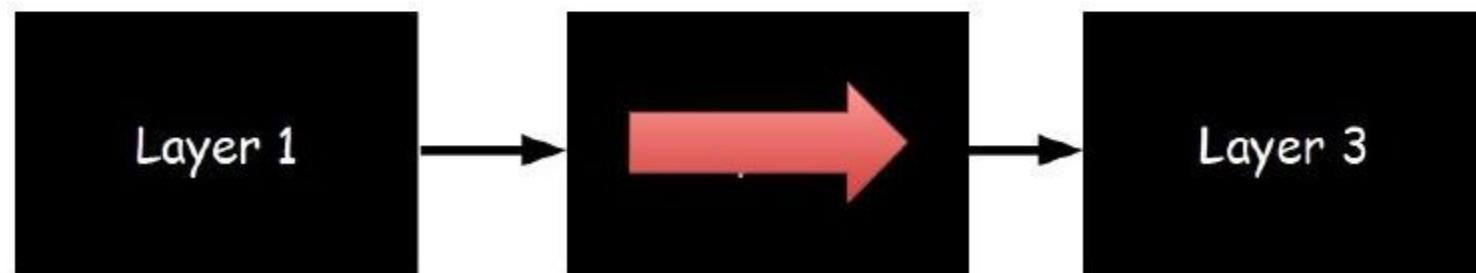
Neural Network Training – Forward Pass

- Step 1: Compute Loss on mini-batch [F-Pass]



Neural Network Training – Forward Pass

- Step 1: Compute Loss on mini-batch [F-Pass]

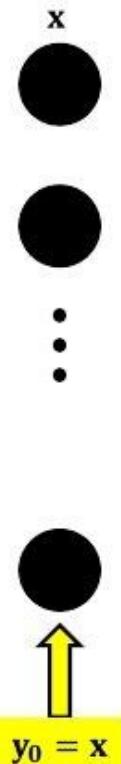


Neural Network Training – Forward Pass

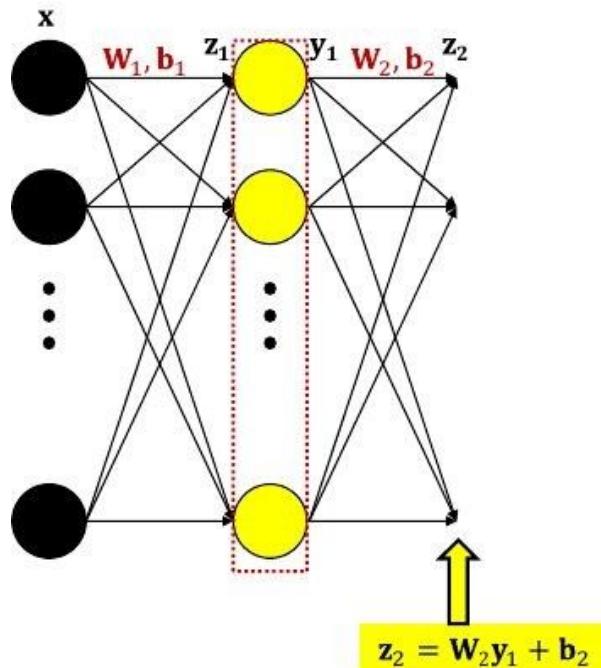
- Step 1: Compute Loss on mini-batch [F-Pass]



The forward pass: Evaluating the network



The forward pass

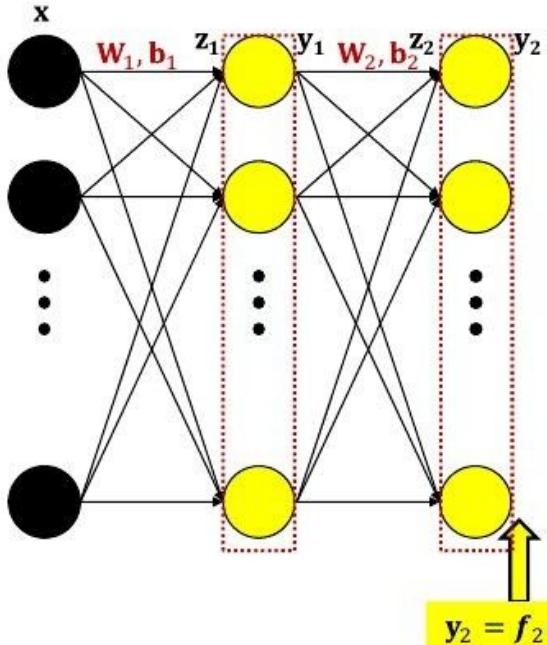


The Complete computation

$$y_1 = f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$



The forward pass

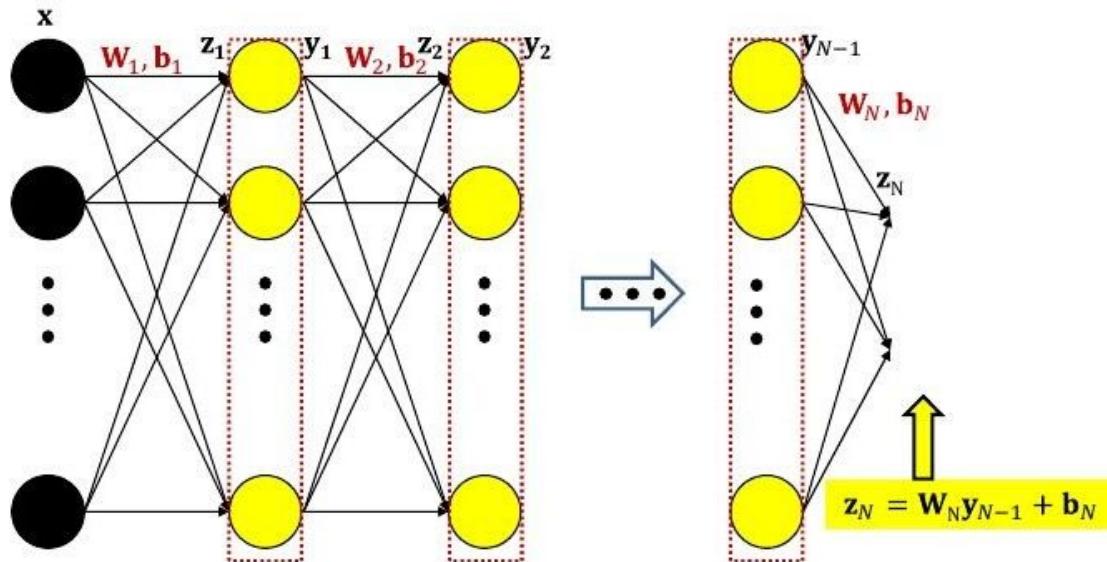


The Complete computation

$$y_2 = f_2(W_2 f_1(W_1 x + b_1) + b_2)$$



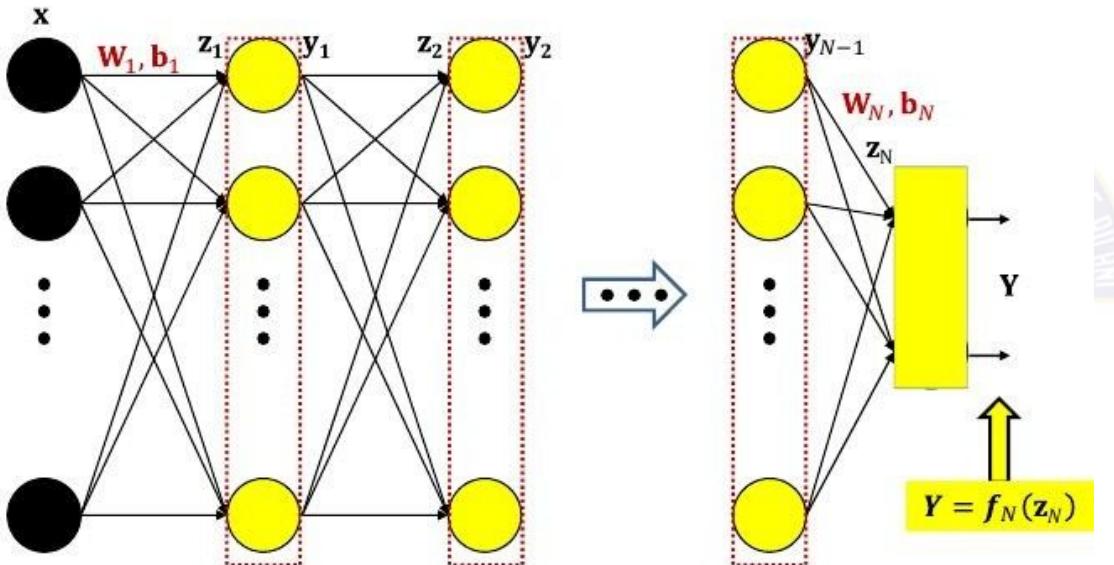
The forward pass



The Complete computation

$$y_2 = f_2(W_2 f_1(W_1 x + b_1) + b_2)$$

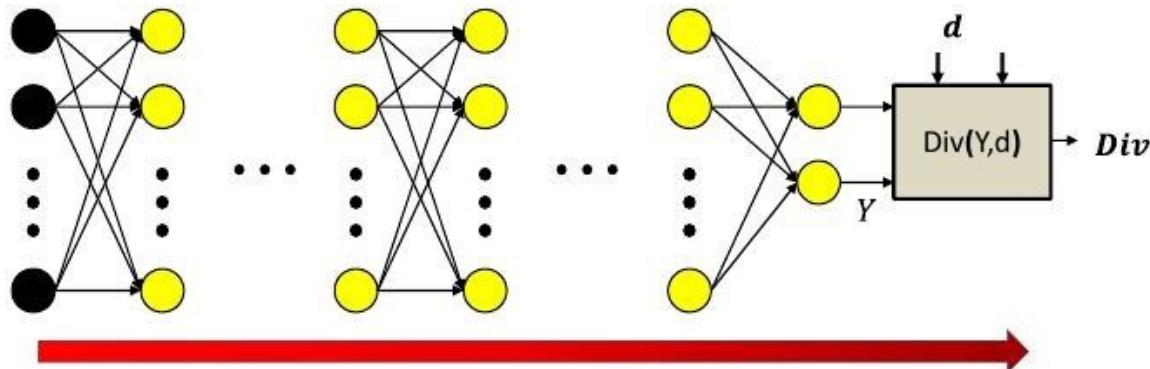
The forward pass



The Complete computation

$$\mathbf{Y} = f_N(\mathbf{W}_N f_{N-1}(\dots f_2(\mathbf{W}_2 f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) \dots) + \mathbf{b}_N)$$

Forward pass



Forward pass:

Initialize

$$\mathbf{y}_0 = \mathbf{x}$$

For $k = 1$ to N :

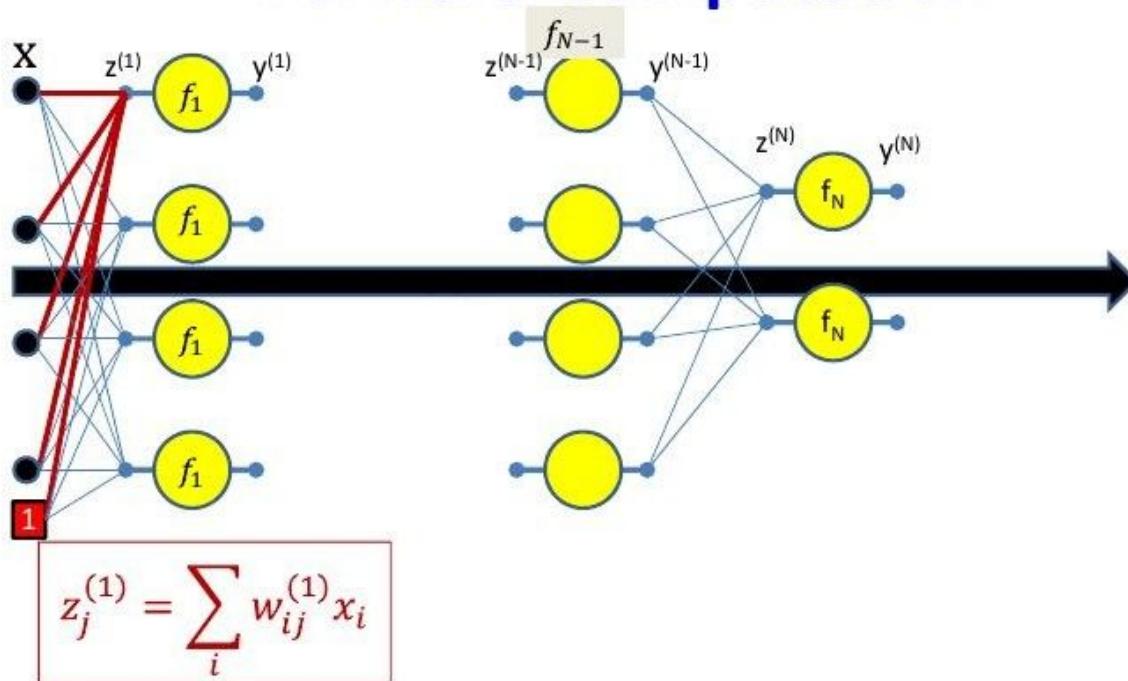
$$\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k$$

$$\mathbf{y}_k = f_k(\mathbf{z}_k)$$

Output

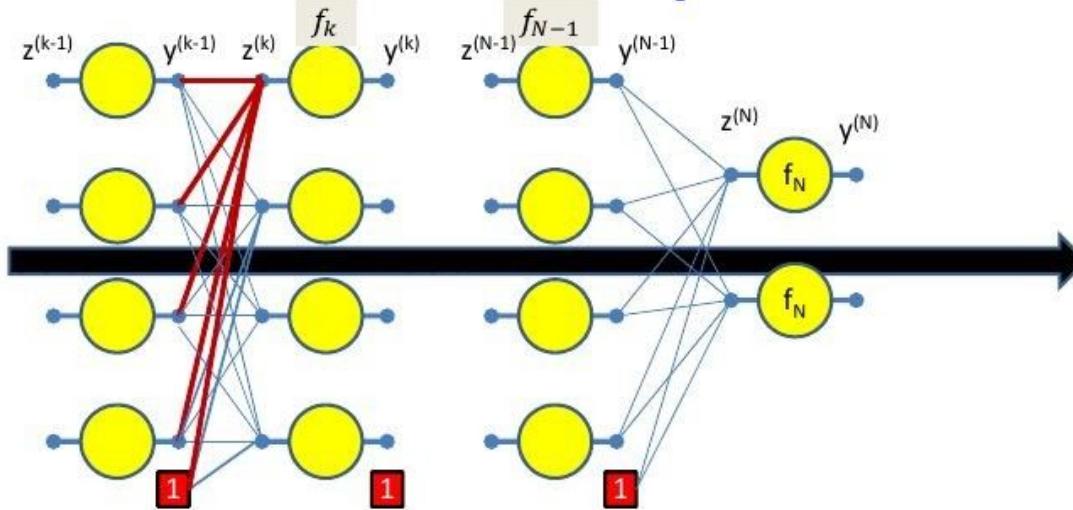
$$\mathbf{Y} = \mathbf{y}_N$$

Forward Computation



Assuming $w_{0j}^{(1)} = b_j^{(1)}$ and $x_0 = 1$

Forward Computation

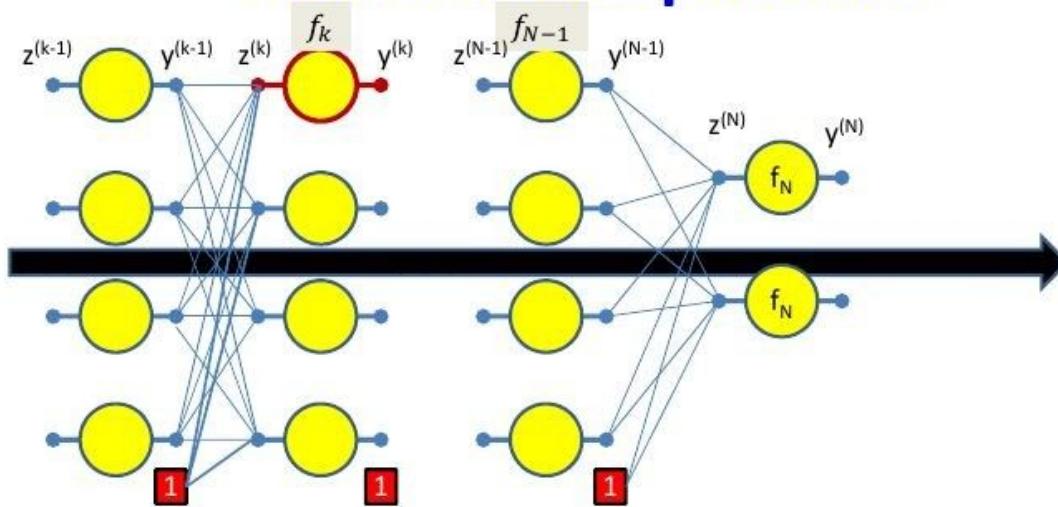


$$z_j^{(1)} = \sum_i w_{ij}^{(1)} x_i$$

$$z_j^{(k)} = \sum_i w_{ij}^{(k)} y_j^{(k-1)}$$

Assuming $w_{0j}^{(k)} = b_j^{(k)}$ and $y_0^{(k-1)} = 1$

Forward Computation

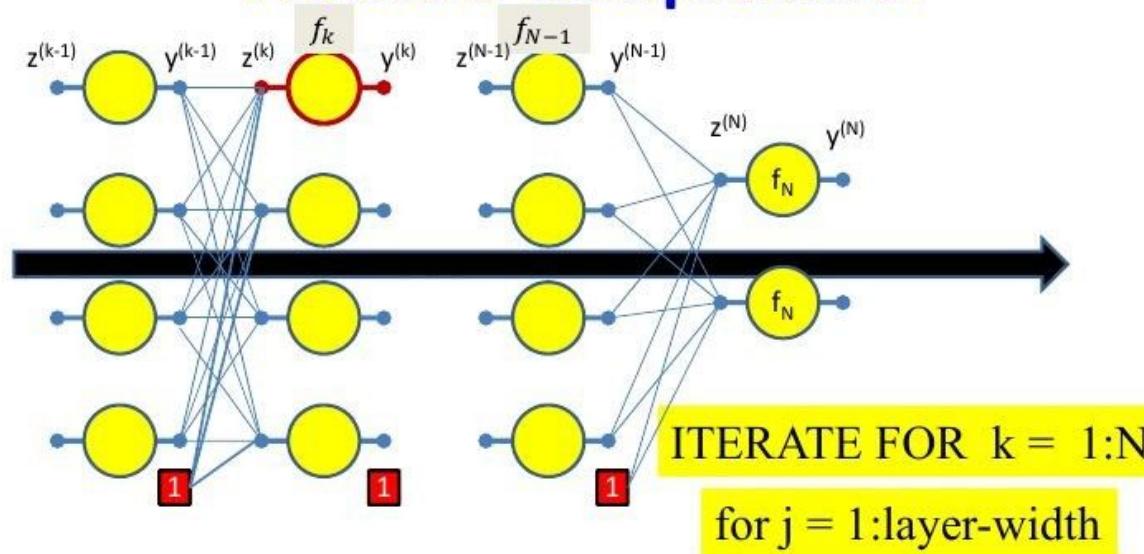


$$z_j^{(1)} = \sum_i w_{ij}^{(1)} x_i$$

$$z_j^{(k)} = \sum_i w_{ij}^{(k)} y_j^{(k-1)}$$

$$y_j^{(k)} = f_k(z_j^{(k)})$$

Forward Computation



$$y_i^{(0)} = x_i$$

$$z_j^{(k)} = \sum_i w_{ij}^{(k)} y_i^{(k-1)}$$

$$y_j^{(k)} = f_k(z_j^{(k)})$$

Forward Propagation Algorithm

Algorithm 1: FORWARD PROPAGATION

1 Initialize the weights and bias randomly.

2 $a^{[0]} = X^{(i)}$

3 $Z^{[1]} = W^{[1]}a^{[0]} + b^{[1]}$

4 $a^{[1]} = \sigma(Z^{[1]})$

5 $Z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$

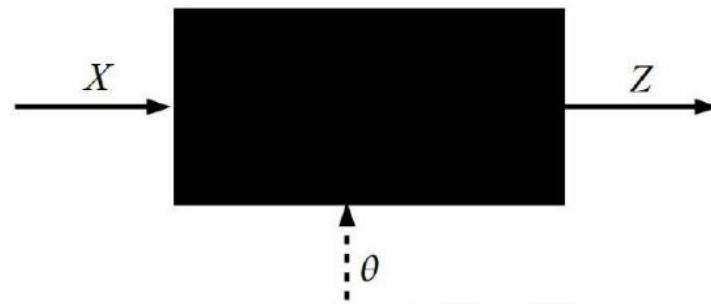
6 $a^{[2]} = \sigma(Z^{[2]})$

7 $\hat{y} = \begin{cases} 1 & \text{if } a^{[2]} > 0.5 \\ 0 & \text{otherwise} \end{cases}$



Note: All training examples are considered in the vectorized form.

Computation Graph for Forward Pass



Cost Function

- The difference between the actual observation y and the computed activation \hat{y} gives the error or the cost function.
- For binary classification,

$$J = \frac{1}{m} \sum_{i=0}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$J = -\frac{1}{m} \sum_{i=0}^m (y^{(i)} \log(a^{[2](i)}) + (1 - y^{(i)}) \log(1 - a^{[2](i)}))$$

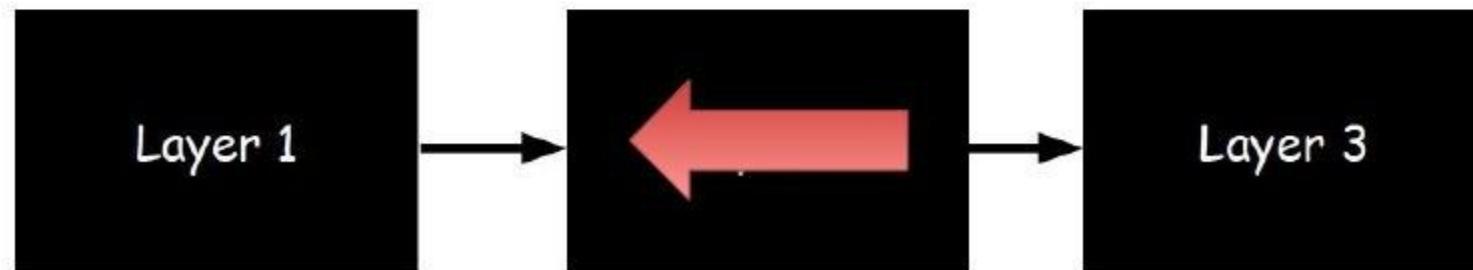
Neural Network Training – Backward Pass

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients w.r.t. parameters [B-Pass]



Neural Network Training – Backward Pass

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients w.r.t. parameters [B-Pass]

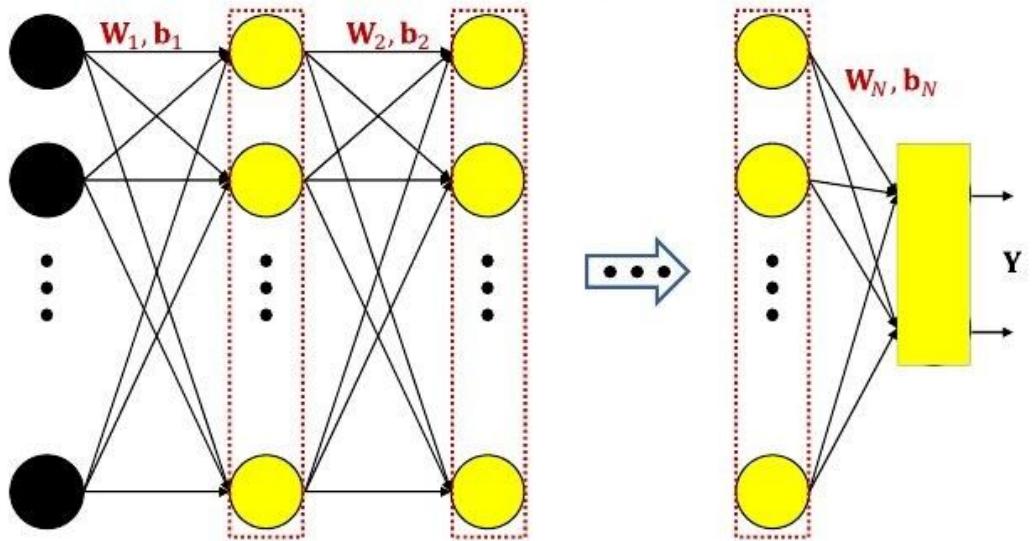


Neural Network Training – Backward Pass

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients w.r.t. parameters [B-Pass]



The backward pass



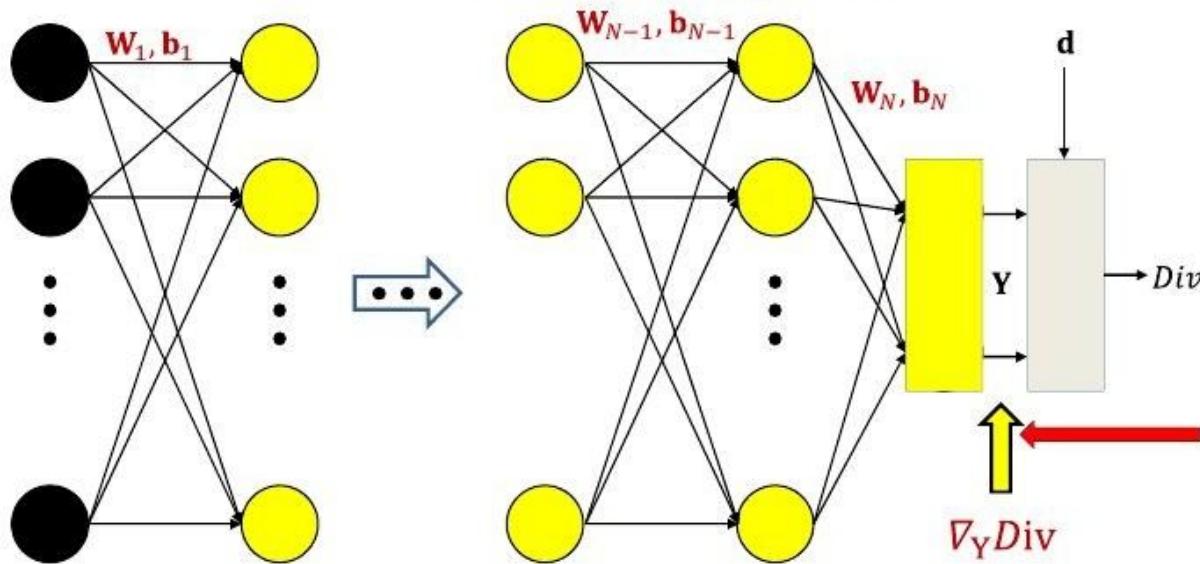
- The network is a nested function

$$Y = f_N(W_N f_{N-1}(\dots f_2(W_2 f_1(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) \dots) + \mathbf{b}_N)$$

- The error for any \mathbf{x} is also a nested function

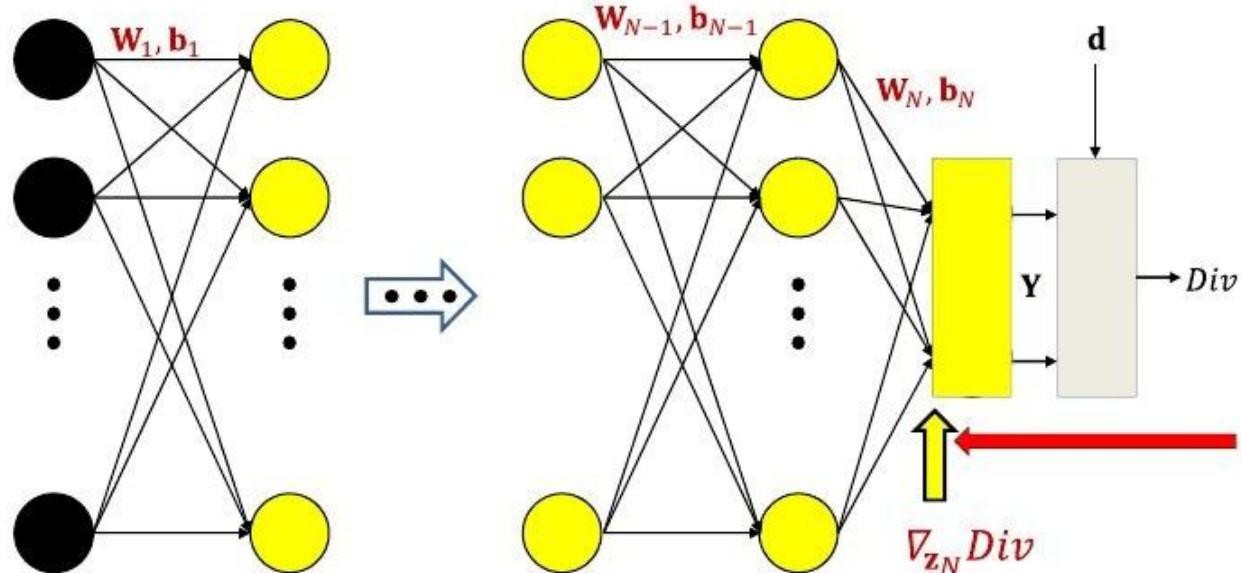
$$\text{Div}(Y, d) = \text{Div}(f_N(W_N f_{N-1}(\dots f_2(W_2 f_1(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) \dots) + \mathbf{b}_N), d)$$

The backward pass



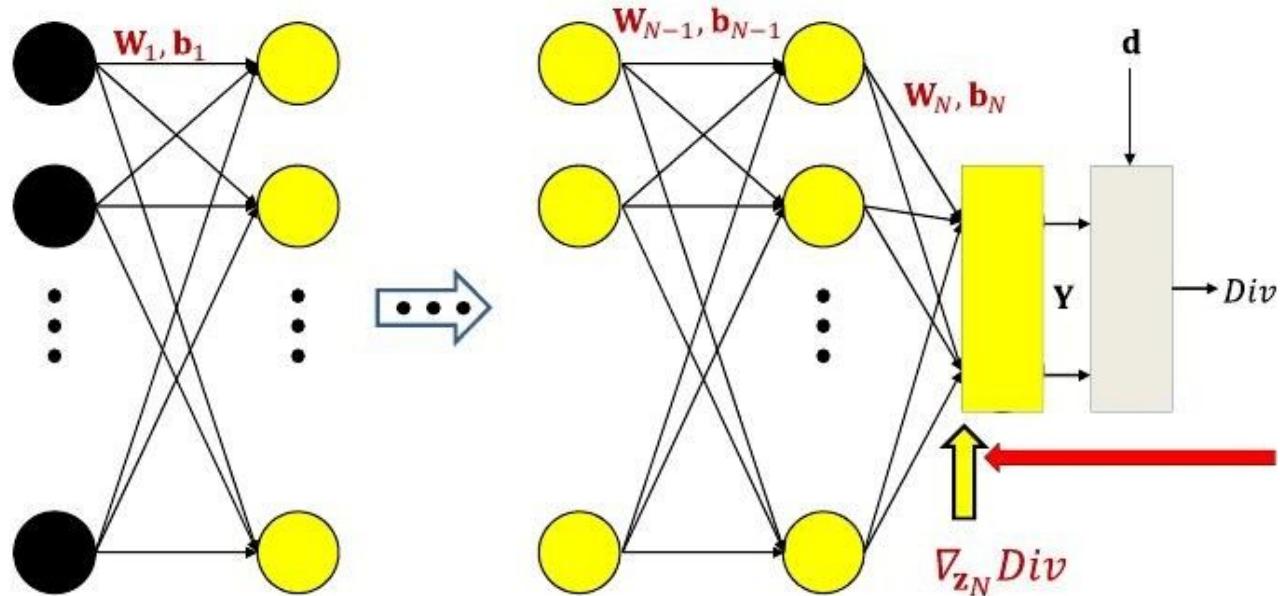
First compute the gradient of the divergence w.r.t. \mathbf{Y} .
The actual gradient depends on the divergence function.

The backward pass



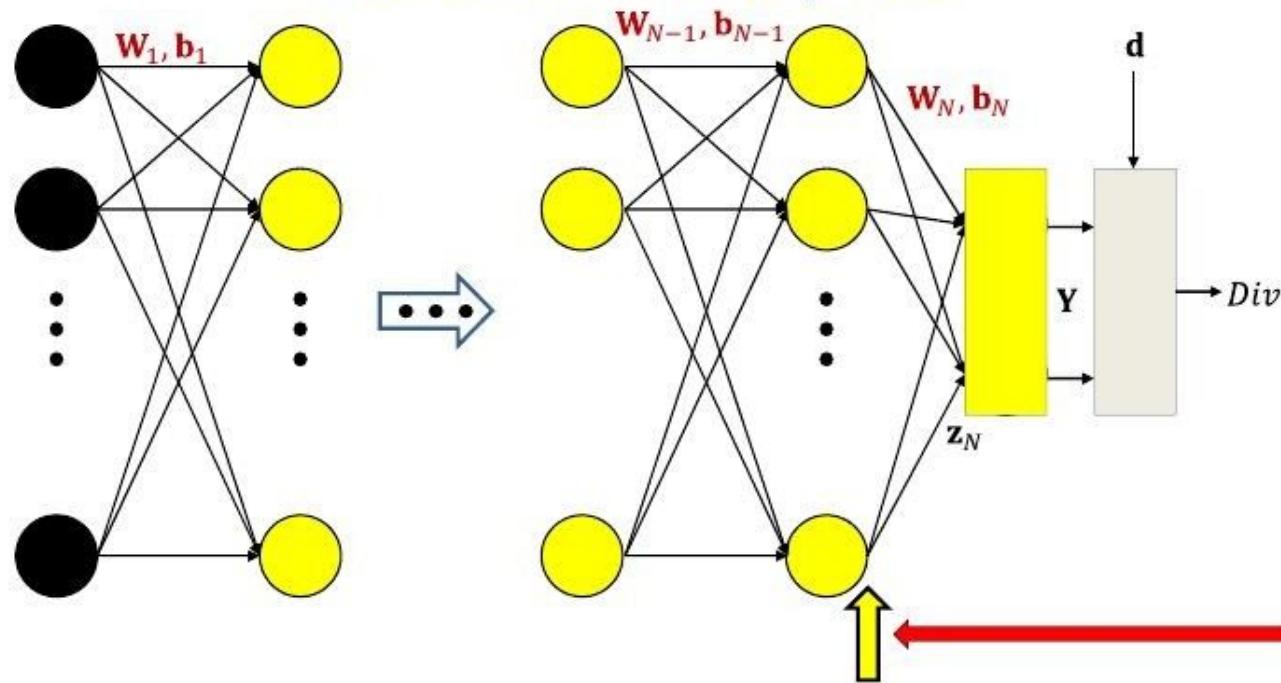
$$\nabla_{\mathbf{z}_N} Div = \nabla_{\mathbf{Y}} Div \cdot \nabla_{\mathbf{z}_N} \mathbf{Y}$$

The backward pass



$$\nabla_{\mathbf{z}_N} Div = \nabla_{\mathbf{Y}} Div J_{\mathbf{Y}}(\mathbf{z}_N)$$

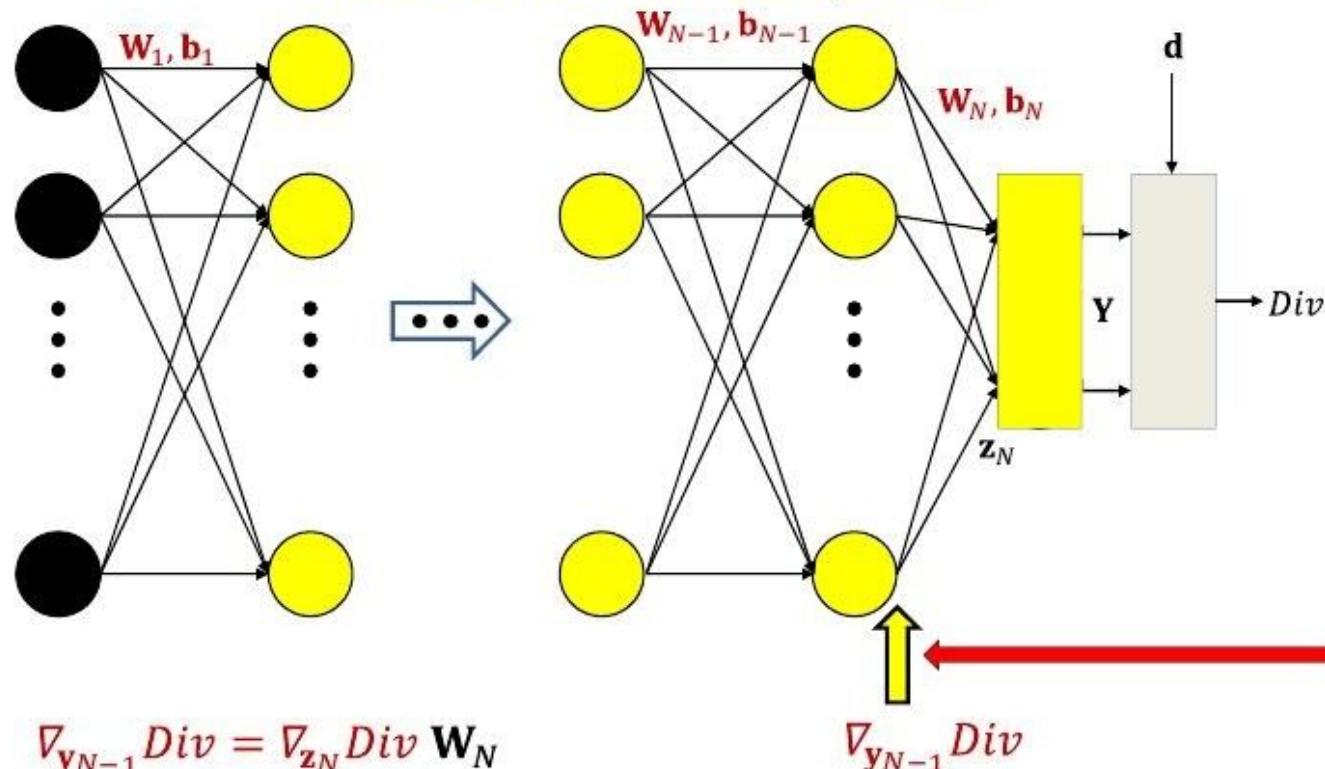
The backward pass



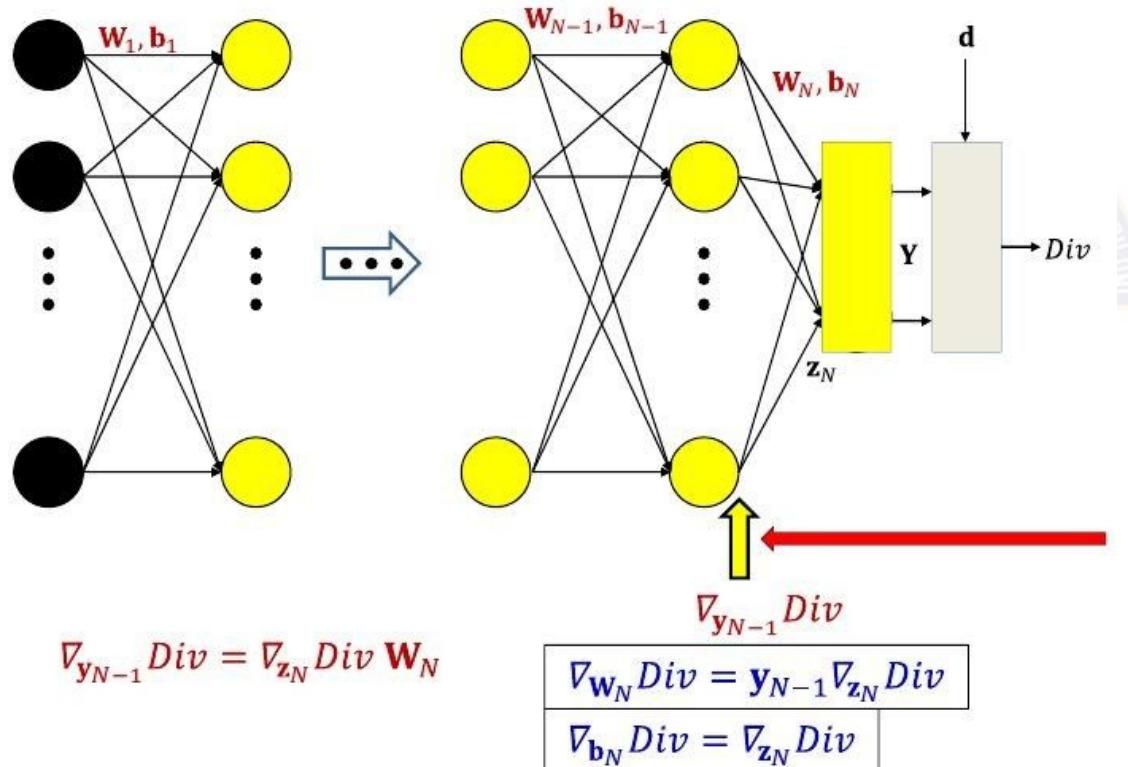
$$\nabla_{\mathbf{y}_{N-1}} Div = \nabla_{\mathbf{z}_N} Div \cdot \nabla_{\mathbf{y}_{N-1}} \mathbf{z}_N$$

$$\nabla_{\mathbf{y}_{N-1}} Div$$

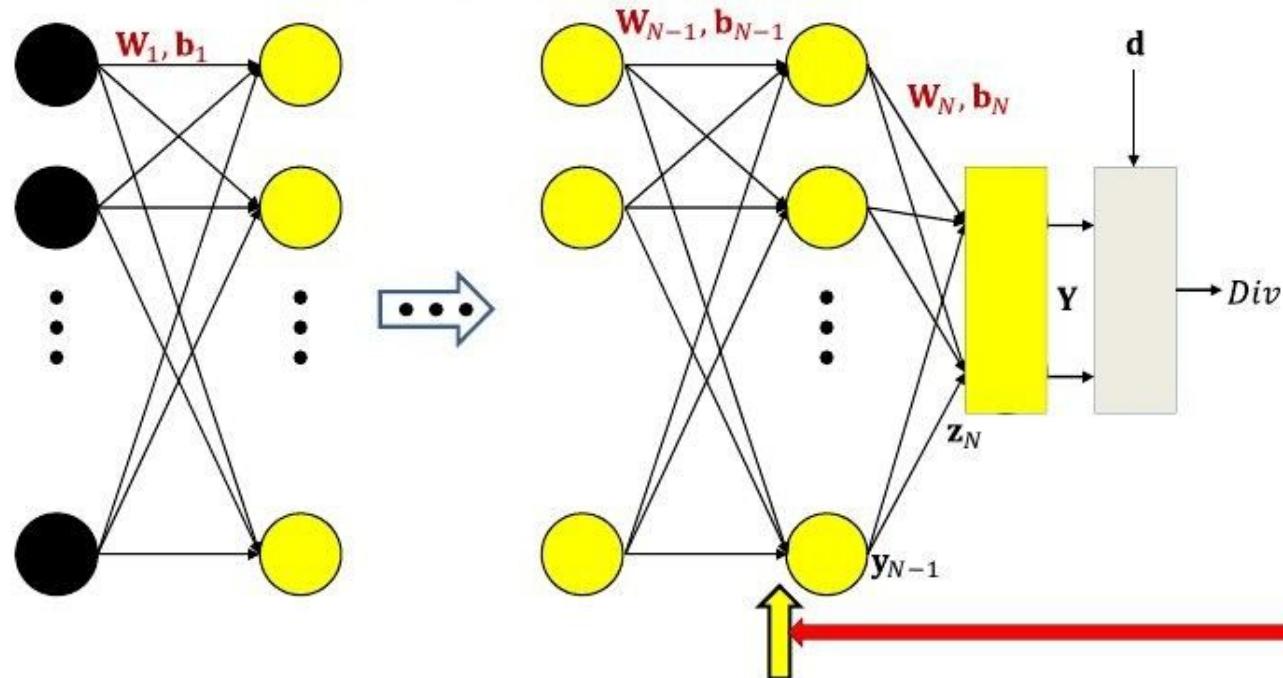
The backward pass



The backward pass



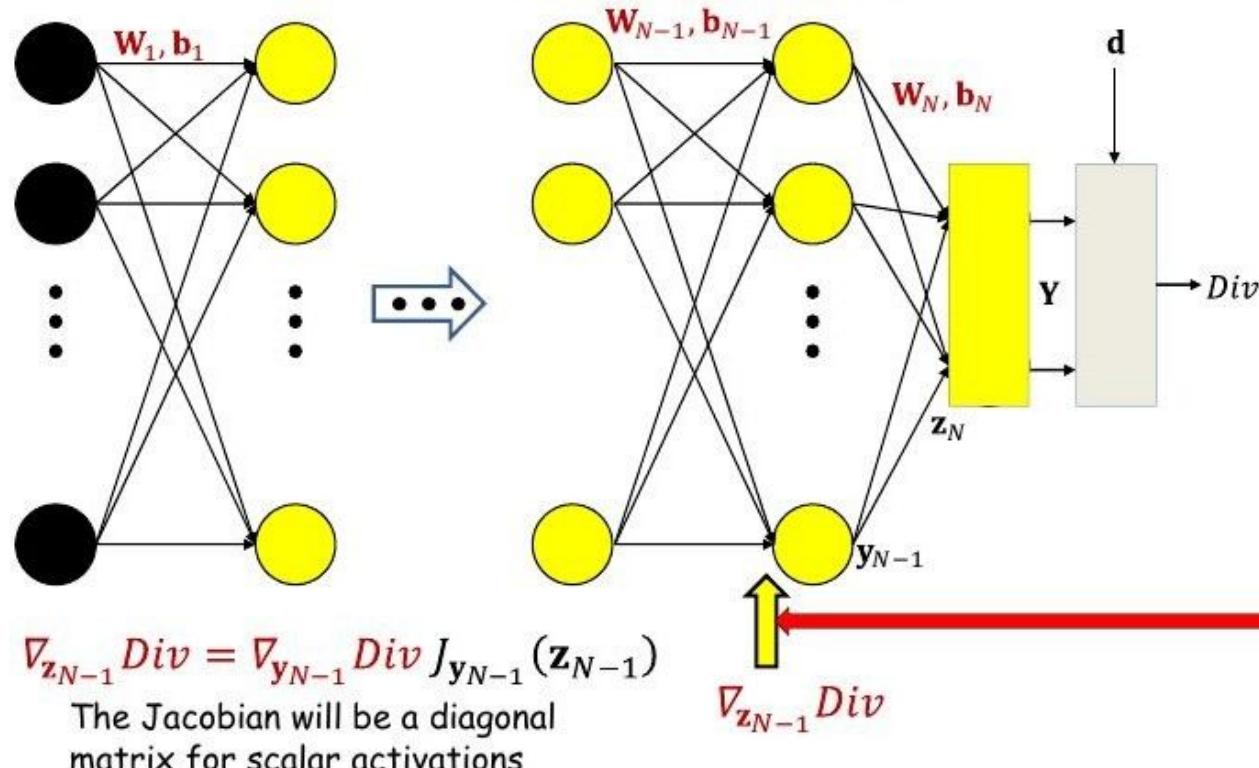
The backward pass



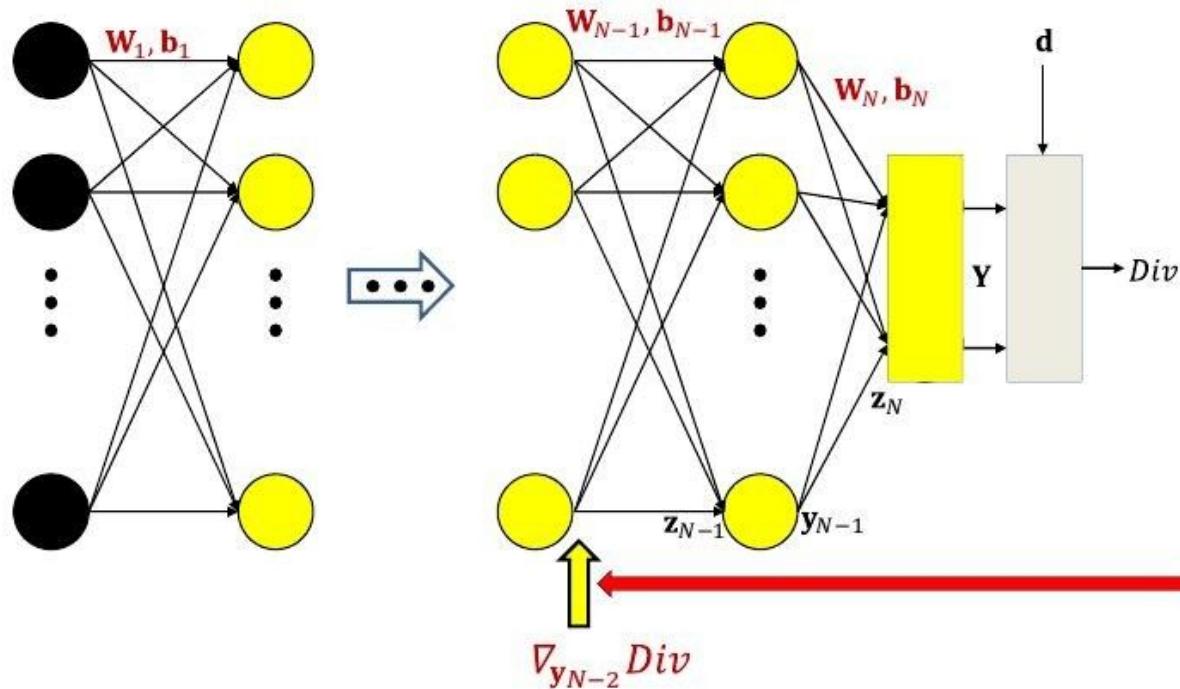
$$\nabla_{z_{N-1}} \text{Div} = \nabla_{y_{N-1}} \text{Div} \cdot \nabla_{z_{N-1}} y_{N-1}$$

$$\nabla_{z_{N-1}} \text{Div}$$

The backward pass

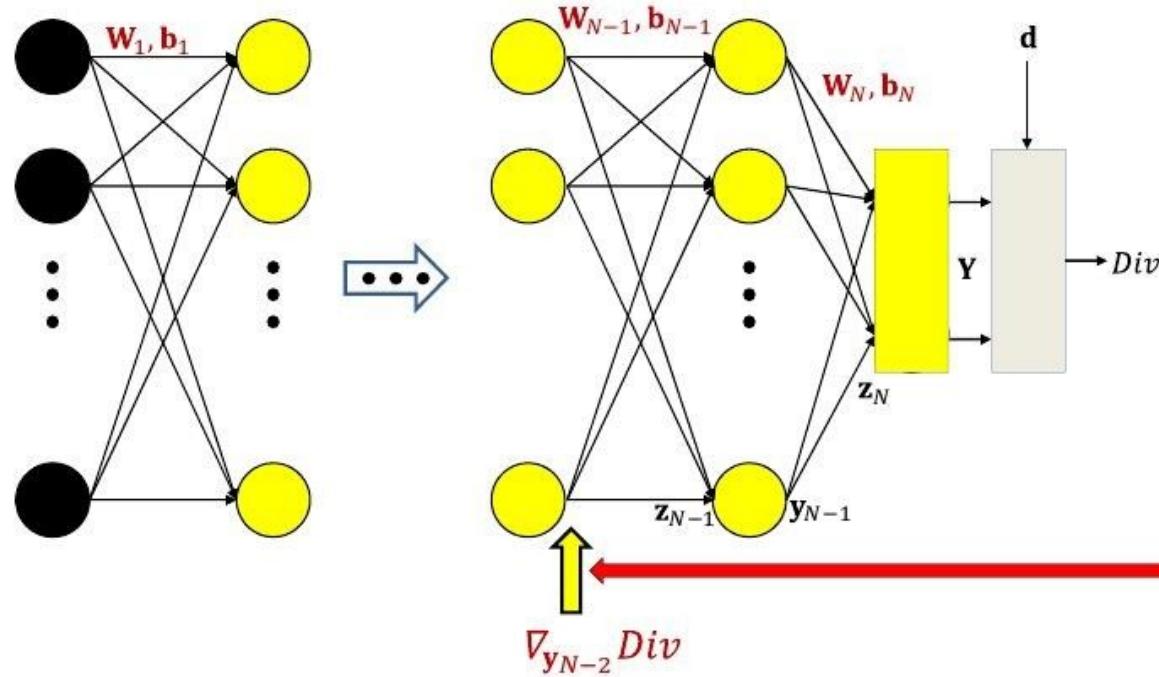


The backward pass



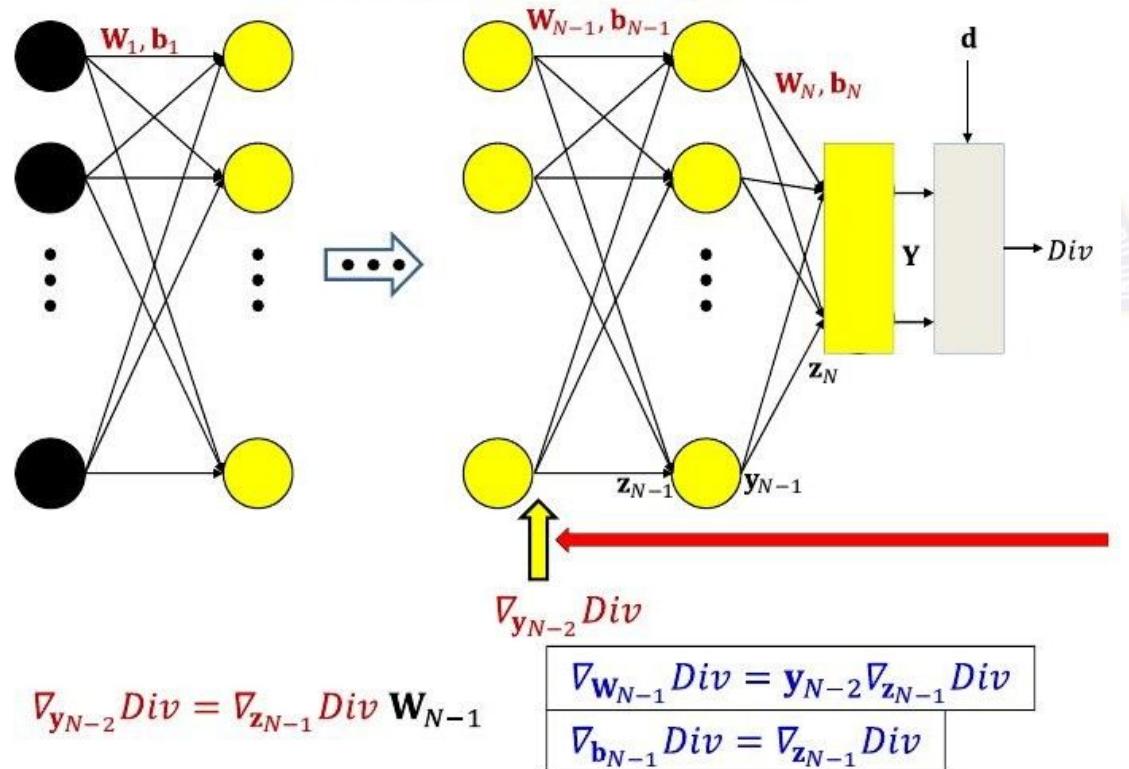
$$\nabla_{y_{N-2}} Div = \nabla_{z_{N-1}} Div \cdot \nabla_{y_{N-2}} z_{N-1}$$

The backward pass

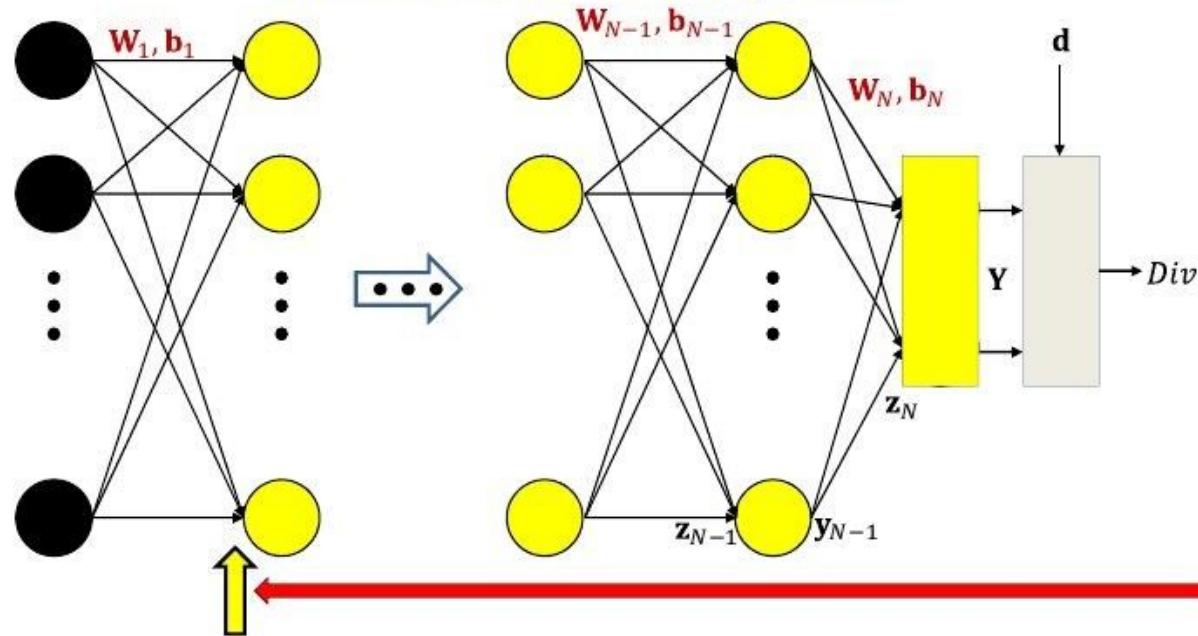


$$\nabla_{y_{N-2}} Div = \nabla_{z_{N-1}} Div W_{N-1}$$

The backward pass

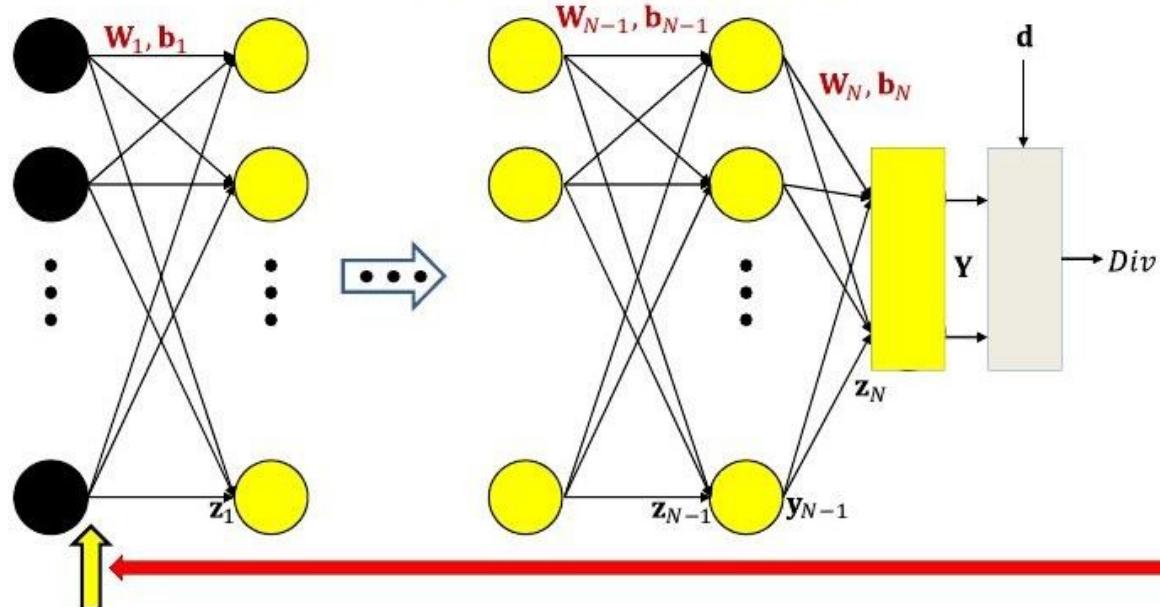


The backward pass



$$\nabla_{\mathbf{z}_1} Div = \nabla_{\mathbf{y}_1} Div J_{\mathbf{y}_1}(\mathbf{z}_1)$$

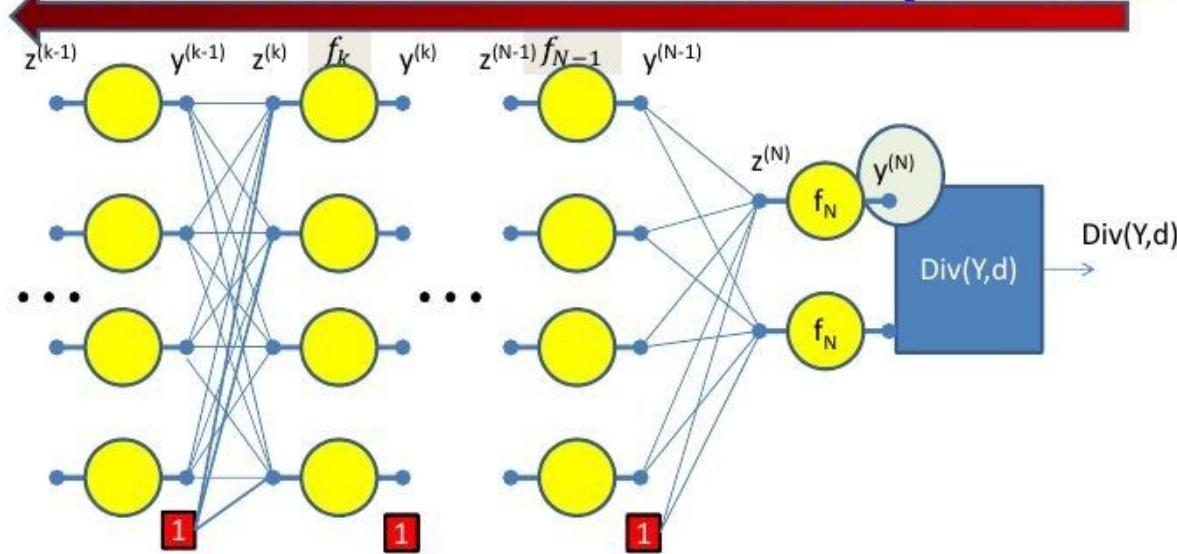
The backward pass



$$\nabla_{\mathbf{W}_1} Div = \mathbf{x} \nabla_{\mathbf{z}_1} Div$$
$$\nabla_{\mathbf{b}_1} Div = \nabla_{\mathbf{z}_1} Div$$

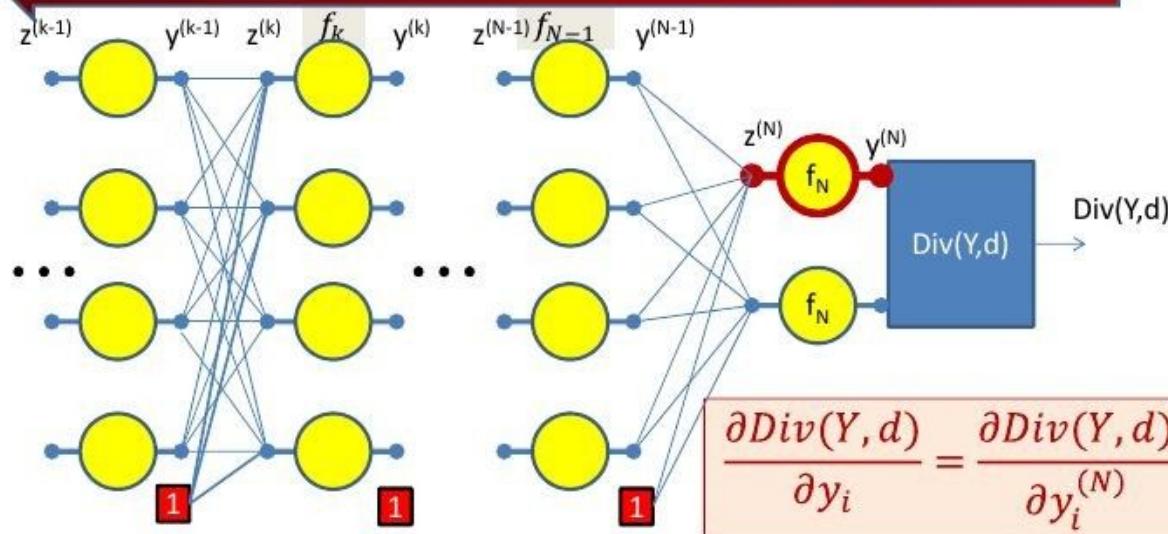
In some problems we will also want to compute the derivative w.r.t. the input

Gradients: Backward Computation



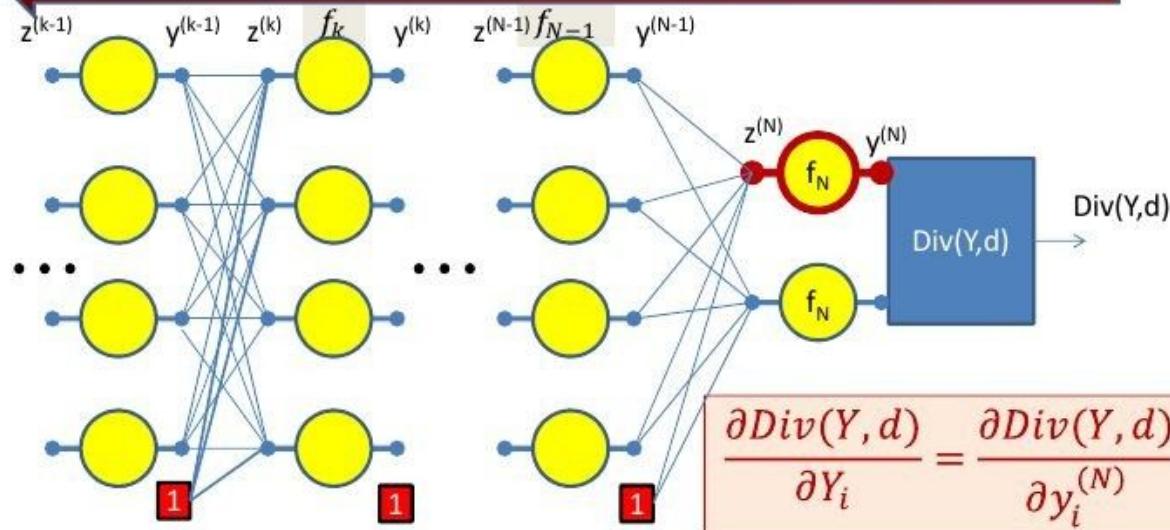
$$\frac{\partial \text{Div}(Y, d)}{\partial y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

Gradients: Backward Computation



$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial \text{Div}}{\partial y_i} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

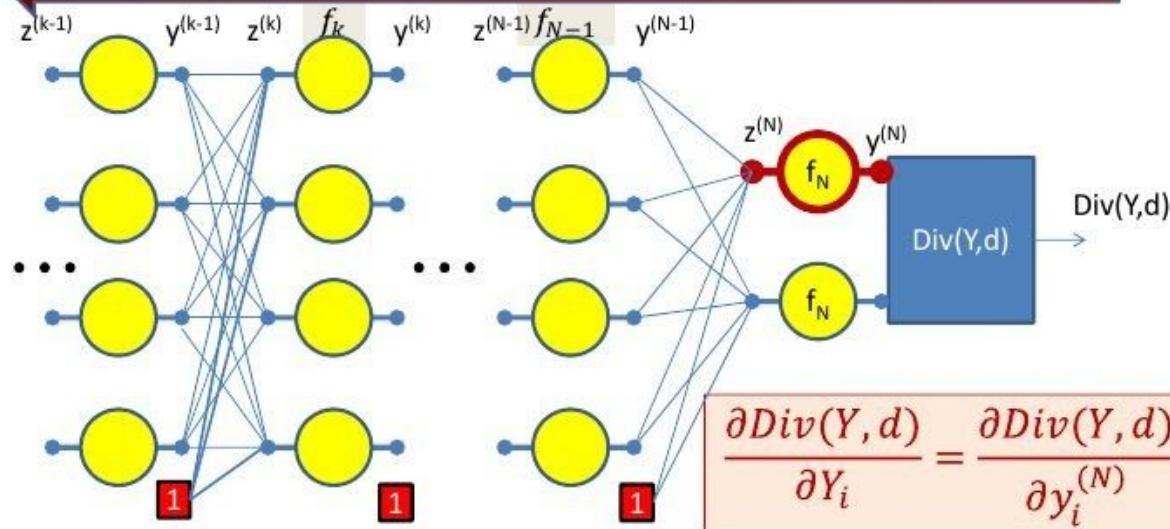
Gradients: Backward Computation



$z_i^{(N)}$ computed during the forward pass

$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial \text{Div}}{\partial Y_i} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

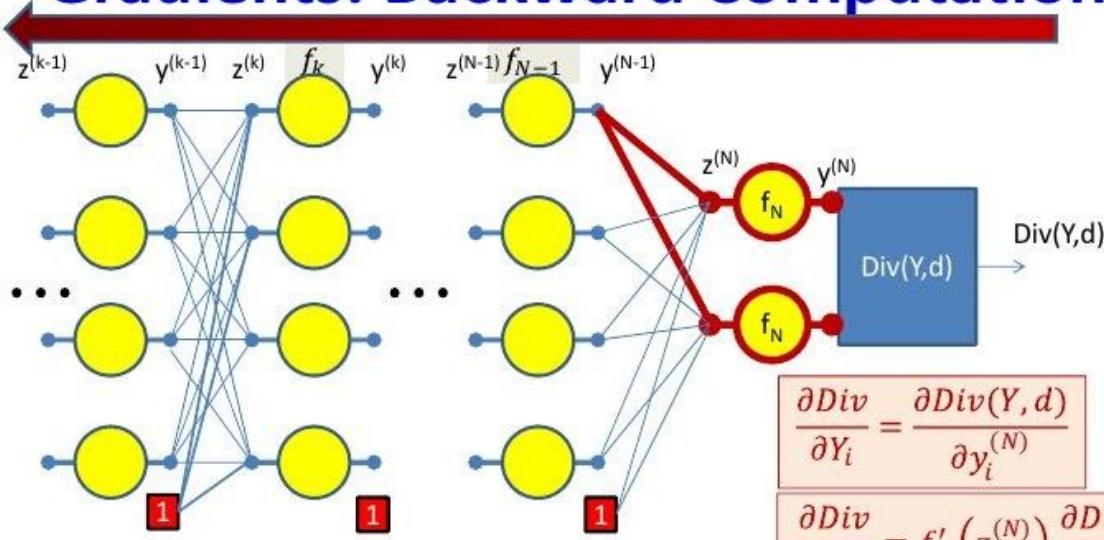
Gradients: Backward Computation



$z_i^{(N)}$ computed during the forward pass

$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}} \frac{\partial \text{Div}}{\partial Y_i} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

Gradients: Backward Computation



$$\frac{\partial \text{Div}}{\partial y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

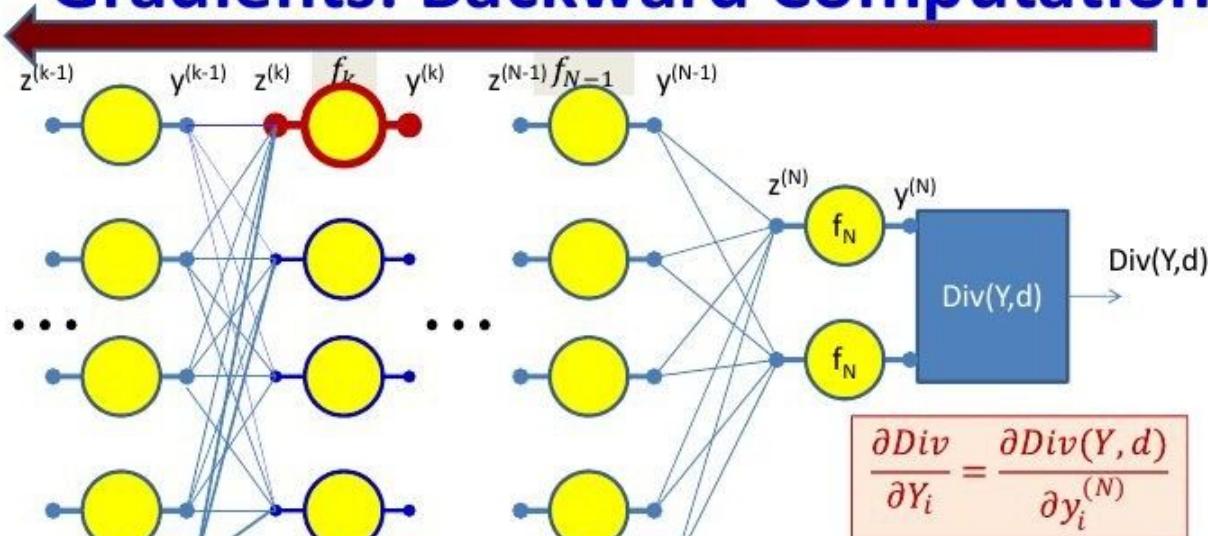
$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} \frac{\partial \text{Div}}{\partial z_j^{(N)}} = \sum_j w_{ij}^{(N)} \frac{\partial \text{Div}}{\partial z_j^{(N)}}$$

Because :

$$\frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}} = w_{ij}^{(N)}$$

Gradients: Backward Computation



$$\frac{\partial \text{Div}}{\partial Y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

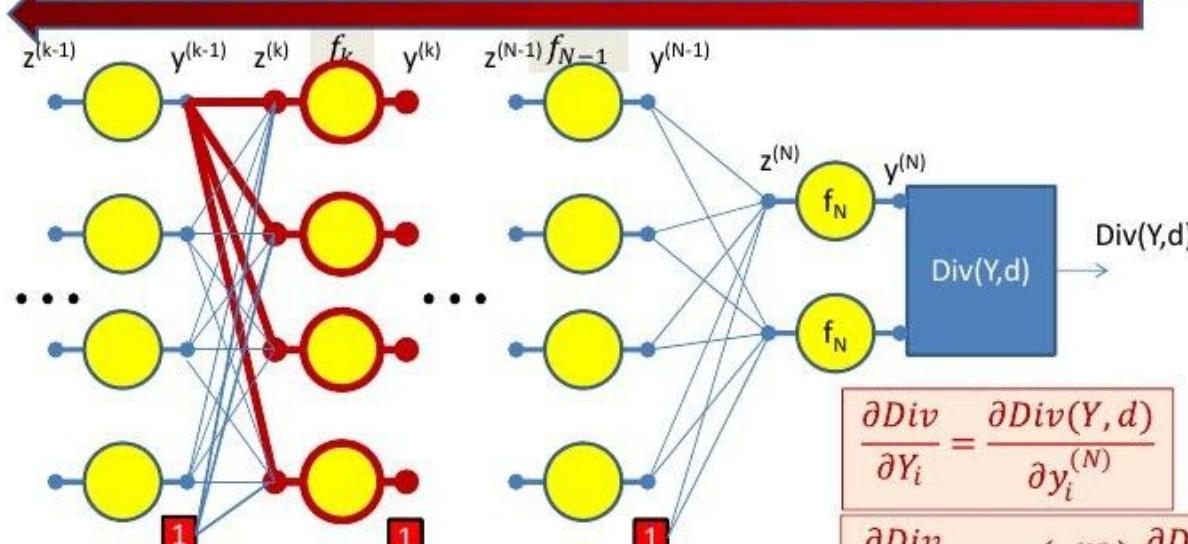
$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

computed during
the forward pass

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial \text{Div}}{\partial y_i^{(k)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(N-1)}} = \sum_j w_{ij}^{(N)} \frac{\partial \text{Div}}{\partial z_j^{(N)}}$$

Gradients: Backward Computation

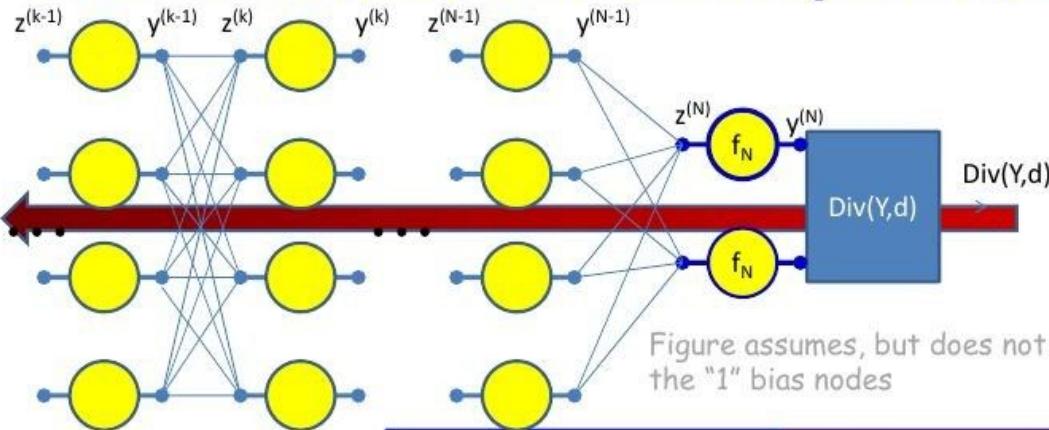


$$\frac{\partial \text{Div}}{\partial Y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(k-1)}} = \sum_j \frac{\partial z_j^{(k)}}{\partial y_i^{(k-1)}} \frac{\partial \text{Div}}{\partial z_j^{(k)}} = \sum_j w_{ij}^{(k)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$

Gradients: Backward Computation



Initialize: Gradient
w.r.t network output

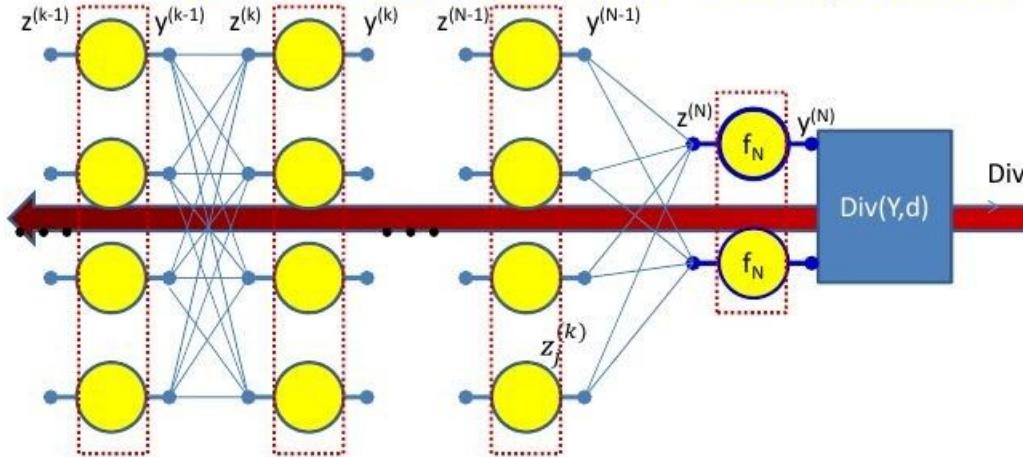
$$\frac{\partial \text{Div}}{\partial y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}}$$

$$\begin{aligned} &\text{For } k = N..1 \\ &\text{For } i = 1:\text{layer_width} \end{aligned} \quad \frac{\partial \text{Div}}{\partial z_i^{(k)}} = f'_k(z_i^{(k)}) \frac{\partial \text{Div}}{\partial y_i^{(k)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$

$$\frac{\partial \text{Div}}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$

Gradients: Backward Computation



For $k = N \dots 1$

For $i = 1 : \text{layer-width}$

If layer has vector activation

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} \frac{\partial y_j^{(k)}}{\partial z_i^{(k)}}$$

Else if activation is scalar

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \frac{\partial \text{Div}}{\partial y_i^{(k)}} \frac{\partial y_i^{(k)}}{\partial z_i^{(k)}}$$

$$\frac{\partial \text{Div}}{\partial y_i^{(k-1)}} = \sum_j w_{ij}^{(k)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$

$$\frac{\partial \text{Div}}{\partial w_{ij}^{(k)}} = y_i^{(k-1)} \frac{\partial \text{Div}}{\partial z_j^{(k)}}$$

Computing the Gradients

For all training examples

For one training example

$$dz^{[2]} = a^{[2]} - y$$

$$dw^{[2]} = dz^{[2]} \cdot a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = (w^{[2]T} dz^{[2]}) * \sigma'(z^{[1]})$$

$$dw^{[1]} = dz^{[1]} \cdot x^T$$

$$db^{[1]} = dz1$$

$$dZ^{[2]} = A^{[2]} - Y$$

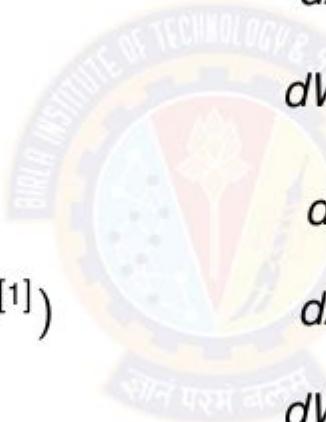
$$dW^{[2]} = \frac{1}{m} (dZ^{[2]} \cdot A^{[1]T})$$

$$db^{[2]} = \frac{1}{m} \sum dZ^{[2]}$$

$$dZ^{[1]} = (W^{[2]T} dZ^{[2]}) * \sigma'(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} (dZ^{[1]} \cdot X^T)$$

$$db^{[1]} = \frac{1}{m} \sum dZ1$$



Backpropagation Algorithm to compute Gradients

Algorithm 2: BACKWARD PROPAGATION

$$1 \quad dZ^{[2]} = A^{[2]} - Y$$

$$2 \quad dW^{[2]} = \frac{1}{m} (dZ^{[2]} \cdot A^{[1]T})$$

$$3 \quad db^{[2]} = \frac{1}{m} \sum dZ^{[2]}$$

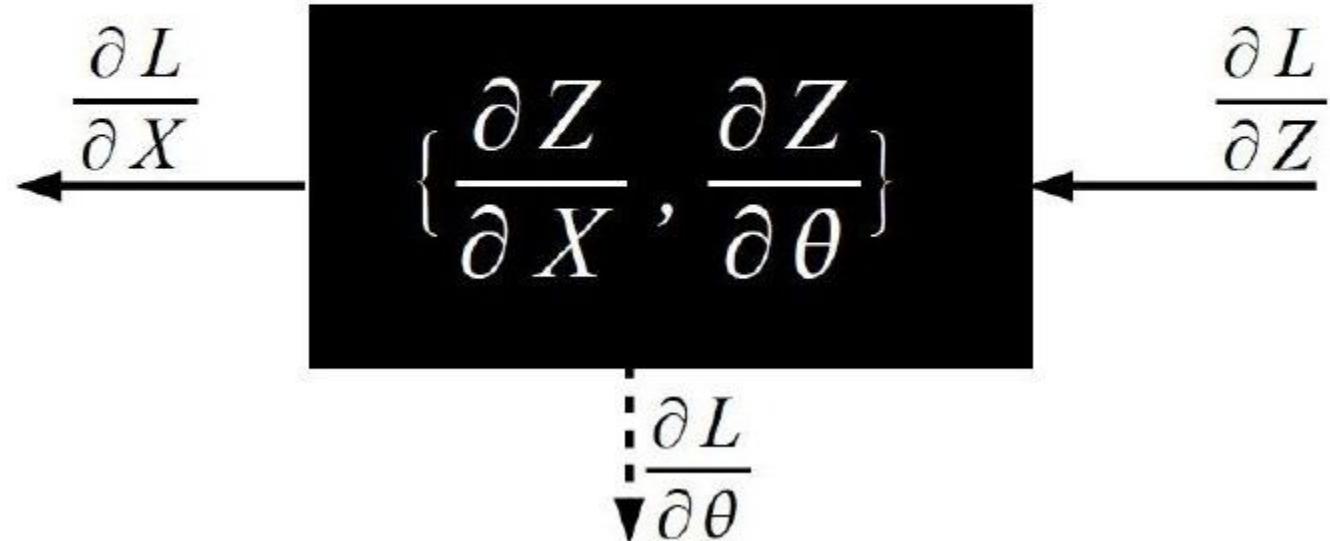
$$4 \quad dZ^{[1]} = (W^{[2]T} dZ^{[2]}) * \sigma'(Z^{[1]})$$

$$5 \quad dW^{[1]} = \frac{1}{m} (dZ^{[1]} \cdot X^T)$$

$$6 \quad db^{[1]} = \frac{1}{m} \sum dZ^{[1]}$$



Computation Graph for BackProp



Neural Network Training – Update Parameters

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients w.r.t. parameters [B-Pass]
- Step 3: Use gradient to update parameters



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

*** θ same as w

Parameter Updation

Algorithm 3: WEIGHT UPDATION

$$1 \quad W^{[1]} = W^{[1]} - \text{learning rate} * dW^{[1]}$$

$$2 \quad b^{[1]} = b^{[1]} - \text{learning rate} * db^{[1]}$$

$$3 \quad W^{[2]} = W^{[2]} - \text{learning rate} * dW^{[2]}$$

$$4 \quad b^{[2]} = b^{[2]} - \text{learning rate} * db^{[2]}$$

Training of MultiLayer Perceptron (MLP)

Requires

1. **Forward Pass** through each layer to compute the output.
2. Compute the deviation or error between the desired output and computed output in the forward pass (first step). This morphs into **objective function**, as we want to minimize this deviation or error.
3. The deviation has to be send back through each layer to compute the delta or change in the parameter values. This is achieved using **back propagation algorithm**.
4. **Update** the parameters.

Neural network training algorithm

- Initialize all weights and biases ($\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_N, \mathbf{b}_N$)
- Do:

- $Err = 0$
 - For all k , initialize $\nabla_{\mathbf{W}_k} Err = 0, \nabla_{\mathbf{b}_k} Err = 0$
 - For all $t = 1:T$
 - Forward pass : Compute
 - Output $\mathbf{Y}(X_t)$
 - Divergence $Div(\mathbf{Y}_t, \mathbf{d}_t)$
 - $Err += Div(\mathbf{Y}_t, \mathbf{d}_t)$
 - Backward pass: For all k compute:
 - $\nabla_{\mathbf{W}_k} Div(\mathbf{Y}_t, \mathbf{d}_t); \nabla_{\mathbf{b}_k} Div(\mathbf{Y}_t, \mathbf{d}_t)$
 - $\nabla_{\mathbf{W}_k} Err += \nabla_{\mathbf{W}_k} Div(\mathbf{Y}_t, \mathbf{d}_t); \nabla_{\mathbf{b}_k} Err += \nabla_{\mathbf{b}_k} Div(\mathbf{Y}_t, \mathbf{d}_t)$
 - For all k , update:

$$\mathbf{W}_k = \mathbf{W}_k - \frac{\eta}{T} (\nabla_{\mathbf{W}_k} Err)^T; \quad \mathbf{b}_k = \mathbf{b}_k - \frac{\eta}{T} (\nabla_{\mathbf{b}_k} Err)^T$$

- Until Err has converged

Scaling up for L layers in MLP

Notations

- L : Number of layers
- $n^{[l]}$: Number of units in layer l
- $w^{[l]}$: Weights for layer l
- $b^{[l]}$: Bias for layer l
- $z^{[l]}$: Hypothesis for layer l
- $g^{[l]}$: Activation Function used for layer l
- $a^{[l]}$: Activation for layer l

Matrix Dimensions

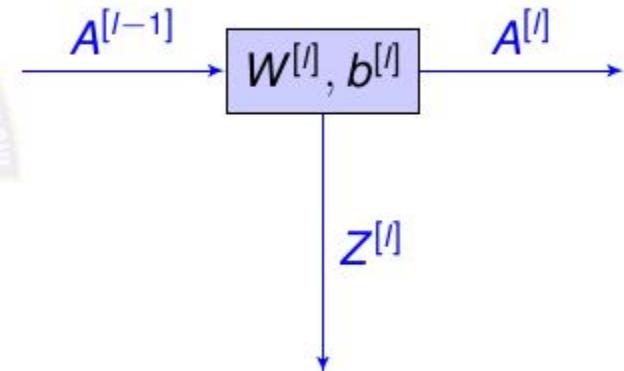
- $W^{[l]} : n^{[l]} \times n^{[l-1]}$
- $b^{[l]} : n^{[l]} \times 1$
- $z^{[l]} : n^{[l]} \times 1$
- $a^{[l]} : n^{[l]} \times 1$
- $dW^{[l]} : n^{[l]} \times n^{[l-1]}$
- $db^{[l]} : n^{[l]} \times 1$
- $dz^{[l]} : n^{[l]} \times 1$
- $da^{[l]} : n^{[l]} \times 1$



Forward Propagation Algorithm

Algorithm 4: Forward Propagation

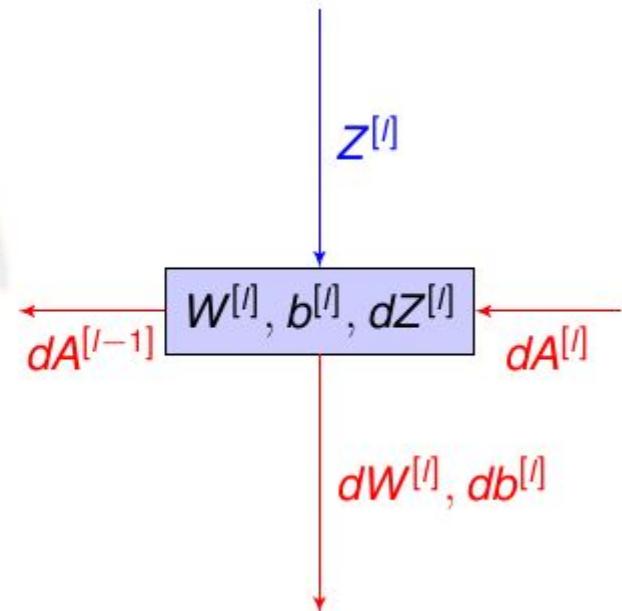
```
1 for  $l$  in range  $(1, L)$  do
2   |  $Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]}$ 
3   |  $A^{[l]} = g^{[l]}(Z^{[l]})$ 
4    $\hat{y} = A^{[L]}$ 
5    $J = L(\hat{y}, y) + \lambda\Omega(\theta)$ 
```



Backward Propagation Algorithm

Algorithm 5: Backward Propagation

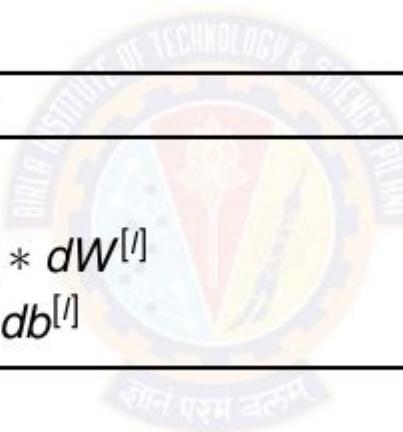
```
1  $da^{[L]} = \frac{-y}{a} + \frac{1-y}{1-a}$ 
2 for  $l$  in range  $(1, L)$  do
3    $dZ^{[l]} = dA^{[l]} * g'^{[l]}(Z^{[l]})$ 
4    $dW^{[l]} = dZ^{[l]} \cdot A^{[l-1]T}$ 
5    $db^{[l]} = dZ^{[l]}$ 
6    $dA^{[l-1]} = W^{[l]T} \cdot dZ^{[l]}$ 
```



Update the Parameters

Algorithm 6: Weight Updation

```
1 for  $l$  in range ( $1, L$ ) do
2      $W^{[l]} = W^{[l]} - \text{learning rate} * dW^{[l]}$ 
3      $b^{[l]} = b^{[l]} - \text{learning rate} * db^{[l]}$ 
```



Ref:
Chapter 3 and 4 of T1



Next Session: Power of MLP





BITS Pilani
Pilani | Dubai | Goa | Hyderabad

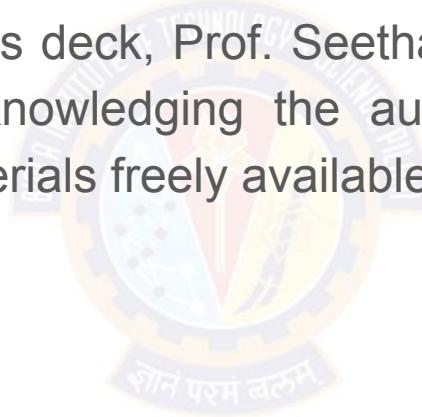
Deep Learning

DSE Module 2

Seetha Parameswaran

BITS Pilani

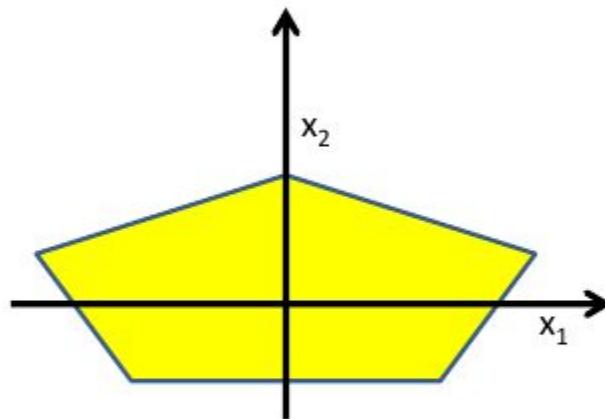
The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Complex Boundaries Using MLP



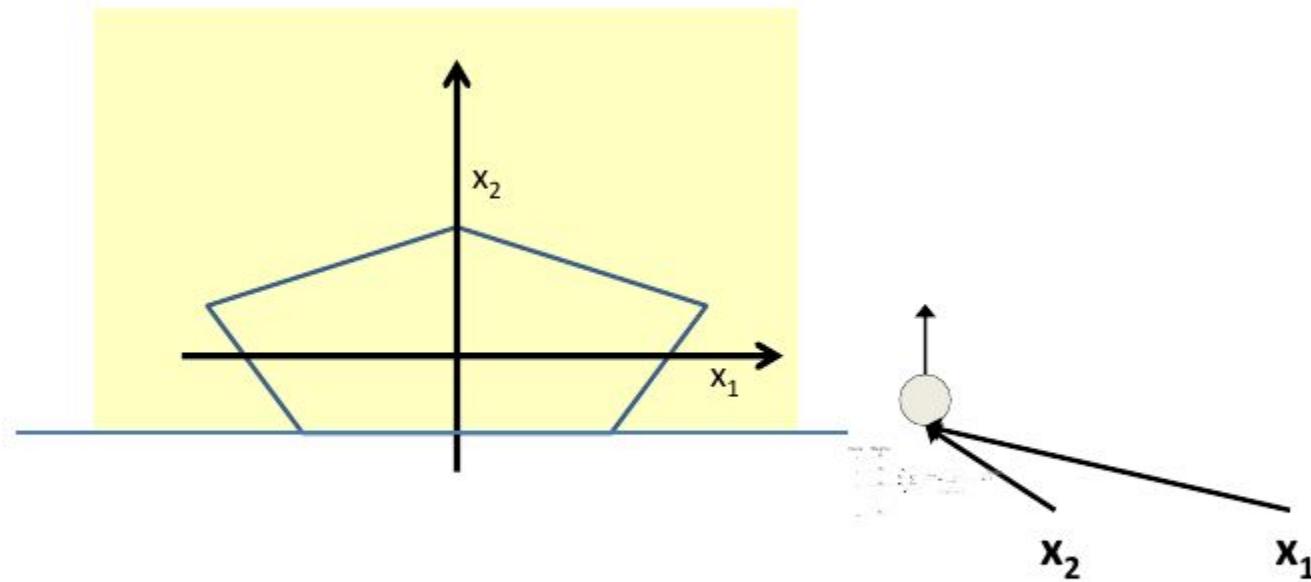
Composing Complex Decision Boundaries



Can now be composed into "networks" to compute arbitrary classification "boundaries"

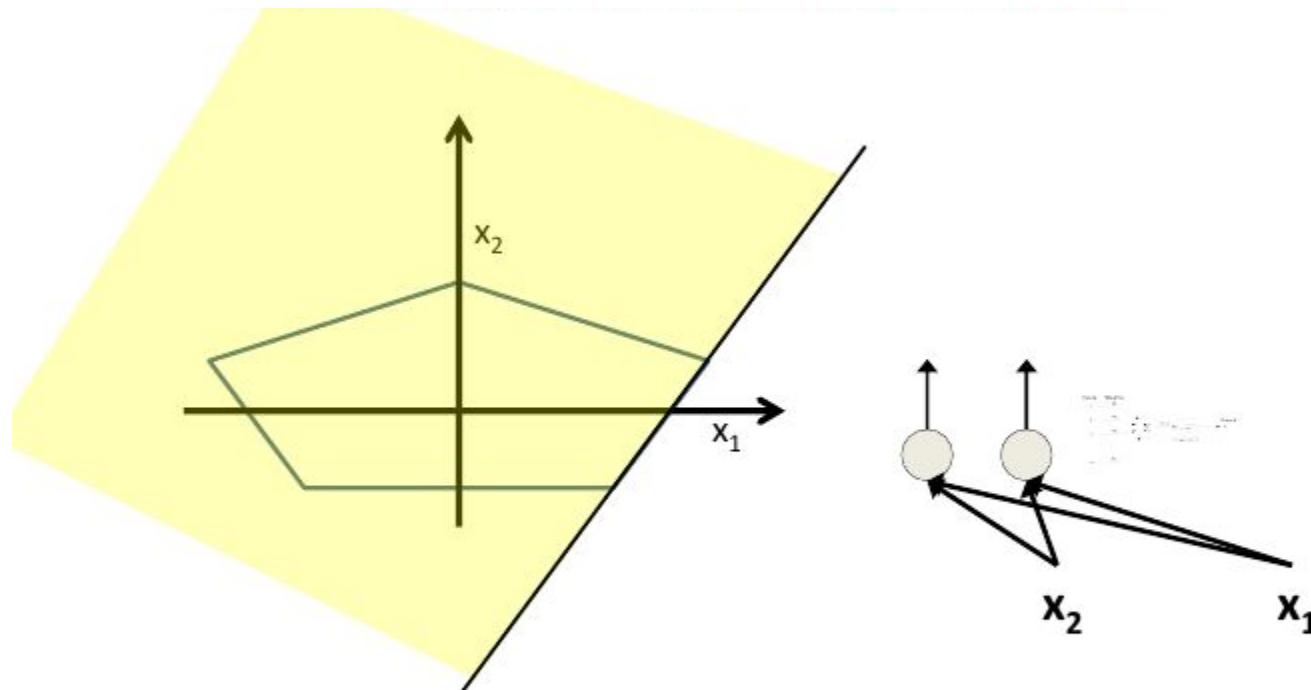
- Build a network of units with a single output that fires if the input is in the coloured area

Boolean over the Reals



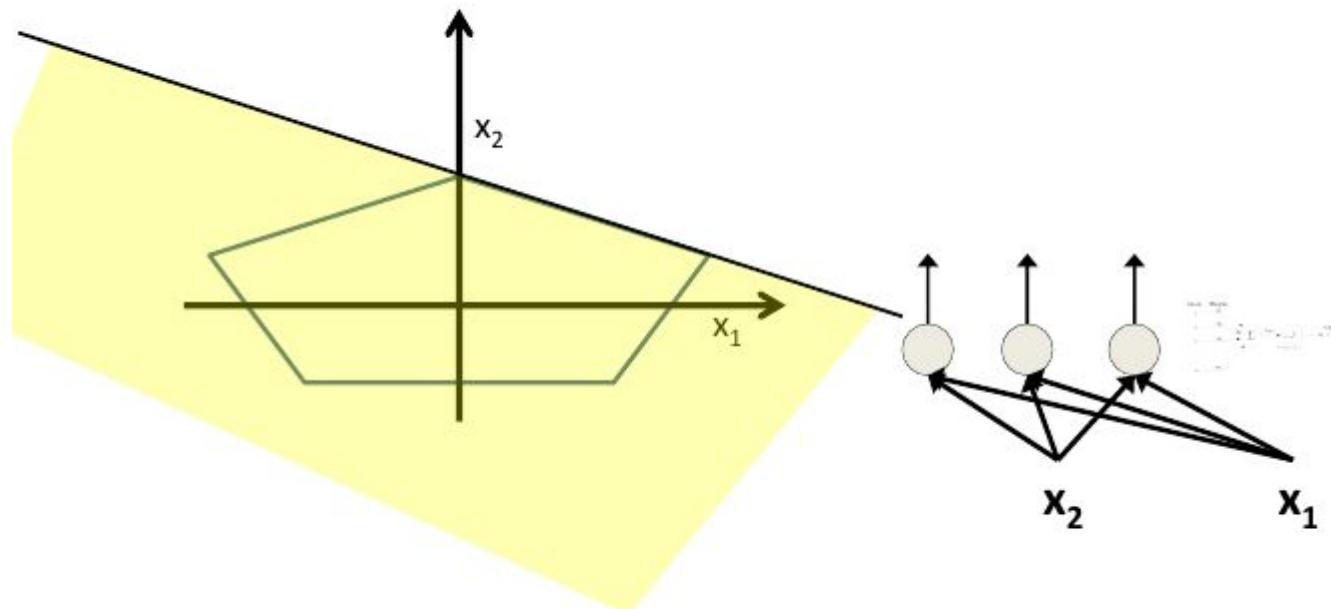
- The network must fire if the input is in the coloured area

Boolean over the Reals



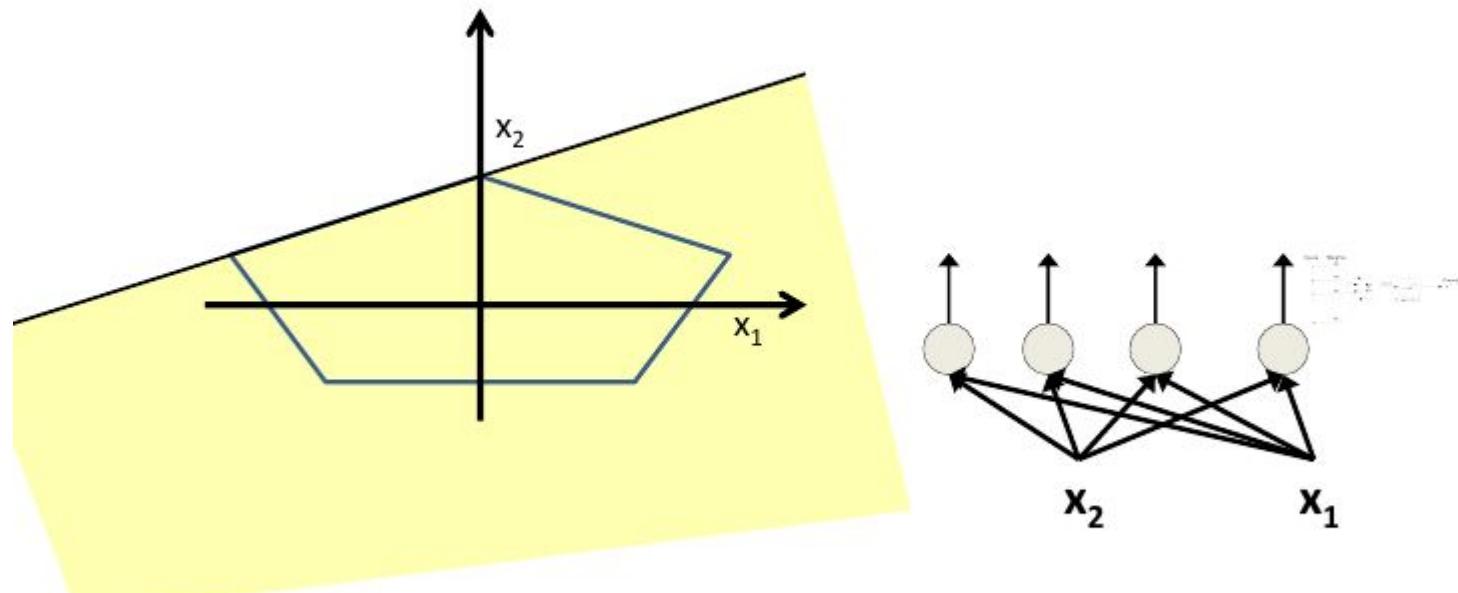
- The network must fire if the input is in the coloured area

Boolean over the Reals



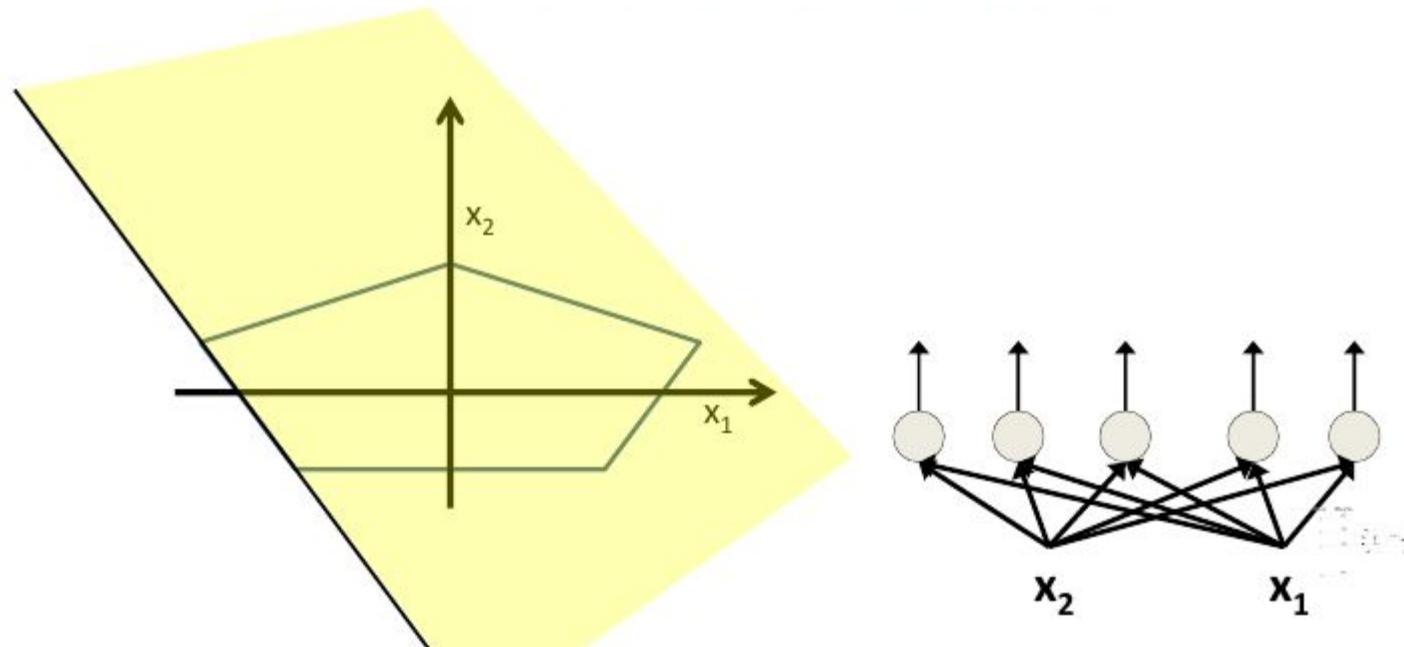
- The network must fire if the input is in the coloured area

Boolean over the Reals



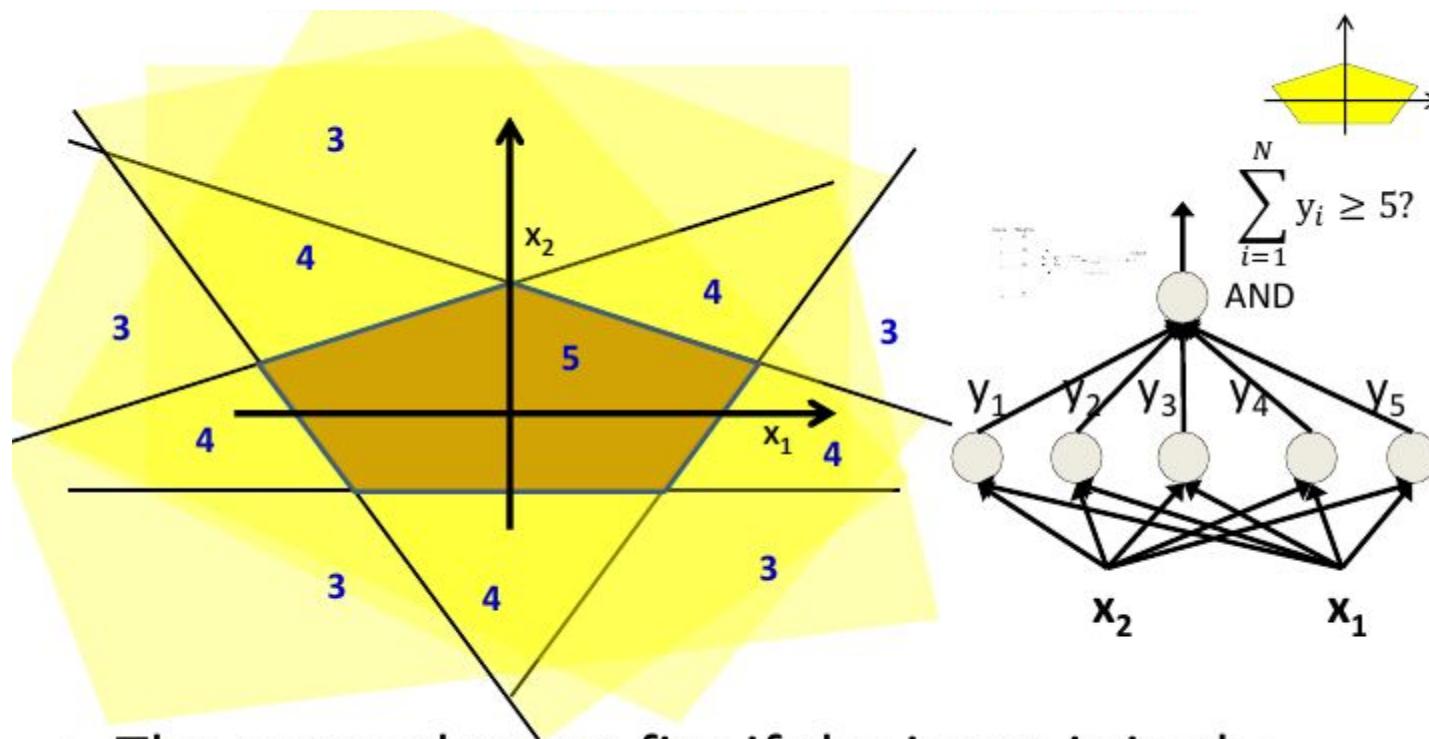
- The network must fire if the input is in the coloured area

Boolean over the Reals



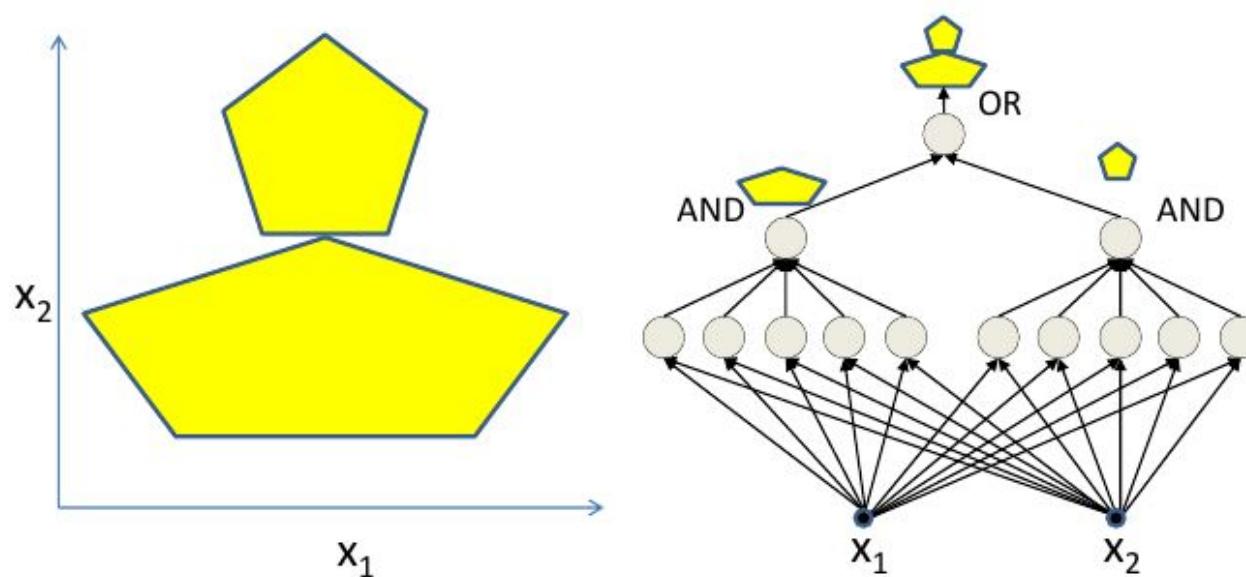
- The network must fire if the input is in the coloured area

Boolean over the Reals



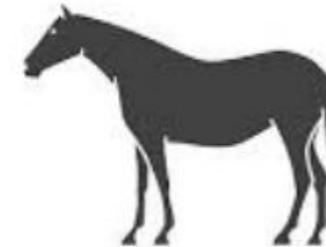
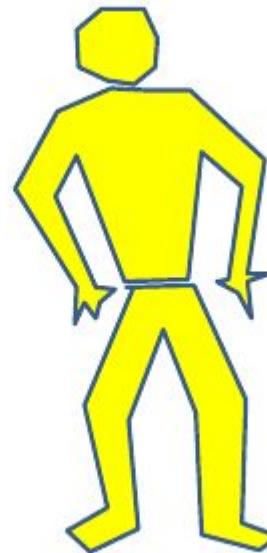
- The network must fire if the input is in the coloured area

More Complex Boundaries



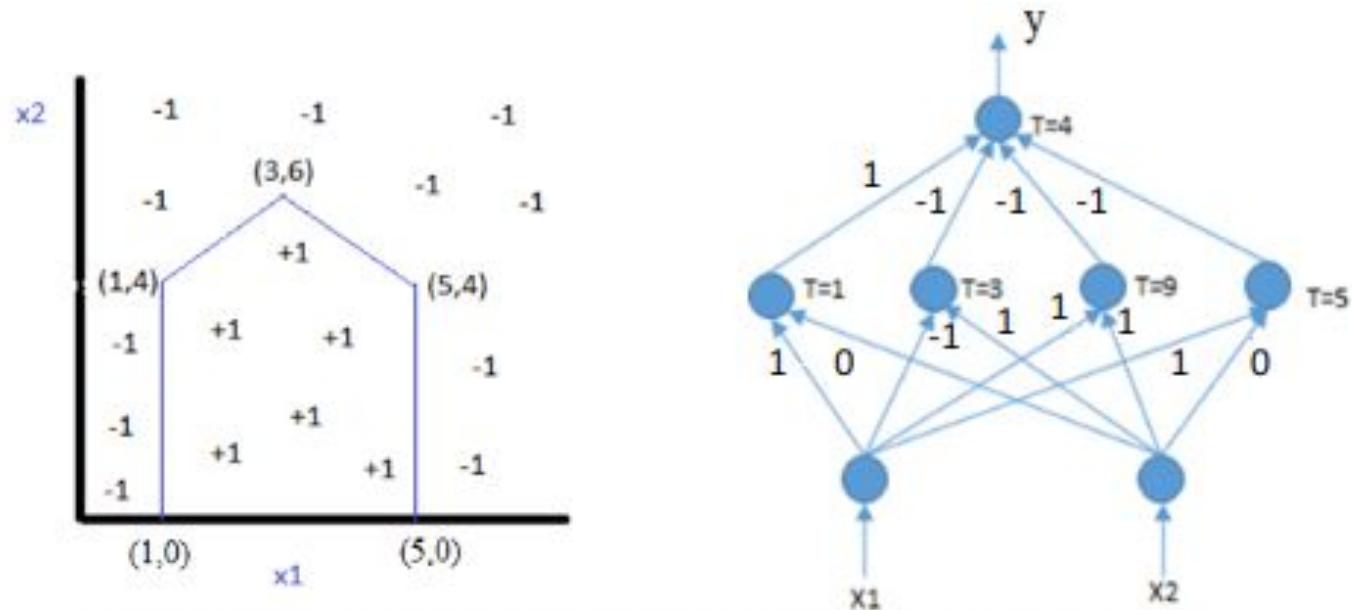
- Network to fire if the input is in the yellow area
 - “OR” two polygons
 - A third layer is required

More Complex Boundaries



- Can compose very complex decision boundaries

Example



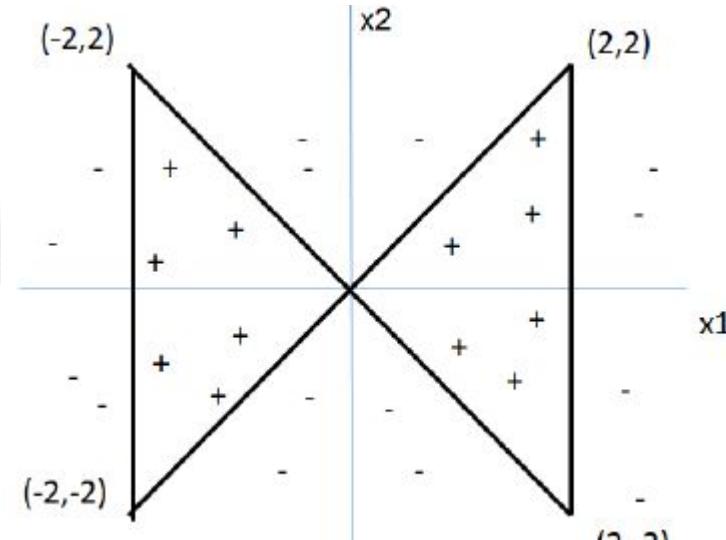
- Note: Different Choices of weights and bias are possible.
- Left hidden node implements $x_1=1$ line
- Right hidden node implements $x_1=5$ line
- 2nd node from left implements $x_2=x_1+3$ and 3rd from left implements $x_1+x_2=9$
- For + class, output of left hidden node = +1, for other nodes output = -1

Question

(x_1, x_2) are input features and target classes are either +1 or -1 as shown in the figure.

A. What is the minimum number of hidden layers and hidden nodes required to classify the following dataset with 100% accuracy using a fully connected multilayer perceptron network? Step activation functions are used at all nodes, i.e., output=+1 if total weighted input \geq bias b at a node, else output = -1.

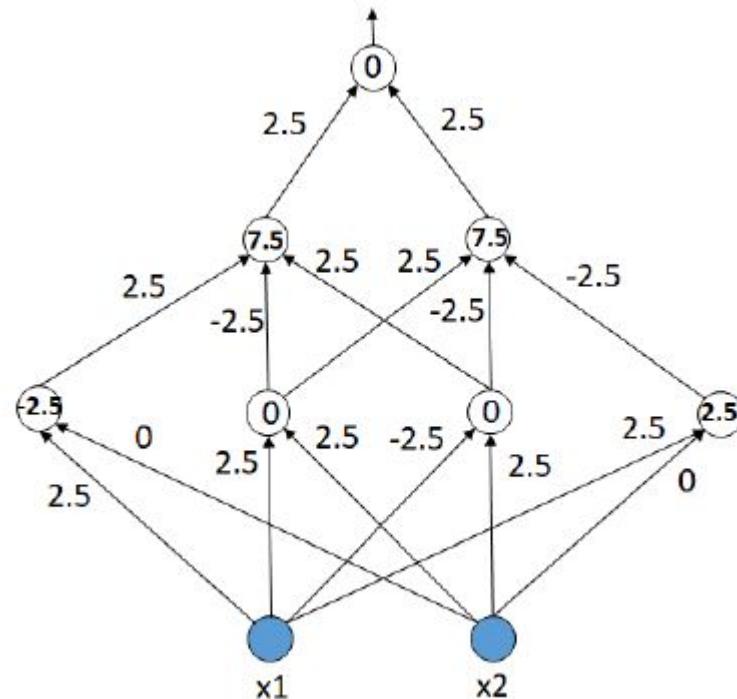
B. Show the minimal network architecture by organizing the nodes in each layer horizontally. Show the node representing x_1 at the left on the input layer. Organize the hidden nodes in ascending order of bias at that node. Specify all weights and bias values at all nodes. Weights can be only -2.5, 2.5 or 0, and bias +ve/-ve multiples of 2.5.



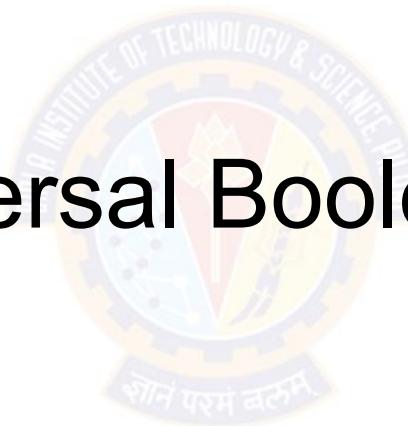
Solution

A. 2 hidden layers, 4 nodes in first hidden layer and 2 nodes in second hidden layer needed.

B.

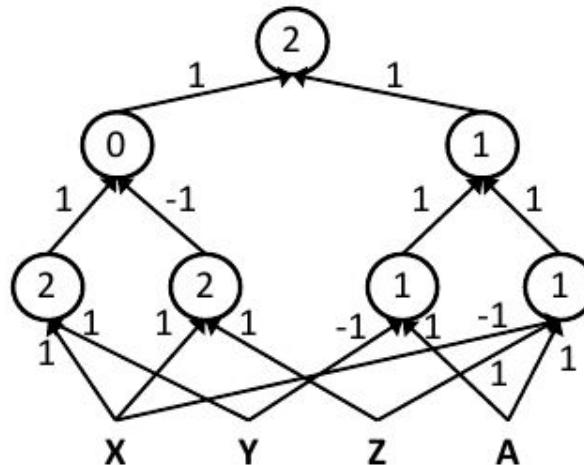


MLP as Universal Boolean Functions



Multi Layered Perceptrons for Boolean Functions

$$((A \& \bar{X} \& Z) | (A \& \bar{Y})) \& ((X \& Y) | (\bar{X} \& Z))$$



- MLPs can compute more complex Boolean functions
- MLPs can compute *any* Boolean function
 - Since they can emulate individual gates
- **MLPs are *universal Boolean functions***

How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows *all* input combinations
for which output is 1

- Expressed in disjunctive normal form

How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows *all* input combinations
for which output is 1

$$Y = \bar{X}_1\bar{X}_2X_3X_4\bar{X}_5 + \bar{X}_1X_2\bar{X}_3X_4X_5 + \bar{X}_1X_2X_3\bar{X}_4\bar{X}_5 + X_1\bar{X}_2\bar{X}_3\bar{X}_4X_5 + X_1\bar{X}_2X_3X_4X_5 + X_1X_2\bar{X}_3\bar{X}_4X_5$$

- Expressed in disjunctive normal form

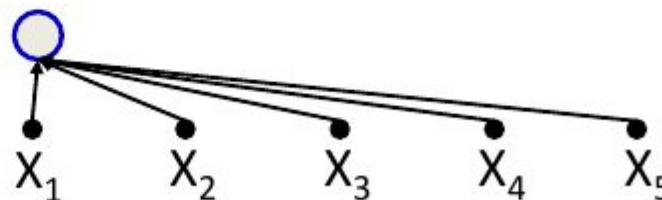
How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows *all* input combinations
for which output is 1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \bar{X}_1 X_2 \bar{X}_3 X_4 X_5 + \bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5 + X_1 X_2 \bar{X}_3 X_4 X_5 + X_1 \bar{X}_2 X_3 X_4 X_5 + X_1 X_2 \bar{X}_3 \bar{X}_4 X_5$$



- Expressed in disjunctive normal form

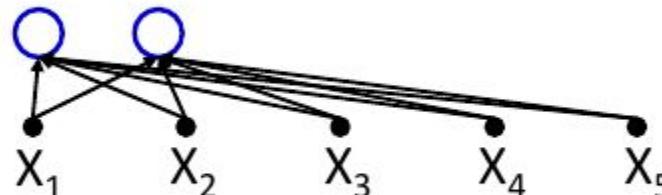
How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows all input combinations for which output is 1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \textcircled{X}_1 X_2 \bar{X}_3 X_4 X_5 + \bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5 + X_1 \bar{X}_2 \bar{X}_3 \bar{X}_4 X_5 + X_1 X_2 X_3 X_4 X_5 + X_1 X_2 \bar{X}_3 \bar{X}_4 X_5$$



- Expressed in disjunctive normal form

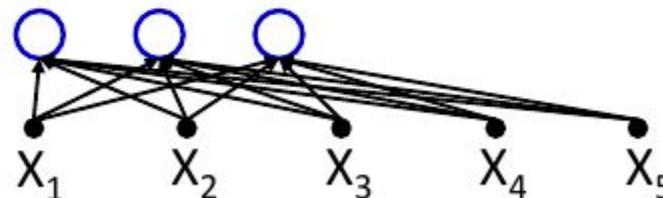
How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows all input combinations for which output is 1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \bar{X}_1 X_2 \bar{X}_3 X_4 X_5 + \textcircled{\bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5} + \\ X_1 \bar{X}_2 \bar{X}_3 \bar{X}_4 X_5 + X_1 \bar{X}_2 X_3 X_4 X_5 + X_1 X_2 X_3 X_4 X_5$$



- Expressed in disjunctive normal form

How many layers for a Boolean MLP?

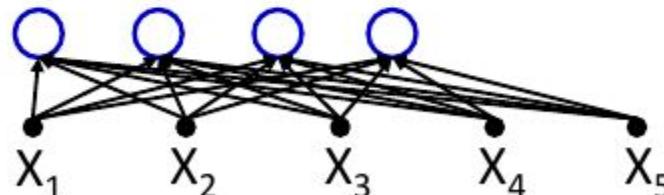
Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows *all* input combinations
for which output is 1

$$Y = \bar{X}_1\bar{X}_2X_3X_4\bar{X}_5 + \bar{X}_1X_2\bar{X}_3X_4X_5 + \bar{X}_1X_2X_3\bar{X}_4\bar{X}_5 +$$

$$\textcircled{X}_1\bar{X}_2\bar{X}_3\bar{X}_4X_5 + X_1\bar{X}_2X_3X_4X_5 + X_1X_2\bar{X}_3\bar{X}_4X_5$$



- Expressed in disjunctive normal form

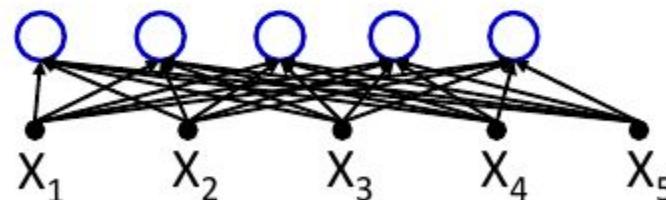
How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows *all* input combinations
for which output is 1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \bar{X}_1 X_2 \bar{X}_3 X_4 X_5 + \bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5 + \\ X_1 \bar{X}_2 \bar{X}_3 \bar{X}_4 X_5 + \textcircled{X}_1 \bar{X}_2 X_3 X_4 X_5 + X_1 X_2 \bar{X}_3 \bar{X}_4 X_5$$



- Expressed in disjunctive normal form

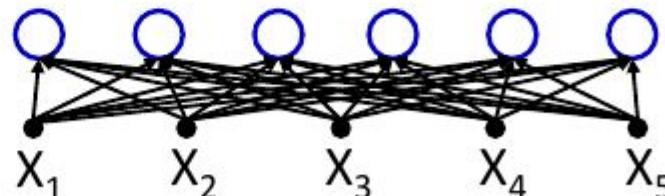
How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows all input combinations for which output is 1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \bar{X}_1 X_2 \bar{X}_3 X_4 X_5 + \bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5 + \\ X_1 \bar{X}_2 \bar{X}_3 \bar{X}_4 X_5 + X_1 \bar{X}_2 X_3 X_4 X_5 + \textcircled{X}_1 X_2 \bar{X}_3 \bar{X}_4 X_5$$



- Expressed in disjunctive normal form

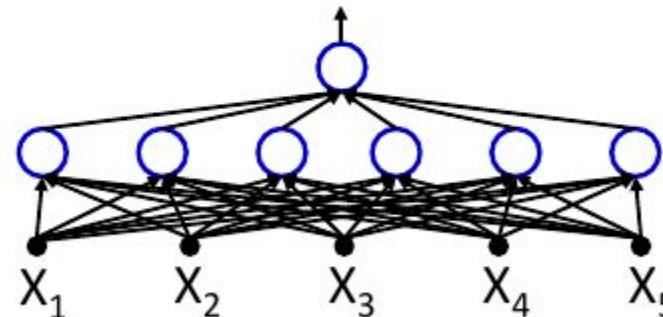
How many layers for a Boolean MLP?

Truth Table

X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

Truth table shows all input combinations for which output is 1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \bar{X}_1 X_2 \bar{X}_3 X_4 X_5 + \bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5 + X_1 \bar{X}_2 \bar{X}_3 \bar{X}_4 X_5 + X_1 \bar{X}_2 X_3 X_4 X_5 + X_1 X_2 \bar{X}_3 \bar{X}_4 X_5$$



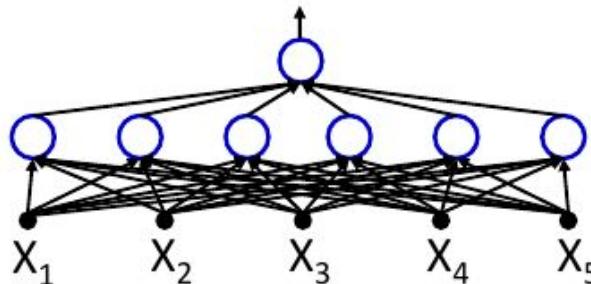
- Expressed in disjunctive normal form

How many layers for a Boolean MLP?

Truth table shows all input combinations for which output is 1

Truth Table					
X_1	X_2	X_3	X_4	X_5	Y
0	0	1	1	0	1
0	1	0	1	1	1
0	1	1	0	0	1
1	0	0	0	1	1
1	0	1	1	1	1
1	1	0	0	1	1

$$Y = \bar{X}_1 \bar{X}_2 X_3 X_4 \bar{X}_5 + \bar{X}_1 X_2 \bar{X}_3 X_4 X_5 + \bar{X}_1 X_2 X_3 \bar{X}_4 \bar{X}_5 + X_1 \bar{X}_2 \bar{X}_3 \bar{X}_4 X_5 + X_1 \bar{X}_2 X_3 X_4 X_5 + X_1 X_2 \bar{X}_3 \bar{X}_4 X_5$$



- Any truth table can be expressed in this manner!
- A one-hidden-layer MLP is a Universal Boolean Function

But what is the largest number of perceptrons required in the single hidden layer for an N-input-variable function?

Reducing a Boolean Function

WX \ YZ	00	01	11	10
00	1	0	0	1
01	1	1	0	0
11	1	0	0	0
10	1	0	0	1

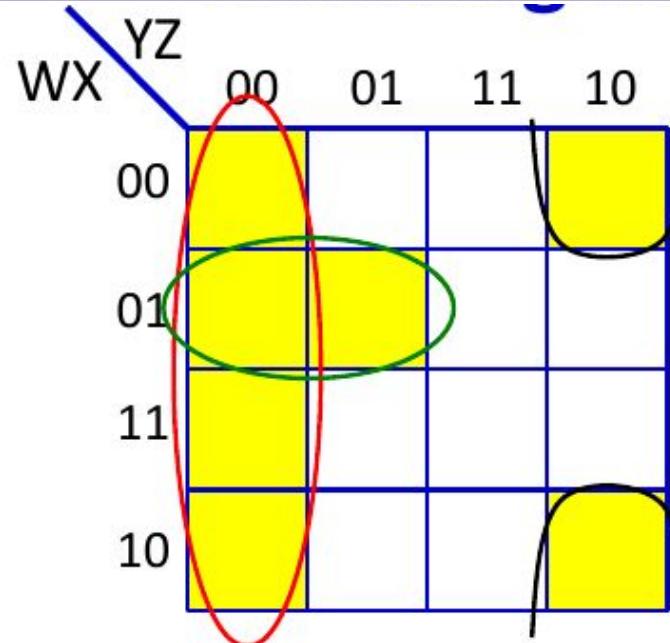
This is a "Karnaugh Map"

It represents a truth table as a grid
Filled boxes represent input combinations
for which output is 1; blank boxes have
output 0

Adjacent boxes can be "grouped" to
reduce the complexity of the DNF formula
for the table

- DNF form:
 - Find groups
 - Express as reduced DNF

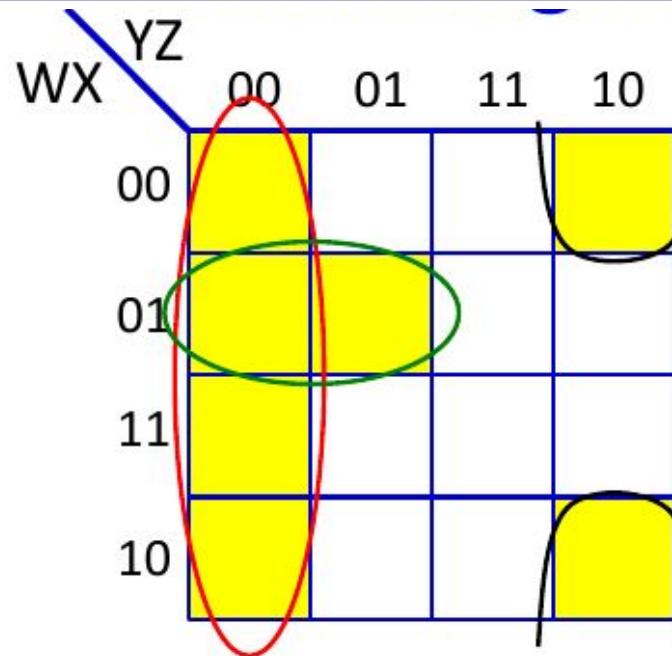
Reducing a Boolean Function



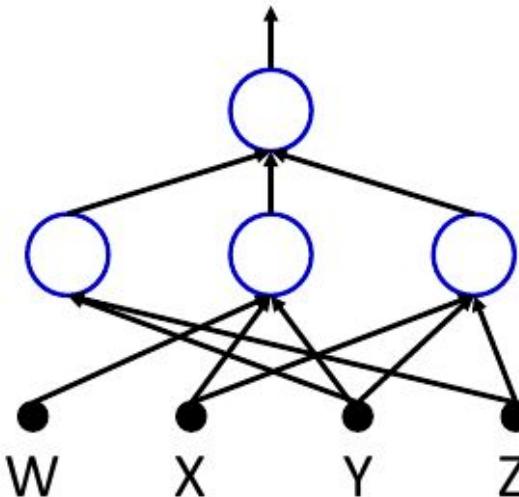
$$O = \bar{Y}\bar{Z} + \bar{W}X\bar{Y} + \bar{X}YZ$$

- *Reduced DNF form:*
 - Find groups
 - Express as reduced DNF

Reducing a Boolean Function

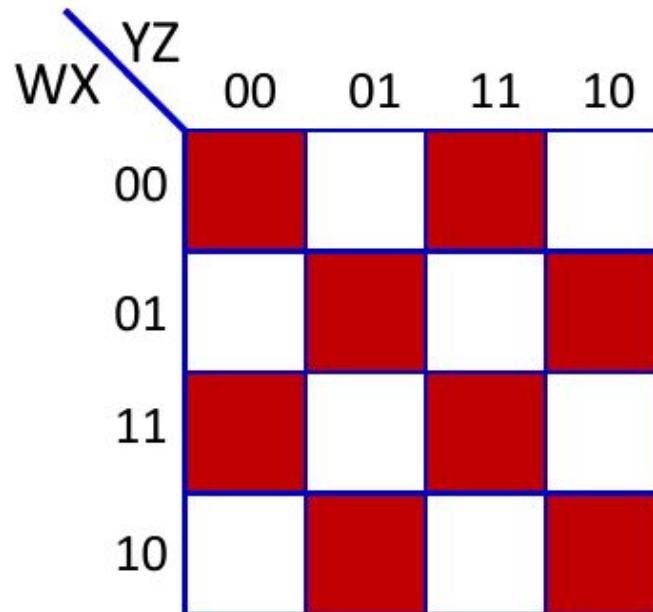


$$O = \bar{Y}\bar{Z} + \bar{W}X\bar{Y} + \bar{X}YZ$$



- *Reduced DNF form:*
 - Find groups
 - Express as reduced DNF

Largest irreducible DNF?



- What arrangement of ones and zeros simply cannot be reduced further?

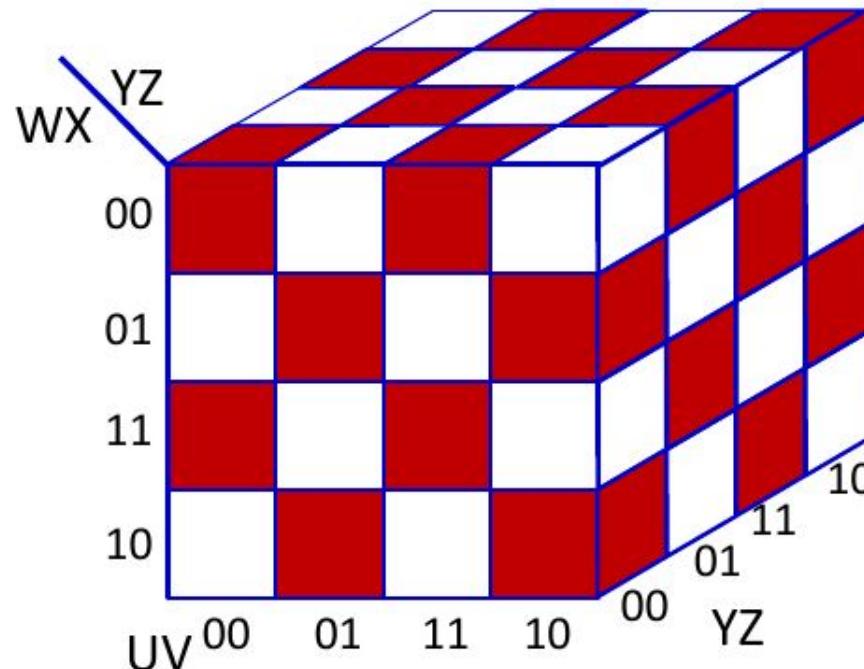
Largest irreducible DNF?

WX \ YZ	00	01	11	10
00	1	0	1	0
01	0	1	0	1
11	1	0	1	0
10	0	1	0	1

How many neurons
in a DNF (one-hidden-layer) MLP
for this Boolean
function?

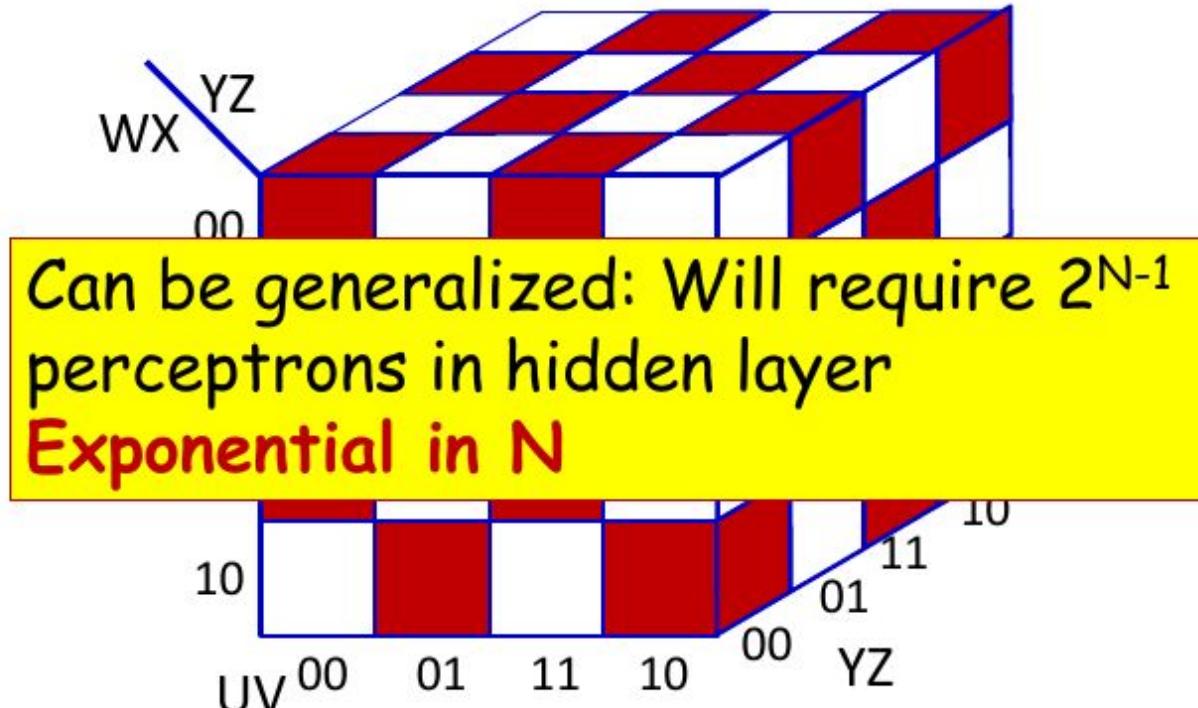
- What arrangement of ones and zeros simply cannot be reduced further?

Width of a Single Layer Boolean MLP



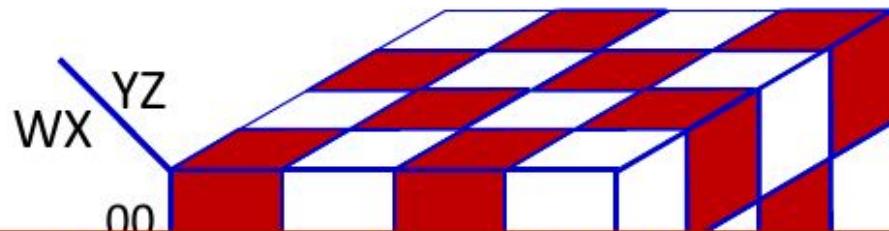
- How many neurons in a DNF (one-hidden-layer) MLP for this Boolean function of 6 variables?

Width of a Single Layer Boolean MLP



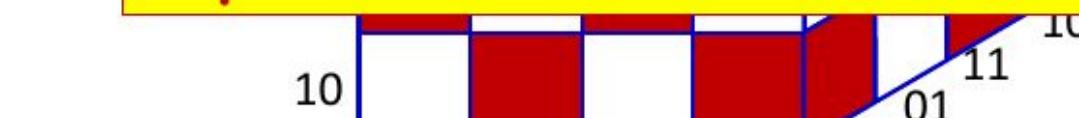
- How many neurons in a DNF (one-hidden-layer) MLP for this Boolean function

Width of a Single Layer Boolean MLP



Can be generalized: Will require 2^{N-1} perceptrons in hidden layer

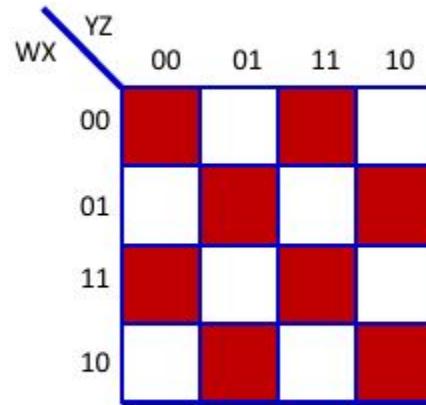
Exponential in N



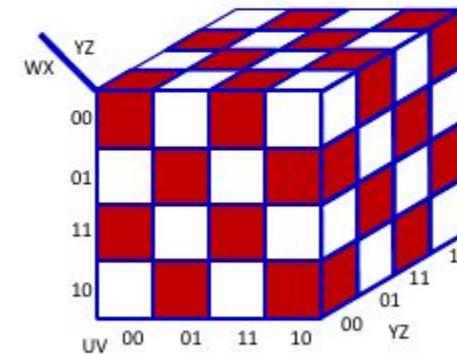
How many units if we use multiple layers?

- How many neurons in a DNF (one-hidden-layer) MLP for this Boolean function

Width of deep MLP

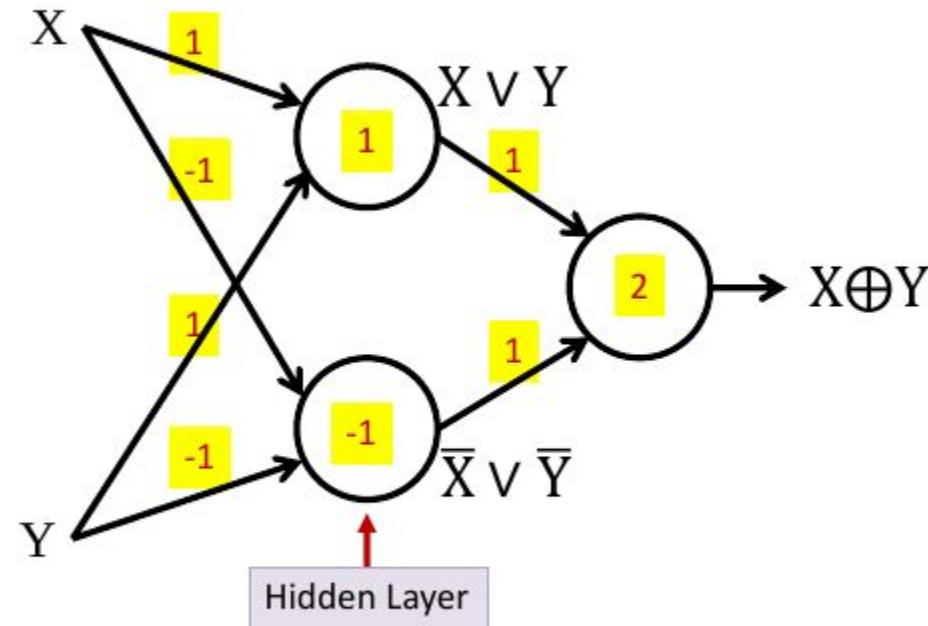


$$O = W \oplus X \oplus Y \oplus Z$$



$$O = U \oplus V \oplus W \oplus X \oplus Y \oplus Z$$

MLP XOR

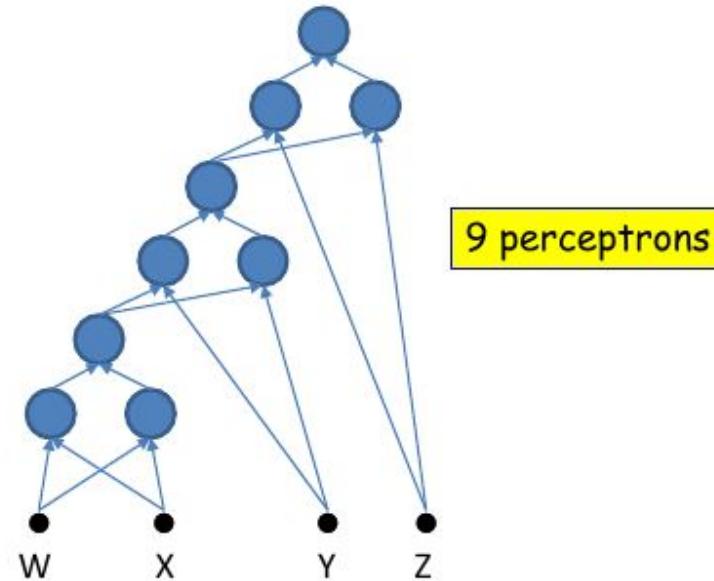


- An XOR takes three perceptrons

Width of deep MLP

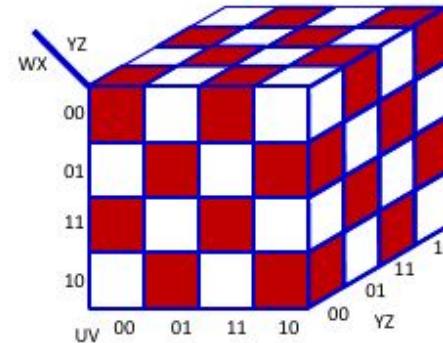
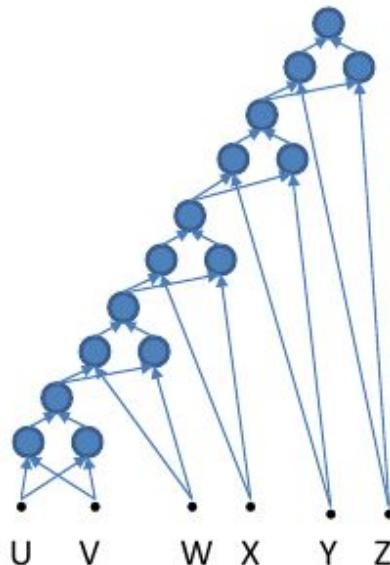
WX	YZ	00	01	11	10
00		Red	White	Red	White
01		White	Red	White	Red
11		Red	White	Red	White
10		White	Red	White	Red

$$O = W \oplus X \oplus Y \oplus Z$$



- An XOR needs 3 perceptrons
- This network will require $3 \times 3 = 9$ perceptrons

Width of deep MLP

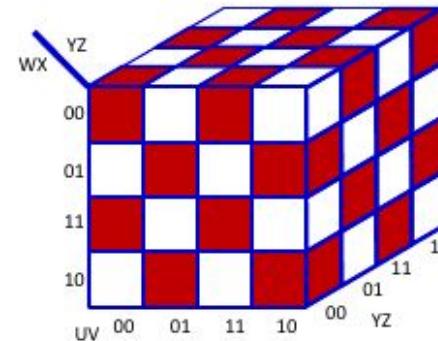
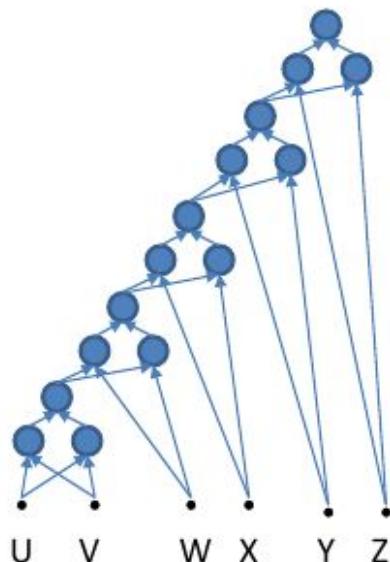


$$O = U \oplus V \oplus W \oplus X \oplus Y \oplus Z$$

15 perceptrons

- An XOR needs 3 perceptrons
- This network will require $3 \times 5 = 15$ perceptrons

Width of deep MLP

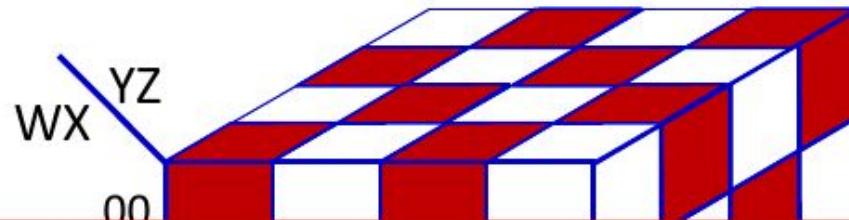


$$O = U \oplus V \oplus W \oplus X \oplus Y \oplus Z$$

More generally, the XOR of N variables will require $3(N-1)$ perceptrons!!

- An XOR needs 3 perceptrons
- This network will require $3 \times 5 = 15$ perceptrons

Width of Single Layer Boolean MLP



Single hidden layer: Will require $2^{N-1}+1$ perceptrons in all (including output unit)
Exponential in N

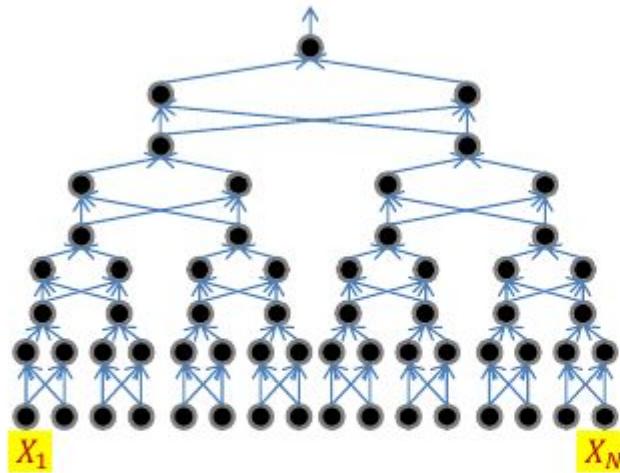


Will require $3(N-1)$ perceptrons in a deep network

Linear in N!!!

Can be arranged in only $2\log_2(N)$ layers

A better representation...

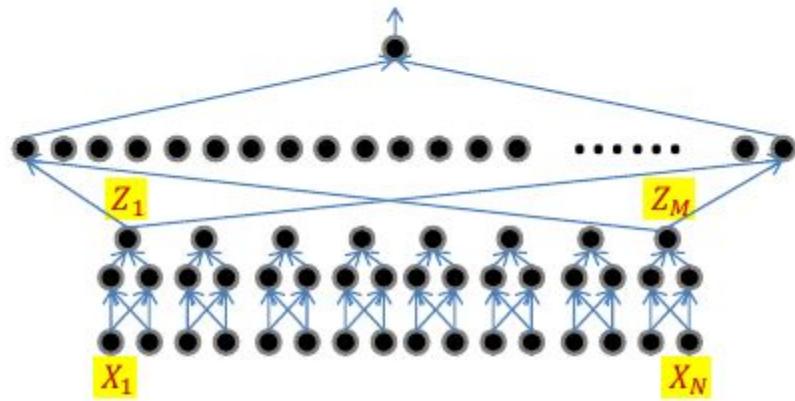


$$O = X_1 \oplus X_2 \oplus \cdots \oplus X_N$$

- Only $2 \log_2 N$ layers
 - By pairing terms
 - 2 layers per XOR

$$O = (((((X_1 \oplus X_2) \oplus (X_3 \oplus X_4)) \oplus ((X_5 \oplus X_6) \oplus (X_7 \oplus X_8))) \oplus (((X_9 \oplus X_{10}) \oplus (X_{11} \oplus X_{12})) \oplus ((X_{13} \oplus X_{14}) \oplus (X_{15} \oplus X_{16})))$$

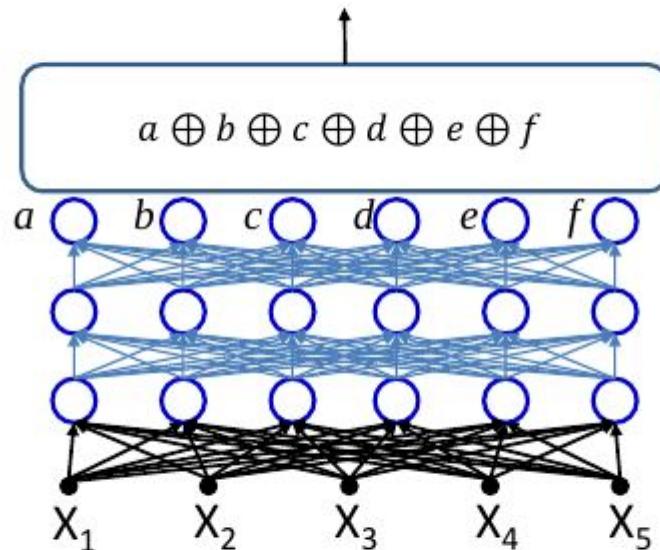
Challenge of Depth



$$\begin{aligned}O &= X_1 \oplus X_2 \oplus \cdots \oplus X_N \\&= Z_1 \oplus Z_2 \oplus \cdots \oplus Z_M\end{aligned}$$

- Using only K hidden layers will require $O(2^{(N-K/2)})$ neurons in the K th layer
 - Because the output can be shown to be the XOR of all the outputs of the $K-1$ th hidden layer
 - **I.e. reducing the number of layers below the minimum will result in an exponentially sized network to express the function fully**
 - **A network with fewer than the required number of neurons cannot model the function**

Need of Depth



- The wide function can happen at any layer
- Having a few extra layers can greatly reduce network size

Network Size

- An MLP is a universal Boolean function
- But can represent a given function only if
 - It is sufficiently wide
 - It is sufficiently deep
 - Depth can be traded off for (sometimes) exponential growth of the width of the network
- Optimal width and depth depend on the number of variables and the complexity of the Boolean function
 - Complexity: *minimal* number of terms in DNF formula to represent it

Depth: Summary

- The number of neurons required in a shallow network is
 - Polynomial in the number of basic patterns
 - Exponential in the dimensionality input
 - (this is the worst case)
 - Alternately, exponential in the number of statistically independent features

Question

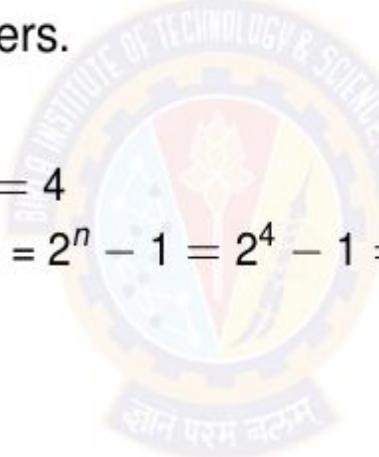
How many perceptrons are required to represent $W \oplus X \oplus Y \oplus Z$?



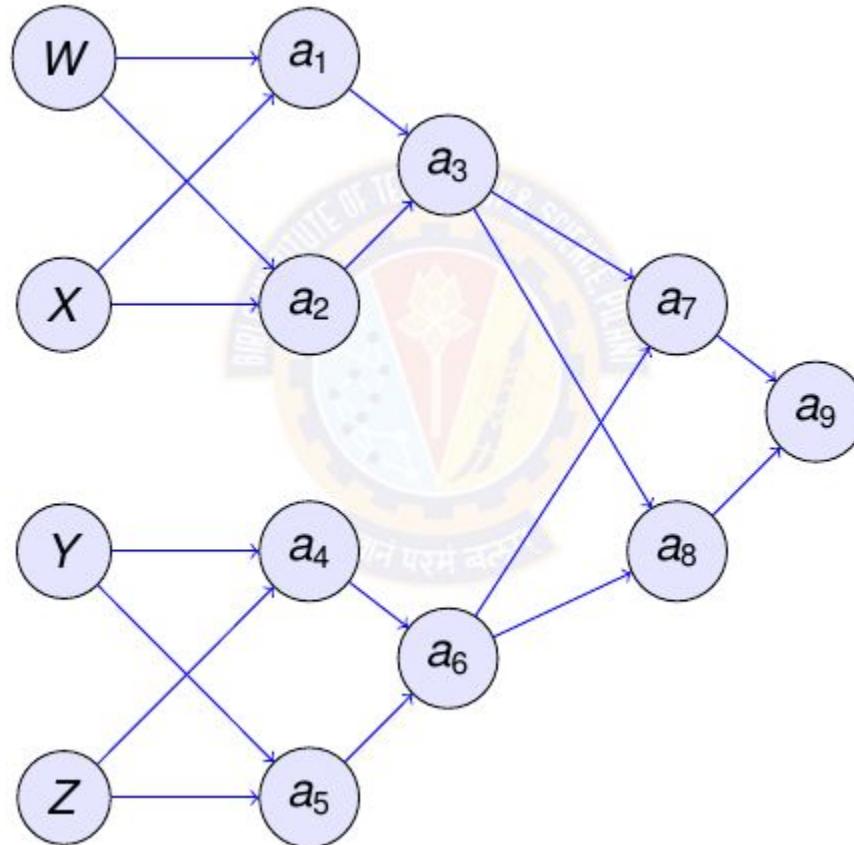
Solution

An XOR cannot be represented using a single perceptron as the XOR has non-linear decision boundary. To represent non-linear decision boundary we need a multi-layer perceptron with one hidden layers.

- Traditional method:
 - ▶ Number of variables $n = 4$
 - ▶ Number of perceptrons $= 2^n - 1 = 2^4 - 1 = 15$
- Alternate method:
 - ▶ Layer 1: $A = W \oplus X$
 - ▶ Layer 2: $B = Y \oplus Z$
 - ▶ Layer 3: $C = A \oplus B$
 - ▶ Number of perceptrons $= 3(n - 1) = 3 * (4 - 1) = 9$
 - ▶ Number of layers = Width $= 2 \log_2 n = 2 \log_2 4 = 4$



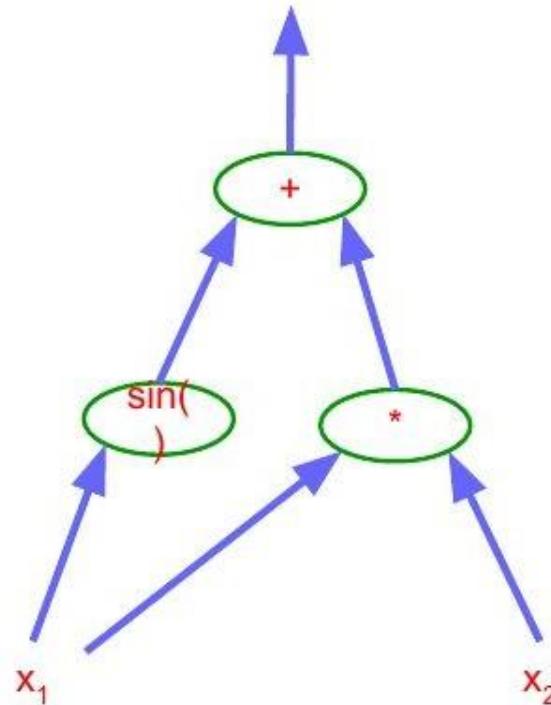
Solution



Computational Graph

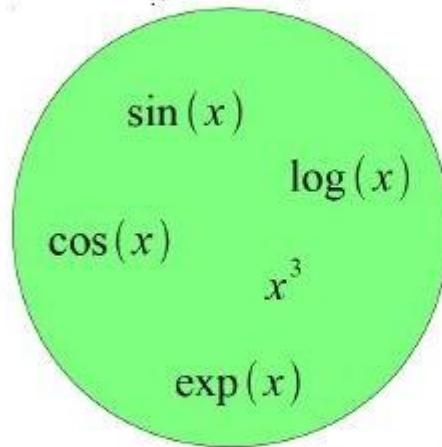
Computational Graph: Example

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



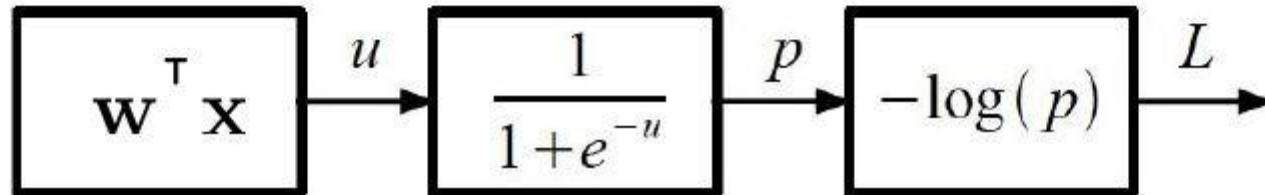
Computational Graph for Logistic Regression

Given a library of simple functions



Compose into a
complicate function

$$-\log \left(\frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}} \right)$$



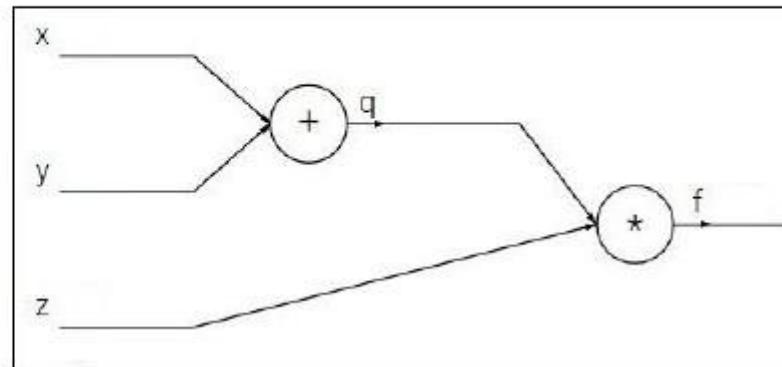
Computational Graph for Back Propagation

$$f(x, y, z) = (x + y)z$$



Computational Graph for Back Propagation

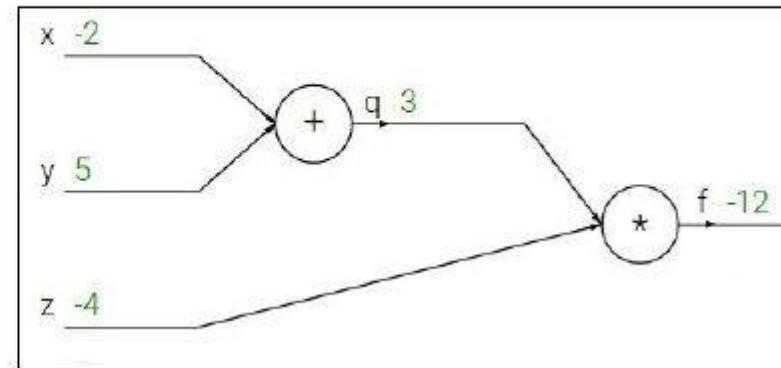
$$f(x, y, z) = (x + y)z$$



Computational Graph for Back Propagation

$$f(x, y, z) = (x + y)z$$

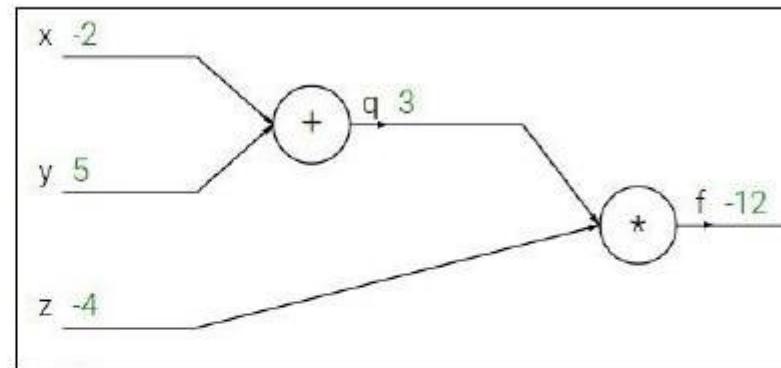
e.g. $x = -2$, $y = 5$, $z = -4$



Computational Graph for Back Propagation

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$



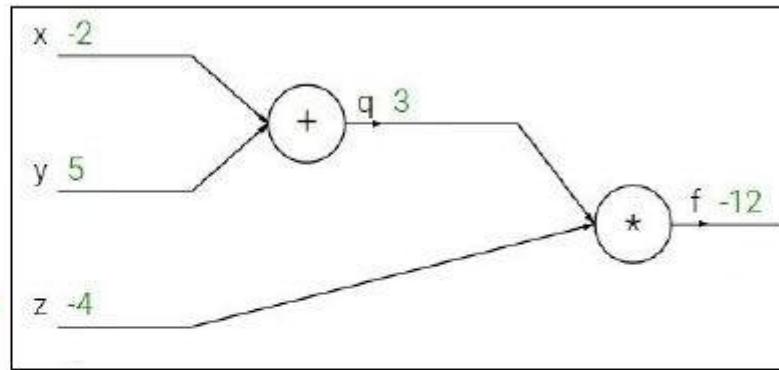
Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Computational Graph for Back Propagation

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$



Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Computational Graph for Back Propagation

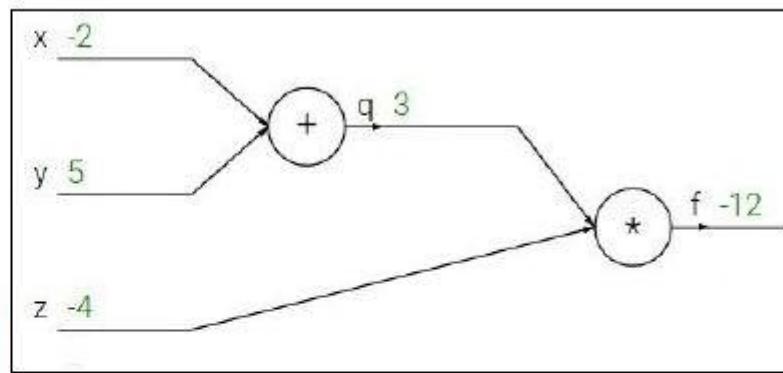
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Computational Graph for Back Propagation

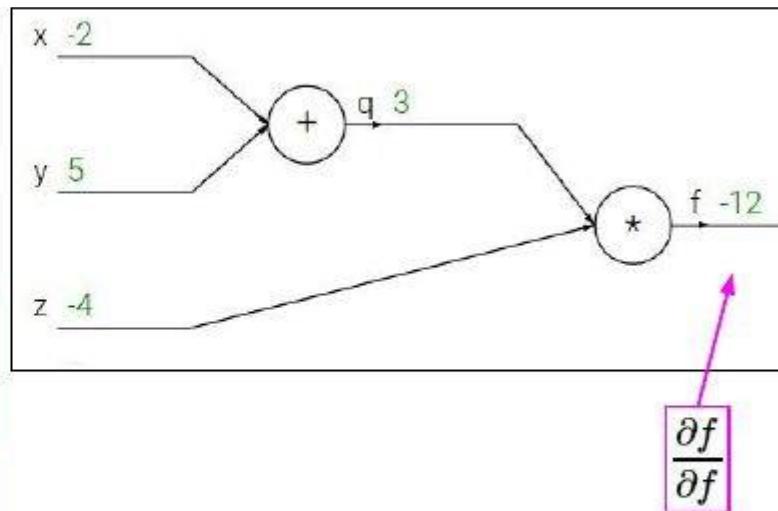
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Computational Graph for Back Propagation

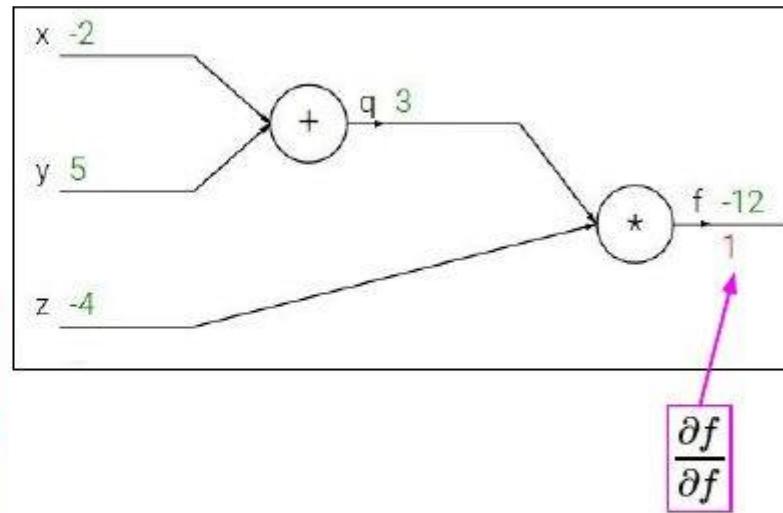
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Computational Graph for Back Propagation

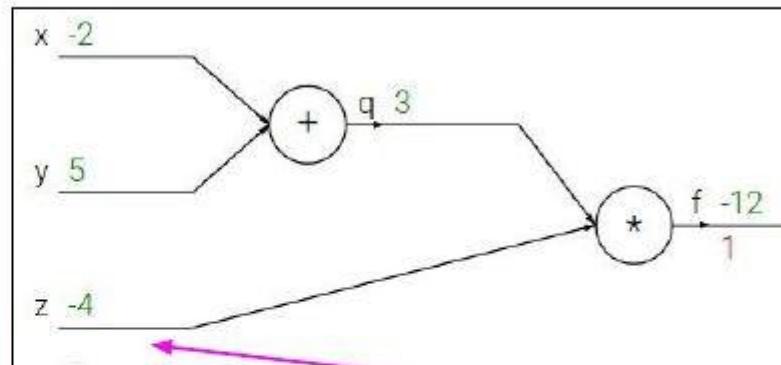
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

Computational Graph for Back Propagation

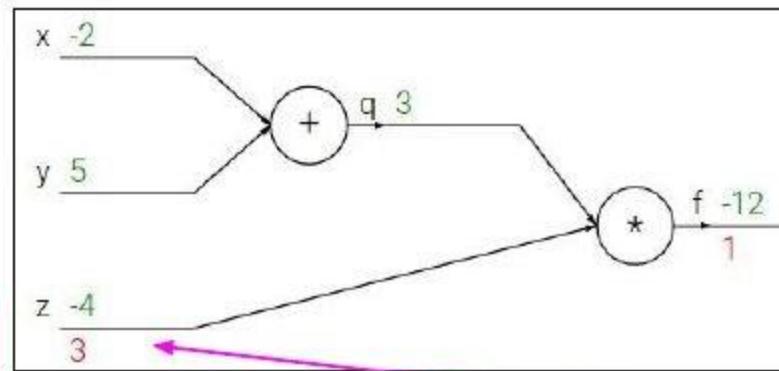
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Computational Graph for Back Propagation

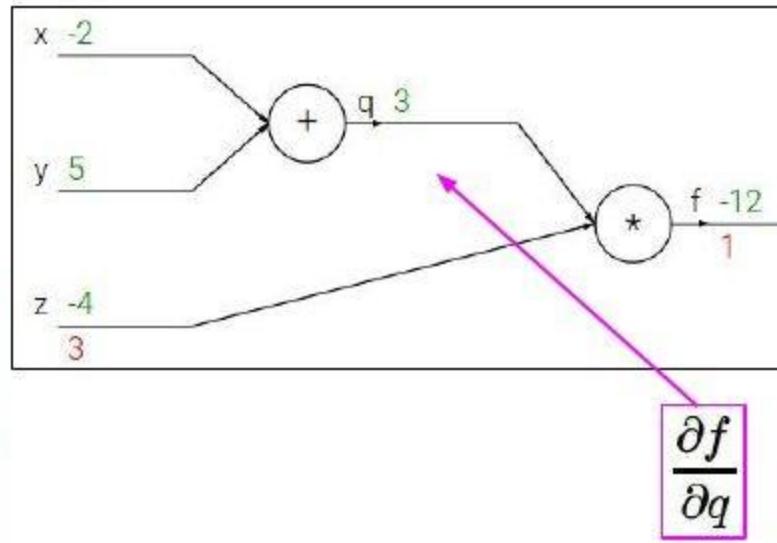
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Computational Graph for Back Propagation

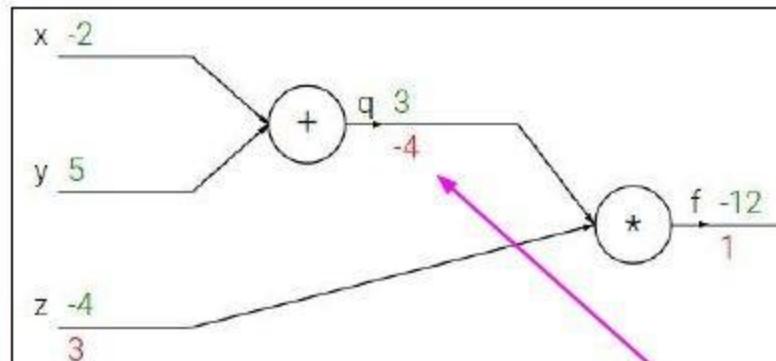
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

Computational Graph for Back Propagation

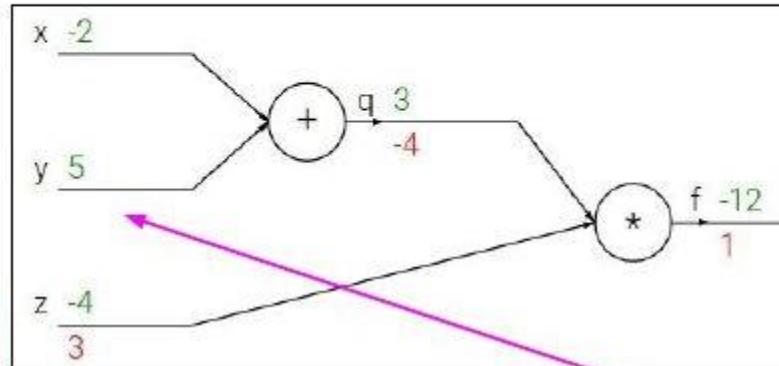
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream
gradient

Local
gradient

$$\frac{\partial f}{\partial y}$$

Computational Graph for Back Propagation

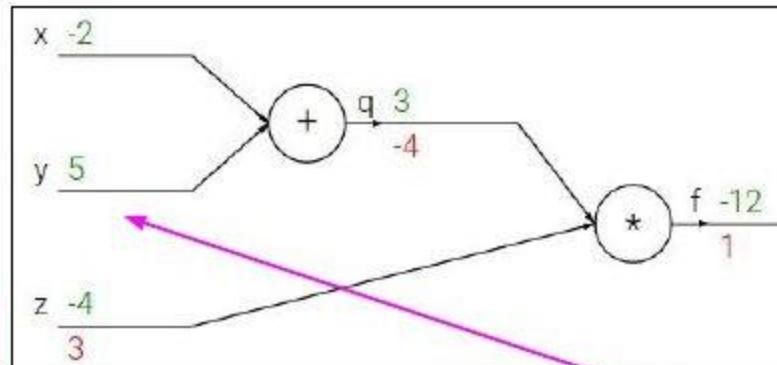
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream
gradient

Local
gradient

$$\frac{\partial f}{\partial y}$$

Computational Graph for Back Propagation

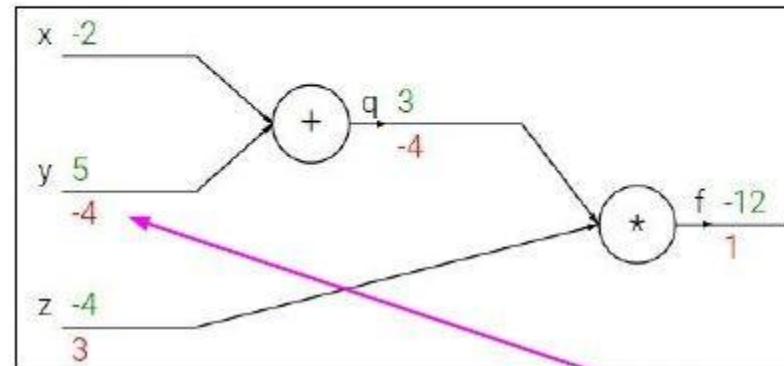
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2$, $y = 5$, $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream
gradient

Local
gradient

$$\frac{\partial f}{\partial y}$$

Computational Graph for Back Propagation

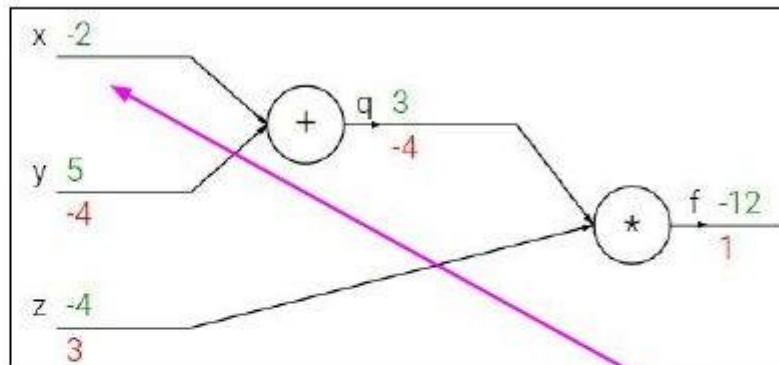
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream gradient

Local gradient

$$\frac{\partial f}{\partial x}$$

Computational Graph for Back Propagation

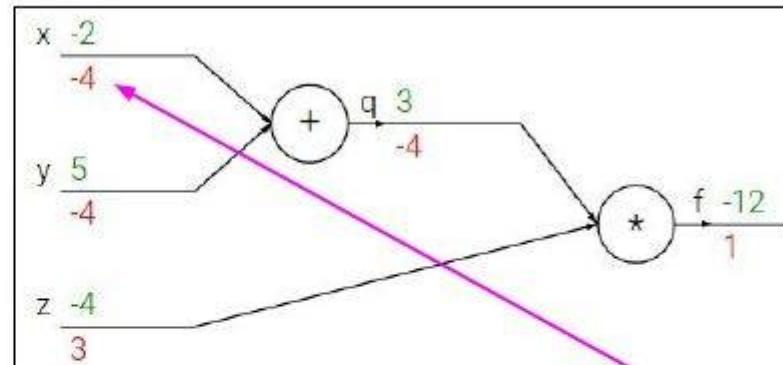
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream
gradient

Local
gradient

$$\frac{\partial f}{\partial x}$$

Question 8

Draw the computational graph for the equation

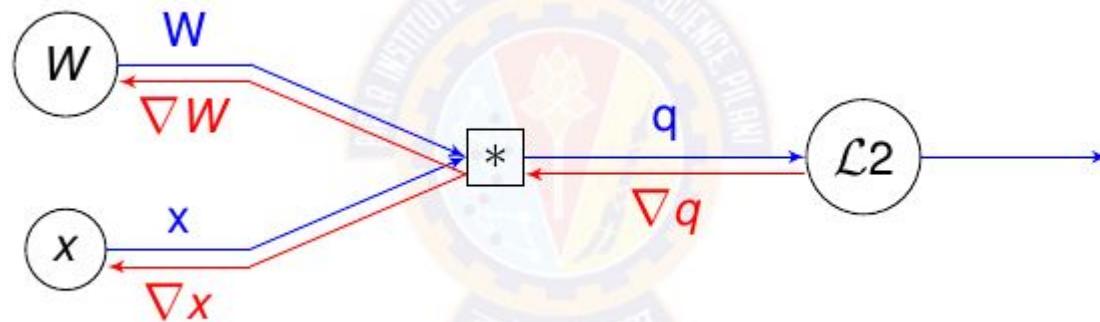
$$f(x, w) = \| W \cdot x \|^2 = \sum_{i=1}^n (W \cdot x)_i^2$$

Using the computation graph show the computation of the gradients also. Assume that

$$W = \begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}$$

$$x = \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

Solution for Q8



Solution for Q8

$$W = \begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix}$$

$$x = \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}$$

$$q = W \cdot x = \begin{bmatrix} 0.1 * 0.2 + 0.5 * 0.4 \\ -0.3 * 0.2 + 0.8 * 0.4 \end{bmatrix} = \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$L2 = \sum_{i=1} q_i^2 = 0.22^2 + 0.26^2 = 0.116$$

$$\nabla q = \nabla_q f = 2q$$

$$= 2 \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix} = \begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

$$\nabla W = \nabla_w f = 2q \cdot x^\top$$

$$= 2 \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix} [0.2 \quad 0.4] = \begin{bmatrix} 0.088 & 0.176 \\ 0.105 & 0.208 \end{bmatrix}$$

$$\nabla x = \nabla_x f = 2W^\top \cdot q$$

$$= 2 \begin{bmatrix} 0.1 & -0.3 \\ 0.5 & 0.8 \end{bmatrix} \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix} = \begin{bmatrix} -0.112 \\ 0.636 \end{bmatrix}$$

Question 7

Draw the computational graph for Sigmoid function and show the gradient computation also. Use generic equations.



Demo of DNN

1. XOR Implementation

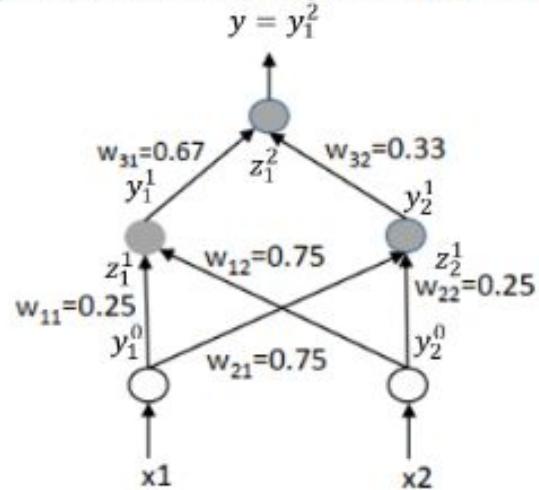
<https://colab.research.google.com/drive/1xVVpeU3q4bIOexV0J3NhYbLVCHwaOI6R#scrollTo=GRaiuHtKI1Sq>

2. DNN using Tensorflow

<https://colab.research.google.com/drive/1307HTGrHmMMQaiBInm10vWHQcUSNUDTR#scrollTo=cPZ81DJmK-W3>

Example with Relu

Weight Updates in Deep Networks



- Training Input: $(x_1, x_2) = (1, 1)$ Output: $d = 1.0$
- *div* function: square error
- Activation function $f()$: ReLU
- Bias: 0 at all nodes
- $\eta = 0.1$
- What are the values of w_{31} and w_{12} in next iteration?

Example with Relu

Weight Updates

$$z_1^1 = 0.25*1 + 0.75*1 = 1.0$$

$$z_2^1 = 0.75*1 + 0.25*1 = 1.0$$

$$y_1^1 = \text{ReLU}(z_1^1) = 1 \quad y_2^1 = \text{ReLU}(z_2^1) = 1$$

$$z_1^2 = 0.67*1 + 0.33*1 = 1.0$$

$$y_1^2 = \text{ReLU}(z_1^2) = 1$$

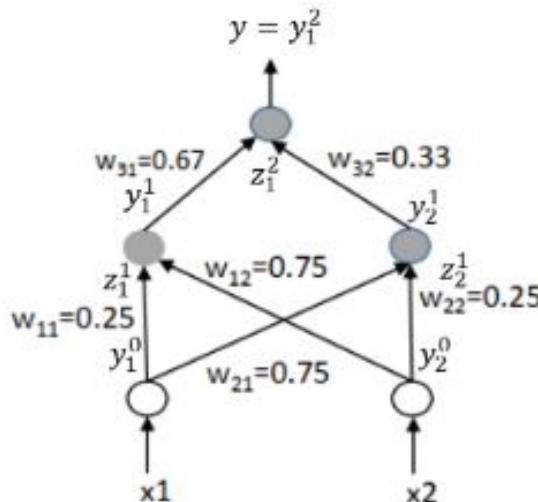
$$\text{div} = 0.5 * (0 - y_1^2)^2$$

$$\frac{\delta \text{div}}{\delta y_1^2} = y_1^2 = 1$$

$$\frac{\delta \text{div}}{\delta z_1^2} = \frac{\delta \text{div}}{\delta y_1^2} * \frac{\delta y_1^2}{\delta z_1^2} = 1 * 1 = 1$$

$$\frac{\delta \text{div}}{\delta w_{11}} = \frac{\delta \text{div}}{\delta z_1^2} * \frac{\delta z_1^2}{\delta w_{11}} = 1 * y_1^1 = 1$$

$$w_{31} = w_{31} - \eta * \frac{\delta \text{div}}{\delta w_{31}} = 0.67 - 0.1 = 0.57$$



$$\frac{\delta \text{div}}{\delta y_1^1} = \frac{\delta \text{div}}{\delta z_1^2} * \frac{\delta z_1^2}{\delta y_1^1} = 1 * w_{31}$$

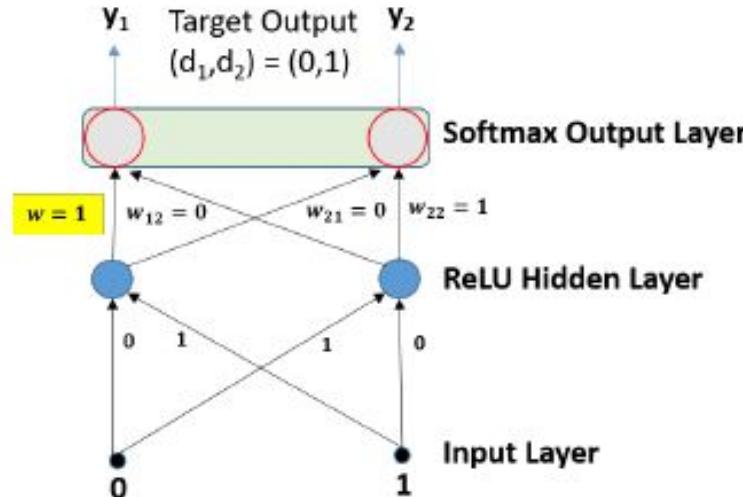
$$\frac{\delta \text{div}}{\delta z_1^1} = \frac{\delta \text{div}}{\delta y_1^1} * \frac{\delta y_1^1}{\delta z_1^1} = 0.67 * 1 = 0.67$$

$$\frac{\delta \text{div}}{\delta w_{12}} = \frac{\delta \text{div}}{\delta z_1^1} * \frac{\delta z_1^1}{\delta w_{12}} = 0.67 * y_2^0 = 0.67$$

$$w_{12} = w_{12} - \eta * \frac{\delta \text{div}}{\delta w_{12}} = 0.75 - 0.067 = 0.683$$

Question

Refer to the Figure below. Hidden nodes use ReLU activation function. Output nodes are softmax. Bias at all nodes is 0.



Calculate the absolute value of the cross-entropy loss for the given input data, bias and weights.

Solution

For the given input and input-hidden weights, $(h_1, h_2) = (1, 0)$ and $(d_1, d_2) = (0, 1)$.
Input to softmax node $z_1 = w_{11} = 1; z_2 = 0$.

$$y_1 = \text{softmax}[\max(1(0 * 0 + 1 * 1), 0)]$$

$$y_1 = e/(1 + e)$$

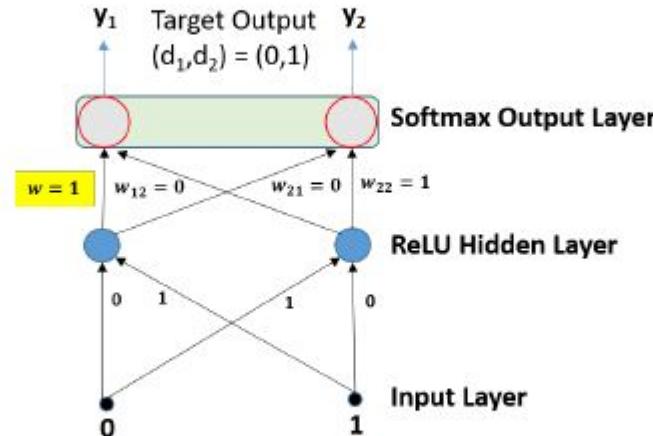
$$y_2 = \text{softmax}[\max(0, 1(0 * 1 + 0 * 1))]$$

$$y_2 = 1/(1 + e)$$

$$\begin{aligned}\text{Cross - entropy Loss} &= -d_1 \log(y_1) - d_2 \log(y_2) \\ &= -\log(y_2) = \log(1 + e) = 1.313\end{aligned}$$

Question

Refer to the Figure below. Hidden nodes use ReLU activation function. Output nodes are softmax. Backpropagation is used to update the weights. Learning rate is 0.1. Bias at all nodes is 0.



Calculate the absolute value of change in weight w (marked in yellow) for the given input data, weights and target outputs (d_1, d_2) . (y_1, y_2) are actual outputs from the network.

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} y_i^{(k)} \left(\delta_{ij} - y_j^{(k)} \right)$$

For the given input and input-hidden weights, $(h_1, h_2) = (1, 0)$ and $(d_1, d_2) = (0, 1)$.
 Change in w_{11}

$$\begin{aligned}\delta w_{11} &= -0.1 * \frac{\partial \text{Div}}{\partial w_{11}} \\ &= -0.1 * \frac{\partial \text{Div}}{\partial z_1} * \frac{\partial z_1}{\partial w_{11}} = -0.1 * \frac{\partial \text{Div}}{\partial z_1} * h_1\end{aligned}$$

$$\frac{\partial \text{Div}}{\partial y_1} = -\frac{d_1}{y_1} = 0$$

$$\frac{\partial \text{Div}}{\partial y_2} = -\frac{d_2}{y_2} = -\frac{1}{y_2}$$

$$\frac{\partial \text{Div}}{\partial z_1} = \frac{\partial \text{Div}}{\partial y_1} * y_1(1 - y_1) + \frac{\partial \text{Div}}{\partial y_2} * y_1(-y_2) = -y_1 y_2 / y_2 = -y_1$$

Input to softmax node

$$z_1 = w_{11} = 1$$

$$z_2 = 0$$

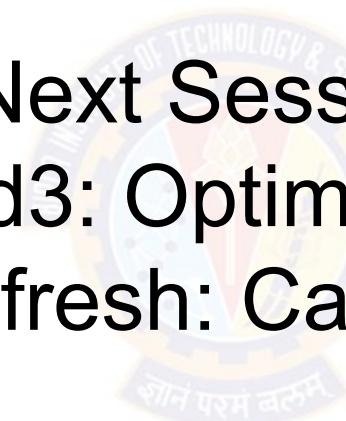
$$y_1 = e/(1 + e)$$

$$\text{abs}(\delta w_{11}) = 0.1 * e/(1 + e) = 0.0731$$

Ref:

- http://mlsp.cs.cmu.edu/people/rsingh/docs/Chapter1_Introduction.pdf
- http://mlsp.cs.cmu.edu/people/rsingh/docs/Chapter2_UniversalApproximators.pdf

Next Session:
Mod3: Optimization
Refresh: Calculus





BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Neural Network

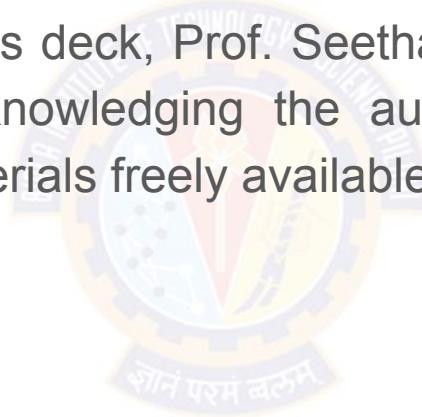
AIML Module 3

Seetha Parameswaran

BITS Pilani



The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.

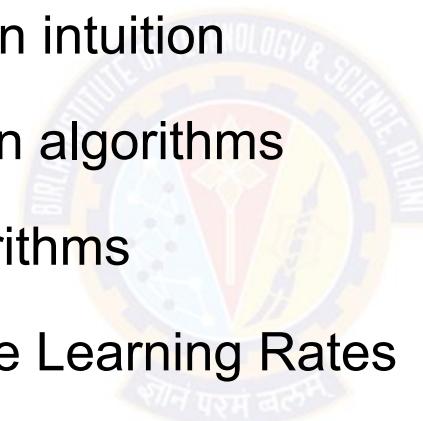


Optimization



What we Learn....

- 3.1 Challenges in Neural Network Optimization – saddle points and plateau
- 3.2 Non-convex optimization intuition
- 3.3 Overview of optimization algorithms
- 3.4 Momentum based algorithms
- 3.5 Algorithms with Adaptive Learning Rates



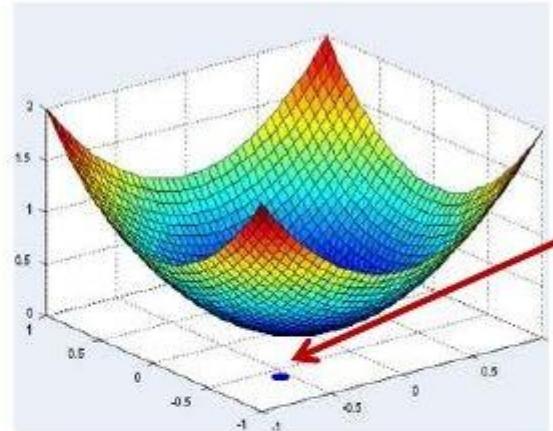
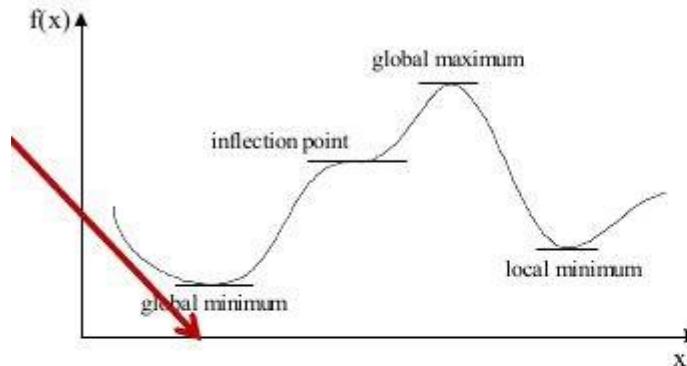
Optimization Algorithms



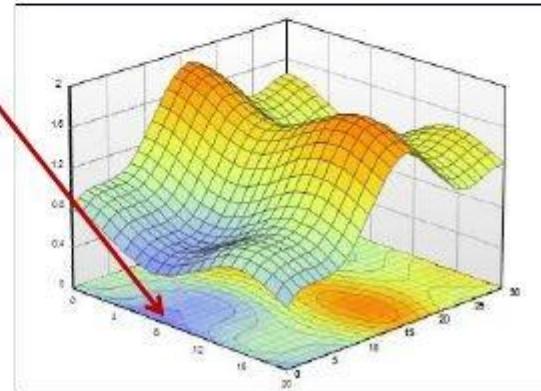
Optimization Algorithm

- Optimization algorithms train deep learning models.
- Optimization algorithms are the tools that allow
 - continue updating model parameters
 - to **minimize the value of the loss function**, as evaluated on the training set.
- In optimization, a loss function is often referred to as the **objective function** of the optimization problem.
- By tradition and convention most optimization algorithms are concerned with **minimization**.

Optimization



- General problem of optimization: find the value of x where $f(x)$ is minimum



Why Optimization Algorithm?

- The **performance** of the optimization algorithm directly affects the model's training efficiency.
- Understanding the principles of different optimization algorithms and the role of their hyperparameters will enable us to **tune the hyperparameters in a targeted manner** to improve the performance of deep learning models.
- **The goal of optimization is to reduce the training error.** The goal of deep learning is to reduce the generalization error, this requires reduction in overfitting also.

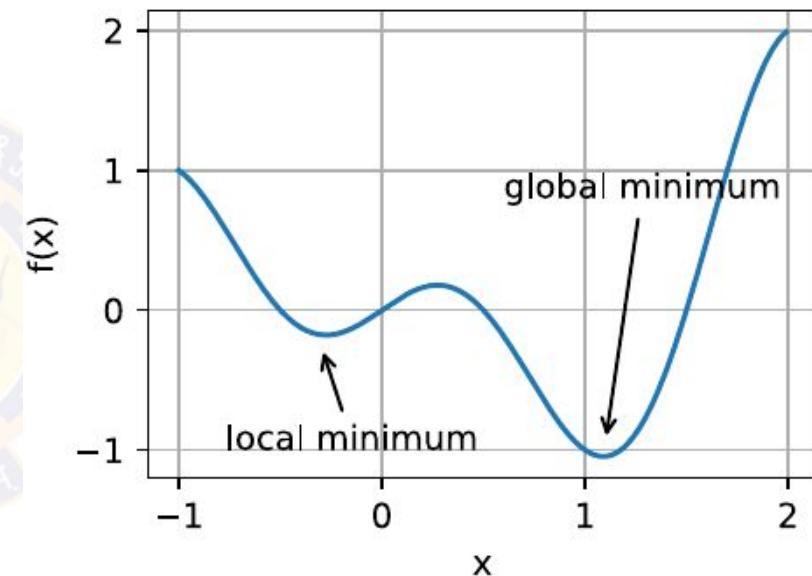
Optimization Challenges in Deep Learning

- Local minima
- Saddle points
- Vanishing gradients

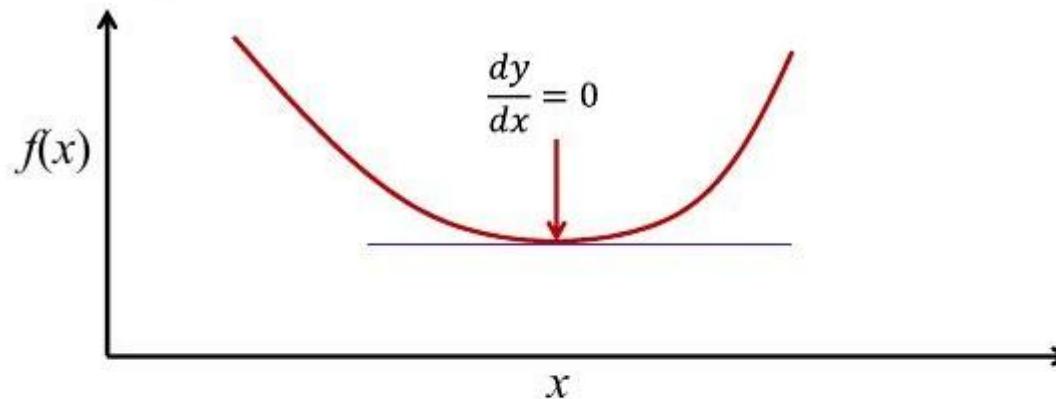


Local Minima

- For any objective function $f(x)$, if the value of $f(x)$ at x is smaller than the values of $f(x)$ at any other points in the vicinity of x , then $f(x)$ could be a **local minimum**.
- If the value of $f(x)$ at x is the minimum of the objective function over the entire domain, then $f(x)$ is the **global minimum**.
- In minibatch stochastic gradient descent, the natural variation of gradients over minibatches is able to dislodge the parameters from local minima.

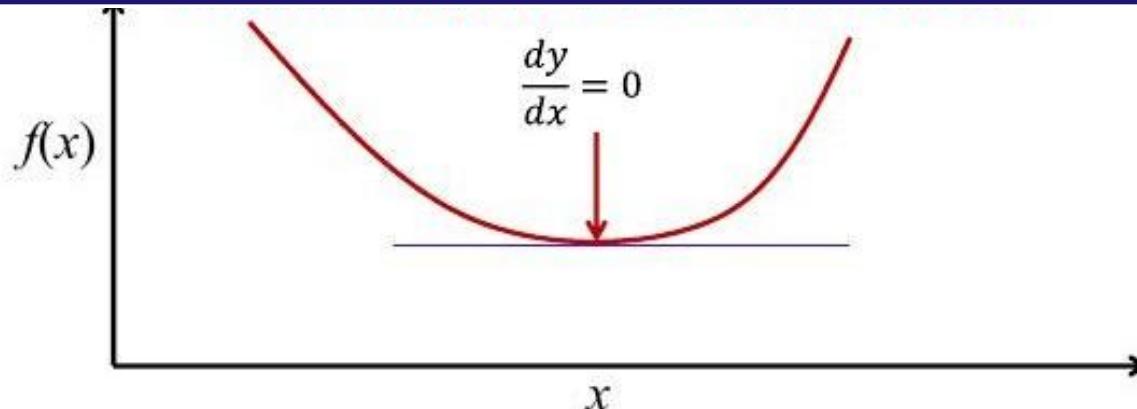


Finding Minimum of a Function



- Find the value x at which $f'(x) = 0$
 - Solve
- The solution is a “turning point”
 - Derivatives go from positive to negative or vice versa at this point
- But is it a minimum?

Finding Minimum of a Function



- Find the value x at which $f'(x) = 0$: Solve

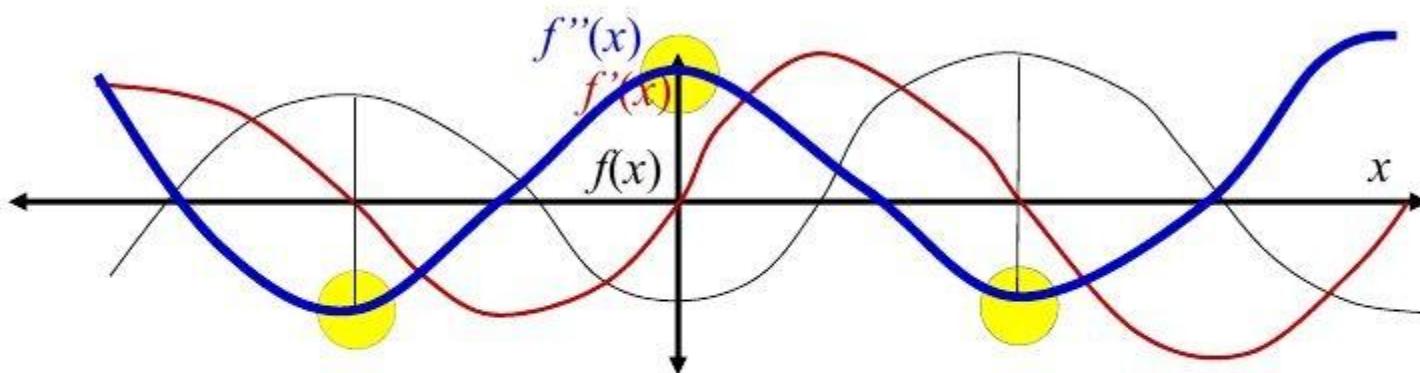
$$\frac{df(x)}{dx} = 0$$

- The solution x_{soln} is a turning point
- Check the double derivative at x_{soln} : compute

$$f''(x_{soln}) = \frac{df'(x_{soln})}{dx}$$

- If $f''(x_{soln})$ is positive x_{soln} is a minimum, otherwise it is a maximum

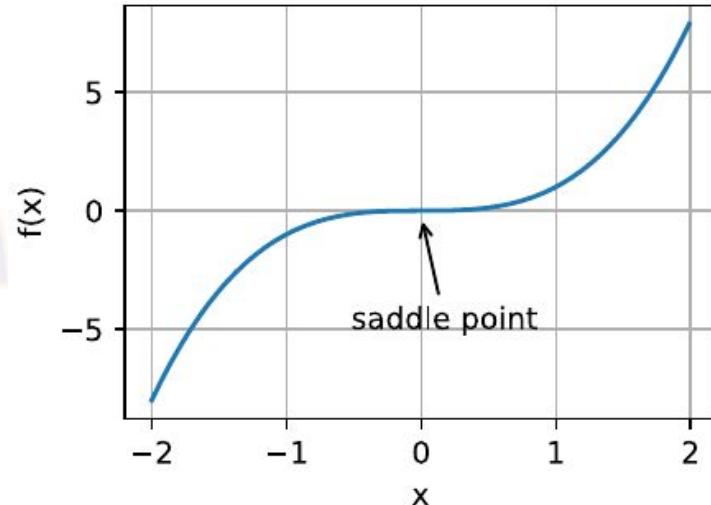
Derivatives of a Function



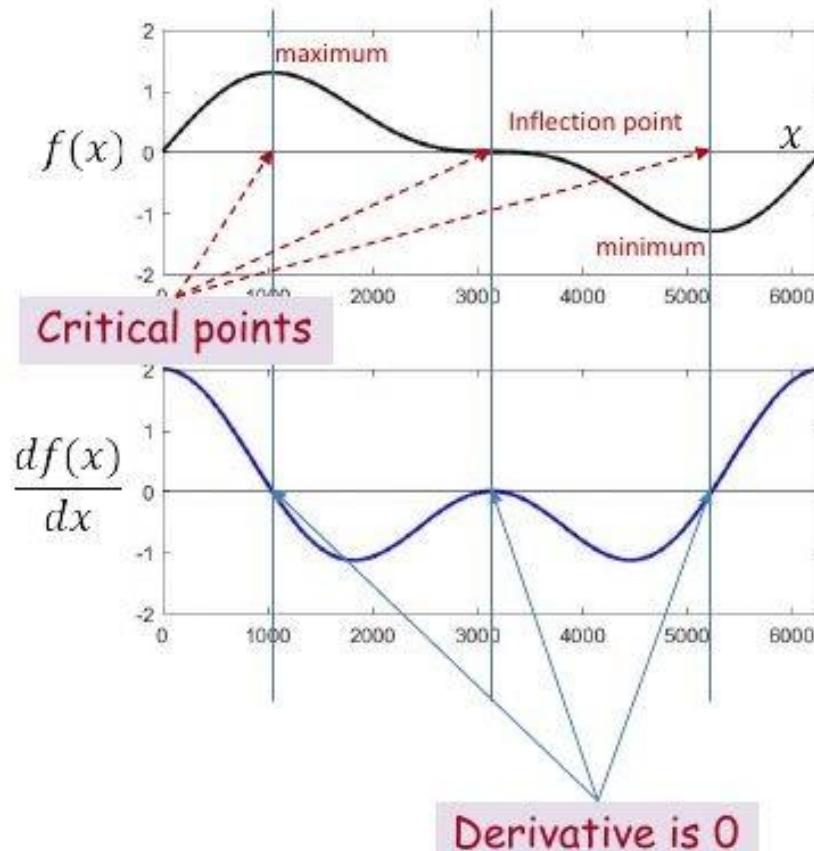
- Both *maxima* and *minima* are turning points
- Both *maxima* and *minima* have zero derivative
- The *second derivative* $f''(x)$ is $-ve$ at maxima and $+ve$ at minima!

Saddle points

- A saddle point is any location where all gradients of a function vanish but which is neither a global nor a local minimum.
- Eg: $f(x, y) = x^2 - y^2$
 - saddle point at $(0, 0)$
 - Maximum wrt y
 - Minimum wrt x

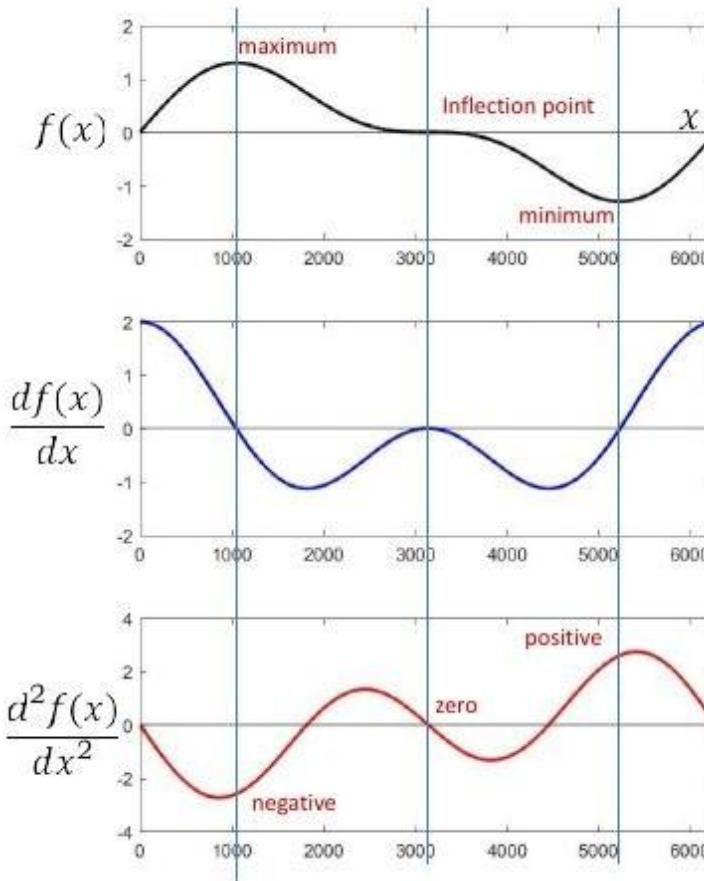


Derivatives at Saddle (Inflection) Point



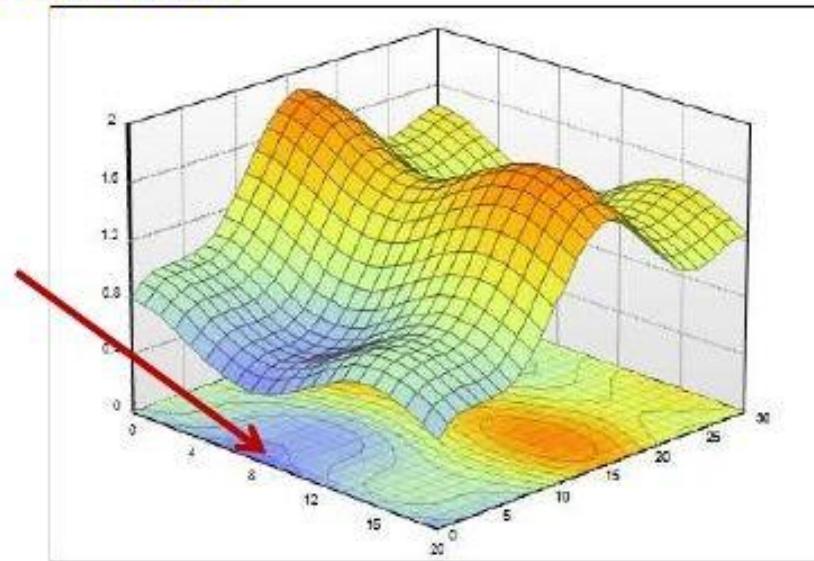
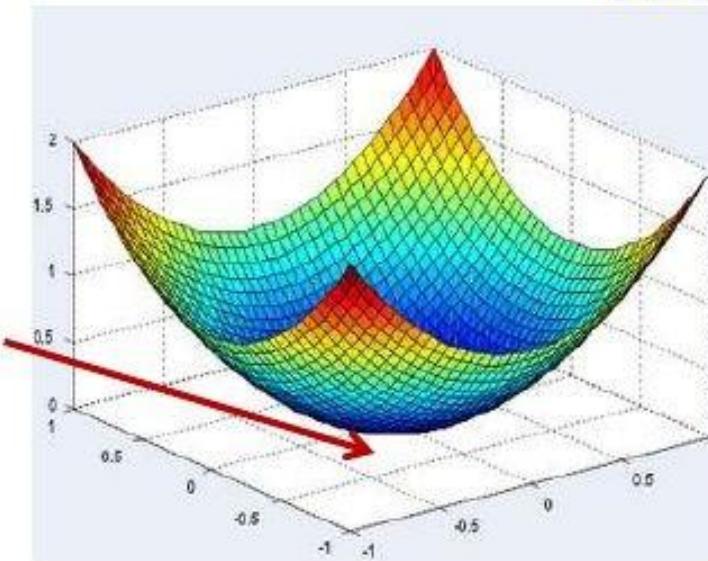
- All locations with zero derivative are *critical points*
 - These can be local maxima, local minima, or inflection points

Derivatives at Saddle (Inflection) Point



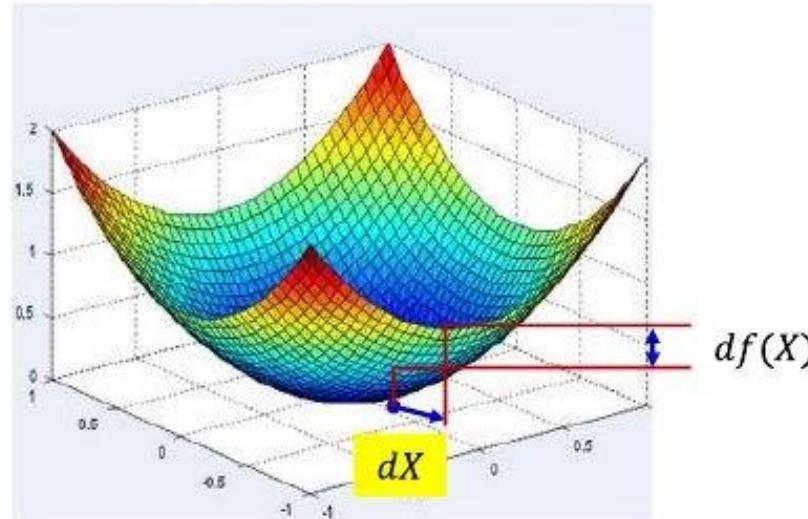
- All locations with zero derivative are *critical points*
 - These can be local maxima, local minima, or inflection points
- The second derivative is
 - ≥ 0 at minima
 - ≤ 0 at maxima
 - Zero at inflection points
- It's a little more complicated for functions of multiple variables..

Functions of Multiple Variables



- The optimum point is still “turning” point
 - Shifting in any direction will increase the value
 - For smooth functions, minuscule shifts will not result in any change at all
- We must find a point where shifting in any direction by a microscopic amount will not change the value of the function

Gradient of a Scalar Function

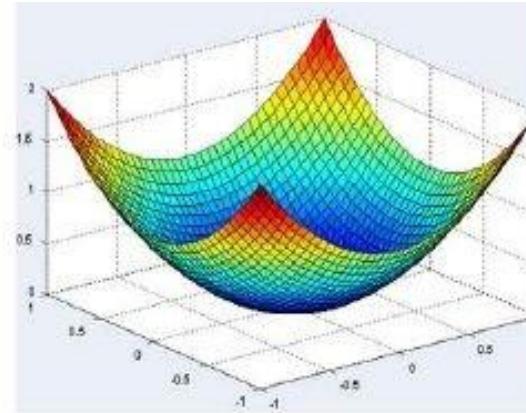


- The *Gradient* $\nabla f(X)$ of a scalar function $f(X)$ of a multi-variate input X is a multiplicative factor that gives us the change in $f(X)$ for tiny variations in X

$$df(X) = \nabla f(X) dX$$

Gradient of a Scalar Function with multiple variables

- Consider $f(X) = f(x_1, x_2, \dots, x_n)$



$$\nabla f(X) = \begin{bmatrix} \frac{\partial f(X)}{\partial x_1} & \frac{\partial f(X)}{\partial x_2} & \dots & \frac{\partial f(X)}{\partial x_n} \end{bmatrix}$$

- Check:

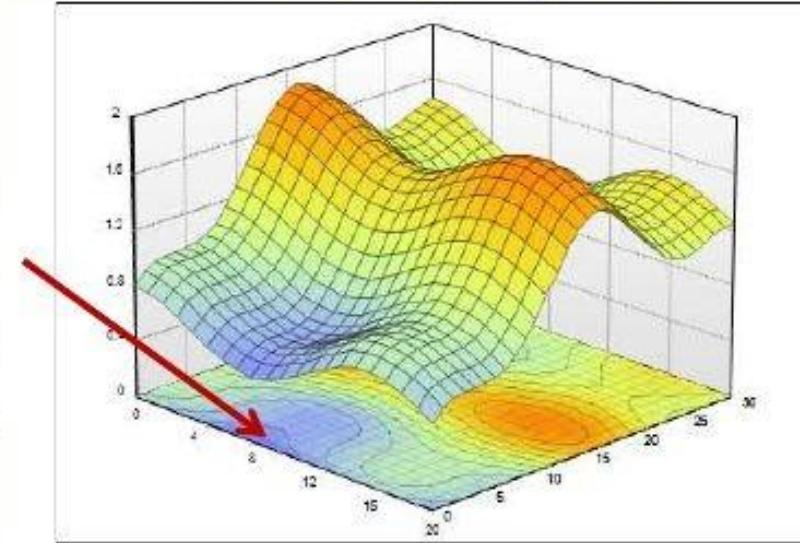
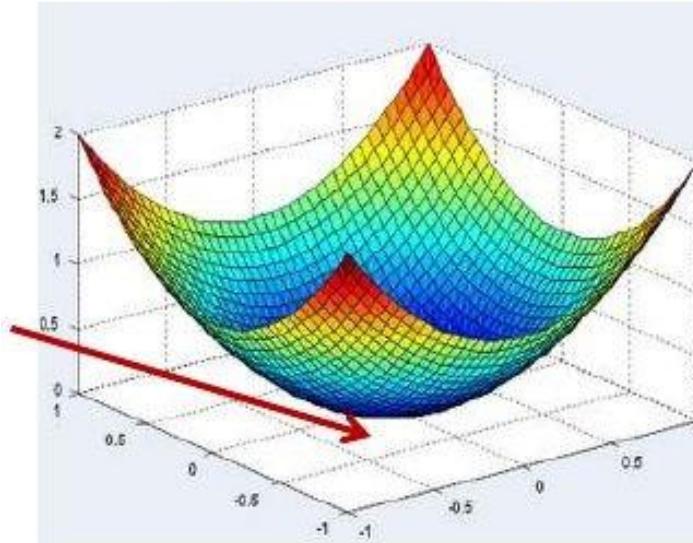
$$\begin{aligned} df(X) &= \nabla f(X) dX \\ &= \frac{\partial f(X)}{\partial x_1} dx_1 + \frac{\partial f(X)}{\partial x_2} dx_2 + \dots + \frac{\partial f(X)}{\partial x_n} dx_n \end{aligned}$$

Hessian

- The Hessian of a function $f(x_1, x_2, \dots, x_n)$ is given by the second derivative

$$\nabla^2 f(x_1, \dots, x_n) := \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdot & \cdot & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdot & \cdot & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdot & \cdot & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

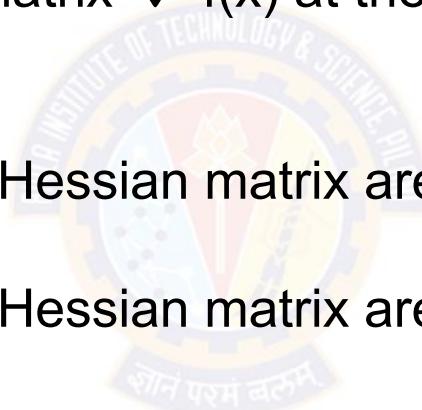
Gradient of a Scalar Function with multiple variables



- The optimum point is a turning point – the gradient will be 0

Solution of Unconstrained Minimization

1. Solve for X where the gradient equation equal to zero.
 $\nabla f(x) = 0$
2. Compute the Hessian matrix $\nabla^2 f(x)$ at the candidate solution and verify that
 - Local Minimum
 - Eigenvalues of Hessian matrix are all positive
 - Local Maximum
 - Eigenvalues of Hessian matrix are all negative
 - Saddle Point
 - Eigenvalues of Hessian matrix at the zero-gradient position are negative and positive



Example

- Minimize

$$f(x_1, x_2, x_3) = (x_1)^2 + x_1(1 - x_2) - (x_2)^2 - x_2x_3 + (x_3)^2 + x_3$$

- Gradient

$$\nabla f = \begin{bmatrix} 2x_1 + 1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 + 2x_3 + 1 \end{bmatrix}^T$$

Example

- Set the gradient to null

$$\nabla f = 0 \Rightarrow \begin{bmatrix} 2x_1 + 1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 + 2x_3 + 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

- Solving the 3 equations system with 3 unknowns

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$$

Example

- Compute the Hessian matrix $\nabla^2 f = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$

- Evaluate the eigenvalues of the Hessian matrix

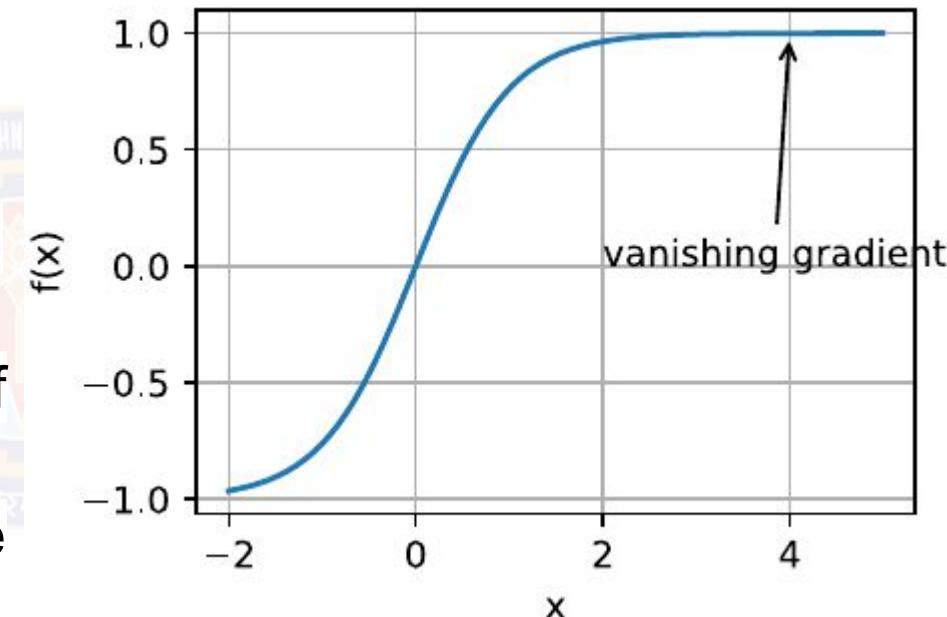
$$\lambda_1 = 3.414, \quad \lambda_2 = 0.586, \quad \lambda_3 = 2$$

- All the eigenvalues are positives \Rightarrow the Hessian matrix is positive definite

- The point $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$ is a minimum

Vanishing gradients

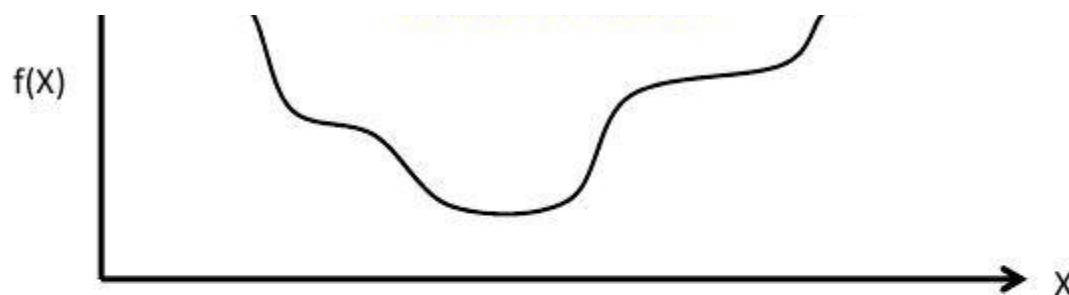
- Function $f(x) = \tanh(x)$
- $f'(x) = 1 - \tanh^2(x)$
 - $f'(4) = 0.0013.$
- Gradient of f is close to nil.
- Vanishing gradients can cause optimization to stall.
 - Reparameterization of the problem helps.
 - Good initialization of the parameters





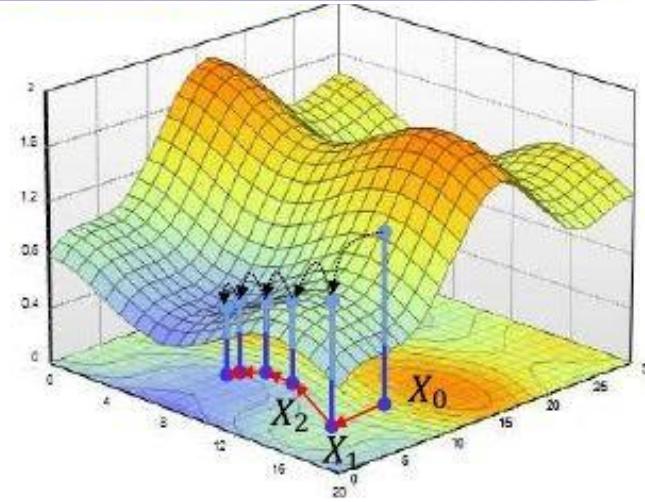
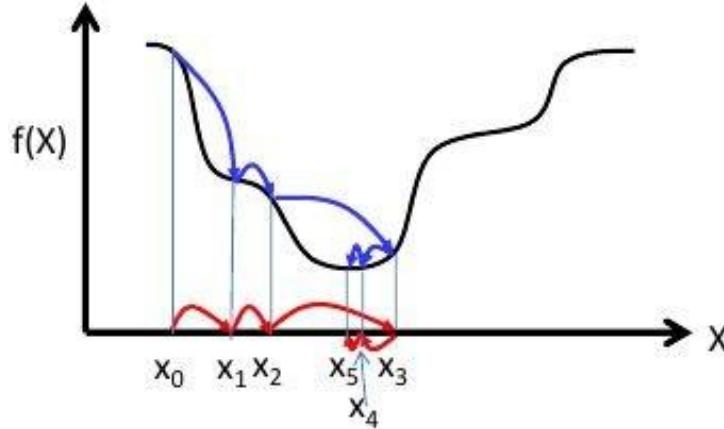
Gradient Descent

How to find Global Minima?



- Often it is not possible to simply solve $\nabla f(X) = 0$
 - The function to minimize/maximize may have an intractable form
- In these situations, iterative solutions are used
 - Begin with a “guess” for the optimal X and refine it iteratively until the correct value is obtained

Find Global Minima Iteratively



- Iterative solutions
 - Start from an initial guess X_0 for the optimal X
 - Update the guess towards a (hopefully) “better” value of $f(X)$
 - Stop when $f(X)$ no longer decreases
- Problems:
 - Which direction to step in
 - How big must the steps be

Approach of Gradient Descent



- Iterative solution:
 - Start at some point
 - Find direction in which to shift this point to decrease error
 - This can be found from the derivative of the function
 - A positive derivative → moving left decreases error
 - A negative derivative → moving right decreases error
 - Shift point in this direction

Approach of Gradient Descent



- Iterative solution: Trivial algorithm
 - Initialize x^0
 - While $f'(x^k) \neq 0$
 - If $\text{sign}(f'(x^k))$ is positive:
 - $x^{k+1} = x^k - \text{step}$
 - Else
 - $x^{k+1} = x^k + \text{step}$
 - What must step be to ensure we actually get to the optimum?

Approach of Gradient Descent



- Iterative solution: Trivial algorithm
 - Initialize x_0
 - While $f'(x^k) \neq 0$
 - $x^{k+1} = x^k - \eta^k f'(x^k)$
 - η^k is the “step size”

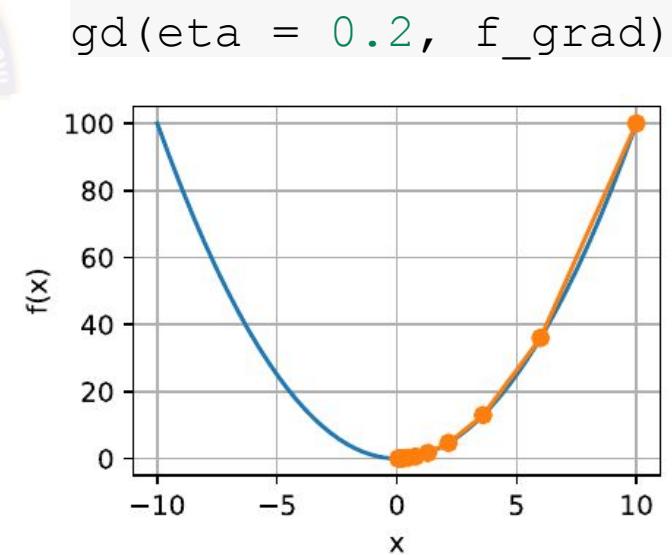
gd(eta, grad)

new = old - eta * grad

Approach of Gradient Descent

- First order Gradient Descent algos consider the first order derivatives to get the value (magnitude) and direction of update.

```
def gd(eta, f_grad):  
    x = 10.0  
    results = [x]  
    for i in range(10):  
        x -= eta * f_grad(x)  
        results.append(float(x))  
    return results
```

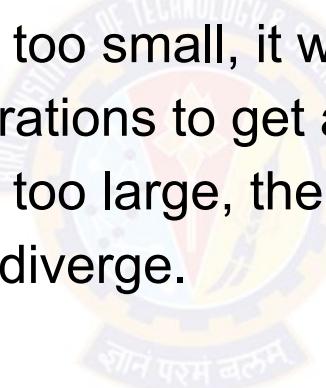


Effect of Learning Rate on Gradient Descent



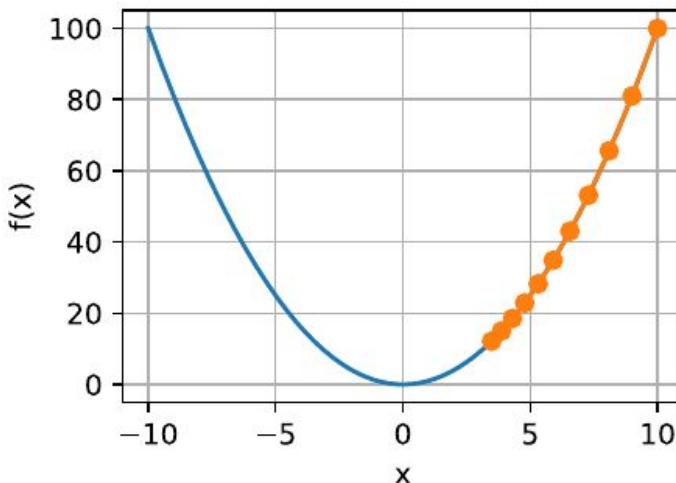
Learning Rate

- The role of the learning rate is to moderate the degree to which weights are changed at each step.
- Learning rate η is set by the algorithm designer.
- If the learning rate that is too small, it will cause x to update very slowly, requiring more iterations to get a better solution.
- If the learning rate that is too large, the solution oscillates and in the worst case it might even diverge.

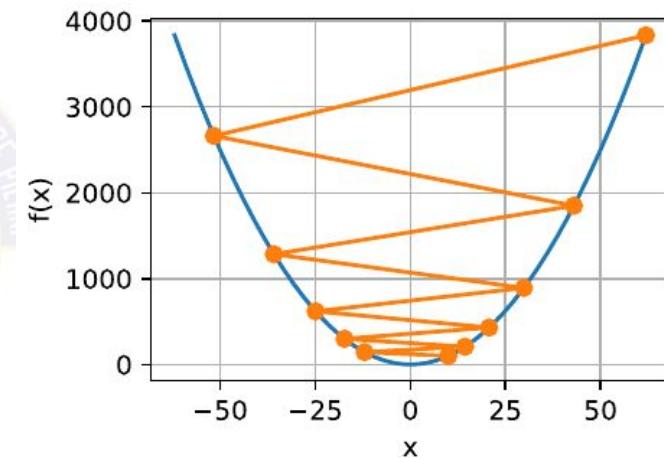


Learning rate and Gradient Descent

```
gd(eta = 0.05, f_grad)
```



```
gd(eta = 1.1, f_grad)
```



Slow Learning

If we pick eta too small, we make little progress.

Oscillations

If we pick eta too large, the solution oscillates and in the worst case it might diverge.

Stochastic Gradient Descent



Stochastic Gradient Descent

- In deep learning, the objective function is the average of the loss functions for each example in the training dataset.
- Given a training dataset of n examples, let $f_i(x)$ is the loss function with respect to the training example of index i , where x is the parameter vector.
- The objective function

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}).$$

- The gradient of the objective function at x

$$\nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}).$$

- Update x as

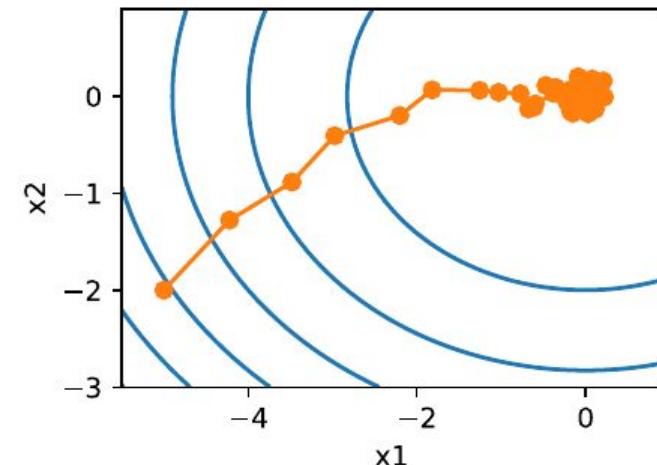
$$\mathbf{x} \leftarrow \mathbf{x} - \eta \nabla f_i(\mathbf{x}),$$

- Computational cost of each iteration is $O(1)$.

SGD with Constant Learning Rate

- Trajectory of the variables in the stochastic gradient descent is much more noisy. This is due to the stochastic nature of the gradient. Even near the minimum, uncertainty is injected by the instantaneous gradient via $\eta \nabla f_i(x)$.

```
def sgd(x1, x2, s1, s2, f_grad):  
    g1, g2 = f_grad(x1, x2)  
    # Simulate noisy gradient  
    g1 += np.random.normal(0.0, 1, (1,))  
    g2 += np.random.normal(0.0, 1, (1,))  
    eta_t = eta * lr()  
    return (x1 - eta_t * g1, x2 - eta_t * g2, 0, 0)
```



Dynamic Learning Rate



Dynamic Learning Rate

- Replace η with a time-dependent learning rate $\eta(t)$
 - adds to the complexity of controlling convergence of an optimization algorithm.
- A few basic strategies that adjust η over time.

$$\eta(t) = \eta_i \text{ if } t_i \leq t \leq t_{i+1} \quad \text{piecewise constant}$$

$$\eta(t) = \eta_0 \cdot e^{-\lambda t} \quad \text{exponential decay}$$

$$\eta(t) = \eta_0 \cdot (\beta t + 1)^{-\alpha} \quad \text{polynomial decay}$$

1. Piecewise constant
 - a. decrease the learning rate, e.g., whenever progress in optimization stalls.
 - b. This is a common strategy for training deep networks. A
2. Exponential decay
 - a. Leads to premature stopping before the algorithm has converged.
3. Polynomial decay with $\alpha = 0.5$.

Dynamic Learning Rate

- Replace η with a time-dependent learning rate $\eta(t)$
 - adds to the complexity of controlling convergence of an optimization algorithm.
- A few basic strategies that adjust η over time.

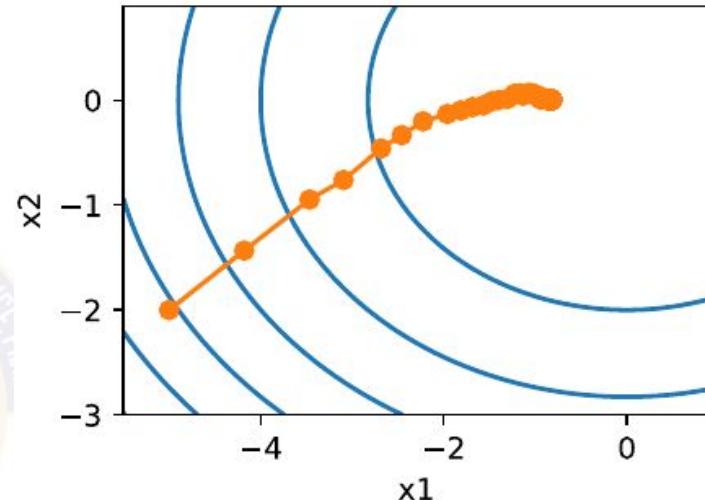
$$\eta(t) = \eta_i \text{ if } t_i \leq t \leq t_{i+1} \quad \text{piecewise constant}$$

$$\eta(t) = \eta_0 \cdot e^{-\lambda t} \quad \text{exponential decay}$$

$$\eta(t) = \eta_0 \cdot (\beta t + 1)^{-\alpha} \quad \text{polynomial decay}$$

Exponential Decay

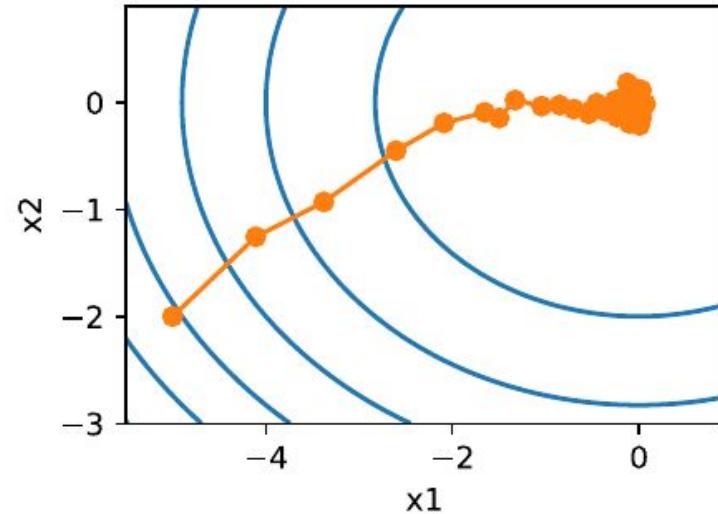
```
def exponential_lr():
    global t
    t += 1
    return math.exp(-0.1 * t)
```



- Variance in the parameters is significantly reduced.
- The algorithm fails to converge at all.

Polynomial decay

```
def polynomial_lr():  
    global t  
    t += 1  
    return (1 + 0.1 * t) ** (-0.5)
```



- Use a polynomial decay
 - Learning rate decays with the inverse square root of the number of steps
 - Convergence gets better after only 50 steps.
-

Review

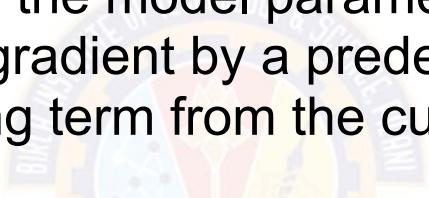
- Gradient descent
 - Uses the full dataset to compute gradients and to update parameters, one pass at a time.
 - Gradient Descent is not particularly data efficient whenever data is very similar.
- Stochastic Gradient descent
 - Processes one observation at a time to make progress.
 - Stochastic Gradient Descent is not particularly computationally efficient since CPUs and GPUs cannot exploit the full power of vectorization.
 - For noisy gradients, choice of the learning rate is critical.
 - If we decrease it too rapidly, convergence stalls.
 - If we are too lenient, we fail to converge to a good enough solution since noise keeps on driving us away from optimality.
- Minibatch SGD
 - Accelerate computation, or better or computational efficiency.
 - Averaging gradients reduce the amount of variance.

Minibatch Stochastic Gradient Descent



Minibatch Stochastic Gradient Descent

- In each iteration, we first randomly sample a minibatch B consisting of a fixed number of training examples.
- We then compute the derivative (gradient) of the average loss on the minibatch with regard to the model parameters.
- Finally, we multiply the gradient by a predetermined positive value η and subtract the resulting term from the current parameter values.


$$(\mathbf{w}, b) \leftarrow (\mathbf{w}, b) - \frac{\eta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \partial_{(\mathbf{w}, b)} l^{(i)}(\mathbf{w}, b).$$

- $|\mathcal{B}|$ represents the number of examples in each minibatch (the batch size) and η denotes the learning rate.

Minibatch Stochastic Gradient Descent Algorithm

- Gradients at time t is calculated as

$$\mathbf{g}_{t,t-1} = \partial_{\mathbf{w}} \frac{1}{|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} f(\mathbf{x}_i, \mathbf{w}_{t-1}) = \frac{1}{|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} \mathbf{h}_{i,t-1}.$$

- $|\mathcal{B}|$ represents the number of examples in each minibatch (the batch size) and η denotes the learning rate.
- $\mathbf{h}_{i,t-1} = \partial_{\mathbf{w}} f(\mathbf{x}_i, \mathbf{w}_{t-1})$ is the stochastic gradient descent for sample i using the weights updated at time $t - 1$.

SGD Algorithm

Algorithm 1: STOCHASTIC GRADIENT DESCENT ALGORITHM

Data: Learning Rate α

Data: Initial Parameters θ

```
1 while stopping criterion not met do
2   for  $t$  in range ( $1, m/batch\_size$ ) do
3     Perform Forward prop on  $X^{\{t\}}$  to compute  $\hat{y}$  and cost  $J(\theta)$ 
4     Perform Backprop on  $(X^{\{t\}}, Y^{\{t\}})$  to compute gradient  $g$ 
5     Update parameters as  $\theta \leftarrow \theta - \alpha g$ 
```

Momentum



Leaky Average in Minibatch SGD

- Replace the gradient computation by a “leaky average “ for better variance reduction. $\beta \in (0, 1)$.

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + \mathbf{g}_{t,t-1}$$

- This effectively replaces the instantaneous gradient by one that's been averaged over multiple past gradients.
- \mathbf{v} is called momentum. It accumulates past gradients.
- Large β amounts to a long-range average and small β amounts to only a slight correction relative to a gradient method.
- The new gradient replacement no longer points into the direction of steepest descent on a particular instance any longer but rather in the direction of a weighted average of past gradients.

Momentum Algorithm

Algorithm 2: STOCHASTIC GRADIENT DESCENT WITH MOMENTUM

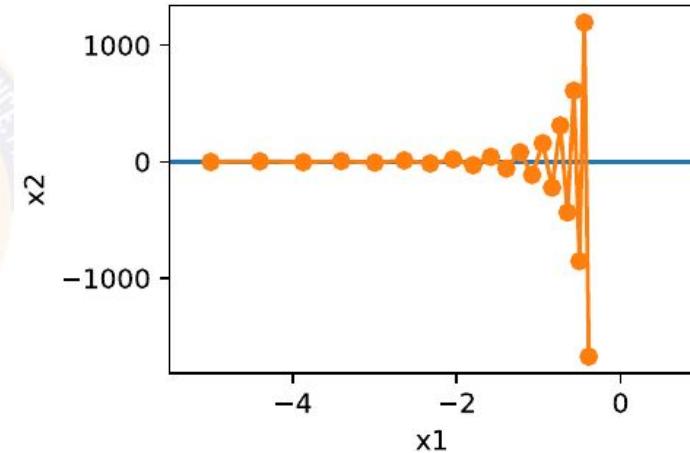
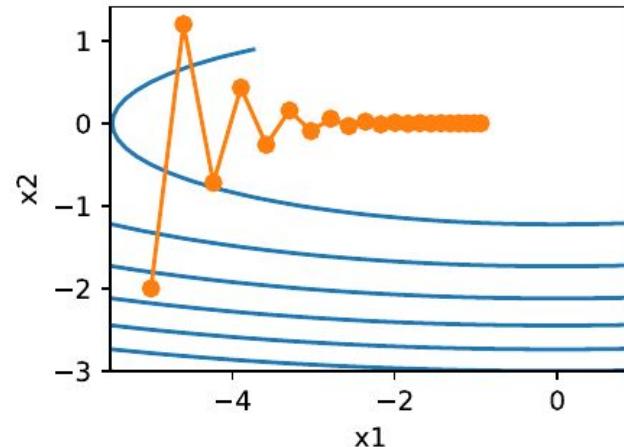
Data: Learning Rate α , Momentum Parameter β

Data: Initial Parameters θ , Inital velocities v

```
1 while stopping criterion not met do
2   for  $t$  in range ( $1, m/batch\_size$ ) do
3     Perform Forward prop on  $X^{\{t\}}$  to compute  $\hat{y}$  and cost  $J(\theta)$ 
4     Perform Backprop on  $(X^{\{t\}}, Y^{\{t\}})$  to compute gradient  $g$ 
5     Update velocities as  $v = \beta v + (1 - \beta)g$ 
6     Update parameters as  $\theta \leftarrow \theta - \alpha v$ 
```

Momentum Method Example

- Consider a moderately distorted ellipsoid objective $f(\mathbf{x}) = 0.1x_1^2 + 2x_2^2$.
- f has its minimum at $(0, 0)$. This function is very flat in x_1 direction.



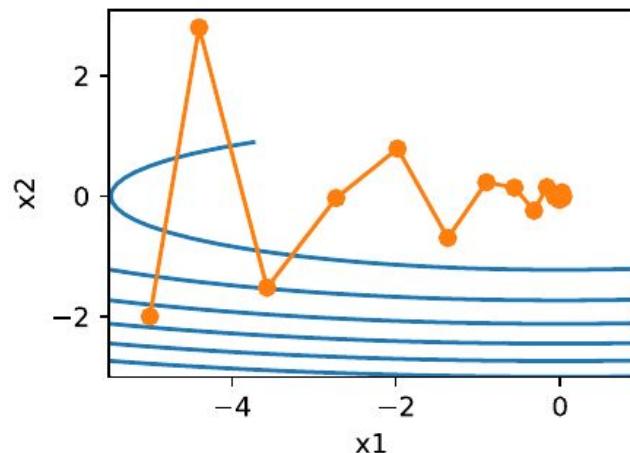
- For $\eta = 0.4$. Without momentum
- The gradient in the x_2 direction oscillates than in the horizontal x_1 direction.

- For $\eta = 0.6$. Without momentum
- Convergence in the x_1 direction improves but the overall solution quality is diverging.

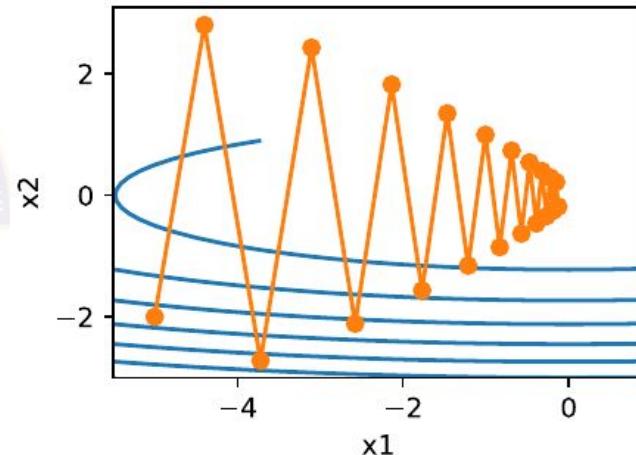
Momentum Method Example

- Consider a moderately distorted ellipsoid objective
- Apply momentum for eta = 0.6

$$f(\mathbf{x}) = 0.1x_1^2 + 2x_2^2$$



- For beta = 0.5.
- Converges well. Lesser oscillations.
Larger steps in x_1 direction.



- For beta = 0.25.
- Reduced convergence. More oscillations. Larger magnitude of oscillations.

Momentum method: Summary

- Momentum replaces gradients with a leaky average over past gradients. This accelerates convergence significantly.
- Momentum prevents stalling of the optimization process that is much more likely to occur for stochastic gradient descent.
- The effective number of gradients is given by $1 / (1 - \beta)$ due to exponentiated downweighting of past data.
- Implementation is quite straightforward but it requires us to store an additional state vector (momentum v).



Adagrad

Adagrad

- Used for features that occur infrequently (sparse features)
- Adagrad uses aggregate of the squares of previously observed gradients.



Adagrad Algorithm

- Variable s_t to accumulate past gradient variance.

$$\mathbf{g}_t = \partial_{\mathbf{w}} l(y_t, f(\mathbf{x}_t, \mathbf{w})),$$

$$\mathbf{s}_t = \mathbf{s}_{t-1} + \mathbf{g}_t^2,$$

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \cdot \mathbf{g}_t.$$

- Operations are applied coordinate wise. $\sqrt{1/v}$ has entries $\sqrt{1/v_i}$ and $u \cdot v$ has entries $u_i v_i$. η is the learning rate and ϵ is an additive constant that ensures that we do not divide by 0.
- Initialize $s_0 = 0$.

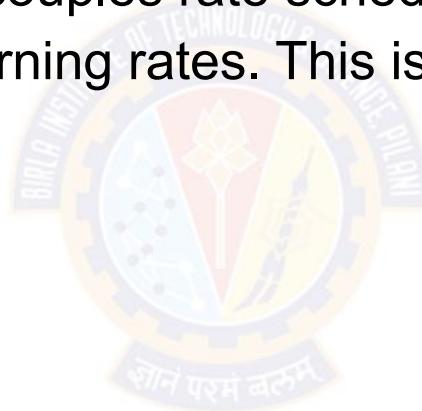
Adagrad: Summary

- Adagrad decreases the learning rate dynamically on a per-coordinate basis.
- It uses the magnitude of the gradient as a means of adjusting how quickly progress is achieved - coordinates with large gradients are compensated with a smaller learning rate.
- If the optimization problem has a rather uneven structure Adagrad can help mitigate the distortion.
- Adagrad is particularly effective for sparse features where the learning rate needs to decrease more slowly for infrequently occurring terms.
- On deep learning problems Adagrad can sometimes be too aggressive in reducing learning rates.



RMSProp

- Adagrad uses learning rate that decreases at a predefined schedule of effectively $O(t^{-1/2})$.
- RMSProp algorithm decouples rate scheduling from coordinate-adaptive learning rates. This is essential for non-convex optimization.



RMSProp Algorithm

- Use leaky average to accumulate past gradient variance.

$$\mathbf{s}_t \leftarrow \gamma \mathbf{s}_{t-1} + (1 - \gamma) \mathbf{g}_t^2,$$
$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t.$$

- Parameter gamma > 0.
- The constant $\epsilon > 0$ is set to 10^{-6} to ensure that we do not suffer from division by zero or overly large step sizes.

RMSProp Algorithm

Algorithm 3: RMSPROP

Data: Learning Rate α , Decay Rate ρ

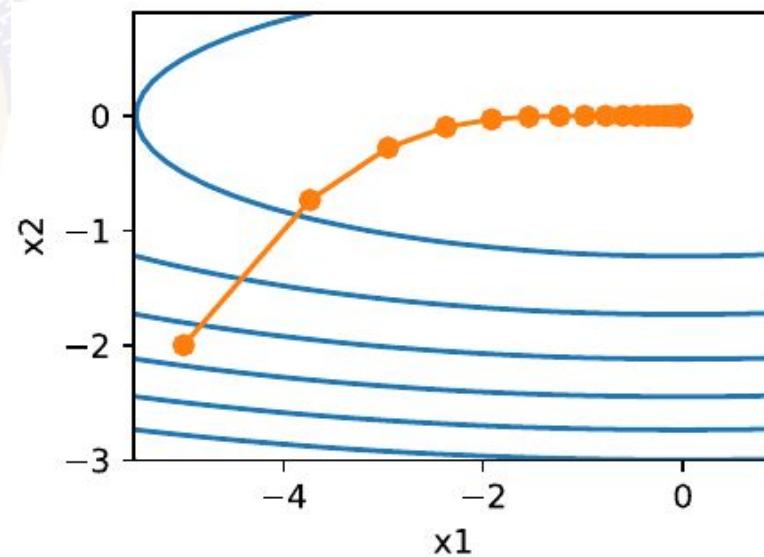
Data: Small constant $\epsilon = 10^{-7}$ to stabilize division

Data: Initial Parameters θ

- 1 Initialize accumulated gradient $r = 0$
 - 2 **while** stopping criterion not met **do**
 - 3 **for** t in range $(1, m/batch_size)$ **do**
 - 4 Perform Forward prop on $X^{\{t\}}$ to compute \hat{y} and cost $J(\theta)$
 - 5 Perform Backprop on $(X^{\{t\}}, Y^{\{t\}})$ to compute gradient g
 - 6 Compute accumulated squared gradient $r = \rho r + (1 - \rho)g \odot g$
 - 7 Compute update parameter as $\Delta\theta = \frac{-\alpha}{\sqrt{r+\epsilon}} \odot g$
 - 8 Update parameters as $\theta \leftarrow \theta + \Delta\theta$
-

RMSProp Example

```
ef rmsprop_2d(x1, x2, s1, s2):  
    g1, g2, eps = 0.2 * x1, 4 * x2, 1e-6  
    s1 = gamma * s1 + (1 - gamma) * g1 ** 2  
    s2 = gamma * s2 + (1 - gamma) * g2 ** 2  
    x1 -= eta / math.sqrt(s1 + eps) * g1  
    x2 -= eta / math.sqrt(s2 + eps) * g2  
    return x1, x2, s1, s2
```



Adagrad: Summary

- RMSProp is very similar to Adagrad as both use the square of the gradient to scale coefficients.
- RMSProp shares with momentum the leaky averaging. However, RMSProp uses the technique to adjust the coefficient-wise preconditioner.
- The learning rate needs to be scheduled by the experimenter in practice.
- The coefficient γ (gamma) determines how long the history is when adjusting the per-coordinate scale.

Adam



Review of techniques learned so far

1. Stochastic gradient descent
 - o more effective than Gradient Descent when solving optimization problems, e.g., due to its inherent resilience to redundant data.
2. Minibatch Stochastic gradient descent
 - o affords significant additional efficiency arising from vectorization, using larger sets of observations in one minibatch. This is the key to efficient multi-machine, multi-GPU and overall parallel processing.
3. Momentum
 - o added a mechanism for aggregating a history of past gradients to accelerate convergence.
4. Adagrad
 - o used per-coordinate scaling to allow for a computationally efficient preconditioner.
5. RMSProp
 - o decoupled per-coordinate scaling from a learning rate adjustment.

Adam

- Adam combines all these techniques into one efficient learning algorithm.
- More robust and effective optimization algorithms to use in deep learning.
- Adam can diverge due to poor variance control. (disadvantage)
- Adam uses exponential weighted moving averages (also known as leaky averaging) to obtain an estimate of both the momentum and also the second moment of the gradient.

Adam Algorithm

- State variables $\mathbf{v}_t \leftarrow \beta_1 \mathbf{v}_{t-1} + (1 - \beta_1) \mathbf{g}_t,$
 $\mathbf{s}_t \leftarrow \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2.$
- β_1 and β_2 are nonnegative weighting parameters. Common choices for them are $\beta_1 = 0.9$ and $\beta_2 = 0.999$. That is, the variance estimate moves much more slowly than the momentum term.
- Initialize $v_0 = s_0 = 0$
- Normalize the state variables

$$\hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_1^t} \text{ and } \hat{\mathbf{s}}_t = \frac{\mathbf{s}_t}{1 - \beta_2^t}.$$

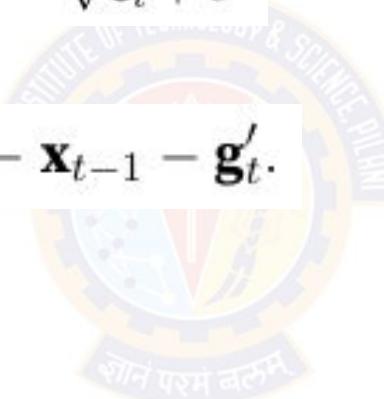
Adam Algorithm

- Rescale the gradient

$$\mathbf{g}'_t = \frac{\eta \hat{\mathbf{v}}_t}{\sqrt{\hat{\mathbf{s}}_t} + \epsilon}.$$

- Compute updates

$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \mathbf{g}'_t.$$



Adam Algorithm

Algorithm 4: ADAM

Data: Learning Rate $\alpha = 0.001$, Momentum parameter $\beta = 0.9$, Decay Rate $\rho = 0.999$, Small constant $\epsilon = 10^{-7}$ to stabilize division

Data: Initial Parameters θ

Data: Initialize first and second moments $s = 0, r = 0$

```
1 while stopping criterion not met do
2   for t in range (1, m/batch_size) do
3     Perform Forward prop on  $X^{\{t\}}$  to compute  $\hat{y}$  and cost  $J(\theta)$  ;
4     Perform Backprop on  $(X^{\{t\}}, Y^{\{t\}})$  to compute gradient  $g$  ;
5     Update first moment as  $s \leftarrow \beta s + (1 - \beta)g$  ;
6     Update second moment as  $r = \rho r + (1 - \rho)g \odot g$  ;
7     Perform bias correction as  $\hat{s} = \frac{s}{1-\beta^t}$  and  $\hat{r} = \frac{r}{1-\rho^t}$  ;
8     Compute update parameter as  $\Delta\theta = \frac{-\alpha\hat{s}}{\sqrt{\hat{r}+\epsilon}}$  ;
9     Update parameters as  $\theta \leftarrow \theta + \Delta\theta$  ;
```

Adam: Summary

- Adam combines features of many optimization algorithms into a fairly robust update rule.
- Adam uses bias correction to adjust for a slow startup when estimating momentum and a second moment.
- For gradients with significant variance we may encounter issues with convergence. They can be amended by using larger minibatches or by switching to an improved estimate for state variables. Yogi algorithm offers such an alternative.

Learning Rate: Summary

1. Adjusting the learning rate is often just as important as the actual algorithm.
2. Magnitude of the learning rate matters. If it is too large, optimization diverges, if it is too small, it takes too long to train or we end up with a suboptimal result. Momentum algo helps.
3. The rate of decay is just as important. If the learning rate remains large we may simply end up bouncing around the minimum and thus not reach optimality. We want the rate to decay, but probably more slowly than $O(t^{-1/2})$.
4. Initialization pertains both to how the parameters are set initially and also how they evolve initially. This is known as **warmup**, i.e., how rapidly we start moving towards the solution initially.

Numerical Problems



Question with Solution

Error surface is given by $E(x, y, z) = 3x^2 + 2y^2 + 4z^2 + 6$. Assume gradient descent is used to find the minimum of this error surface. What is the optimal learning rate that leads to fastest convergence to the global minimum?

Explanation:

$$E(x, y, z) = \frac{1}{2} * (6x^2 + 4y^2 + 8z^2) + 6$$

$$\eta_{xopt} = 1/6$$

$$\eta_{yopt} = 1/4$$

$$\eta_{zopt} = 1/8$$

optimal learning rate for convergence = $\min[\eta_{xopt}, \eta_{yopt}, \eta_{zopt}] = 1/8$

Largest learning rate for convergence = $\min[2\eta_{xopt}, 2\eta_{yopt}, 2\eta_{zopt}] = 0.33$

Learning rate for divergence > $2\eta_{opt} = 2 * 1/8 = 0.25$

Question

Consider an error function

$$E(w_1, w_2) = 0.05 + \frac{(w_1 - 3)^2}{4} + \frac{(w_2 - 4)^2}{9} - \frac{(w_1 - 3)(w_2 - 4)}{6}$$

. Different variants of gradient descent algorithm can be used to minimize this error function w.r.t (w_1, w_2) . Assume $(w_1, w_2) = (1, 1)$ at time $(t-1)$ and after update $(w_1, w_2) = (1.5, 2.0)$ at time (t) . Assume $\alpha = 1.5, \beta = 0.6, \eta = 0.3$.

1. Compute the value that minimizes (w_1, w_2) . Compute the minimum possible value of error.
2. What will be value of (w_1, w_2) at time $(t+1)$ if standard gradient descent is used?
3. What will be value of (w_1, w_2) at time $(t+1)$ if momentum is used?
4. What will be value of (w_1, w_2) at time $(t+1)$ if RMSProp is used?
5. What will be value of (w_1, w_2) at time $(t+1)$ if Adam is used?

Solution

- (a) Compute the value that minimizes (w_1, w_2) . Compute the minimum possible value of error.

$$\frac{\partial E}{\partial w_1} = \frac{2(w_1 - 3)}{4} - \frac{(w_2 - 4)}{6} = 0$$

$$\frac{\partial E}{\partial w_2} = \frac{2(w_2 - 4)}{9} - \frac{(w_1 - 3)}{6} = 0$$

Solving $w_1 = 3$

$$w_2 = 4$$

Substituting the above $E = 0.05$

Solution

What will be value of (w_1, w_2) at time $(t + 1)$ if standard gradient descent is used?

$$\frac{\partial E}{\partial w_1} = \frac{2(w_1 - 3)}{4} - \frac{(w_2 - 4)}{6}$$

$$\frac{\partial E}{\partial w_2} = \frac{2(w_2 - 4)}{9} - \frac{(w_1 - 3)}{6}$$

at $(t + 1)$ $w_{1(t+1)} = w_{1(t)} - \eta \frac{\partial E}{\partial w_1} = 1.5 - 0.3 * \left[\frac{(1.5 - 3)}{2} - \frac{(2 - 4)}{6} \right] = 1.625$

$$w_{2(t+1)} = w_{2(t)} - \eta \frac{\partial E}{\partial w_2} = 2 - 0.3 * \left[\frac{2(2 - 4)}{9} - \frac{(1.5 - 3)}{6} \right] = 2.05$$

Solution

What will be value of (w_1, w_2) at time $(t + 1)$ if momentum is used?

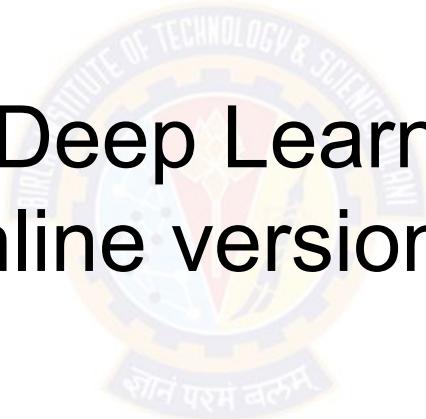
$$\frac{\partial E}{\partial w_1} = \frac{2(w_1 - 3)}{4} - \frac{(w_2 - 4)}{6}$$

$$\frac{\partial E}{\partial w_2} = \frac{2(w_2 - 4)}{9} - \frac{(w_1 - 3)}{6}$$

at $(t + 1)$

$$w_{1(t+1)} = w_{1(t)} - \eta \frac{\partial E}{\partial w_1} + \beta \Delta w_1$$
$$= 1.5 - 0.3 * \left[\frac{(1.5 - 3)}{2} - \frac{(2 - 4)}{6} \right] + 0.9 * (1.5 - 1) = 2.075$$

$$w_{2(t+1)} = w_{2(t)} - \eta \frac{\partial E}{\partial w_2} + \beta \Delta w_1$$
$$= 2 - 0.3 * \left[\frac{2(2 - 4)}{9} - \frac{(1.5 - 3)}{6} \right] + 0.9 * (2 - 1) = 2.958$$



Ref TB Dive into Deep Learning

- Chapter 12 (online version)

Next Session: Regularization





BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Neural Network

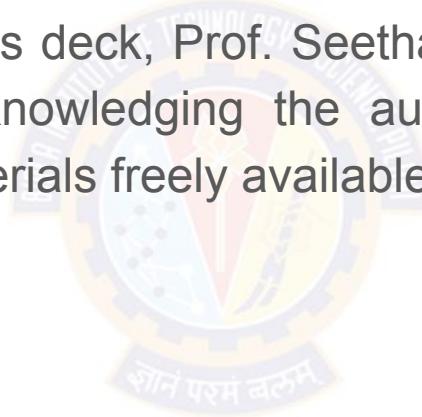
AIML Module 5

Seetha Parameswaran

BITS Pilani



The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Regularization Techniques



What we Learn....

4.1 Model Selection

4.2 Underfitting, and Overfitting

4.3 L1 and L2 Regularization

4.4 Dropout

4.5 Challenges - Vanishing and Exploding Gradients, Covariance shift

4.6 Parameter Initialization

4.7 Batch Normalization



Generalization in DNN



Generalization

- Goal is to discover patterns that generalize.
 - The goal is to discover patterns that capture regularities in the underlying population from which our training set was drawn.
 - Models are trained on a sample of data.
 - When working with finite samples, we run the risk that we might discover apparent associations that turn out not to hold up when we collect more data or on newer samples.
- The trained model should predict for newer or unseen data. This problem is called generalization.

Training Error and Generalization Error

- **Training error** is the error of our model as calculated on the training dataset.
 - Obtained while training the model.
- **Generalization error** is the expectation of our model's error, if an infinite stream of additional data examples drawn from the same underlying data distribution as the original sample were applied on the model.
 - Cannot be computed, but estimated.
 - Estimate the generalization error by applying the model to an independent test set, constituted of a random selection of data examples that were withheld from the training set.

Model Complexity

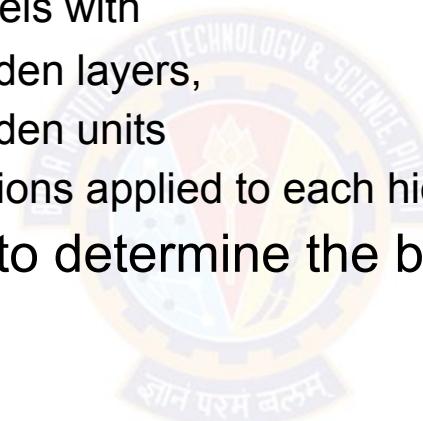
- Simple models and abundant data
 - Expect the generalization error to resemble the training error.
- More complex models and fewer examples
 - Expect the training error to go down but the generalization gap to grow.
- Model complexity
 - A model with more parameters might be considered more complex.
 - A model whose parameters can take a wider range of values might be more complex.
 - A neural network model that takes more training iterations are more complex, and one subject to early stopping (fewer training iterations) are less complex.
-

Factors that influence the generalizability of a model

1. The number of tunable parameters.
 - When the number of tunable parameters, called the degrees of freedom, is large, models tend to be more susceptible to overfitting.
2. The values taken by the parameters.
 - When weights can take a wider range of values, models can be more susceptible to overfitting.
3. The number of training examples.
 - It is trivially easy to overfit a dataset containing only one or two examples even if your model is simple. But overfitting a dataset with millions of examples requires an extremely flexible model.
- 4.

Model Selection

- Model selection is the process of selecting the final model after evaluating several candidate models.
- With MLPs, compare models with
 - different numbers of hidden layers,
 - different numbers of hidden units
 - different activation functions applied to each hidden layer.
- Use Validation dataset to determine the best among our candidate models.

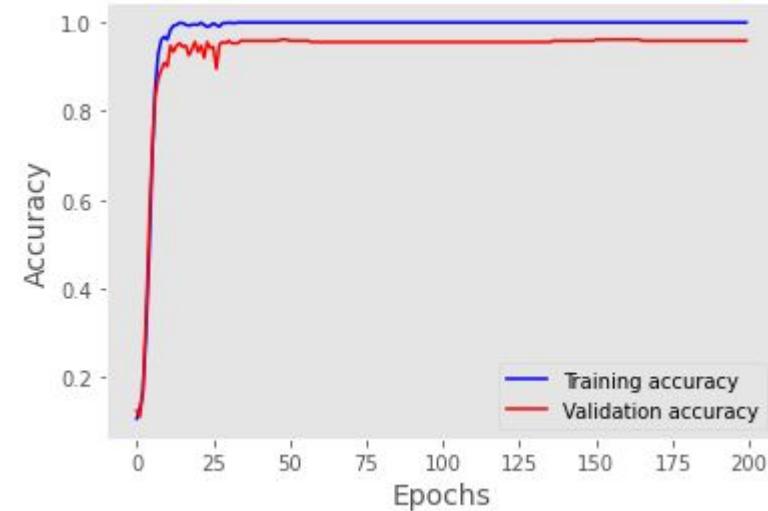
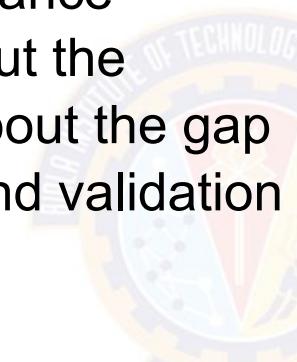


Validation dataset

- Never rely on the test data for model selection.
 - Risk of overfit the test data
- Do not rely solely on the training data for model selection
 - We cannot estimate the generalization error on the very data that we use to train the model.
- Split the data three ways, incorporating a validation dataset (or validation set) in addition to the training and test datasets.
- In deep learning, with millions of data available, the split is generally
 - Training = 98-99 % of the original dataset
 - Validation = 1-2 % of training dataset
 - Testing = 1-2 % of the original dataset

Just Right Model

- High Training accuracy
- High Validation accuracy
- Low Bias and Low Variance
- Usually care more about the validation error than about the gap between the training and validation errors.



```
844/844 [=====] - 3s 3ms/step - loss: 0.0848 - accuracy: 0.9750 - val_loss: 0.0714 - val_accuracy: 0.9790
Epoch 2/25
844/844 [=====] - 2s 2ms/step - loss: 0.0821 - accuracy: 0.9761 - val_loss: 0.0730 - val_accuracy: 0.9802
Epoch 3/25
844/844 [=====] - 2s 2ms/step - loss: 0.0793 - accuracy: 0.9769 - val_loss: 0.0706 - val_accuracy: 0.9798
Epoch 4/25
844/844 [=====] - 2s 3ms/step - loss: 0.0769 - accuracy: 0.9771 - val_loss: 0.0753 - val_accuracy: 0.9795
```

Underfitting

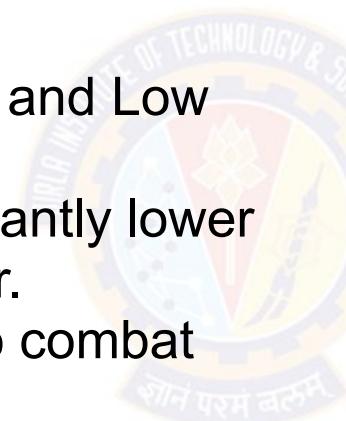
- Low Training accuracy and Low Validation accuracy.

```
Train on 2594 samples, validate on 495 samples
Epoch 1/200
2594/2594 [=====] - 6s 2ms/step - loss: 2.8738 - acc: 0.1191 - val_loss: 2.2041 - val_a
cc: 0.1232
Epoch 2/200
2594/2594 [=====] - 5s 2ms/step - loss: 2.1731 - acc: 0.1415 - val_loss: 2.1776 - val_a
cc: 0.1192
Epoch 3/200
2594/2594 [=====] - 5s 2ms/step - loss: 2.0862 - acc: 0.1974 - val_loss: 2.1261 - val_a
cc: 0.1596
Epoch 4/200
2594/2594 [=====] - 5s 2ms/step - loss: 1.9268 - acc: 0.2787 - val_loss: 2.0487 - val_a
cc: 0.2525
```

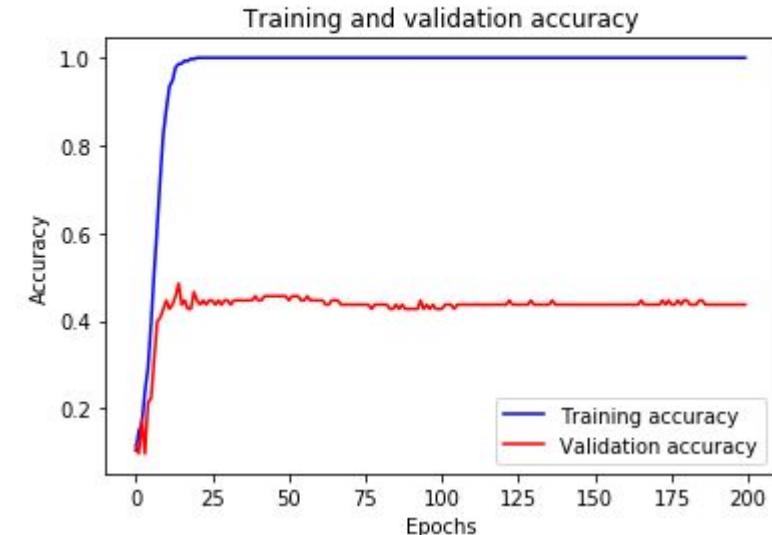
- Training error and validation error are both substantial but there is a little gap between them.
- The model is too simple (insufficiently expressive) to capture the pattern that we are trying to model.
- If generalization gap between our training and validation errors is small, a more complex model may be better.

Overfitting

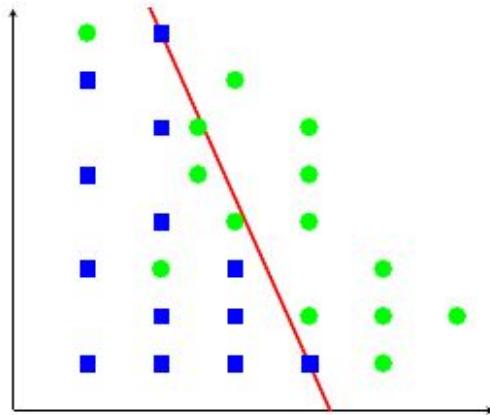
- The phenomenon of fitting the training data more closely than fit the underlying distribution is called **overfitting**.
- High Training accuracy and Low Validation accuracy
- Training error is significantly lower than the validation error.
- The techniques used to combat overfitting are called **regularization**.



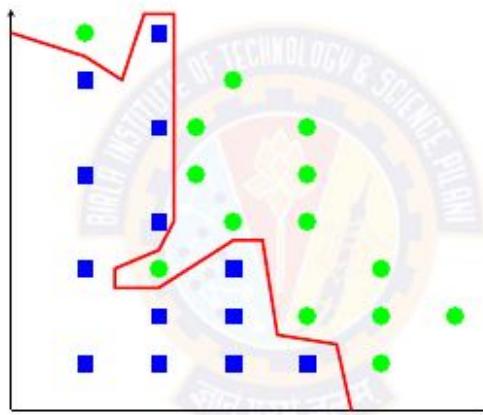
Training accuracy = 1.0
Validation accuracy = 0.43689320475152393
Testing accuracy = 0.5614035087719298



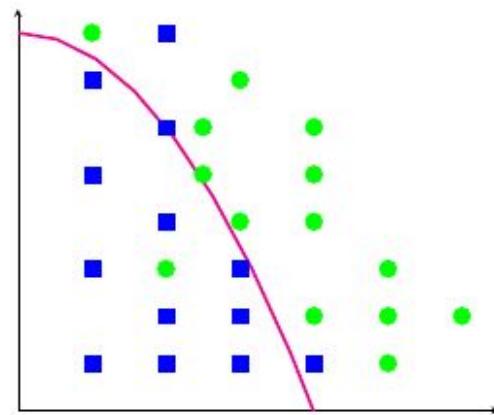
Underfitting or Overfitting?



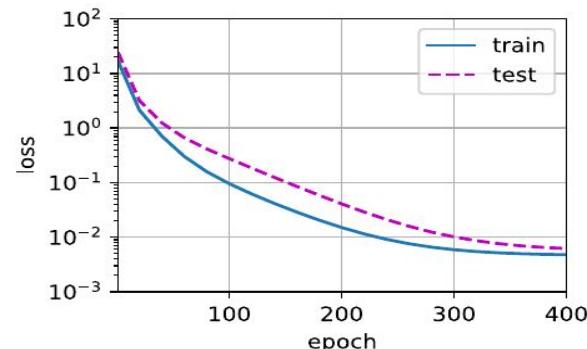
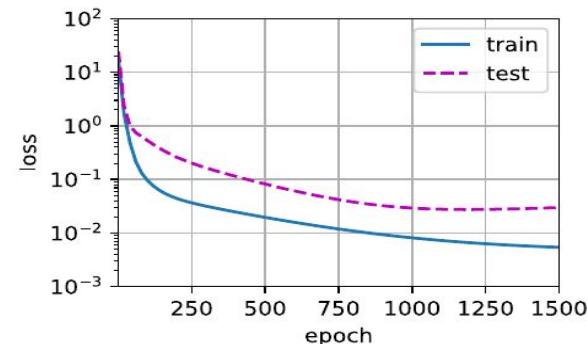
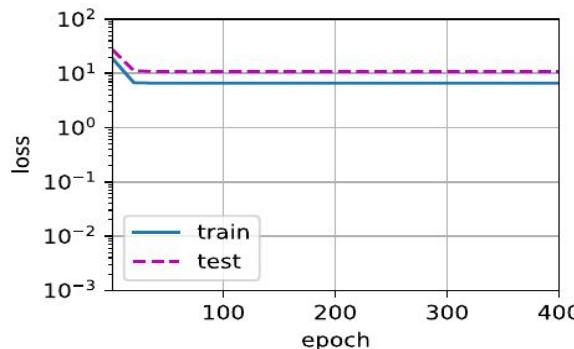
Simple Model
Underfitting



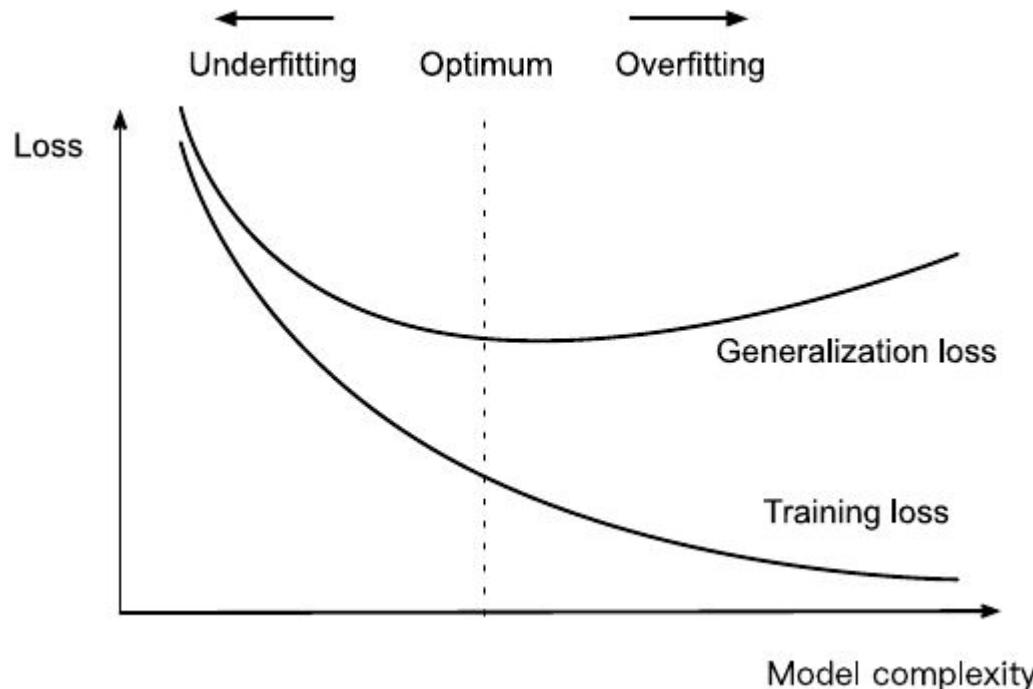
Complex Model
Overfitting



Just right model

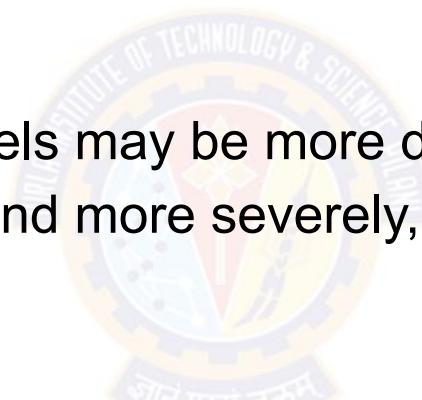


Polynomial degree and underfitting vs. overfitting



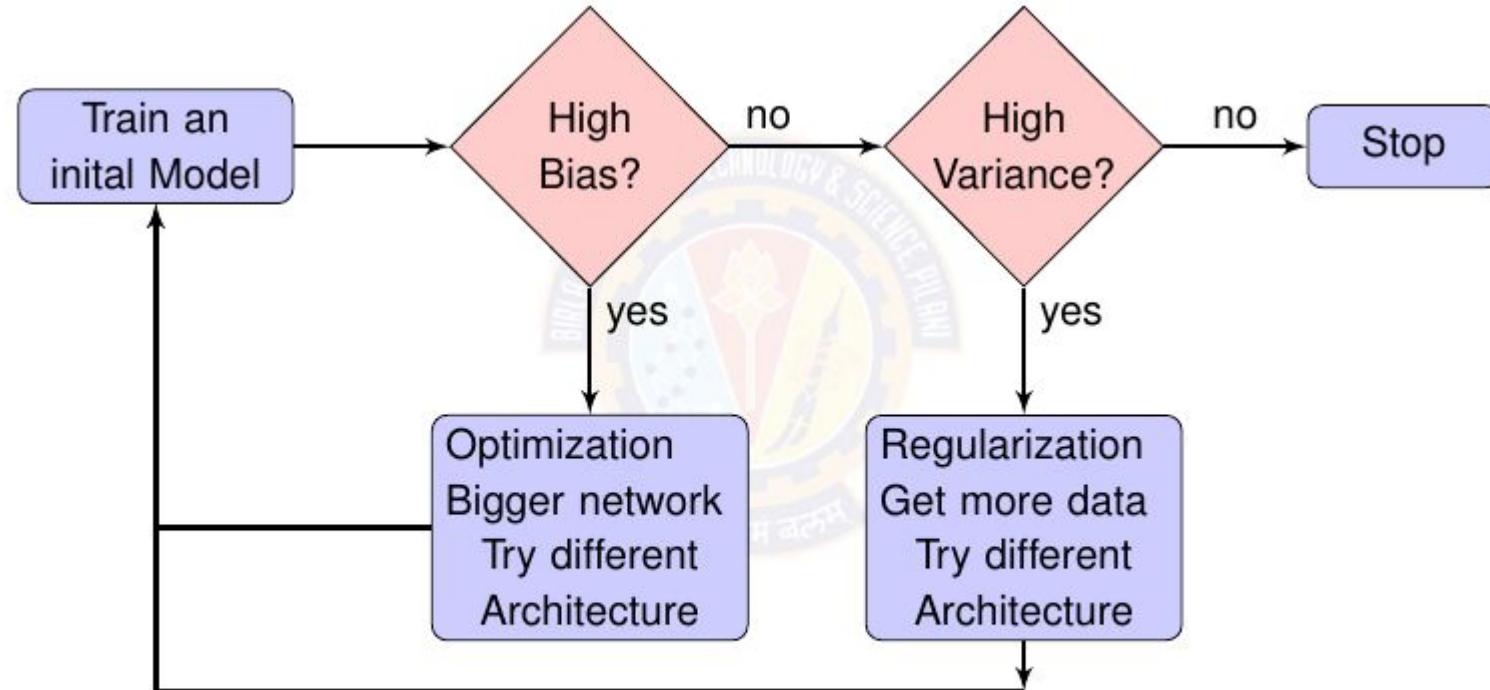
Model complexity and dataset size

- More data, fit a more complex model.
 - More data, the generalization error typically decreases.
-
- Less data, simpler models may be more difficult to beat.
 - Less data, more likely and more severely, models may over-fit.



The current success of deep learning owes to the current abundance of massive datasets due to Internet companies, cheap storage, connected devices, and the broad digitization of the economy.

Deep Learning Model Selection



Credit: Andrew Ng

Regularization



Regularization Techniques

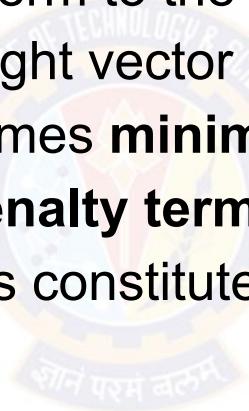
- Weight Decay (L2 regularization)
- Dropout
- Early Stopping



Weight Decay

L2 Regularization

- Measure the complexity of a linear function $f(x) = w^T x$ by some norm of its weight vector, e.g., $\|w\|^2$.
- Add the norm as a penalty term to the problem of minimizing the loss. This will ensure that the weight vector is small.
- The objective function becomes **minimizing the sum of the prediction loss and the penalty term.**
- L2-regularized linear models constitute the ridge regression algorithm.



L2 Regularization

- The trade off between standard loss and the additive penalty is given by regularization constant λ , a non-negative hyperparameter.

$$\text{Objective function} = \mathcal{L}(w, b) + \frac{\lambda}{2} \| w \|^2$$

$$\text{Weight updation} \quad w \leftarrow (1 - \eta\lambda)w - \eta g$$

- For $\lambda = 0$, we recover the original loss function.
- For $\lambda > 0$, we restrict the size of $\| w \|$.
- Smaller values of λ correspond to less constrained w , whereas larger values of λ constrain w more considerably.
- By squaring the L2 norm, we remove the square root, leaving the sum of squares of each component of the weight vector.
- The derivative of a quadratic function ensure that the 2 and 1/2 cancel out.

L1 Regularization

- L1-regularized linear regression is known as lasso regression.
- L1 penalties lead to models that concentrate weights on a small set of features by clearing the other weights to zero. This is called feature selection.

$$\text{Objective function} = \mathcal{L}(w, b) + \lambda \sum_i |w_i|$$

$$\text{Weight updation} \quad w \leftarrow (1 - \eta\lambda)w - \eta g$$

L2 Regularization

L1 Regularization

Sum of square of weights

Sum of absolute value of weights

Learn complex data patterns

Generates simple and interpretable models

Estimate mean of data

Estimate median of data

Not robust to outliers

Robust to outliers

Shrink coefficients equally

Shrink coefficients to zero

Non sparse solution

Sparse solution

One solution

Multiple solutions

No feature selection

Selects features

Useful for collinear features

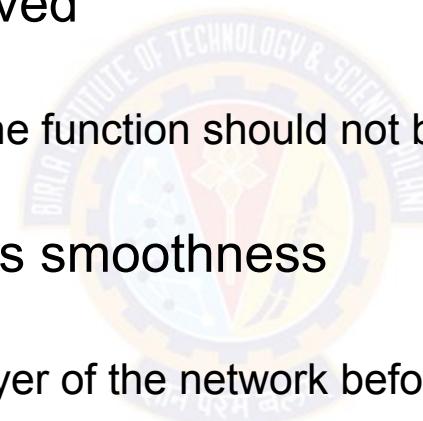
Useful for dimensionality reduction



Dropout

Smoothness

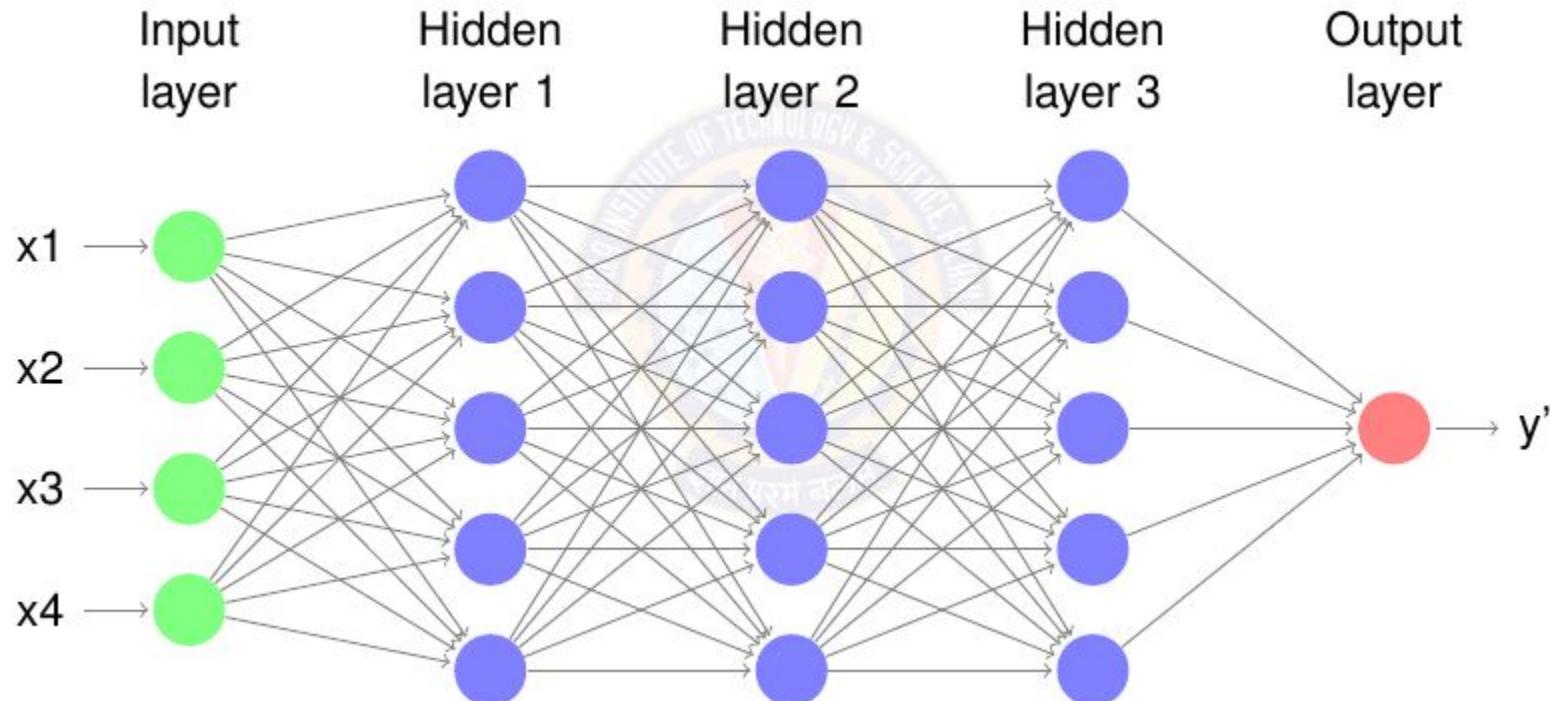
- Classical generalization theory suggests that to close the gap between train and test performance, aim for a simple model.
- Simplicity can be achieved
 - Using weight decay
 - Smoothness, i.e., that the function should not be sensitive to small changes to its inputs.
- Injecting noise enforces smoothness
 - training with input noise
 - inject noise into each layer of the network before calculating the subsequent layer during training.



Dropout

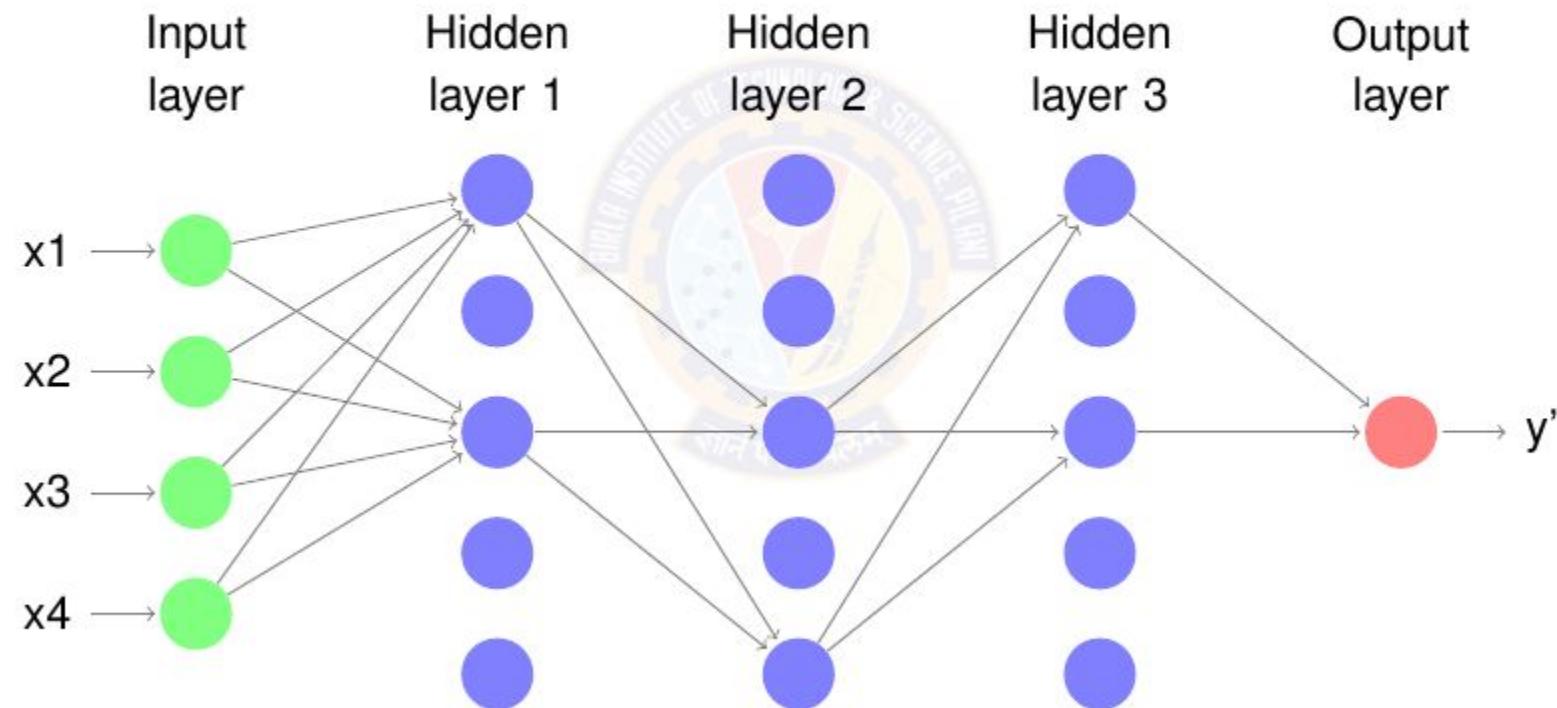
- **Dropout** involves injecting noise while computing each internal layer during forward propagation.
- It has become a **standard** technique for training neural networks.
- The method is called **dropout** because we literally drop out some neurons during training.
- Apply dropout to a hidden layer, zeroing out each hidden unit with probability p .
- The calculation of the outputs no longer depends on dropped out neurons and their respective gradient also vanishes when performing backpropagation. (in that iteration)
- Dropout is disabled at test time.

Without Dropout



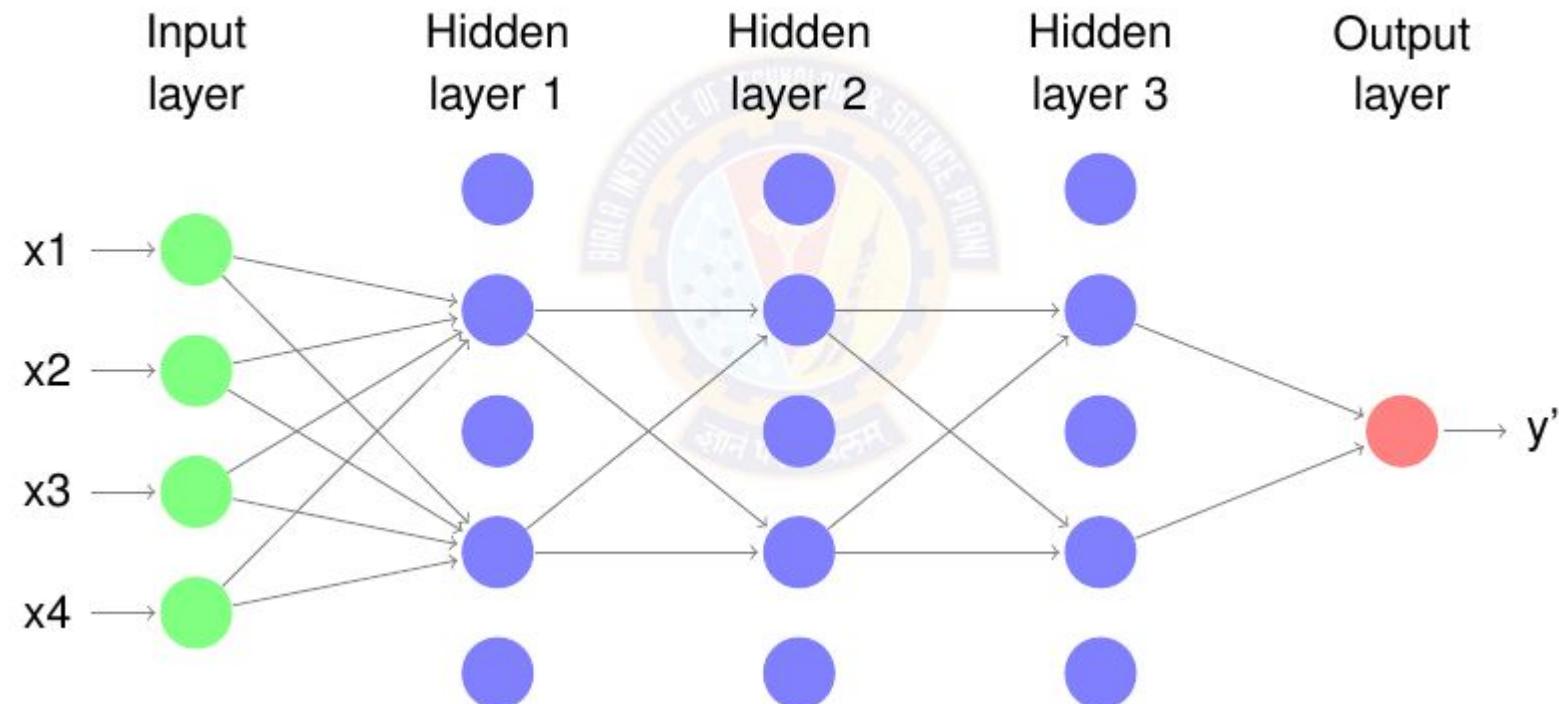
Dropout for first iteration

Keep probability = 0.5



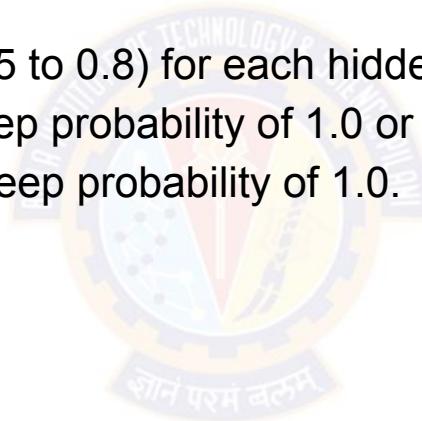
Dropout for second iteration

Keep probability = 0.5



Dropout

- Dropout gives a smaller neural network, giving the effect of regularization.
- In general,
 - Vary keep probability (0.5 to 0.8) for each hidden layer.
 - The input layer has a keep probability of 1.0 or 0.9.
 - The output layer has a keep probability of 1.0.
-

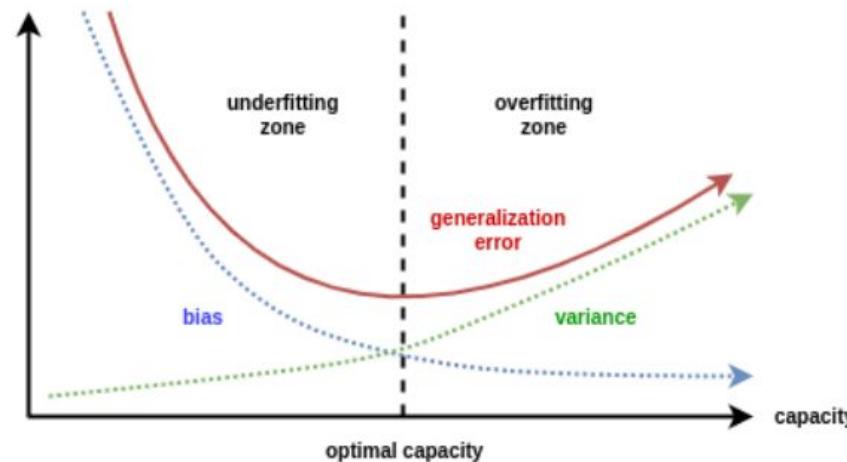


Early Stopping



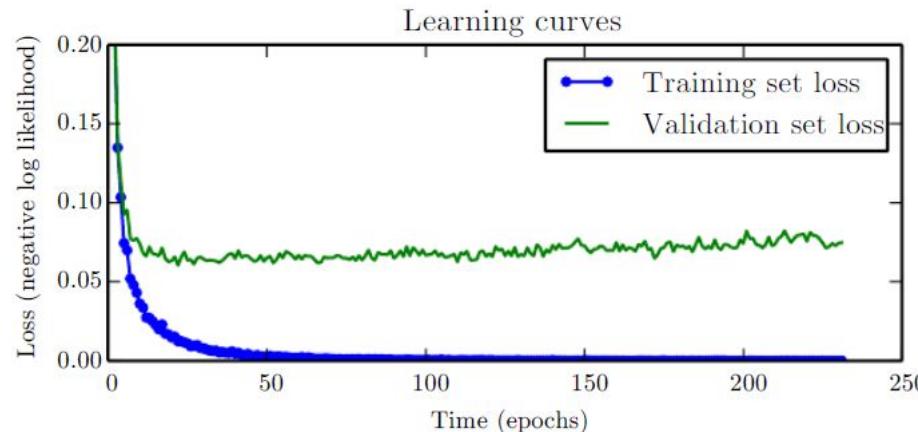
Before Early stopping

- When training large models, training error decreases steadily over time, but validation set error begins to rise again.
- Training objective decreases consistently over time.
- Validation set average loss begins to increase again, forming an asymmetric U-shaped curve.



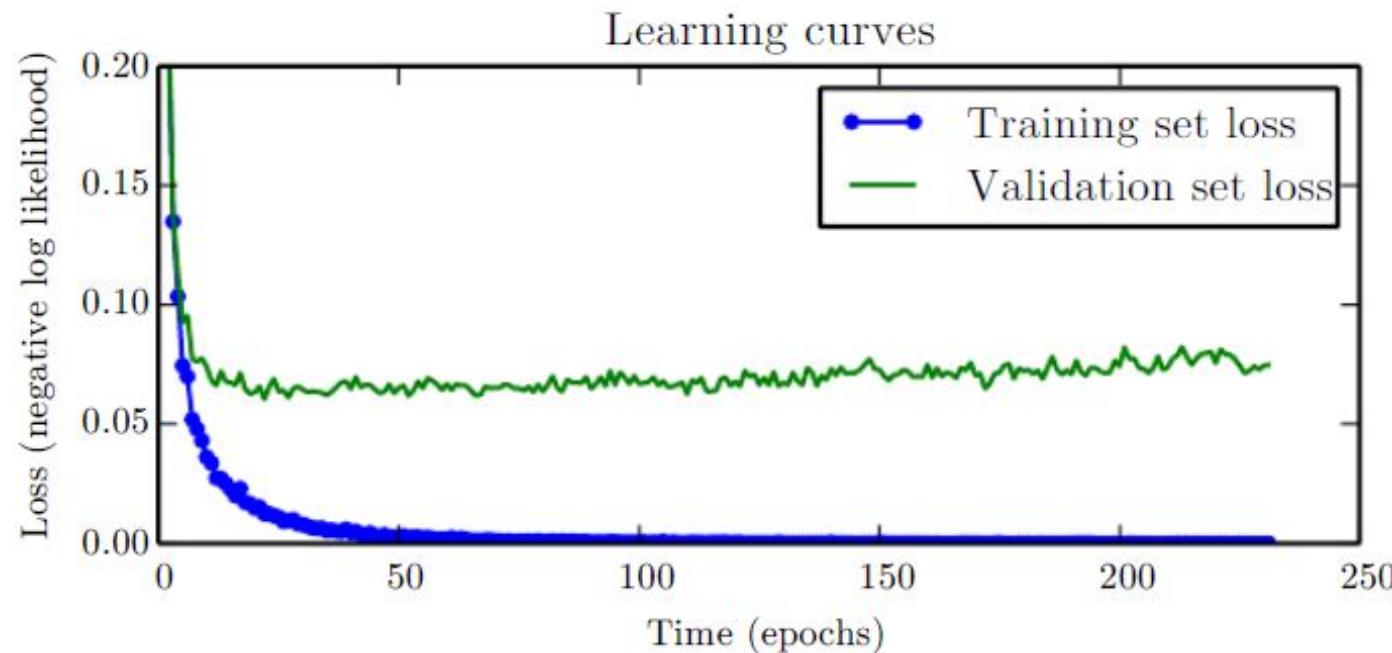
Early stopping

- No longer looking for a local minimum of validation error, while training.
- Train until the validation set error has not improved for some amount of time.
- Every time the error on the validation set improves, store a copy of the model parameters. When the training algorithm terminates, return these parameters.



Early stopping

- Effective and simple form of regularization.
- Trains simpler models



Early Stopping code

```
>>> callback = tf.keras.callbacks.EarlyStopping(monitor='loss', patience=3)
>>> # This callback will stop the training when there is no improvement in
>>> # the loss for three consecutive epochs.
>>> model = tf.keras.models.Sequential([tf.keras.layers.Dense(10)])
>>> model.compile(tf.keras.optimizers.SGD(), loss='mse')
>>> history = model.fit(np.arange(100).reshape(5, 20), np.zeros(5),
...                      epochs=10, batch_size=1, callbacks=[callback],
...                      verbose=0)
>>> len(history.history['loss']) # Only 4 epochs are run.
```

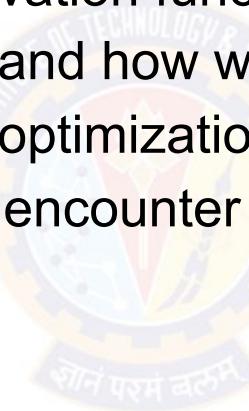
4

Numerical Stability and Initialization



Why Initialization is important?

- The choice of initialization is crucial for maintaining numerical stability.
- The choices of initialization can be tied up in interesting ways with the choice of the nonlinear activation function.
- Which function we choose and how we initialize parameters can determine how quickly our optimization algorithm converges.
- Poor choices can cause to encounter exploding or vanishing gradients while training.



Vanishing and Exploding Gradients

- Consider a deep network with L layers, input \mathbf{x} and output \mathbf{o} . With each layer I defined by a transformation f_I parameterized by weights $\mathbf{W}(I)$, whose hidden variable is $\mathbf{h}(I)$

$$\mathbf{h}^{(l)} = f_l(\mathbf{h}^{(l-1)}) \text{ and thus } \mathbf{o} = f_L \circ \dots \circ f_1(\mathbf{x}).$$

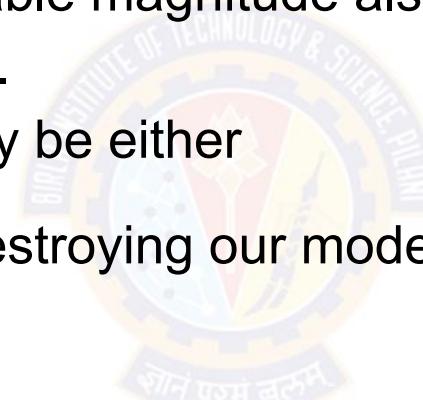
- If all the hidden variables and the input are vectors, then the gradient of \mathbf{o} with respect to any set of parameters $\mathbf{W}(I)$

$$\partial_{\mathbf{W}^{(l)}} \mathbf{o} = \underbrace{\partial_{\mathbf{h}^{(L-1)}} \mathbf{h}^{(L)}}_{\mathbf{M}^{(L)} \stackrel{\text{def}}{=} \dots} \cdot \dots \cdot \underbrace{\partial_{\mathbf{h}^{(l)}} \mathbf{h}^{(l+1)}}_{\mathbf{M}^{(l+1)} \stackrel{\text{def}}{=} \dots} \underbrace{\partial_{\mathbf{W}^{(l)}} \mathbf{h}^{(l)}}_{\mathbf{v}^{(l)} \stackrel{\text{def}}{=}}.$$

- Gradient is the product of $L - 1$ matrices $\mathbf{M}(L) \cdot \dots \cdot \mathbf{M}(l+1)$ and the gradient vector $\mathbf{v}(l)$.

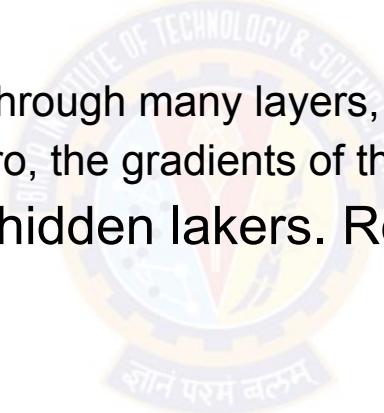
Vanishing and Exploding Gradients

- The matrices $M(I)$ may have a wide variety of eigenvalues. They might be small or large, and their product might be very large or very small.
- Gradients of unpredictable magnitude also threaten the stability of the optimization algorithms.
- Parameter updates may be either
 - (i) excessively large, destroying our model (the **exploding gradient** problem);
 - (ii) excessively small (the **vanishing gradient** problem), rendering learning impossible as parameters hardly move on each update.



Vanishing Gradients

- Activation function sigmoid σ can cause the vanishing gradient problem.
 - The sigmoid's gradient vanishes both when its inputs are large and when they are small.
 - When backpropagating through many layers, where the inputs to many of the sigmoids are close to zero, the gradients of the overall product may vanish.
- Solution: Use ReLU for hidden layers. ReLU is more stable.



Parameter Initialization

1. Default Initialization
 - Used a normal distribution to initialize the values of the parameters.
2. Xavier Initialization
 - samples weights from a Gaussian distribution with zero mean and variance

$$\sigma = \sqrt{\frac{2}{n_{\text{in}} + n_{\text{out}}}}.$$

- now-standard and practically beneficial

Batch Normalization



Why Batch Normalization?

1. Standardize the input features to each have a mean of zero and variance of one. This standardization puts the parameters a priori at a similar scale. Better optimization.
2. In a MLP or CNN, as we train, the variables in intermediate layers may take values with widely varying magnitudes: both along the layers from the input to the output, across units in the same layer, and over time due to our updates to the model parameters. This drift in the distribution of such variables could hamper the convergence of the network.
3. Deeper networks are complex and easily capable of overfitting. This means that regularization becomes more critical.

Batch Normalization

- Batch normalization is a popular and effective technique that consistently accelerates the convergence of deep networks.
- Batch normalization is applied to individual layers.
- It works as follows:
 - In each training iteration, first normalize the inputs (of batch normalization) by subtracting their mean and dividing by their standard deviation, where both are estimated based on the statistics of the current minibatch.
 - Next, apply a scale coefficient and a scale offset.
- Due to the normalization based on batch statistics that batch normalization derives its name.
- Batch normalization works best for moderate minibatches sizes in the 50 to 100 range.

Batch Normalization

- Denote by $\mathbf{x} \in \mathcal{B}$ an input to batch normalization (BN) that is from a minibatch \mathcal{B} , batch normalization transforms \mathbf{x} as

$$\text{BN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \hat{\mu}_{\mathcal{B}}}{\hat{\sigma}_{\mathcal{B}}} + \beta.$$

- $\hat{\mu}_{\mathcal{B}}$ is the sample mean and $\hat{\sigma}_{\mathcal{B}}$ is the sample standard deviation of the minibatch \mathcal{B} .
- After applying standardization, the resulting minibatch has zero mean and unit variance.

Batch Normalization

- Denote by $\mathbf{x} \in \mathcal{B}$ an input to batch normalization (BN) that is from a minibatch \mathcal{B} , batch normalization transforms \mathbf{x} as

$$\text{BN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \hat{\mu}_{\mathcal{B}}}{\hat{\sigma}_{\mathcal{B}}} + \beta.$$

- $\hat{\mu}_{\mathcal{B}}$ is the sample mean and $\hat{\sigma}_{\mathcal{B}}$ is the sample standard deviation of the minibatch \mathcal{B} .
- After applying standardization, the resulting minibatch has zero mean and unit variance.

Batch Normalization

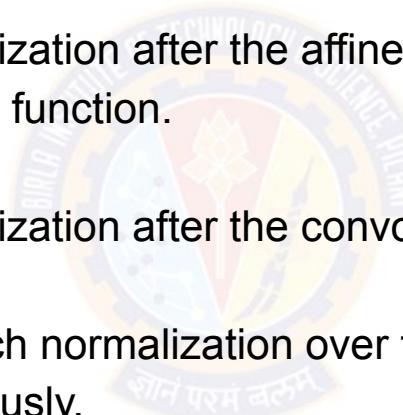
- Elementwise scale parameter γ and shift parameter β that have the same shape as x . γ and β are parameters are learned jointly with the other model parameters.
- Batch normalization actively centers and rescales the inputs to each layer back to a given mean and size.
- Calculate $\hat{\mu}_B$ and $\hat{\sigma}_B^2$

$$\hat{\mu}_B = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x} \in \mathcal{B}} \mathbf{x},$$

$$\hat{\sigma}_B^2 = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x} \in \mathcal{B}} (\mathbf{x} - \hat{\mu}_B)^2 + \epsilon.$$

Batch Normalization Layers

- Batch normalization implementations for fully-connected layers and convolutional layers are slightly different.
 - Fully-Connected Layers
 - Insert batch normalization after the affine transformation and before the nonlinear activation function.
 - Convolutional Layers
 - Apply batch normalization after the convolution and before the nonlinear activation function.
 - Carry out each batch normalization over the $m \cdot p \cdot q$ elements per output channel simultaneously.
- It operates on a full minibatch at a time.

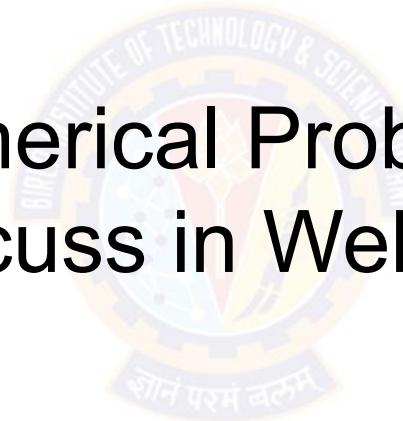


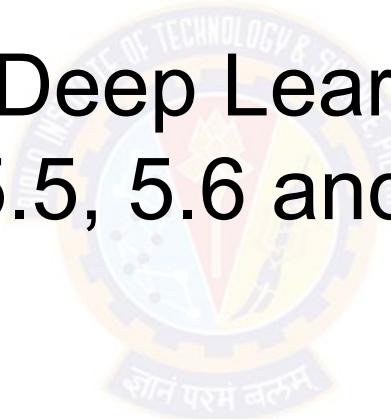
Batch Normalization During Prediction

- After training, use the entire dataset to compute stable estimates of the variable statistics and then fix them at prediction time.



Numerical Problems (discuss in Webinar)





Ref TB Dive into Deep Learning

- Sections 5.4, 5.5, 5.6 and 8,5 (online version)

Next Session: CNN





BITS Pilani
Pilani | Dubai | Goa | Hyderabad

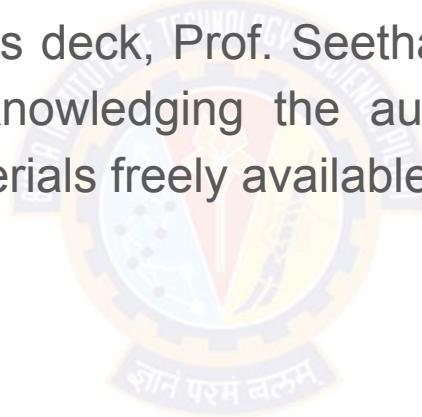
Deep Learning

DSE Review Session 8

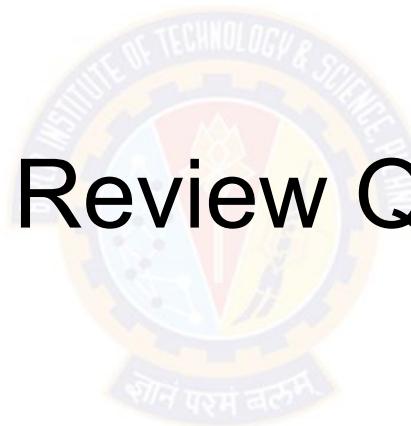
Seetha Parameswaran

BITS Pilani

The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Midsem Review Questions



Question

For a newsfeed classification system with 46 topics, the following two models are tried. 10,000 unique words are used as input in both models. 20% of the training data is used as the validation set, `partial_x_train`.

Network Model 1

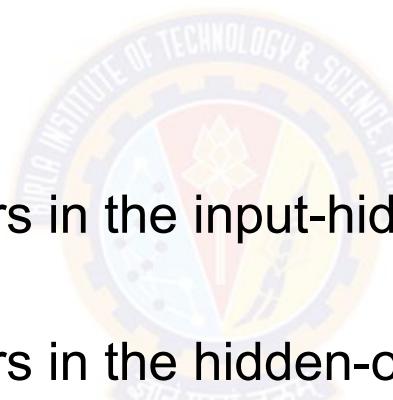
```
model = models.Sequential()  
  
model.add(layers.Dense(64, activation = 'relu',  
input_shape =(10000,)))  
  
model.add(layers.Dense(64, activation = 'relu'))  
  
model.add(layers.Dense(46, activation = 'softmax'))  
  
model.compile(optimizer='rmsprop', loss='categorica  
l_crossentropy ', metrics =['accuracy'])  
  
model.fit(partial_x_train, partial_y_train,  
epochs=9, batch_size=512, validation_data =(x_val,  
y_val))
```

Network Model 2

```
model = models.Sequential()  
  
model.add(layers.Dense(64, activation = 'relu',  
input_shape =(10000,)))  
  
model.add(layers.Dense(4, activation = 'relu'))  
  
model.add(layers.Dense(46, activation  
='softmax'))  
  
model.compile(optimizer='rmsprop', loss='categori  
cal_crossentropy ', metrics =['accuracy'])  
  
model.fit(partial_x_train, partial_y_train,  
epochs=20, batch_size=128,  
validation_data=(x_val, y_val))
```

Solution

What is the number of trainable parameters in the input-hidden layer and hidden-output layer in Model 2? Which of the networks will give lower validation error and why?

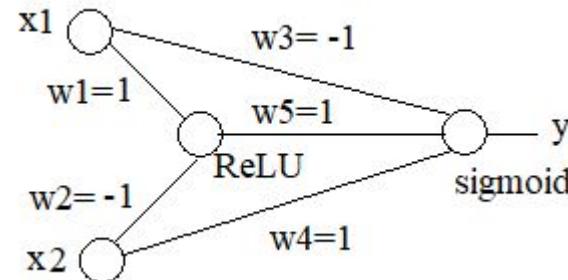


- # of trainable parameters in the input-hidden layer:
 $10000*64+64=640064$
- # of trainable parameters in the hidden-output layer: $4*46+46=230$
- Model 2 will give less validation error since hidden layer 2 is too small leading to information loss and reduced generalization ability.

Question

Hidden node and output node use, respectively, ReLU and sigmoid activation functions. Bias values at hidden and output nodes are zero. Weights for the current iteration are given in the above figure. Target output d is specified as 0. Learning rate is 0.3.

- A. Calculate the actual output y for the current iteration with input $(x_1, x_2) = (1, 1)$.
- B. Calculate the binary cross-entropy error for the current iteration.
- C. Assuming L1 regularization constant = 0.2, calculate the new w_3 for the next iteration.
- D. Assuming L2 regularization constant = 0.2, calculate the w_1 for the next iteration.
- E. Assuming both L1 (with regularization constant=0.2) and L2 (regularization constant=0.2) are applied, calculate the value of w_5 in next iteration.



Question

A. $y=1/(1+1)=1/2$

B. $L= - (1-\log(1/2))=-1-\log(2)$

C. Let z be the total input to the output node.

$$y = 1/(1+\exp(-z))$$

$$w3(t+1) = w3(t) - 0.3 * dL/dw3 - 0.2 * 0.3 * \text{sign}(w3)$$

$$= -0.94 - 0.3 * dL/dz * dz/dw3$$

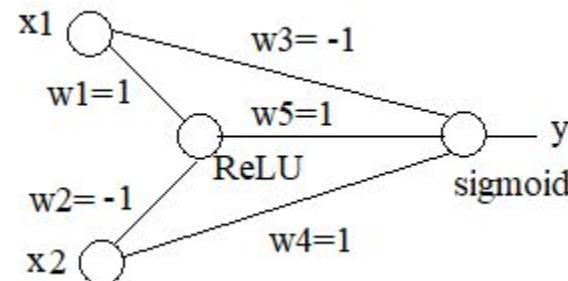
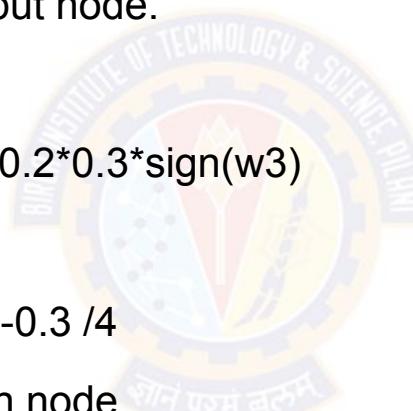
$$= -0.94 - 0.3 * x1 * 1/2 * 1/2 = -0.94 - 0.3 / 4$$

D. Let $h1$ be the output of the hidden node.

$$w1(t+1) = w1(t) - 0.3 * dL/dh1 * dh1/dw1 - 0.3 * 0.2 * w1$$

$$= 0.94, \text{ since } dh1/dw1 = \text{ReLU}'(0) = 0$$

E. $w5(t+1) = w5(t) - 0.3 * dL/dz * dz/dw5 - 0.3 * 0.2 * w5 - 0.3 * 0.2 * \text{sign}(w5) = 0.88$



Question

- A. What is a saddle point? What is the advantage/disadvantage of using Stochastic Gradient Descent in dealing with saddle points?
- B. What are five strategies to prevent overfitting in deep networks, when used for classification, say of images?
- C. What is the difference between kernel regularizers and activity regularizers?
- D.



Answer

- A. In a multi-dimensional surface, points where all the partial first derivatives are 0 but some of the partial double derivatives are +ve and some are –ve are known as saddle points. Or in other words, Hessian is –ve definitive. At Saddle points, the surface has one or more minima along some directions and maxima along other directions.
- B. L1/L2 regularization, dropout, data augmentation, early stopping, adding noise to input/target output.
- C. kernel_regularizer: Regularizer to apply a penalty on the layer's kernel. activity_regularizer: Regularizer to apply a penalty on the layer's output.

Question

Consider the following DNN for image classification for a dataset that consists of RGB images of size 32x32.

```
model=models.Sequential()
# Layer1
model.add(layers.Dense( 50,activation='relu',input_shape=**A**))
#Layer2
model.add(layers.Dense( 40,activation='relu'))
#Layer3
model.add(layers.Dense( 30,activation='relu'))
#Layer4
model.add(layers.Dense(**B**,activation=**C**))
model.compile(optimizer='sgd',loss=**D**,metrics=[ 'accuracy'])
```

Question

- A. What is the input shape **A** in Layer 1?
- B. What will be the value of **B**, activation function **C** and loss **D** if the total number of classes in the dataset is (a) 2 (b) 10
- C. What will be the total number of parameters in Layer 1, Layer 2 and Layer 3?
- D. If a dropout layer of value 0.5 is added after Layer 2, what will be the change in the number of parameters?

Answer

A. **A** is (32*32*3,) or (3072,)

B. (**B**, **C**, **D**)

(a) 2 is 1, sigmoid, binary_crossentropy

(b) 10 is 10, softmax, categorical_crossentropy

C. What will be the total number of parameters in Layer 1, Layer 2 and Layer 3?

- Layer 1- $3072 \times 50 + 50 = 153,650$
- Layer 2- $50 \times 40 + 40 = 2040$
- Layer 3- $40 \times 30 + 30 = 1230$
- Total = $153650 + 2040 + 1230 = 156,920$

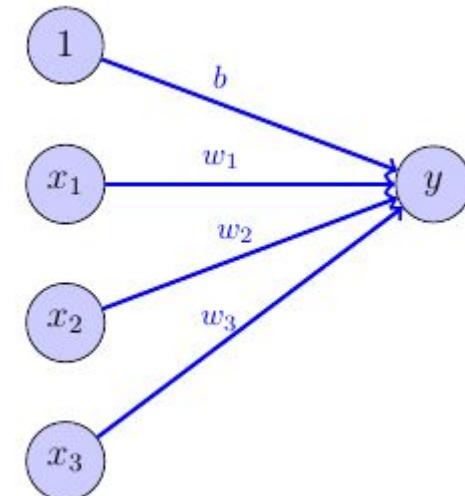
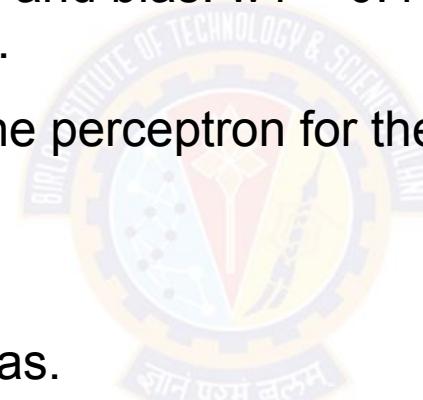
D. No change in the number of parameters if dropout is added.

Question

A perceptron structure and the training data are given below.

Assume the following weights and bias. $w_1 = 0.41$, $w_2 = 0.23$, $w_3 = 0.5$ and $b = 0.01$.

- (a) Compute the output y of the perceptron for the first training example.
- (b) Compute the error.
- (c) Update the weights and bias.
- (d) Using the updated parameters, compute the output for the second training example.



x_1	x_2	x_3	y
0.8	0.2	0.4	A
0.3	0.1	0.9	B
0.2	0.4	0.6	A
0.1	0.7	0.3	B

Solution

- (a) Output of the perceptron for the first training example.

Let $A = +1 \quad B = -1$

$$\hat{y}^{(1)} = \text{sign}(w_1x_1^{(1)} + w_2x_2^{(1)} + w_3x_3^{(1)} + b)$$

- (b) Error

$$\text{Error} = \text{difference}[\hat{y}^{(1)}, d^{(1)}]$$

- (c) Update the weights and bias

If error = 0, no need to update the parameters. Otherwise update weights.

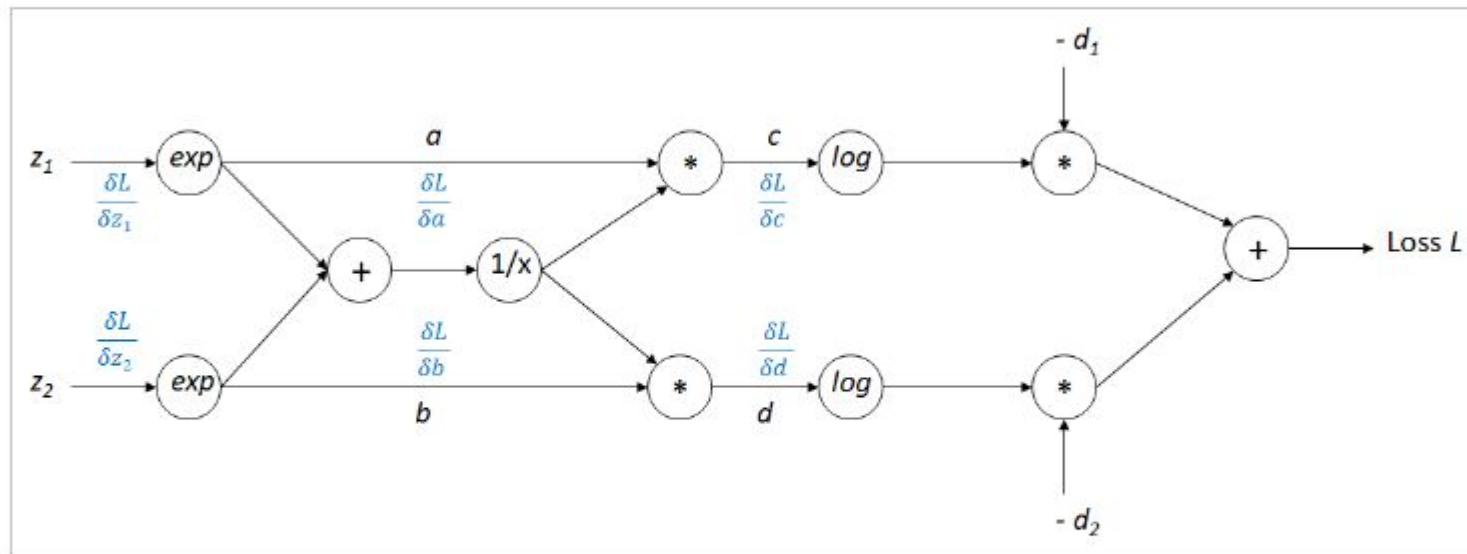
- (d) Output for the second training example.

$$X^{(2)} = \begin{bmatrix} 0.3 & 0.1 & 0.9 \end{bmatrix} \quad d^{(2)} = -1$$

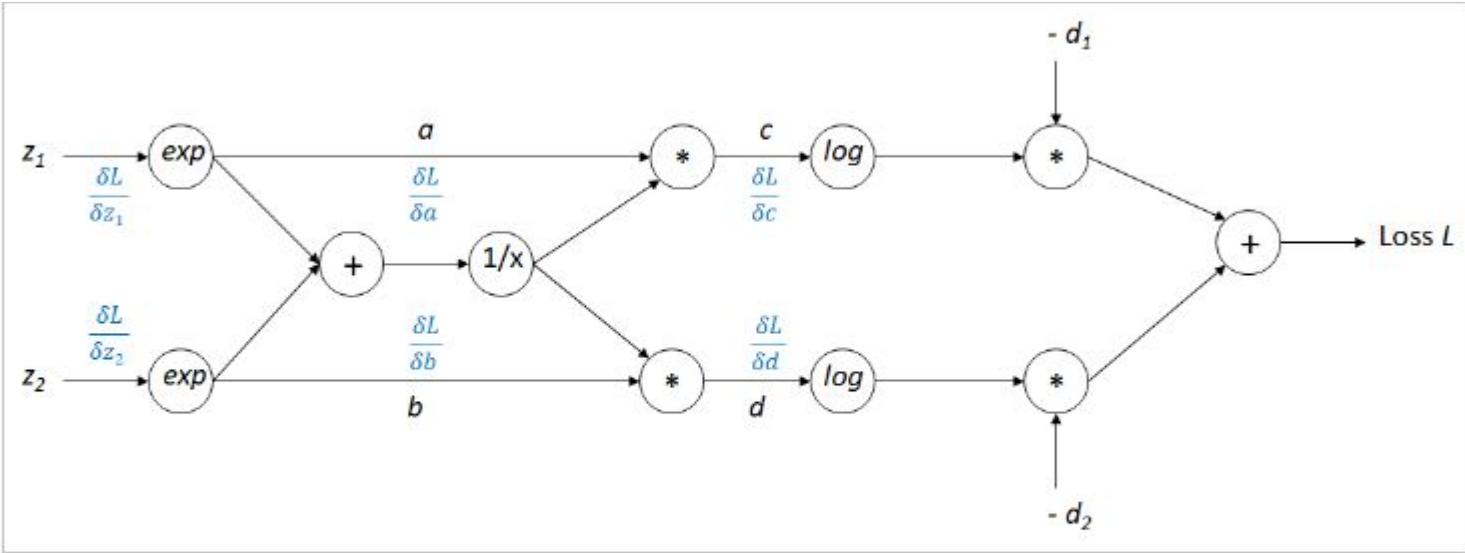
$$\hat{y}^{(2)} = \text{sign}(W^T X)$$

Question

Use the following computation graph to calculate the loss L , $\frac{\delta L}{\delta c}, \frac{\delta L}{\delta d}, \frac{\delta L}{\delta a}, \frac{\delta L}{\delta b}, \frac{\delta L}{\delta z_1}, \frac{\delta L}{\delta z_2}$.



Solution



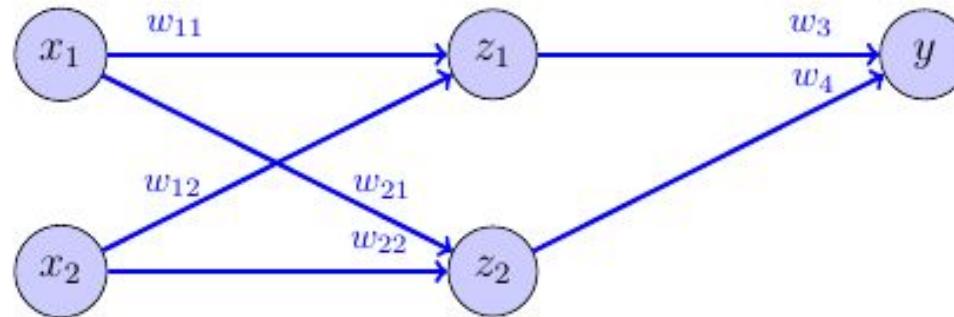
$$L = -\log(1/(1+e)) = \log(1+e),$$

$$a=1, \quad b=e, \quad c=1/(1+e), \quad d=e/(1+e)$$

$$\frac{\delta L}{\delta c} = -(1+e) \quad \frac{\delta L}{\delta d} = 0 \quad \frac{\delta L}{\delta a} = \frac{-e}{1+e} \quad \frac{\delta L}{\delta b} = \frac{1}{1+e} \quad \frac{\delta L}{\delta z_1} = \frac{-e}{1+e} \quad \frac{\delta L}{\delta z_2} = \frac{e}{1+e}$$

Do it yourself :)

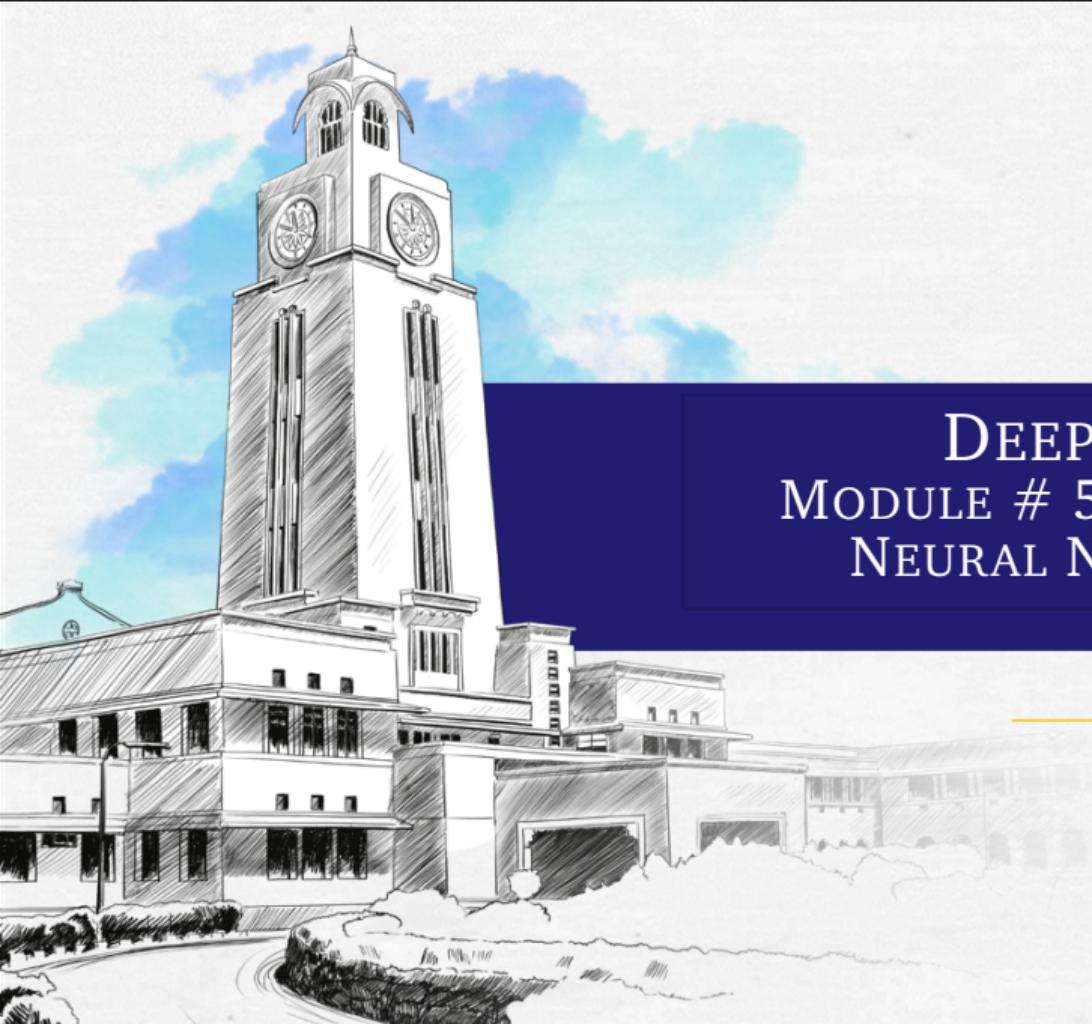
Consider the following network structure and the initial weights and bias a given. $w_{11} = 0.3$, $w_{12} = 0.4$, $w_{21} = -0.3$, $w_{22} = 0.05$, $w_3 = -0.5$, $w_4 = 0.2$ and $\text{bias} = 0$. Given that $\langle x_1, x_2, y \rangle = \langle 1, 0, 0.5 \rangle$.



- Compute the forward propagation and generate the output. Use Sigmoid activation function.
- Compute the RMSE loss.
- Let the given weights be at time $(t-1)$. Compute the weights at time t using SGD.
- Compute the weights at $(t + 1)$ using Momentum. Assume $\alpha = 1.1$ and $\beta = 0.8$.

All the best :)



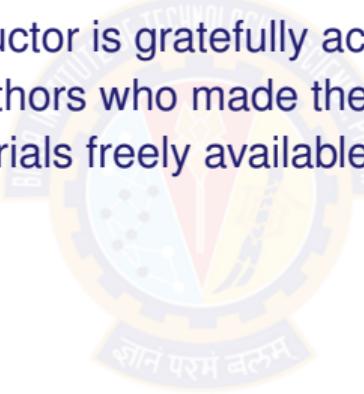


BITS Pilani
Pilani | Dubai | Goa | Hyderabad

DEEP LEARNING MODULE # 5 : CONVOLUTIONAL NEURAL NETWORK [CNN]

Seetha Parameswaran
Asst Prof, BITS Pilani

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.

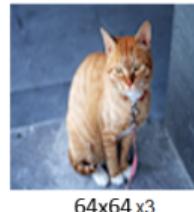


IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



IMAGE CLASSIFICATION



→ Cat? (0/1)



→ Cat? (0/1)



Classify into 3 classes

FACE DETECTION



OBJECT RECOGNITION

"SIFT" & Object Recognition, David Lowe, 1999

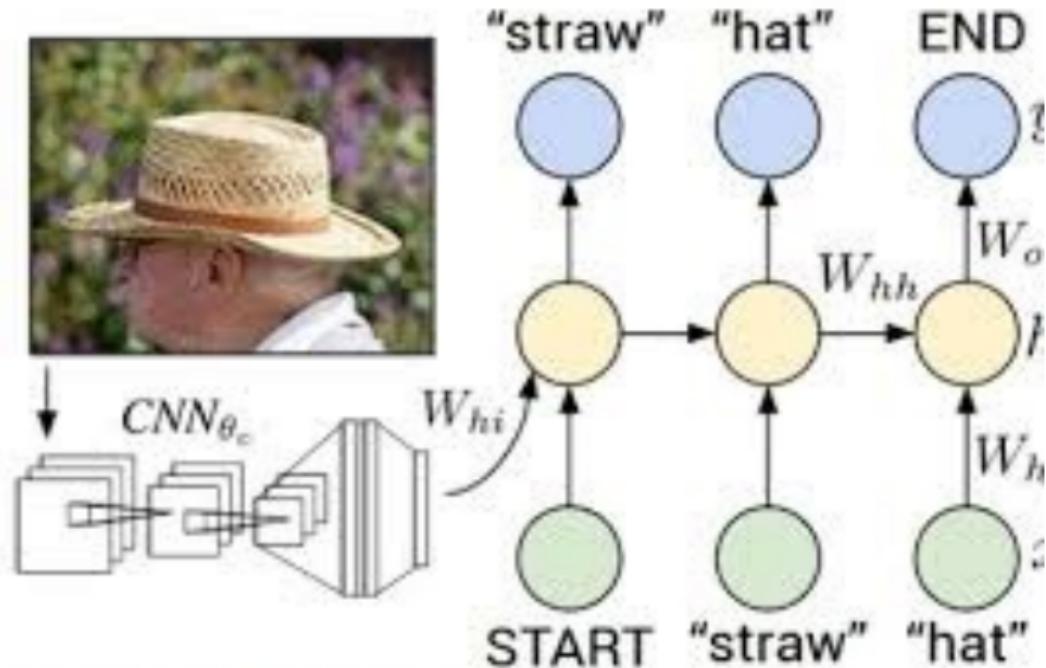


[Image](#) is public domain



[Image](#) is CC BY-SA 2.0

IMAGE CAPTIONING



Andrej Karpathy, Fei-Fei Li

IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



DIGITAL GRayscale IMAGE

Pixel intensities = $I(x, y)$



0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49	
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	21	252	255	248	144	6	0	
0	13	112	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	4	91	255	255	255	248	252	255	244	255	182	10	0	4	
0	22	206	252	246	251	241	104	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	2	62	255	250	125	3	0	
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	249	255	240	255	129	0	5	0	
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1	
0	0	5	5	0	0	0	0	0	14	1	0	6	6	0	0	

0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49	
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	21	252	255	248	144	6	0	
0	13	112	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0	
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4	
0	22	206	252	246	251	241	104	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	2	62	255	250	125	3	0	
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	249	255	240	255	129	0	5	0	
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1	
0	0	5	5	0	0	0	0	0	14	1	0	6	6	0	0	

DIGITAL COLOR IMAGE

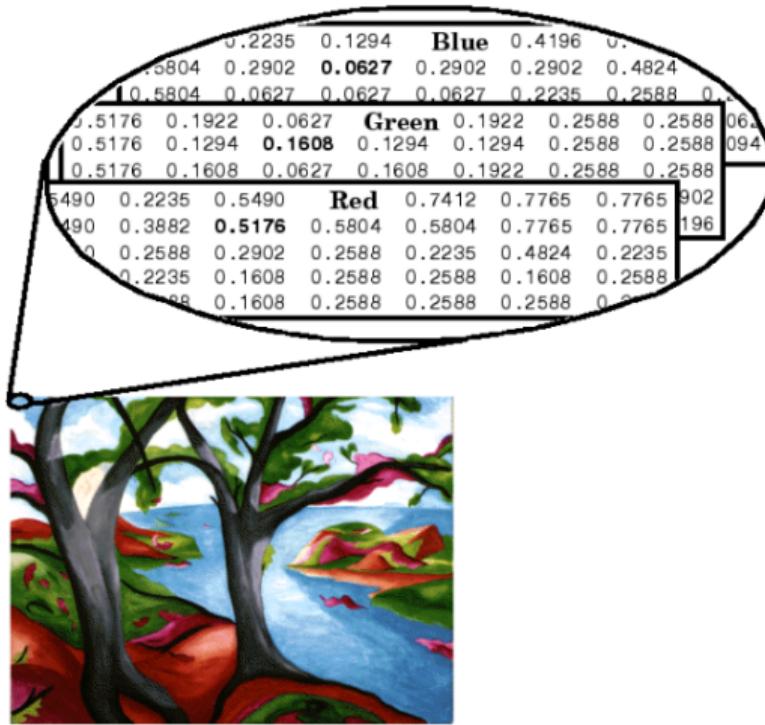


Colour Image

This image is composed of many colors and almost all colors can be generated from the three primary colors- **Red, Green, and Blue**. We can say that each colored image is composed of these three colors or 3 channels- Red, Green, and Blue-



DIGITAL COLOR IMAGE



IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



IMAGE CLASSIFICATION - CHALLENGES

Image Classification: A core task in Computer Vision



This image by Nikita is
licensed under CC-BY 2.0

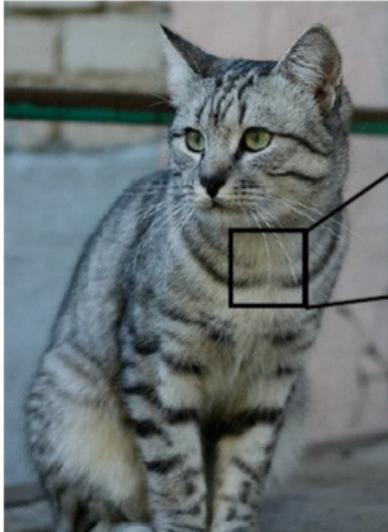
(assume given set of discrete labels)
{dog, cat, truck, plane, ...}



cat

IMAGE CLASSIFICATION - CHALLENGES

The Problem: Semantic Gap



This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)

[1] 185 112 106 111 104 99 106 99 96 103 112 119 104 97 93 87]
[1] 91 98 192 106 104 79 98 183 99 105 123 136 118 105 94 85]
[1] 76 85 98 105 128 105 87 96 95 99 115 112 106 103 99 85]
[1] 99 81 81 93 128 131 127 189 95 98 102 99 96 92 181 94]
[1] 106 91 61 64 69 91 88 85 181 187 109 98 75 84 96 95]
[1] 114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[1] 133 137 147 103 65 81 88 65 52 54 74 84 102 93 85 82]
[1] 128 137 144 148 109 95 86 79 62 65 63 63 68 73 88 181]
[1] 125 133 148 137 119 121 117 94 65 79 86 65 54 64 72 98]
[1] 127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[1] 115 114 109 123 158 148 131 118 113 109 106 92 74 65 72 78]
[1] 89 93 98 97 108 147 131 118 113 114 113 109 106 95 77 88]
[1] 63 77 86 81 77 79 102 123 117 115 117 125 125 138 115 87]
[1] 62 65 82 89 78 71 88 101 124 126 119 101 107 114 131 119]
[1] 63 65 75 88 89 71 62 81 120 138 135 105 81 90 118 118]
[1] 87 85 71 87 106 95 69 45 76 130 126 107 92 94 185 112]
[1] 118 97 82 86 117 122 116 66 41 51 95 93 89 95 182 187]
[1] 164 146 112 88 82 128 124 184 76 48 45 66 80 101 182 189]
[1] 157 179 157 128 93 86 114 132 112 97 69 55 78 82 99 94]
[1] 139 128 134 161 139 108 105 118 121 134 114 87 65 53 69 86]
[1] 128 112 96 117 158 144 128 115 184 197 102 93 87 81 72 79]
[1] 123 187 96 86 83 112 153 149 122 189 194 75 89 107 112 99]
[1] 122 121 192 98 82 85 94 117 145 148 153 102 58 78 92 147]
[1] 122 164 148 143 71 56 78 83 93 103 118 139 102 61 69 84]

What the computer sees

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)

IMAGE CLASSIFICATION - CHALLENGES

Challenges: Viewpoint variation



All pixels change when
the camera moves!

This image by Nikita is
licensed under [CC-BY 2.0](#)

IMAGE CLASSIFICATION - CHALLENGES

Challenges: Illumination



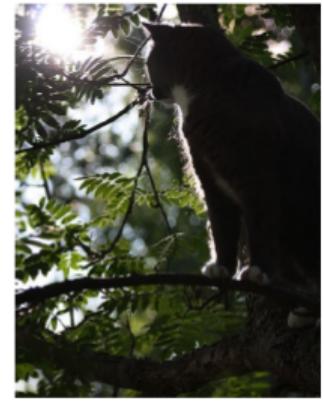
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



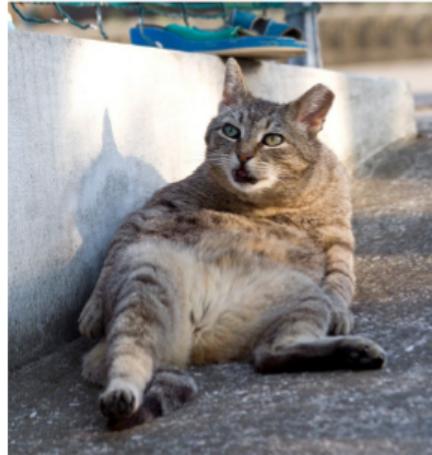
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

IMAGE CLASSIFICATION - CHALLENGES

Challenges: Deformation



This image by Umberto Salvagnin
is licensed under CC-BY 2.0



This image by Umberto Salvagnin
is licensed under CC-BY 2.0



This image by sare bear is
licensed under CC-BY 2.0



This image by Tom Thai is
licensed under CC-BY 2.0

IMAGE CLASSIFICATION - CHALLENGES

Challenges: Occlusion



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)



[This image by jonsson is licensed under CC-BY 2.0](#)

IMAGE CLASSIFICATION - CHALLENGES

Challenges: Background Clutter



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

IMAGE CLASSIFICATION - CHALLENGES

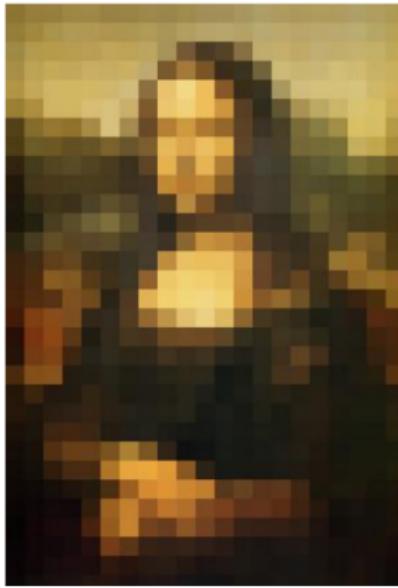
Challenges: Intraclass variation



This image is CC0 1.0 public domain

IMAGE CLASSIFICATION - CHALLENGES

Pixelated Images



Larger Images



- A 50x50 Image requires 2500 input dimensions.
- A HD image 1920x1080 is about 2 million inputs. (2 Mega Pixels)

IMAGE INVARIANCE - CHALLENGES

Object can appear anywhere in the image.

Image A

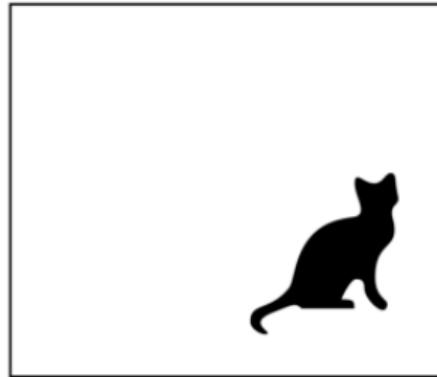
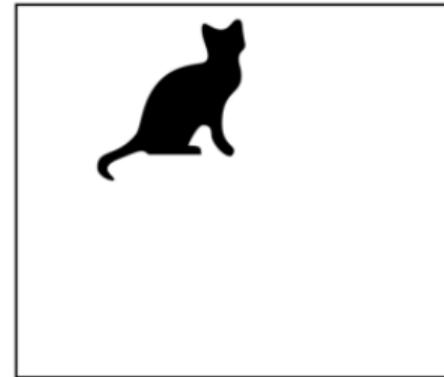


Image B

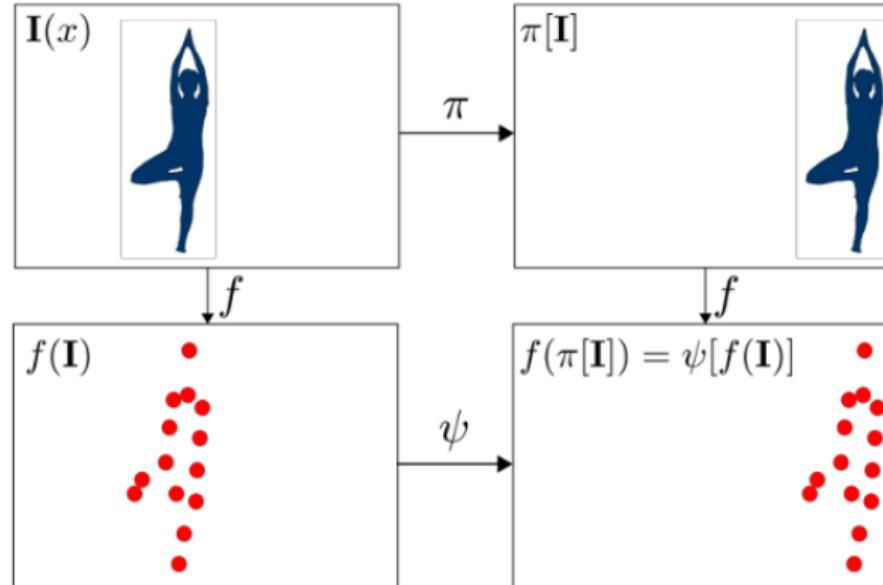


Model predicts “cat” for both images!

Source: <https://penkovsky.com/neural-networks/day5/>

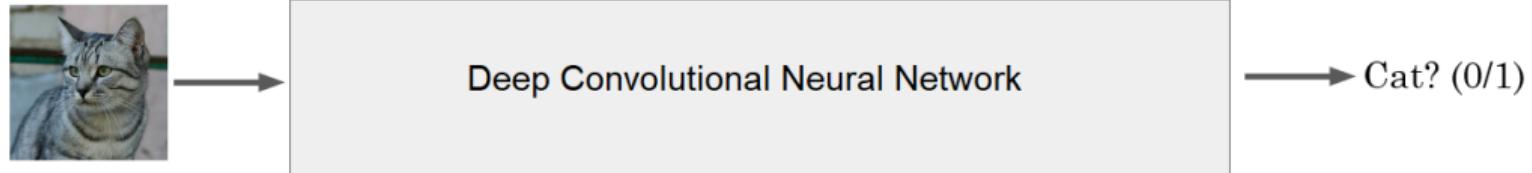
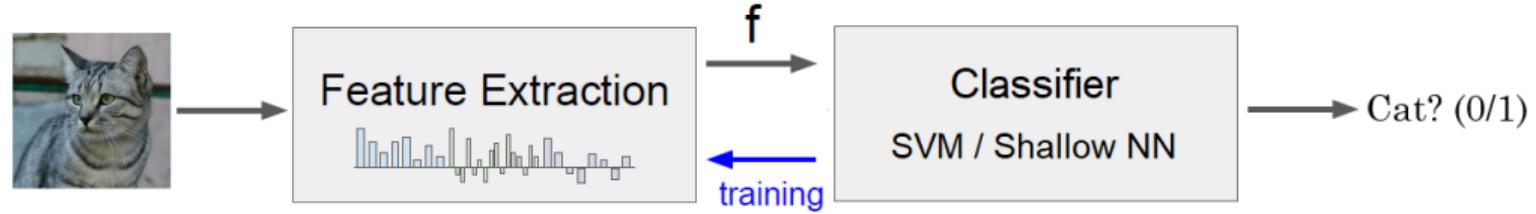
IMAGE EQUIVARIANCE - CHALLENGES

Changes to the input predictably (equivalently) should not affect the output.



Source: <https://towardsdatascience.com/translational-invariance-vs-translational-equivariance-f9fbc8fcfa63a>

IMAGE CLASSIFICATION - SOLUTION



IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



CONVOLUTIONAL NEURAL NETWORKS (CNN)

- Convolutional networks are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers.
- Special kind of neural network for processing data that has a known, grid-like topology.
- Eg: Time-series data – 1D grid taking samples at regular time intervals.
- Eg: Image data – 2D grid of pixels.
- Network employs a mathematical operation called **convolution**.

EXAMPLE OF CNN – LENET-5

- Handwritten Character recognition
- 10-categories classification problem

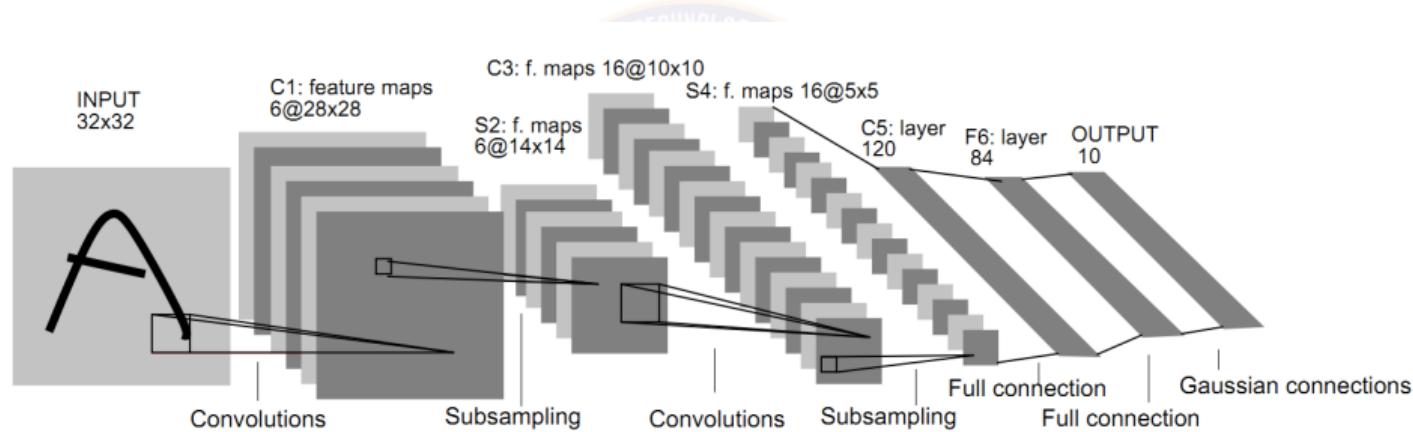


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

WHY CNN FOR IMAGES

- Some patterns are much smaller than the whole image.



WHY CNN FOR IMAGES

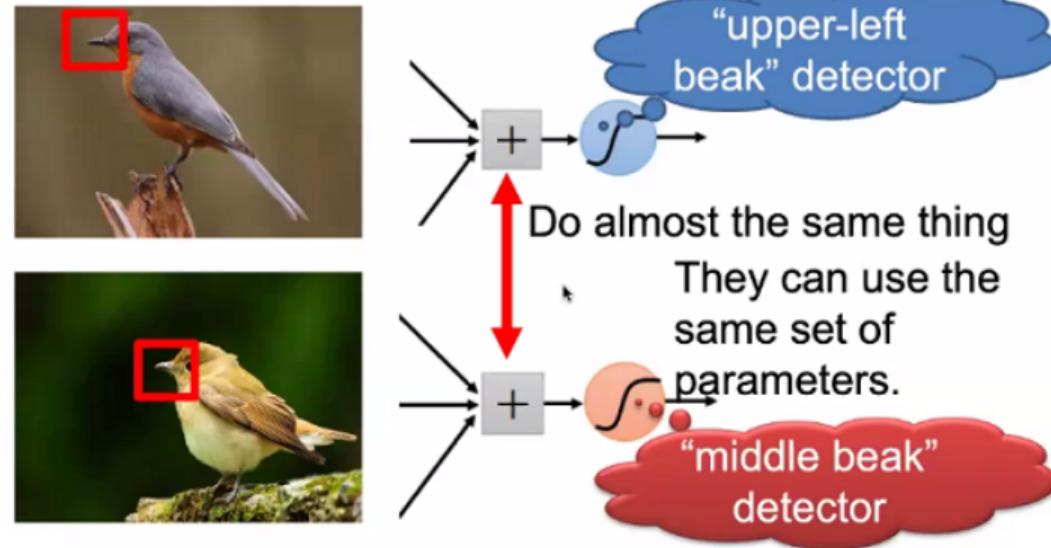
- Some patterns are much smaller than the whole image.

A neuron does not have to see the whole image to discover the pattern.



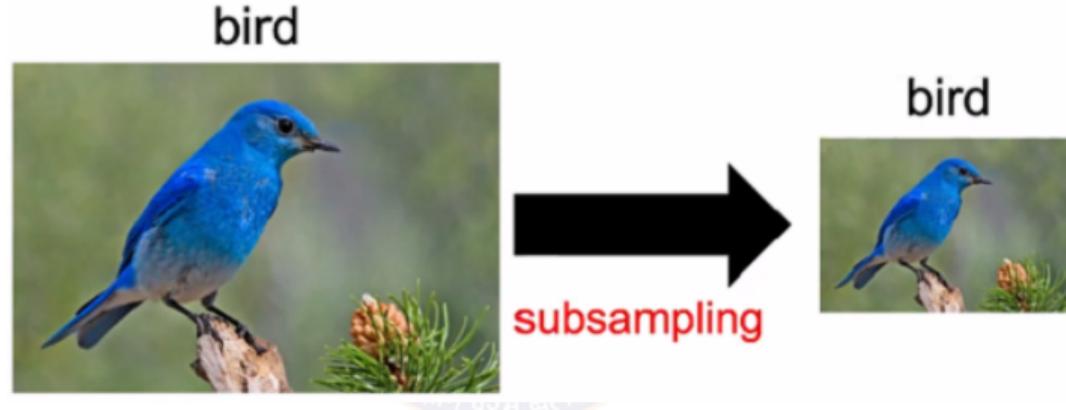
WHY CNN FOR IMAGES

- Same pattern may appear in different regions.



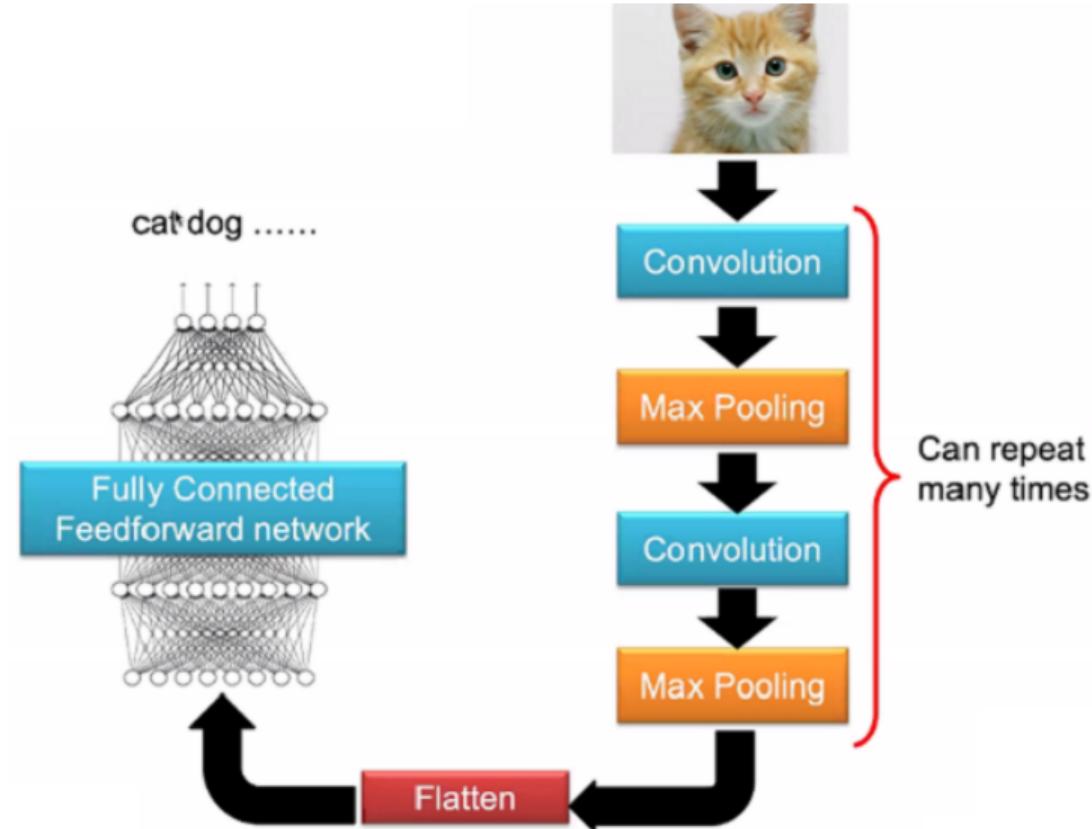
WHY CNN FOR IMAGES

- Sub-sampling the pixels will not change the object.



- We can sub-sample the image to make it smaller.
- Less parameters for the network to process the image.

THE WHOLE CNN



THE WHOLE CNN

Property 1

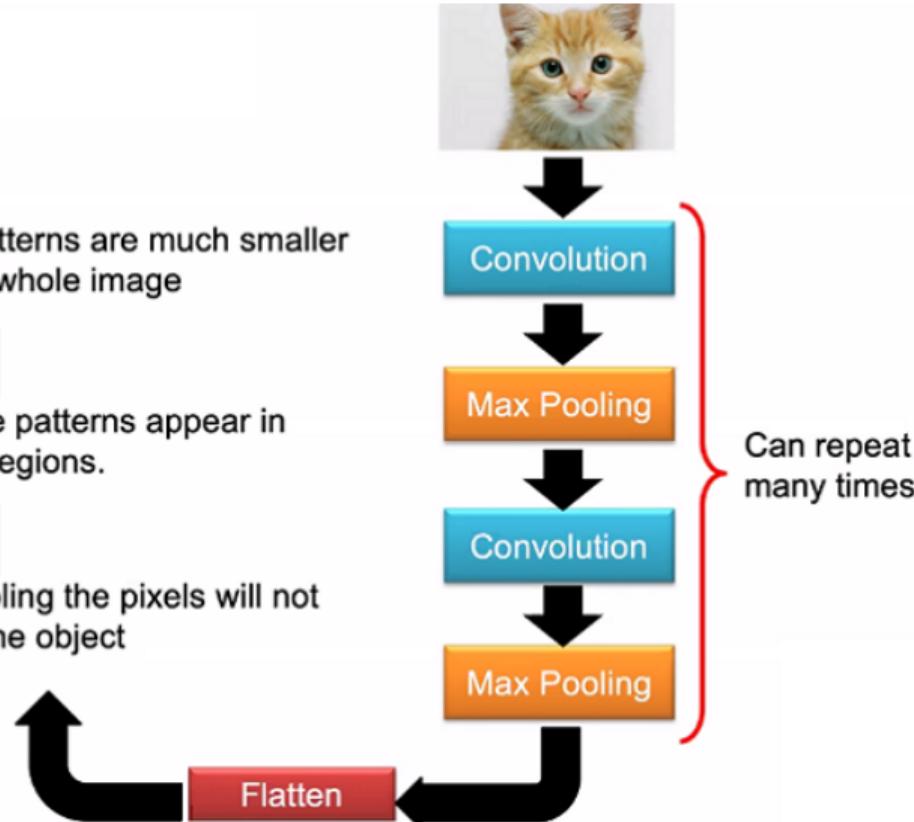
- Some patterns are much smaller than the whole image

Property 2

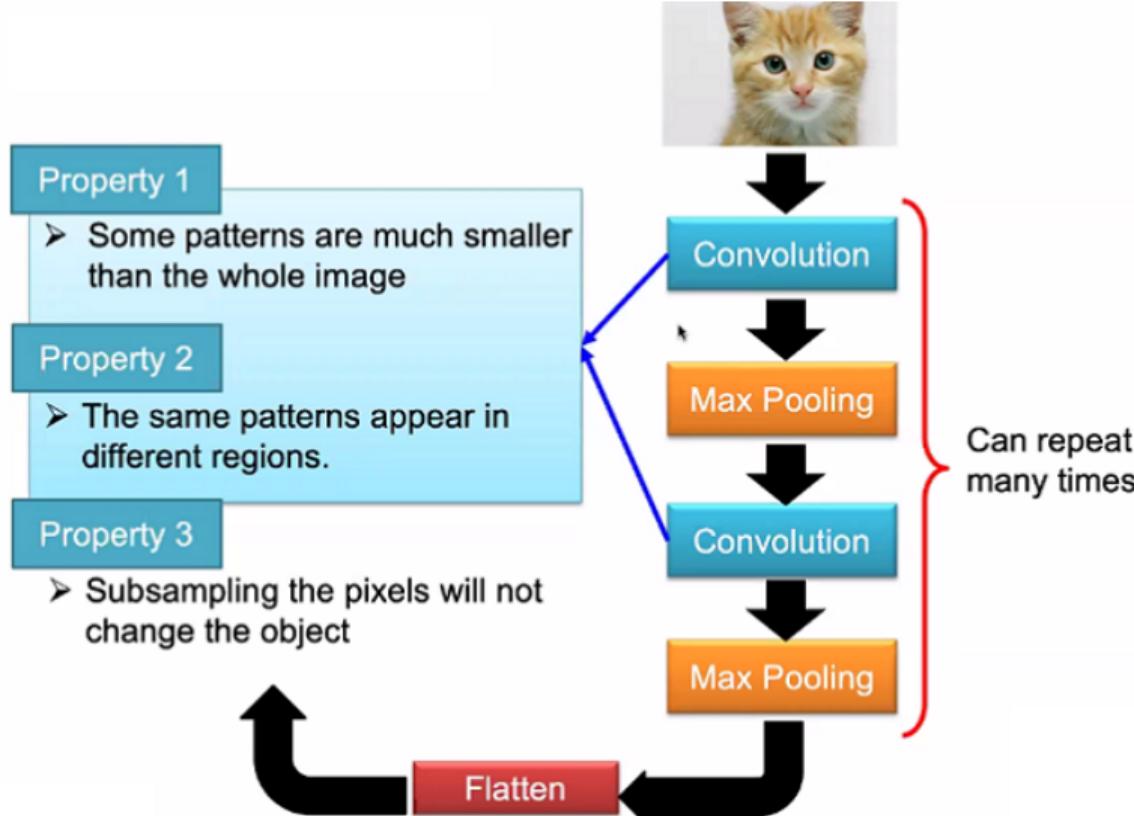
- The same patterns appear in different regions.

Property 3

- Subsampling the pixels will not change the object



THE WHOLE CNN



THE WHOLE CNN

The whole CNN

Property 1

- Some patterns are much smaller than the whole image

Property 2

- The same patterns appear in different regions.

Property 3

- Subsampling the pixels will not change the object



Convolution



Max Pooling



Convolution



Max Pooling



Can repeat
many times

IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



CONVOLUTION LAYER

- Perform Convolution operation.
- In mathematics, convolution is a mathematical operation on two functions that produces a third function that expresses how the shape of one is modified by the other. [Wikipedia]
- The term convolution refers to both the result function and to the process of computing it. [Wikipedia]
- We will assume that the kernel is square matrix.
- We will assume that the kernel is centered on the pixel of interest. (One reason why kernel is odd matrix.)
- Convolution operation

$$S_{ij} = (I * K)_{ij} = \sum_{a=\lfloor -\frac{m}{2} \rfloor}^{\lfloor \frac{m}{2} \rfloor} \sum_{b=\lfloor -\frac{n}{2} \rfloor}^{\lfloor \frac{n}{2} \rfloor} I_{i-a, j-b} K_{\frac{m}{2}+a, \frac{n}{2}+b}$$

1-D CONVOLUTION

- Suppose we track position of an object using a sensor or GPS device at discrete time intervals.
- Assume that the sensor is noisy.
- To obtain a less noisy estimate, we would like to take a weighted average.
- We would only sum over a small window or a weight array (filter) (w).

$$s_t = \sum_{a=0}^k x_{t-a} w_a$$

- Slide the filter over the input and convolve.
- The input and the kernel are one dimensional.

1-D CONVOLUTION

$$S = x_5w_0 + x_4w_1 + x_3w_2 + x_2w_3 + x_1w_4 + x_0w_5$$

X	1.00	1.10	1.20	1.40	1.70	1.80	1.90	2.10	2.20	2.40	2.50
---	------	------	------	------	------	------	------	------	------	------	------

W	0.01	0.02	0.03	0.04	0.4	0.5					
---	------	------	------	------	-----	-----	--	--	--	--	--

S			1.80								
---	--	--	------	--	--	--	--	--	--	--	--

W	0.01	0.02	0.03	0.04	0.4	0.5					
---	------	------	------	------	-----	-----	--	--	--	--	--

S		1.80	1.96								
---	--	------	------	--	--	--	--	--	--	--	--

W	0.01	0.02	0.03	0.04	0.4	0.5					
---	------	------	------	------	-----	-----	--	--	--	--	--

S		1.80	1.96	2.11							
---	--	------	------	------	--	--	--	--	--	--	--

2D CONVOLUTION OPERATION

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

Input Image 6×6

These values inside the matrix are to be learned by the network. We call these as parameters or weights.

1	-1	-1
-1	1	-1
-1	-1	1

Kernel / Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Kernel / Filter 2

Property 1: Each filter detects a small pattern (3×3)

CONVOLUTION OPERATION

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

Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Start at top left corner.

Perform element wise product and add all the values.

Output

3

CONVOLUTION OPERATION

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

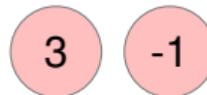
Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Slide by one pixel towards right.
Perform element wise product and
add all the values.

Output



CONVOLUTION OPERATION

stride $s = 1$

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

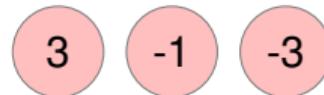
Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Slide by one pixel towards right.
Perform element wise product and
add all the values.

Output



CONVOLUTION OPERATION

stride $s = 1$

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

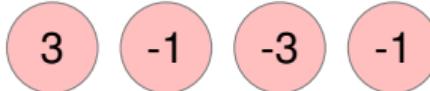
Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Slide by one pixel towards right.
Perform element wise product and
add all the values.

Output



CONVOLUTION OPERATION

stride $s = 1$

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

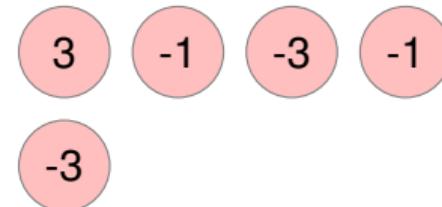
Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Jump to next row, left most pixel.
The filter has to completely
superimpose the pixels of input.
NOT PARTIALLY.

Output



CONVOLUTION OPERATION

stride $s = 1$

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

Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

After sliding through the entire image.

3	-1	-3	-1
-3	-1	0	-3
-3	-2	0	1
3	-2	-2	-1

CONVOLUTION OPERATION

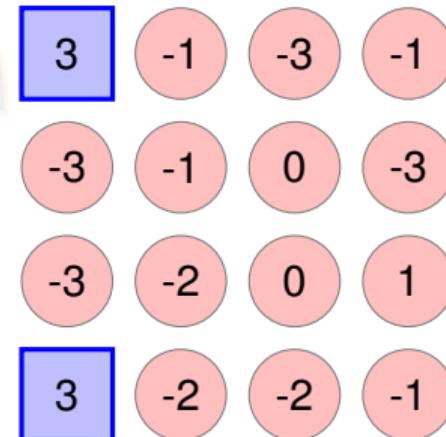
stride $s = 1$

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

Input Image 6×6

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



Property 2 is satisfied.

CONVOLUTION OPERATION

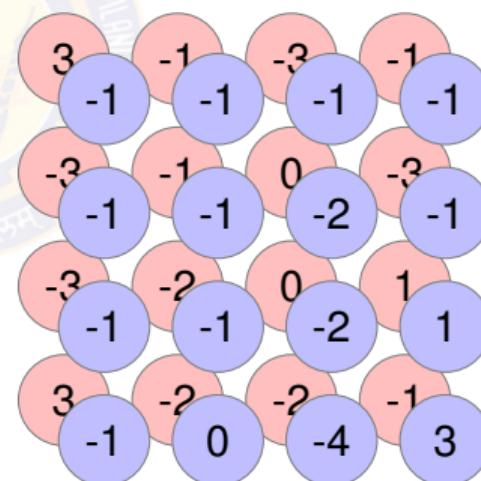
stride $s = 1$

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

Input Image 6×6



Filter 1			Filter 2		
1	-1	-1	-1	1	-1
-1	1	-1	-1	1	-1
-1	-1	1	-1	1	-1



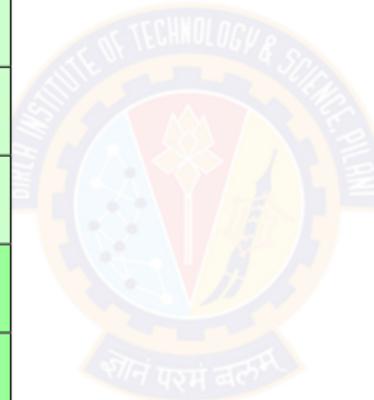
Feature
Maps

CONVOLUTION OPERATION

stride $s = 3$

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

Input Image 6×6



Filter1

1	-1	-1
-1	1	-1
-1	-1	1

Output

3	-1
3	-1

CONVOLUTION WITH STRIDING

- Stride means skip by how many pixels.

Input Size : $n \times n$

Stride : s

Kernel Size : $f \times f$

Output Size : $\left\lfloor \frac{n-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n-f}{s} + 1 \right\rfloor$

- The filter must lie completely inside the input.

PADDING

- When filter is convolved with the input image
 - ▶ Image shrinks – 6×6 image became a 4×4 image.
 - ▶ Pixels in the corners are used only once, when compared the pixels in the middle. The data in the corners are thrown away.
- To resolve, apply padding to the input image before convolution.
- Padding means append zeros around the boundary of the input.

CONVOLUTION WITH PADDING

Input Image with padding 8×8

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

Filter1

$$* \begin{matrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{matrix}$$

PADDING

- Two kinds of padding

- Valid Padding

- Padding is not applied.

$$p = 0$$

- Input size \neq output size.
 - Output shrinks when compared to the input.

- Same Padding

- Padding is applied.

$$p = \frac{f - 1}{2}$$

- Input size = output size.
 - We maintain the size of input and output as the same.

PADDING

- How much to pad?

$$\text{Padding : } p = \frac{f - 1}{2}$$

Kernel $3 \times 3 : p = 1$

Kernel $5 \times 5 : p = 2$

Kernel $7 \times 7 : p = 3$

- Padding preserves the input size.

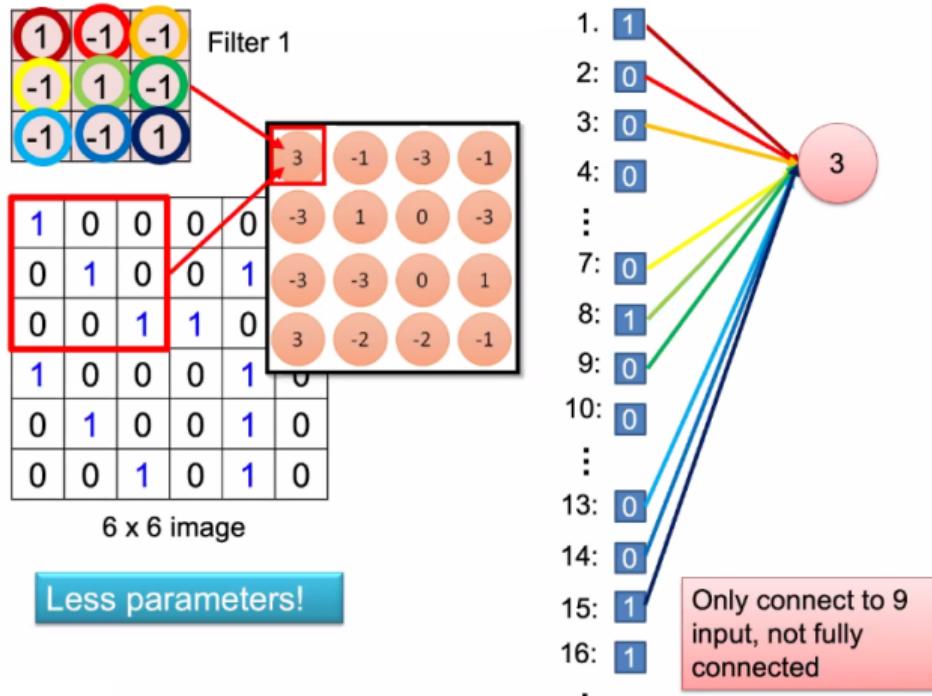
Input Size : $n \times n$

Padding : p

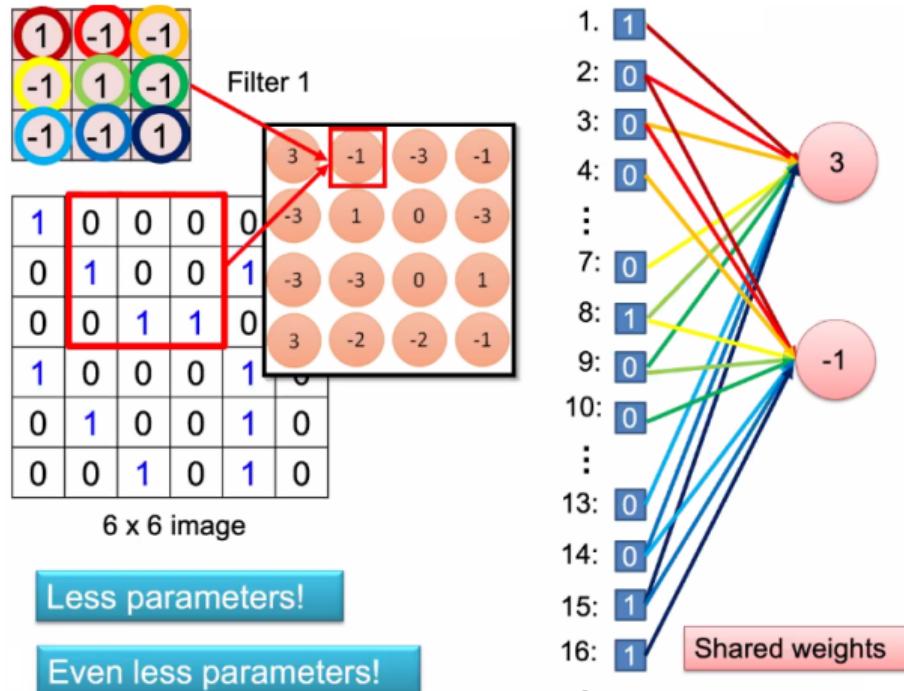
Kernel Size : $f \times f$

Output Size : $(n + 2p - f + 1) \times (n + 2p - f + 1)$

CONVOLUTIONAL LAYER

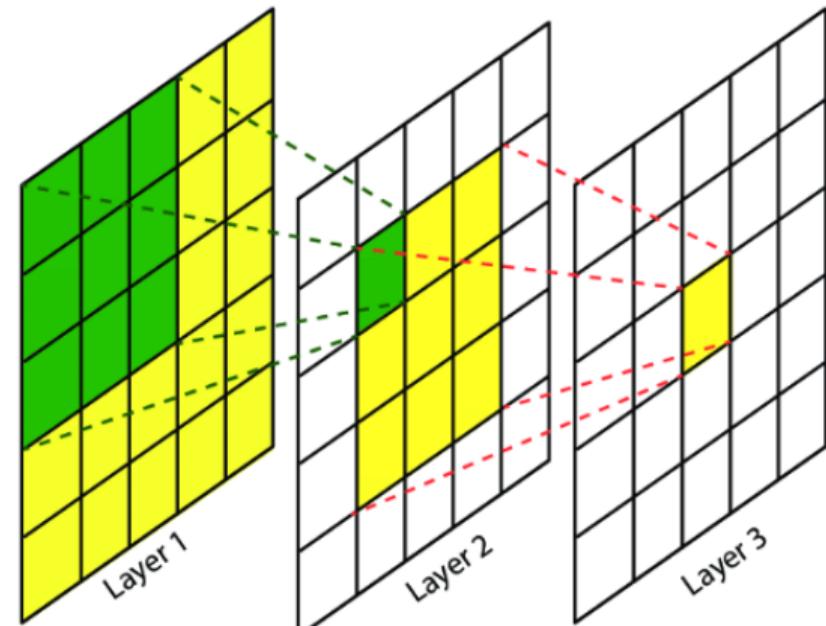


CONVOLUTIONAL LAYER

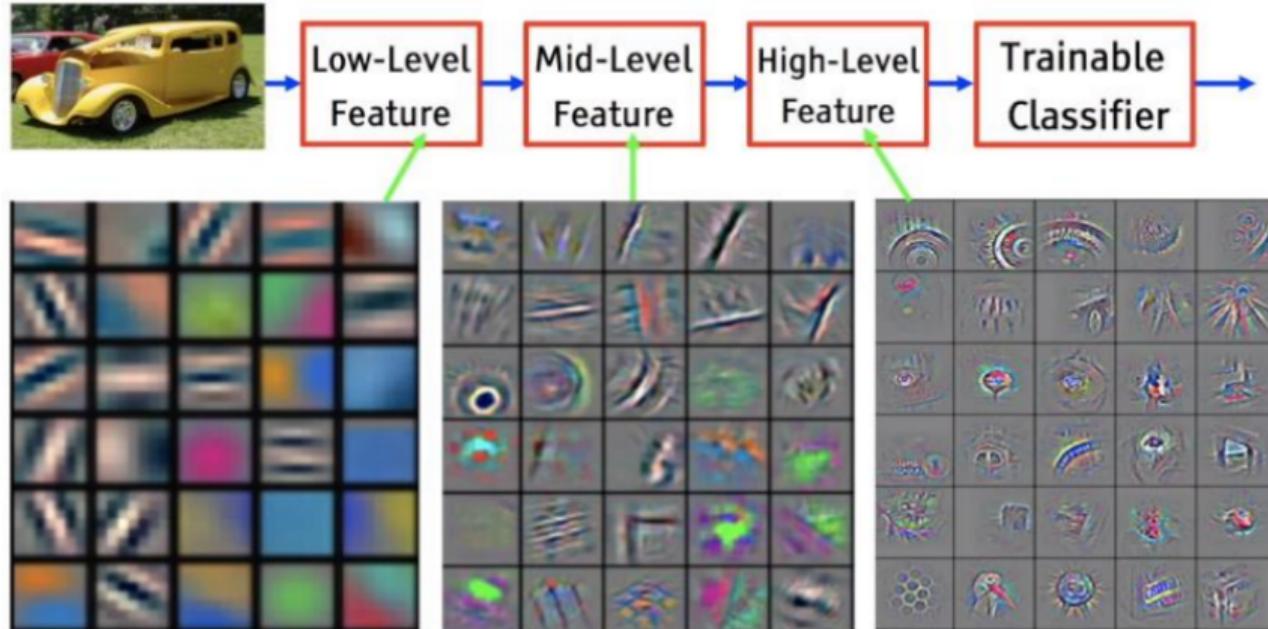


ACTIVATIONS OR RECEPTIVE FIELDS

- Effective area “seen” by a neuron
- What is the effective number of pixels seen at each layer?
 - Input image x
 - $h_1 = \text{Conv2d}(x, k = 3)$
 - $h_2 = \text{Conv2d}(h_1, k = 3)$
 - $h_3 = \text{Conv2d}(h_2, k = 3)$
 - ...
 - $h_n = \text{Conv2d}(h_{n-1}, k = 3)$
- Homework: What happens if the kernel size was 5?



LOCAL TO GLOBAL



PARAMETER SIZE IN CONVOLUTION LAYER

- Input image: $B \times W_{in} \times H_{in} \times C_{in}$
 - ▶ Input data has B samples
 - ▶ Input or Image size is $W_{in} \times H_{in} \times C_{in}$
- Convolution filter
 - ▶ Weight parameters: $C_{in} \times C_{out} \times K_W \times K_H$
 - ▶ Bias parameters: C_{out}
- Output activation: $B \times W_{out} \times H_{out} \times C_{out}$
 - ▶ C_{out} filters of shape $C_{in} \times K_W \times K_H$ operate on each sample

PARAMETER SIZE IN CONVOLUTION LAYER

- Each filter gives one 2D output.
- K filter will give K such 2D outputs.
- The number of parameters depend on the number of kernels and the kernel size.

Number of kernels in layer ($I - 1$) : K_{I-1}

Kernel size : $f \times f$

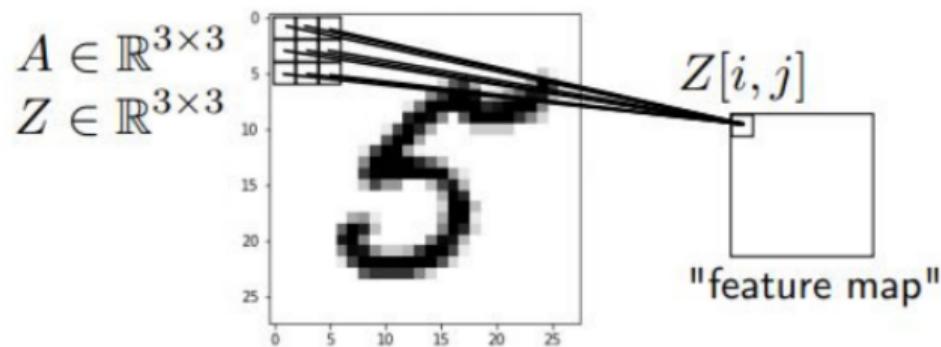
Number of kernels in layer (I) : K_I

Number of Parameters in layer (I) : $K_{I-1} * f * f * K_I$

CROSS-CORRELATION VS CONVOLUTION

- Convolution in Deep Learning is actually cross-correlation.
- Cross-correlation is sliding dot product over the images.
- Cross-correlation

$$Z[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k K[u, v] A[i + u, j + v]$$



CROSS-CORRELATION VS CONVOLUTION

Cross-Correlation

$$Z[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k K[u, v] A[i+u, j+v]$$

(Note the looping direction given in red.)

1) -1,-1	2) -1,0	3) -1,1
4) 0,-1	5) 0,0	6) 0,1
7) 1,-1	8) 1,0	9) 1,1

Convolution

$$Z[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k K[u, v] A[i-u, j-v]$$

(Flipping the kernel.)

9) -1,-1	8) -1,0	7) -1,1
6) 0,-1	5) 0,0	4) 0,1
3) 1,-1	2) 1,0	1) 1,1

CROSS-CORRELATION VS CONVOLUTION

- Convolution in Deep Learning is actually cross-correlation.
- Convolution (as in Signal Processing) has nice association property.

$$(A * B) * C = A * (B * C)$$

- In DL, we usually don't care about that (as opposed to traditional computer vision and signal processing applications.)
- Cross-correlation is easier to implement.
- Maybe the term "convolution" for cross-correlation because popular, because "Cross-correlation Neural Network" sounds weird.

TRANSPOSED CONVOLUTION

Regular Convolution:

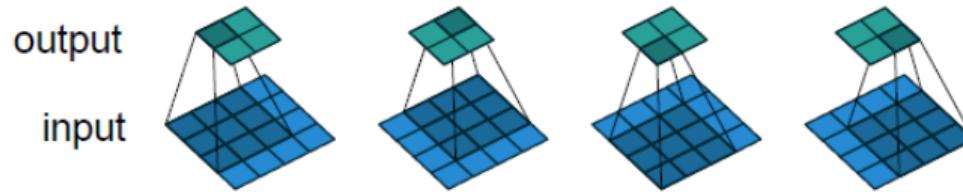
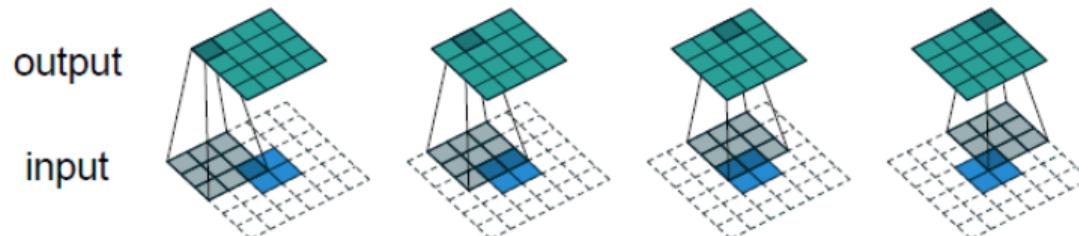


Figure 2.1: (No padding, unit strides) Convolving a 3×3 kernel over a 4×4 input using unit strides (i.e., $i = 4$, $k = 3$, $s = 1$ and $p = 0$).

Transposed Convolution (emulated with direct convolution):



Dumoulin, Vincent, and Francesco Visin. "[A guide to convolution arithmetic for deep learning](#)." *arXiv preprint arXiv:1603.07285* (2016).

DILATED CONVOLUTION

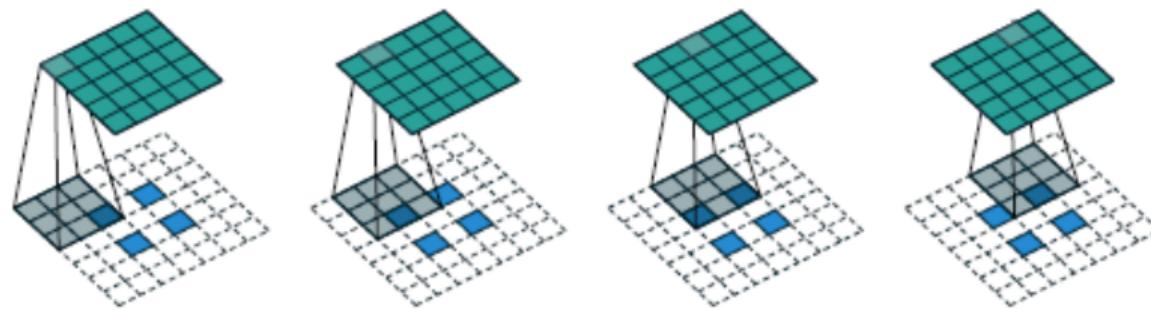


Figure 4.5: The transpose of convolving a 3×3 kernel over a 5×5 input using 2×2 strides (i.e., $i = 5, k = 3, s = 2$ and $p = 0$). It is equivalent to convolving a 3×3 kernel over a 2×2 input (with 1 zero inserted between inputs) padded with a 2×2 border of zeros using unit strides

Dumoulin, Vincent, and Francesco Visin. "[A guide to convolution arithmetic for deep learning](#)." arXiv preprint arXiv:1603.07285 (2016).

REMEMBER IMAGE INVARIANCE CHALLENGE

Object can appear anywhere in the image.

Image A

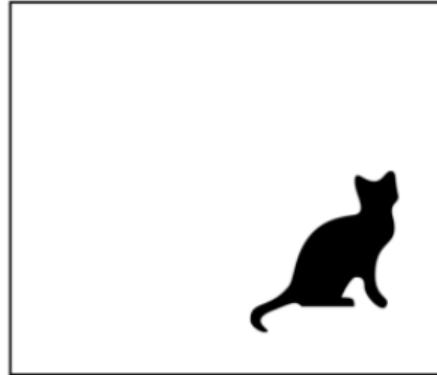
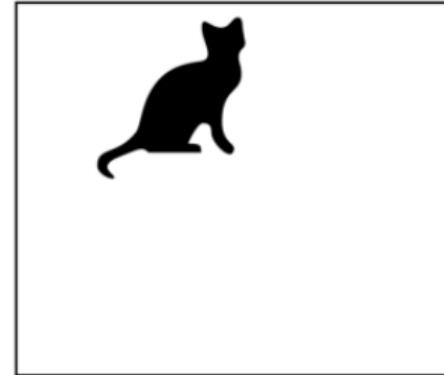


Image B

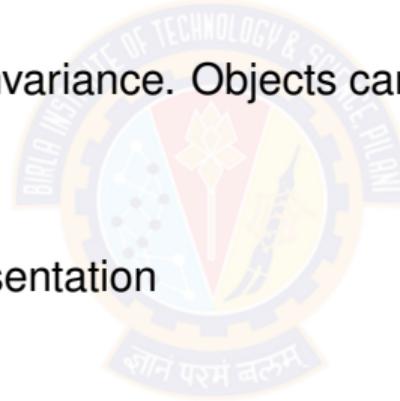


Model predicts “cat” for both images!

Source: <https://penkovsky.com/neural-networks/day5/>

POOLING (POOL) LAYER

- Helps to achieve image invariance. Objects can appear anywhere in the Image and still get detected.
- Used for sub-sampling
- Reduce the size of representation
- Speed up computation



POOLING

- Apply a filter of size f and stride s . Padding is not applied.
- The output dimensions will be

$$\left\lfloor \frac{n-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n-f}{s} + 1 \right\rfloor$$

- Fixed computation
- Gradient descent is not applied as there are not parameters to be learned.
- Pooling is applied on each of the channels.

POOLING LAYER

- Two types
 - ① Max pooling
 - ② Average pooling



MAX POOLING

- Take the maximum of the sub region that is considered.

Input 4×4

2	4	7	4
6	6	9	8
3	4	8	3
7	5	3	6

Max pooling

$$f = 2$$

$$s = 2$$

Output

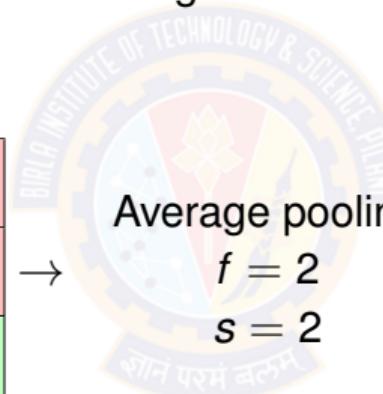
2×2

6	9
7	8

AVERAGE POOLING

- Compute the average of the sub region that is considered.

Input 4×4			
2	4	7	4
6	6	9	8
3	4	8	3
7	5	3	6



Average pooling

$$f = 2$$

$$s = 2$$

Output

$$2 \times 2$$

4.5	7
4.75	5

CONVOLUTION + MAX POOLING

Input Image 6×6

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

Filter 1

1	-1	-1
-1	1	-1
-1	-1	1

Filter 2

-1	1	-1
-1	1	-1
-1	1	-1

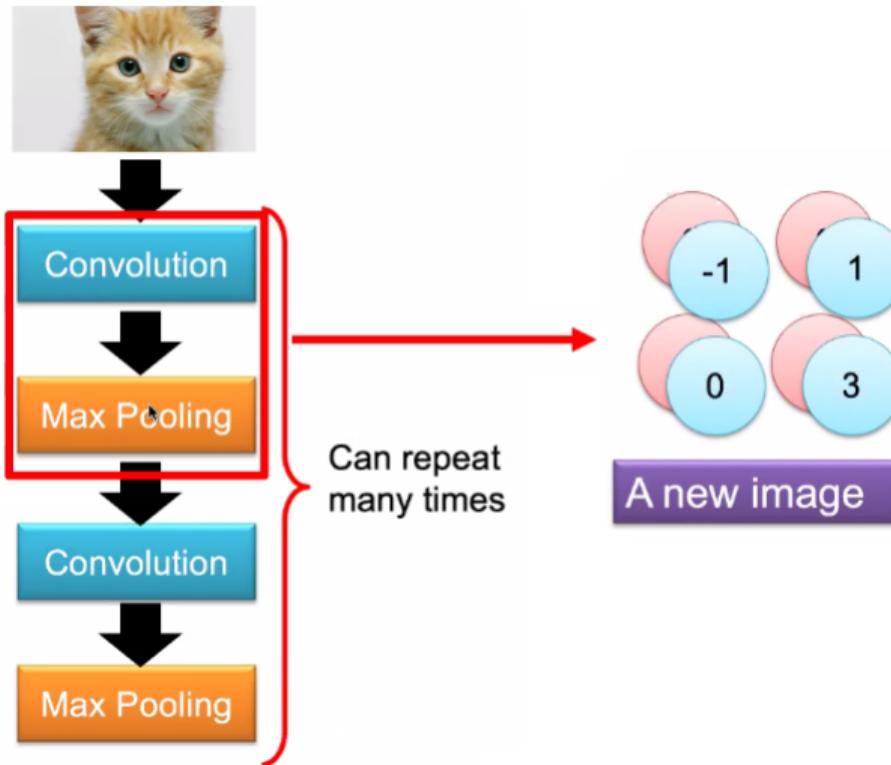
Convolve s=1

3	-1	-3	-1
-3	-1	0	-3
-3	-2	0	1
3	-2	-2	-1
-1	-1	-1	-1
-1	-1	-2	-1
-1	-1	-2	1
-1	0	-4	3

Maxpool f=2, s=2

3	0
3	1
-1	-1
0	3

CONVOLUTION + MAX POOLING



FLATTEN

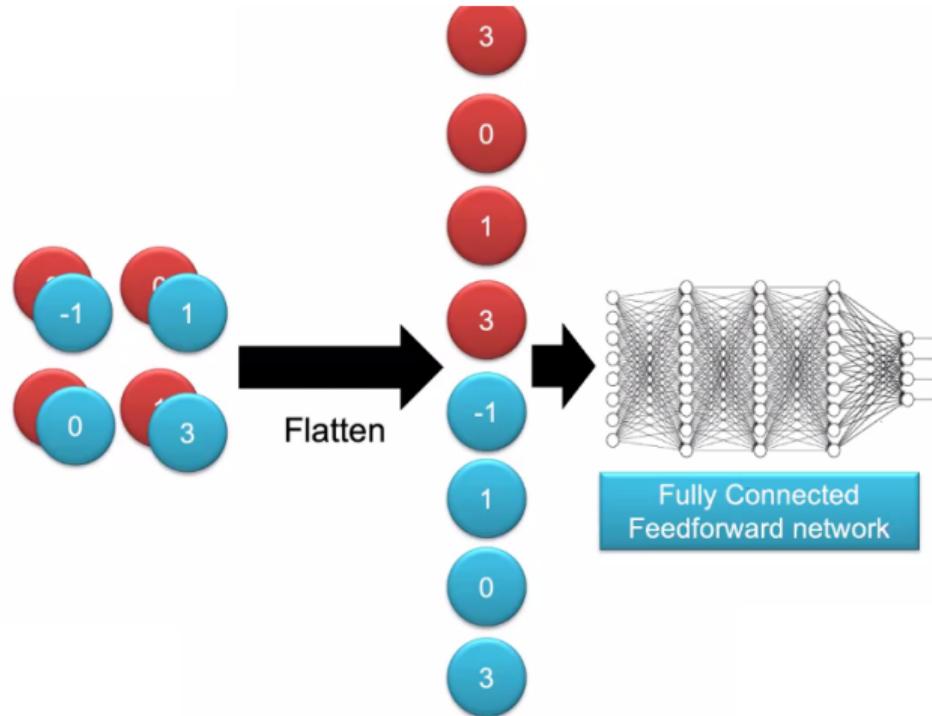


IMAGE IS HIGHLY NON-LINEAR

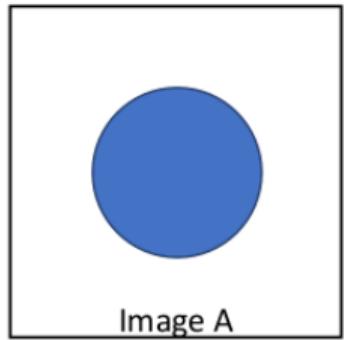


Image A

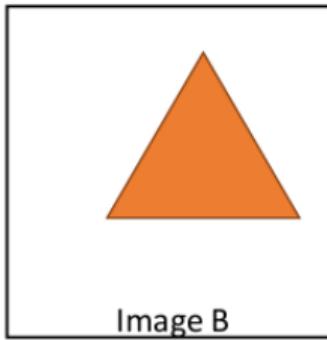
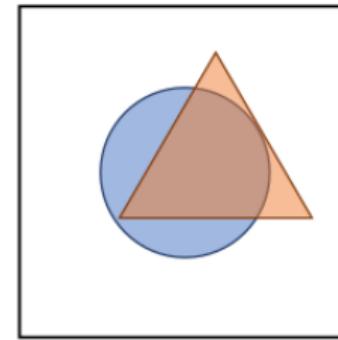
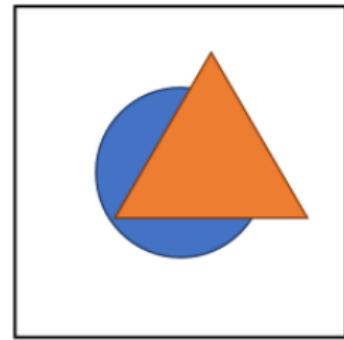


Image B



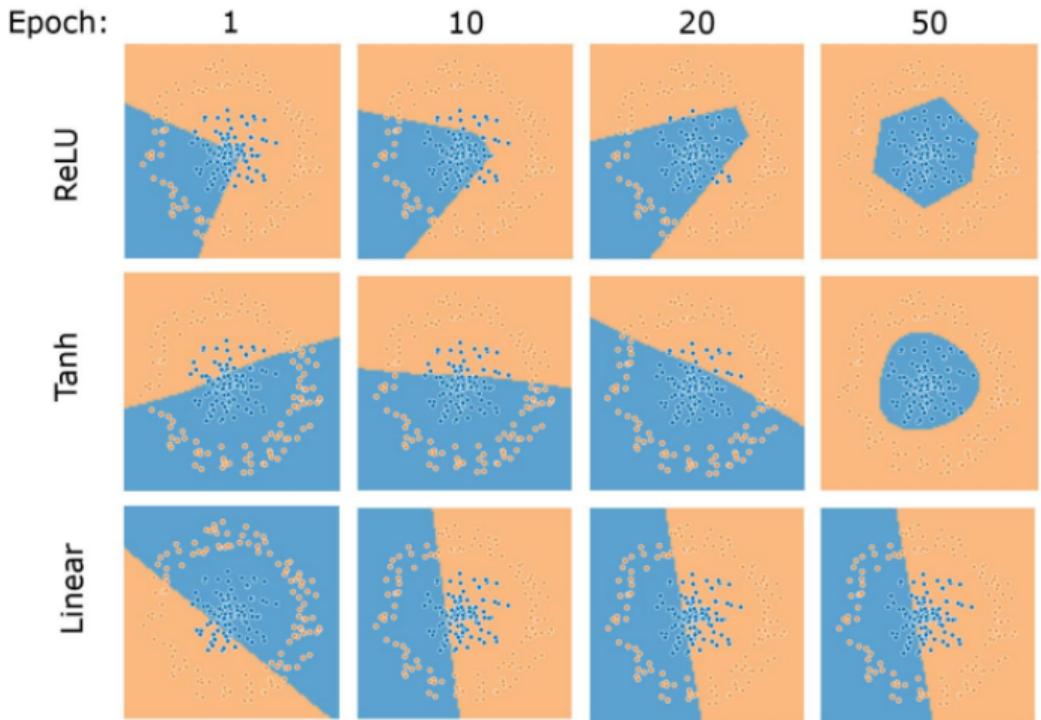
Linear combination
 $A+B$



Actual combination
 $A+B$

RECTIFIED LINEAR UNIT (ReLU)

- Very fast
- Gradients
- Piece-wise linear approximation



RELU IN IMAGES

Ground truth image



Standard fully-connected net

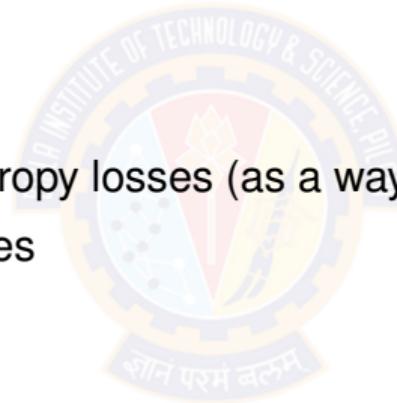


With Positional Encoding



LOSS FOR TRAINING

- Mean Squared Error
- Binary Cross Entropy
- Multi class Cross Entropy
- Class-weighted cross-entropy losses (as a way to counter imbalance)
- L1 (MAE), L2 (MSE) losses
- Focal, Dice, Huber loss
- Margin ranking losses
- Generalized Intersection-over-Union (GloU) loss
- Connectionist Temporal Classification (CTC) loss



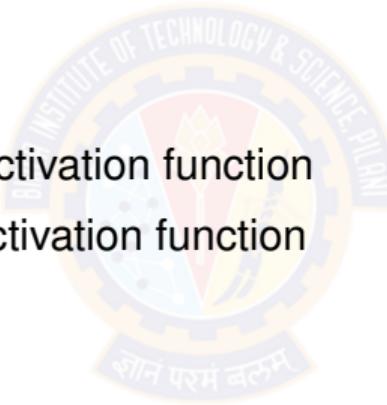
REGULARISATION AND OPTIMIZATION

- Drop out
- Batch Normalisation
- RMSProp, Adam or any other optimizer can be used.



CLASSIFICATION

- $p \in [0, 1]$ - use Sigmoid activation function
- k classes - use softmax activation function



IN THIS SEGMENT

- 1 COMPUTER VISION TASKS
- 2 IMAGE BASICS
- 3 IMAGE CLASSIFICATION CHALLENGES
- 4 CONVOLUTIONAL NEURAL NETWORK
- 5 CONVOLUTION
- 6 PUTTING IT TOGETHER



PUTTING EVERYTHING TOGETHER

- ① Identify the task as classification or regression. This helps in identify the output size.
- ② Collect the N images. This helps in identifying input size as $W \times H \times C$.
- ③ Feature extraction
 - ▶ Decide on the number of Convolution layers. This can be derived by looking at how large or small pattern we are interested in.
 - ▶ Decide on the kernel size and number of kernels in each layer. Smaller patterns in images require smaller kernels. Number of kernels can be derived from the non-linearity.
 - ▶ Use ReLU activation function for better approximation.
- ④ Classification
 - ▶ Use Fully connected layers for classification.
 - ▶ The number of (dense) layers depends on the non-linearity.
 - ▶ The last (output) layer activation is decided based on Step 1.
- ⑤ Train and Test on the available dataset.

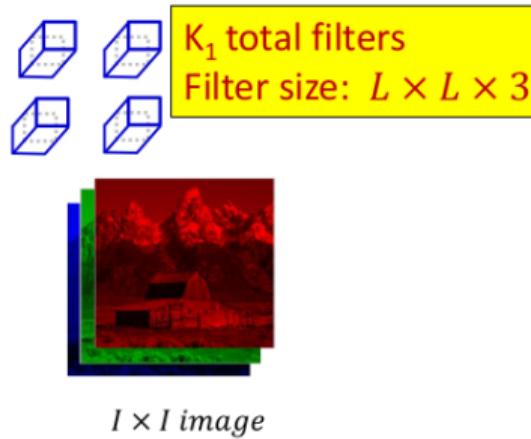
CONVOLUTIONAL NEURAL NETWORK



- Input: 1 or 3 images
 - Grey scale or color
 - Will assume color to be generic

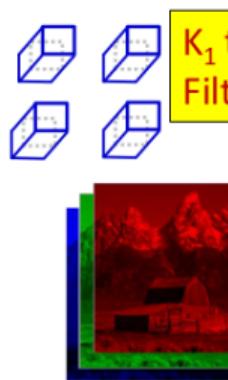
CONVOLUTIONAL NEURAL NETWORK

- Input is convolved with a set of K_1 filters or kernels.
 - ▶ Typically, K_1 is a power of 2 e.g 2,4,8,16,32,...
 - ▶ Filters are typically $5 \times 5 \times 3$, $3 \times 3 \times 3$ or $1 \times 1 \times 3$ (better notation).
 - ▶ Typical stride is 1 or 2.



CONVOLUTIONAL NEURAL NETWORK

- Input is convolved with a set of K_1 filters or kernels.
 - ▶ Typically, K_1 is a power of 2 e.g 2,4,8,16,32,...
 - ▶ Filters are typically 5x5, 3x3 or 1x1.
 - ▶ Filters should be small enough to capture fine features, particularly important for scaled down images.



Parameters to choose: K_1 , L and S

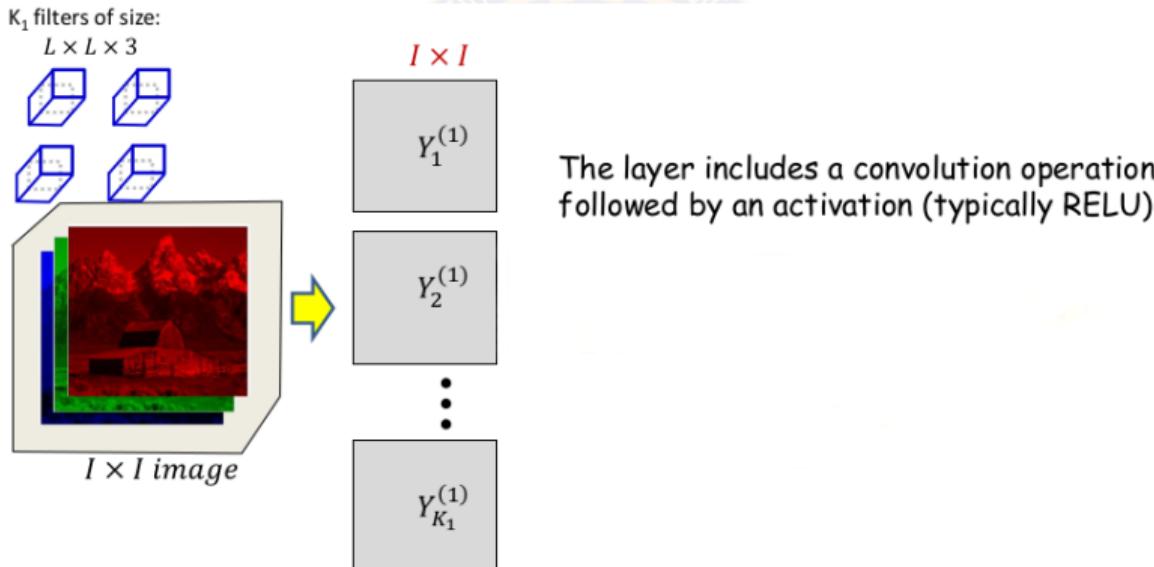
1. Number of filters K_1
2. Size of filters $L \times L \times 3 + \text{bias}$
3. Stride of convolution S

Total number of parameters: $K_1(3L^2 + 1)$

CONVOLUTIONAL NEURAL NETWORK

- First Convolutional layer:

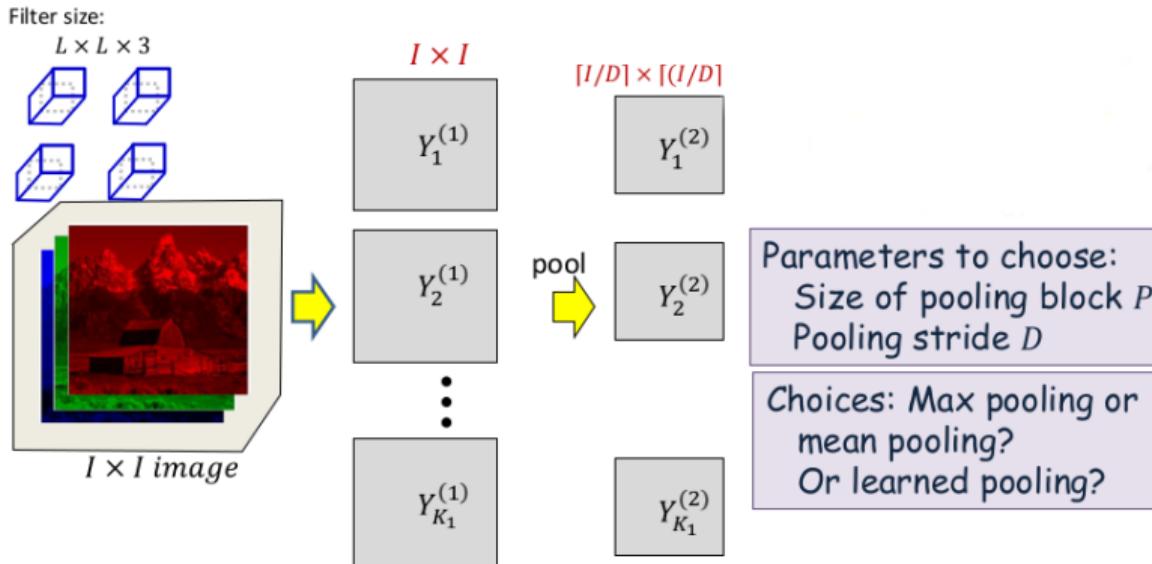
- ▶ Several convolutional filters, which are 3D.
- ▶ Convolution is followed typically by a ReLU activation.
- ▶ Each filter creates a singly 2D feature map.



CONVOLUTIONAL NEURAL NETWORK

- First Pooling layer:

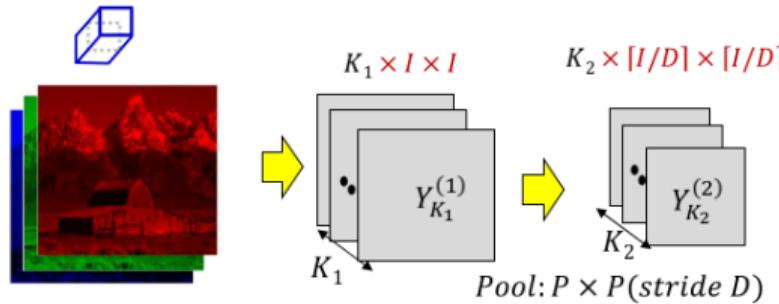
- ▶ For each $P \times P$ block of each map, pool down to a single value.



CONVOLUTIONAL NEURAL NETWORK

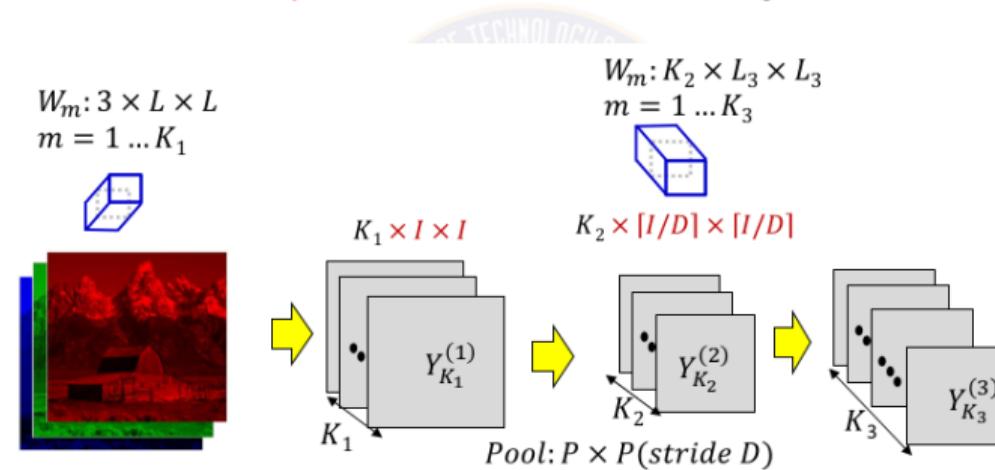
- Better drawing

$$W_m: 3 \times L \times L$$
$$m = 1 \dots K_1$$



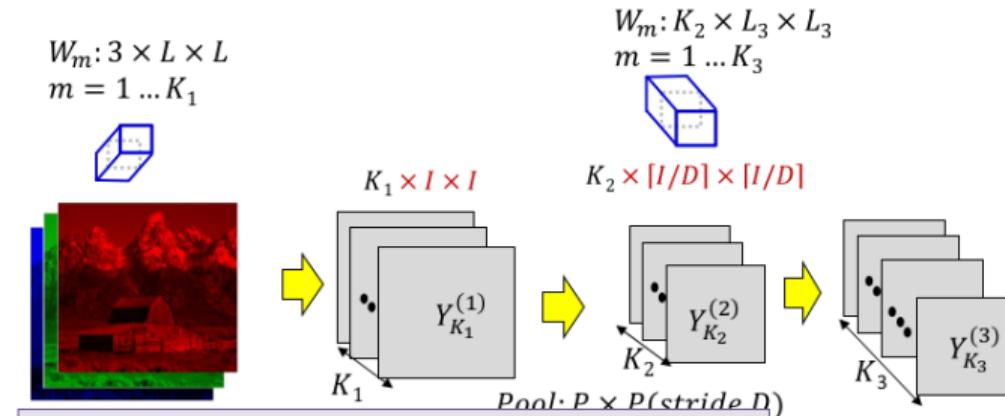
CONVOLUTIONAL NEURAL NETWORK

- Second Convolutional layer" K_2 3D filters resulting in K_2 feature maps.



CONVOLUTIONAL NEURAL NETWORK

- Second Convolutional layer



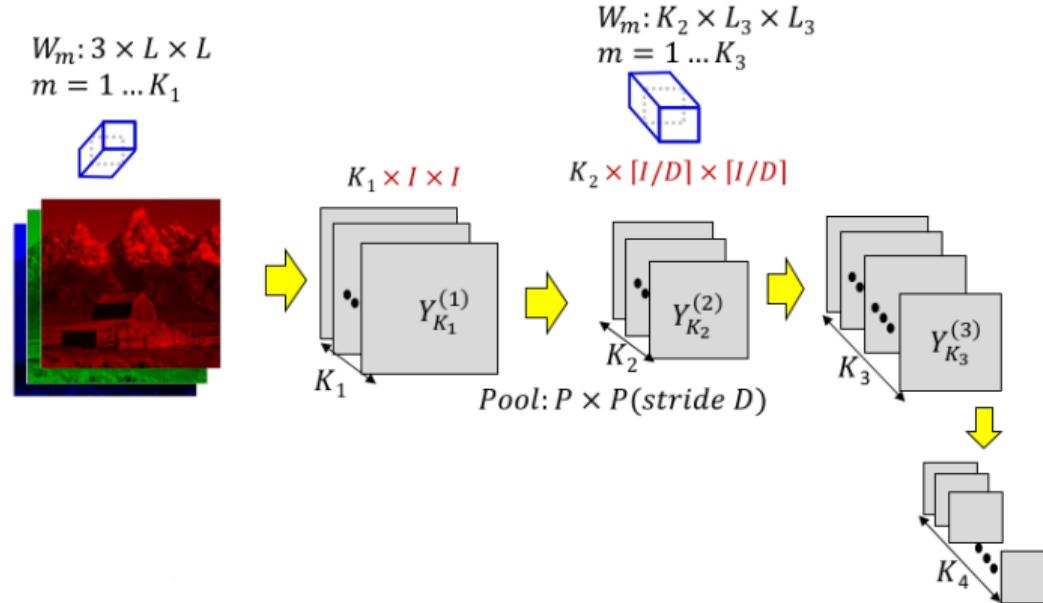
Parameters to choose: K_3, L_3 and S_3

1. Number of filters K_3
2. Size of filters $L_3 \times L_3 \times K_2 + bias$
3. Stride of convolution S_3

Total number of parameters: $K_3(K_2L_3^2 + 1)$
All these parameters must be learned

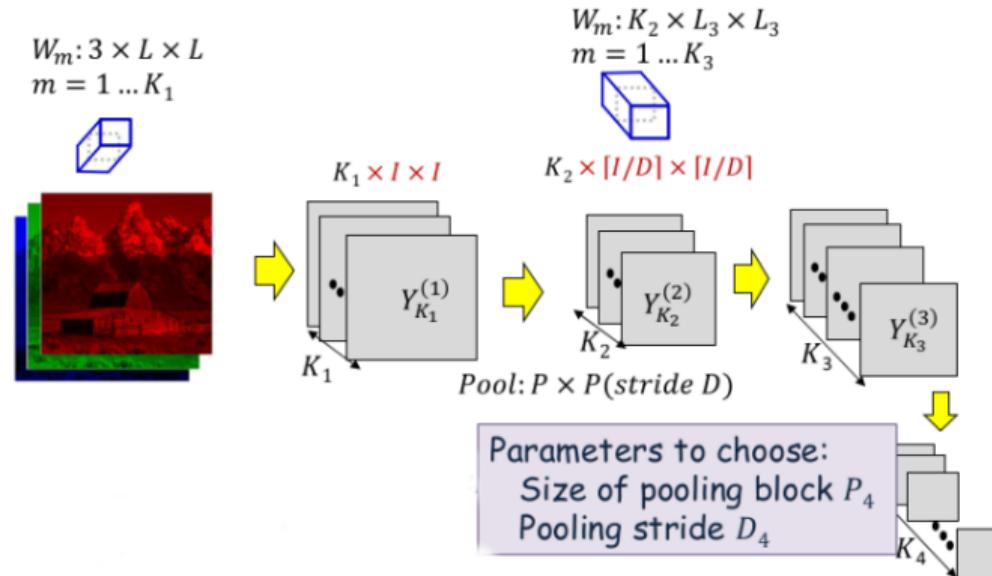
CONVOLUTIONAL NEURAL NETWORK

- **Second Pooling layer:** K_2 pooling operations results in K_2 reduced 2D maps.



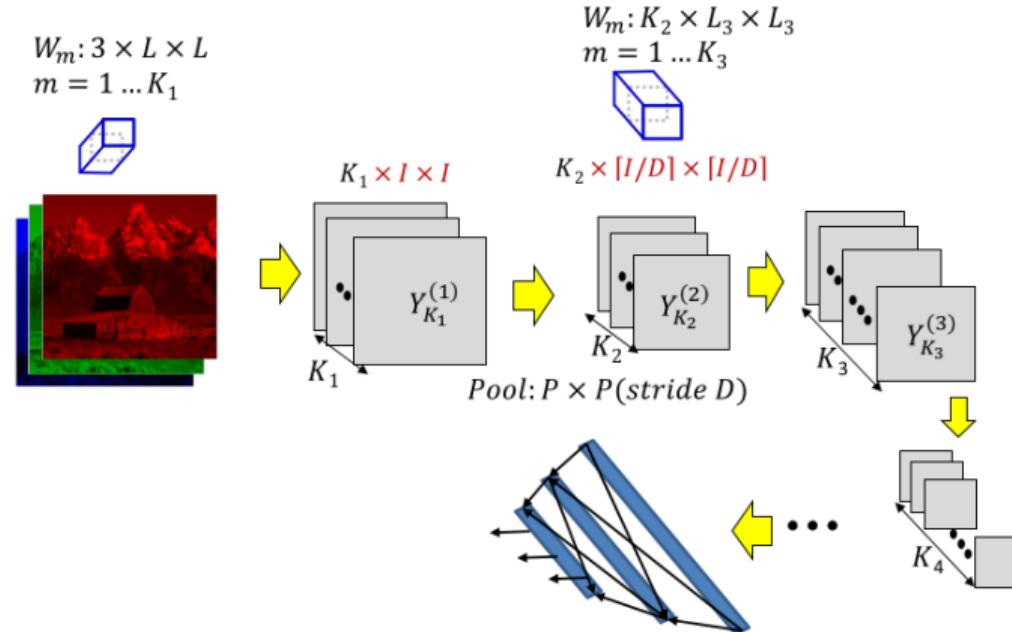
CONVOLUTIONAL NEURAL NETWORK

- Second Pooling layer



CONVOLUTIONAL NEURAL NETWORK

- **Fully connected layer:** This continues for several layers until the final convolved feature maps are fed to MLP.



CONVOLUTIONAL NEURAL NETWORK

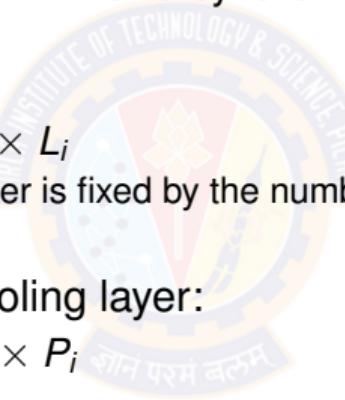
Size of Layers

- Each convolution layer with stride 1 typically maintains the size of the image.
 - ▶ With appropriate zero padding.
 - ▶ If performed without zero padding it will decrease the size of the input.
- Each convolution layer will generally increase the number of maps from the previous layer.
 - ▶ Increasing layers reduces the amount of information lost by subsequent downsampling.
- Each pooling layer with stride D decreases the size of the maps by a factor of D.
- Filters within a layer must all be the same size, but sizes may vary with layer. Similarly for pooling D may vary with layer.
- In general the number of convolutional filters increases with layers.

CONVOLUTIONAL NEURAL NETWORK

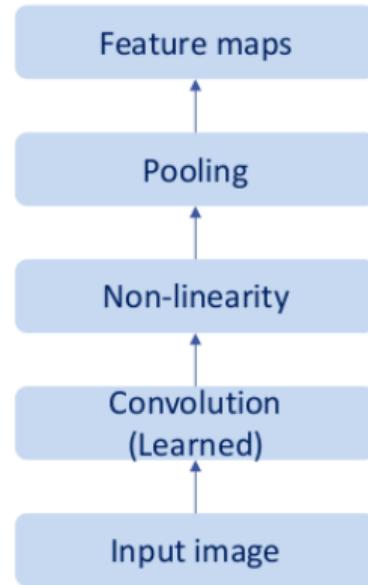
Parameters to choose (design choices)

- Number of convolutional and downsampling layers
 - ▶ And arrangement (order in which they follow one another)
- For each convolution layer:
 - ▶ Number of filters K_i
 - ▶ Spatial extent of filter $L_i \times L_i$
 - ★ The “depth” of the filter is fixed by the number of filters in the previous layer K_{i-1} .
 - ▶ The stride S_i
- For each downsampling/pooling layer:
 - ▶ Spatial extent of filter $P_i \times P_i$
 - ▶ The stride D_i
- For the final MLP:
 - ▶ Number of layers
 - ▶ Number of neurons in each layer



CNN-Recap

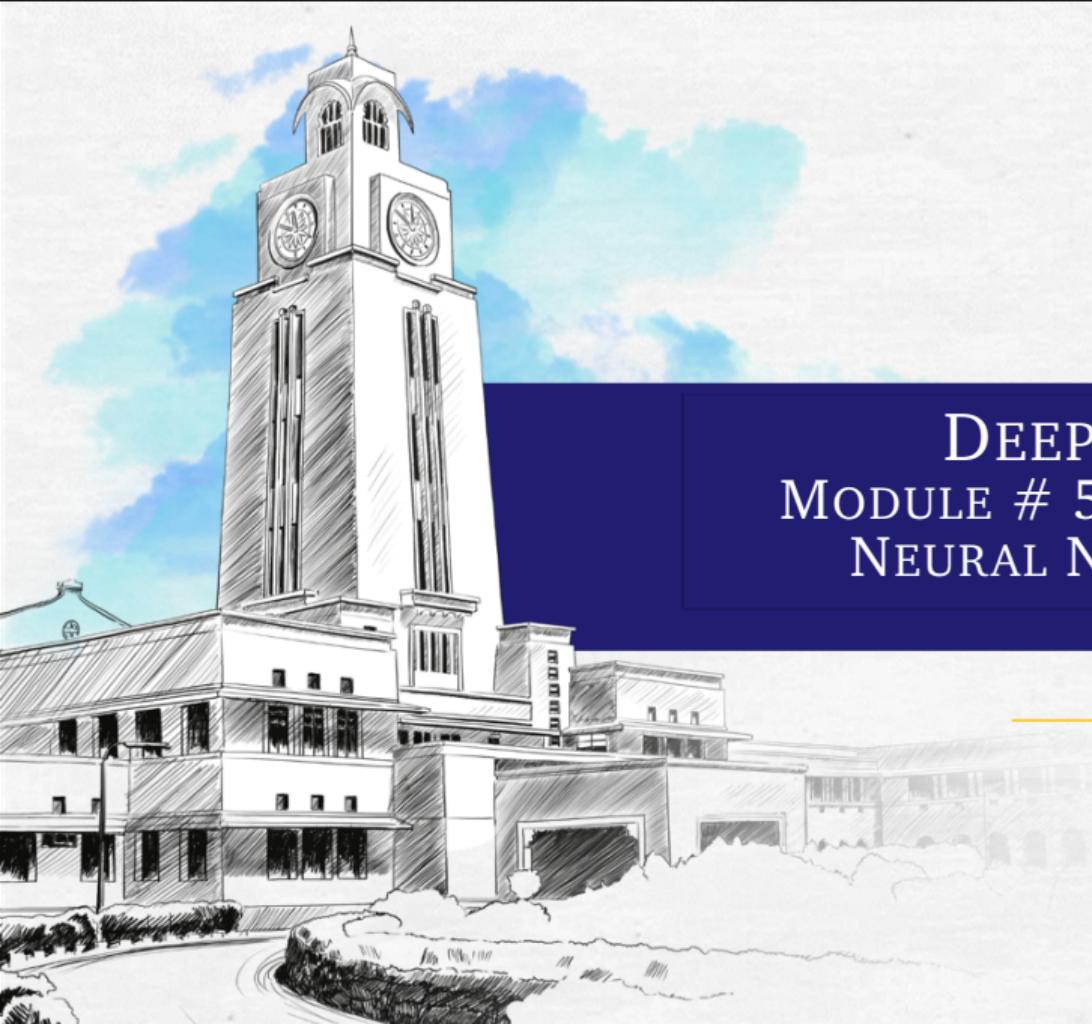
- Neural network with specialized connectivity structure
- Feed-forward:
 - Convolve input
 - Non-linearity (rectified linear)
 - Pooling (local max)
- Supervised training
- Train convolutional filters by back-propagating error
- Convolution over time



References

- ① Deep Learning with Python by Francois Chollet.
<https://livebook.manning.com/book/deep-learning-with-python/>
- ② Deep Learning by Ian Goodfellow, Yoshua Bengio, Aaron Courville
<https://www.deeplearningbook.org/>

Thank You!

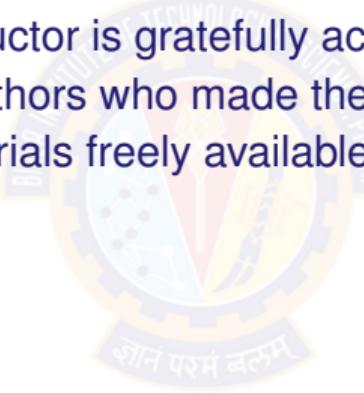


BITS Pilani
Pilani | Dubai | Goa | Hyderabad

DEEP LEARNING MODULE # 5 : CONVOLUTIONAL NEURAL NETWORK [CNN]

Seetha Parameswaran
Asst Prof, BITS Pilani

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.



IN THIS SEGMENT

1 CNN ARCHITECTURES

2 TRANSFER LEARNING

3 FURTHER READING



LENET 5

feature repr.

- Every convolutional layer includes three parts: convolution, pooling, and nonlinear activation functions.
- Using convolution to extract spatial features.
- Conv filters were 5x5, applied at stride 1.
- Subsampling average pooling layer. Subsampling (Pooling) layers were 2x2 applied at stride 2.
- tanh activation function.
- Using MLP as the last classifier.
- Architecture is [CONV-POOL-CONV-POOL-FC-FC]
- Sparse connection between layers to reduce the complexity of computation.

CA CA FC FC

LENET-5

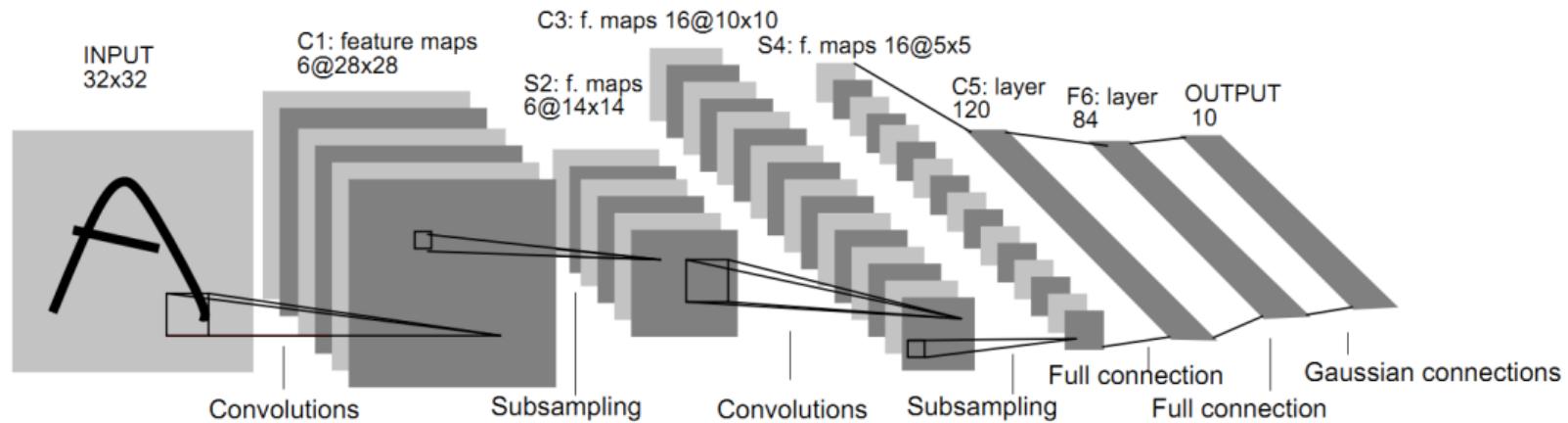
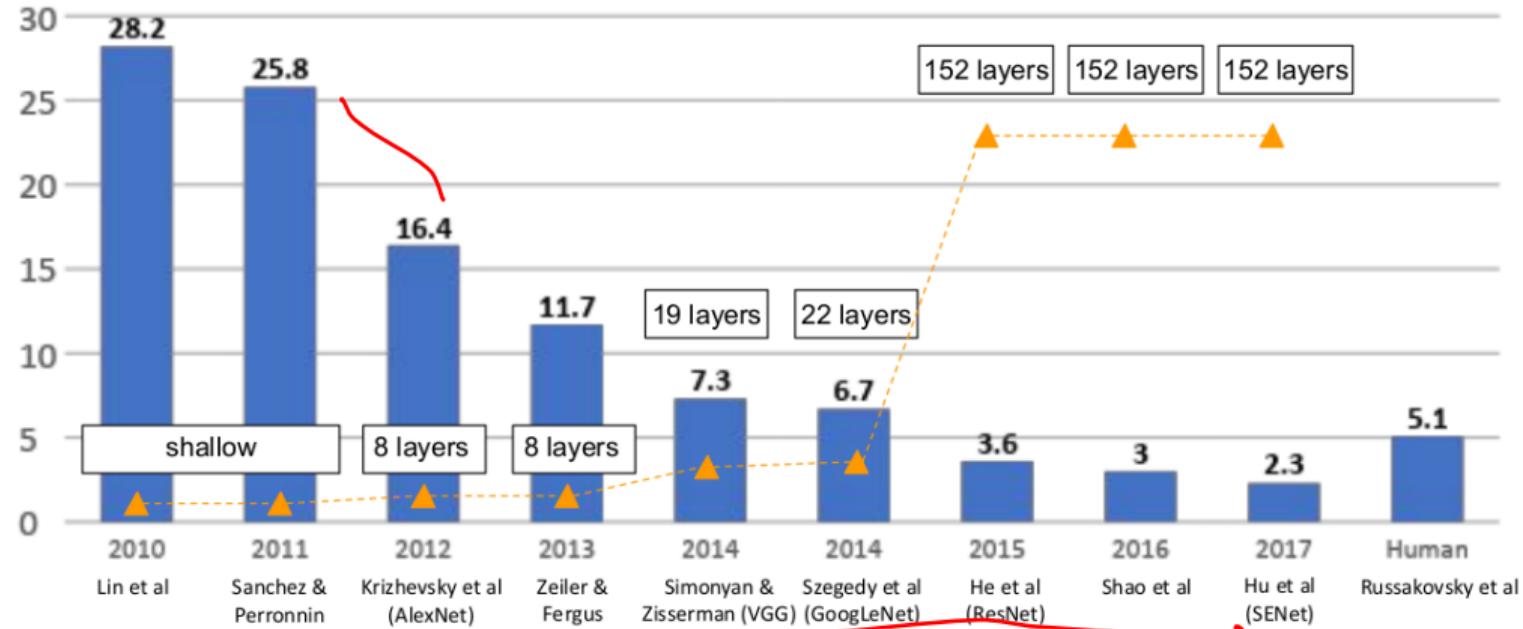


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

IMAGENET LARGE SCALE VISUAL RECOGNITION



IMAGENET DATA

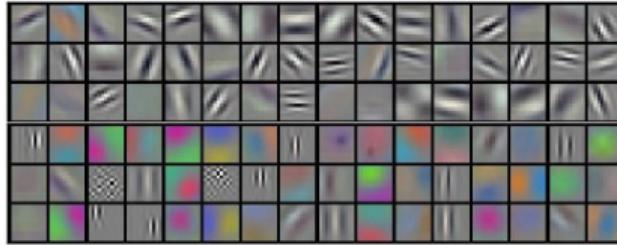


Figure 3: 96 convolutional kernels of size $11 \times 11 \times 3$ learned by the first convolutional layer on the $224 \times 224 \times 3$ input images. The top 48 kernels were learned on GPU 1 while the bottom 48 kernels were learned on GPU

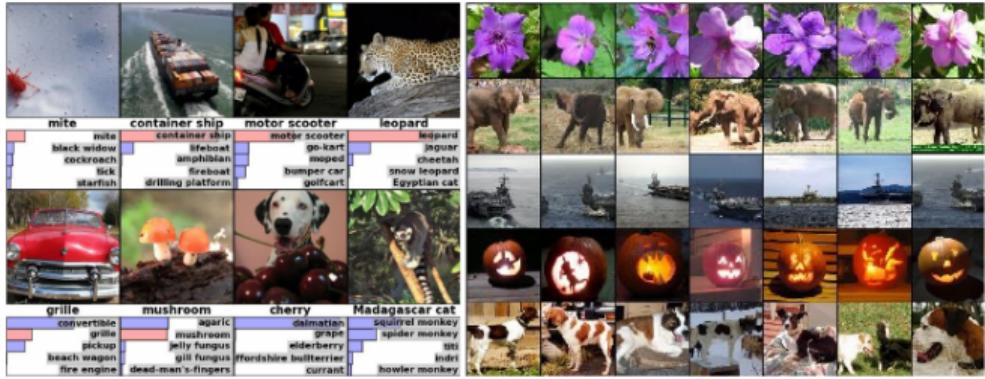


Figure 4: (Left) Eight ILSVRC-2010 test images and the five labels considered most probable by our model. The correct label is written under each image, and the probability assigned to the correct label is also shown with a red bar (if it happens to be in the top 5). (Right) Five ILSVRC-2010 test images in the first column. The remaining columns show the six training images that produce feature vectors in the last hidden layer with the smallest Euclidean distance from the feature vector for the test image.

IMAGENET DATA



The ImageNet set that was used has ~1.2 million images and 1000 classes

Accuracy is measured as top-5 performance:
Correct prediction if the true label matches one of the top 5 predictions of the model

ALEXNET

- AlexNet competed in the ImageNet Large Scale Visual Recognition Challenge on September 30, 2012.[3] The network achieved a top-5 error of 15.3%, more than 10.8 percentage points lower than that of the runner up.
- Details/Retrospectives:
 - ▶ first use of ReLU
 - ▶ used Local Response Normalisation (LRN) layers (not common anymore)
 - ▶ Approx. 60 million parameters are learned.
 - ▶ heavy data augmentation
 - ▶ dropout 0.5
 - ▶ batch size 128
 - ▶ SGD Momentum 0.9
 - ▶ Learning rate $1e - 2$, reduced by 10 manually when val accuracy plateaus
 - ▶ L2 weight decay $5e - 4$
 - ▶ 7 CNN ensemble: 18.2% – > 15.4%

ALEXNET ARCHITECTURE

227x227x3 INPUT

55x55x96 CONV1 : 96 11x11 filters at stride 4, pad 0

27x27x96 MAX POOL1 : 3x3 filters at stride 2

27x27x96 NORM1 : Normalization layer

27x27x256 CONV2 : 256 5x5 filters at stride 1, pad 2

13x13x256 MAX POOL2 : 3x3 filters at stride 2

13x13x256 NORM2 : Normalization layer

13x13x384 CONV3 : 384 3x3 filters at stride 1, pad 1

13x13x384 CONV4 : 384 3x3 filters at stride 1, pad 1

13x13x256 CONV5 : 256 3x3 filters at stride 1, pad 1

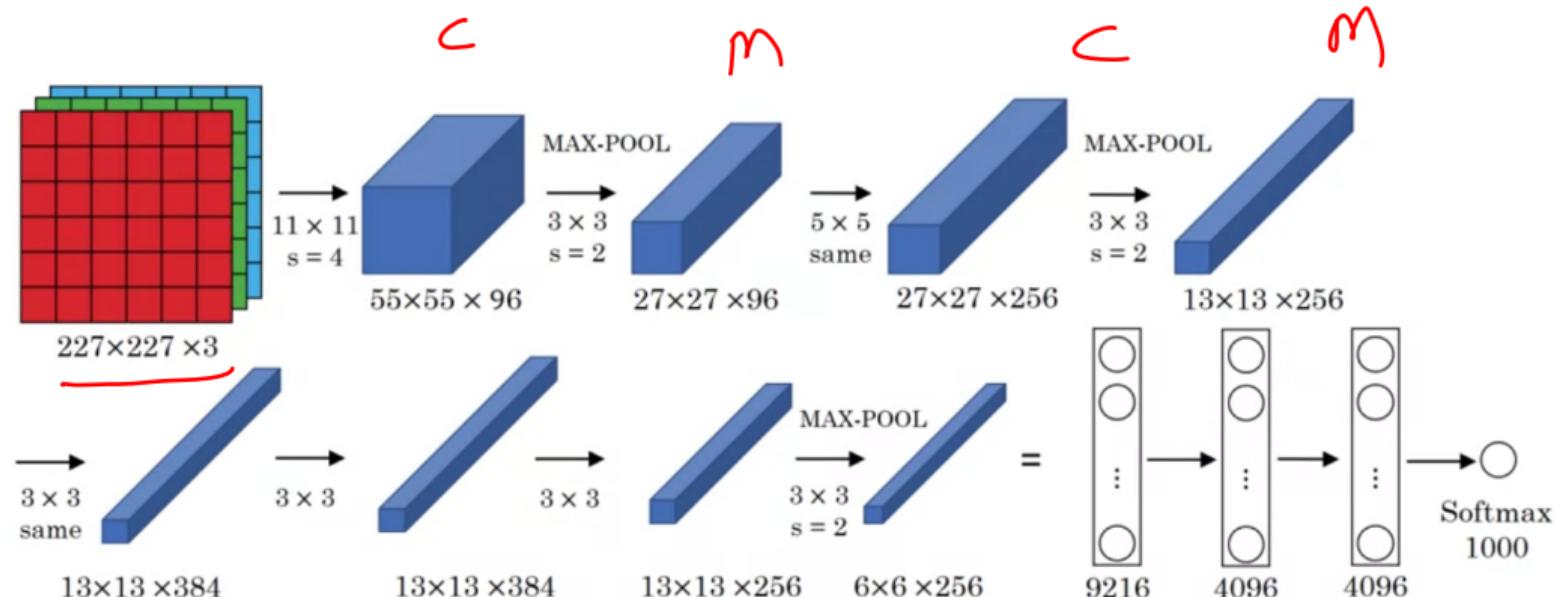
6x6x256 MAX POOL3 : 3x3 filters at stride 2

4096 FC6 : 4096 neurons

4096 FC7 : 4096 neurons

1000 FC8 : 1000 neurons (class scores)

ALEXNET



[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

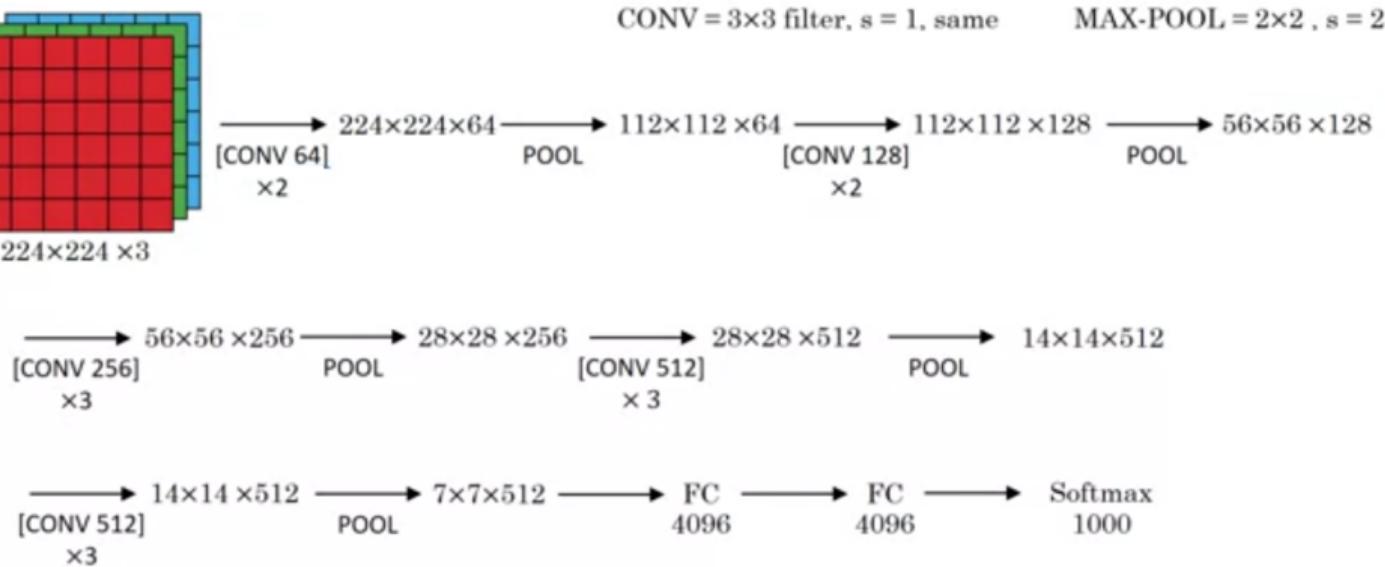
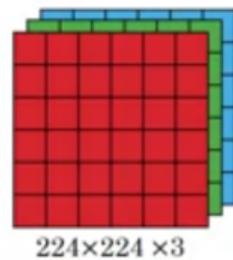
Andrew Ng

VGG 16

- VGG stands for Visual Geometry Group with 16 Layers
- Plug and play in Caffe
- Deeper the better
- Details
 - ▶ ILSVRC'14 2nd in classification, 1st in localization
 - ▶ Similar training procedure as Krizhevsky 2012
 - ▶ Approx. 138 million parameters are learned.
 - ▶ No Local Response Normalisation (LRN)
 - ▶ Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
 - ▶ Use ensembles for best results
 - ▶ All convolutions with a 3×3 kernel
 - ▶ All max-pooling layers with a 2×2 kernel
 - ▶ FC7 features generalize well to other tasks

VGG 16

CCM CCM CCM CCM CCM FC FC FC



[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]

Andrew Ng

NETWORK IN NETWORK

- A convolution kernel can be thought of as a generalized linear model (GLM).
- Using a sophisticated non-linear function approximator like MLP, enhance the abstraction ability of the local model.
- Replace GLM by **micro network**.
- Replace the MLP micro structure via convolutions.
- Replace fully connected layers in the last layers by global average pooling.

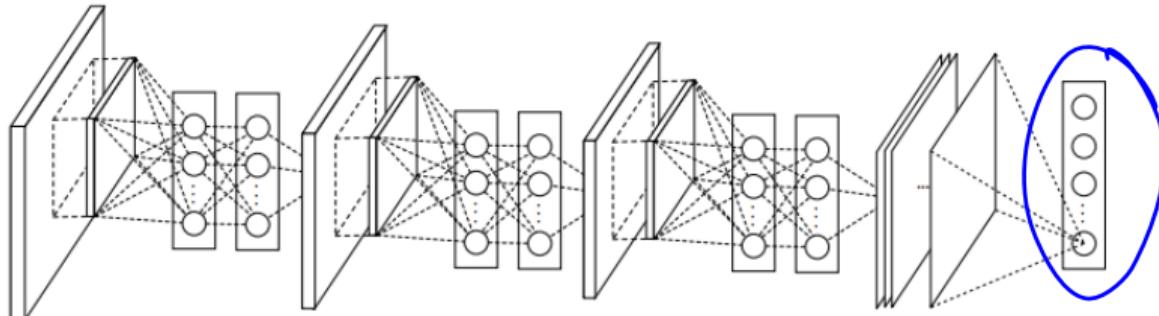
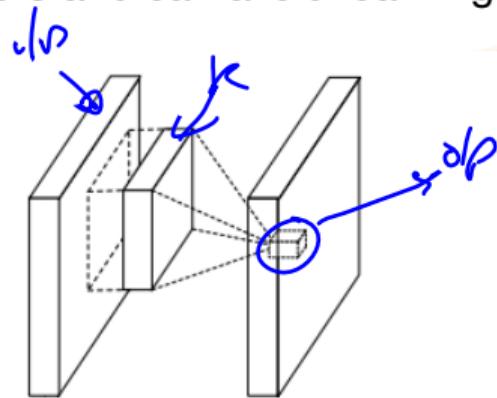


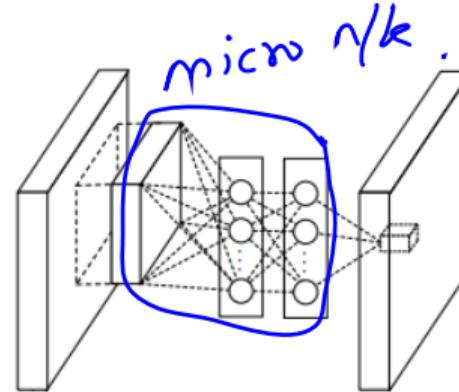
Figure 2: The overall structure of Network In Network. In this paper the NINs include the stacking of three mlpconv layers and one global average pooling layer.

NETWORK IN NETWORK

Using micro-networks allow to extract more sophisticated features and may need fewer extractors and can avoid learning too simple or redundant features.



(a) Linear convolution layer



(b) Mlpconv layer

Figure 1: Comparison of linear convolution layer and mlpconv layer. The linear convolution layer includes a linear filter while the mlpconv layer includes a micro network (we choose the multilayer perceptron in this paper). Both layers map the local receptive field to a confidence value of the latent concept.

NETWORK IN NETWORK

Fully connected layers have a lot of learnable parameters and may cause overfitting.
Replace FC by global average pooling for better generalization.

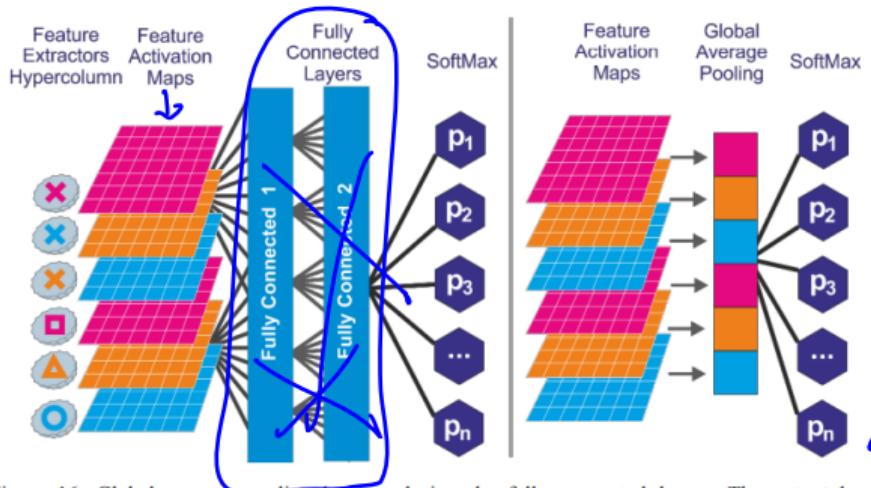


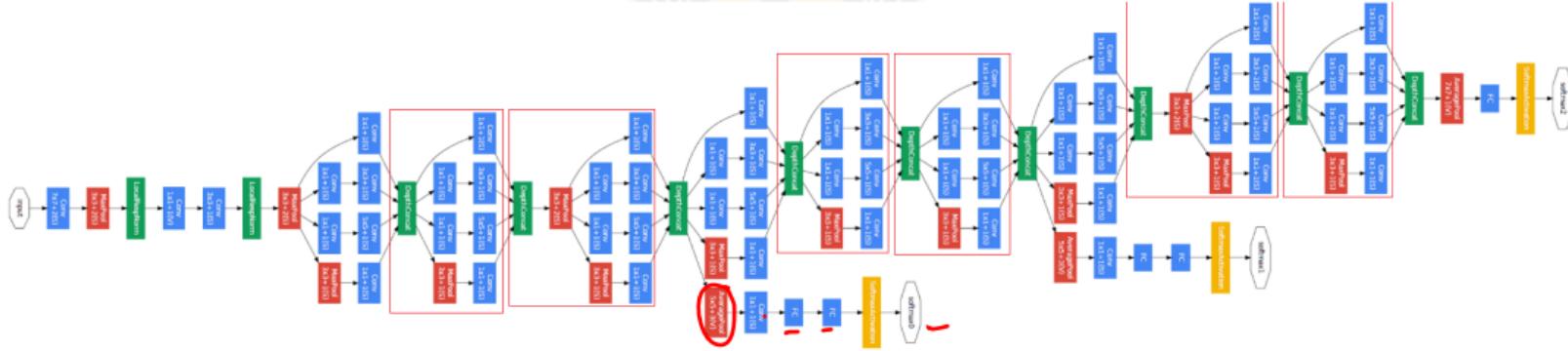
Figure 16: Global average pooling layer replacing the fully connected layers. The output layer implements a Softmax operation with p_1, p_2, \dots, p_n the predicted probabilities for each class.

Figure Source: Singh, Anshuman Vikram. "Content-based image retrieval using deep learning." (2015)

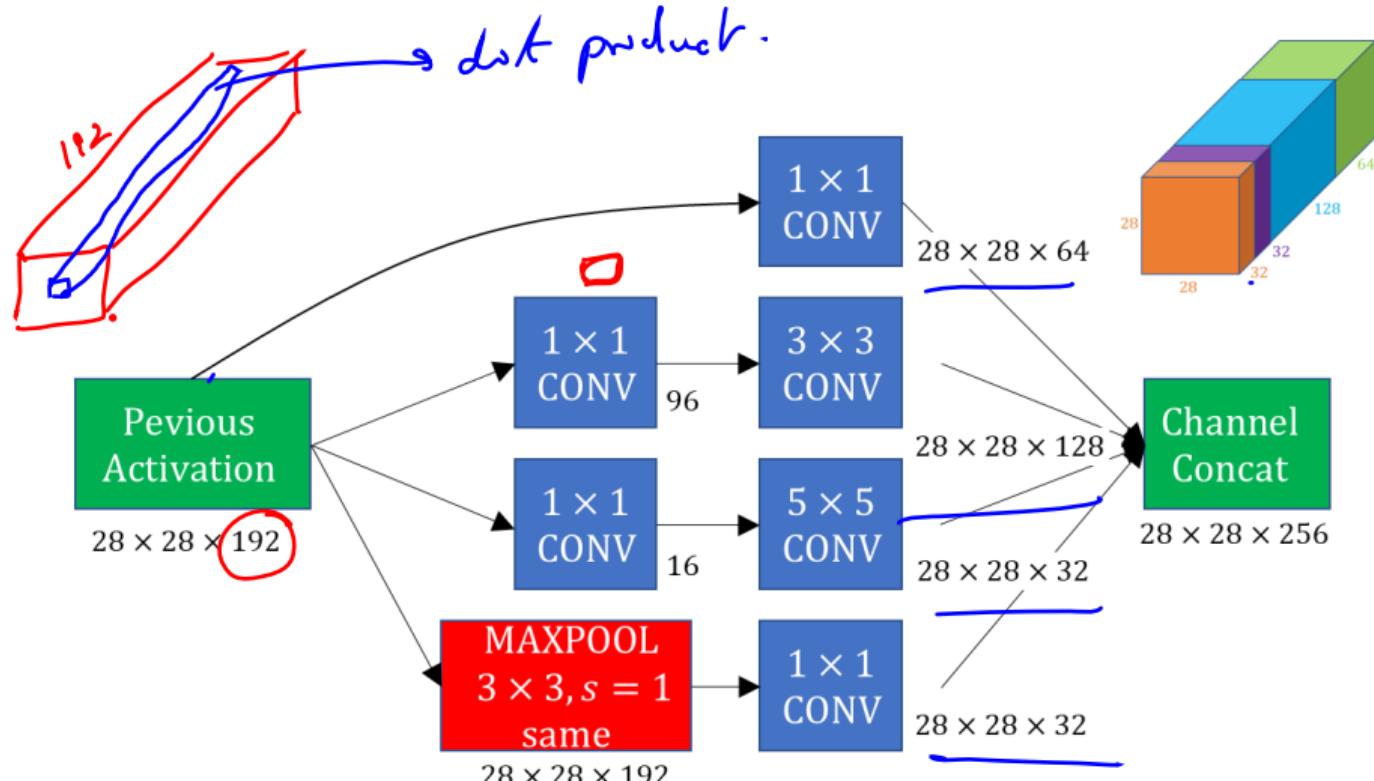
GOOGLENET / INCEPTION NET

- ILSVRC'14 classification winner (6.7% top 5 error)
- 22 layers with weights
- Only 5 million parameters (12x less than AlexNet and 27x less than VGG-16)
- Inception Module - convolutional “blocks” – efficient
 - ▶ Design a good local network topology (network within a network) and then stack these modules on top of each other.
- Linear layers at the end
- Max pooling in between, multiple Conv layers between pooling
- Great ideas for data augmentation
- Deeper networks, with computational efficiency
- No FC layers.
- After the last convolutional layer, a global average pooling layer is used that spatially averages across each feature map, before final FC layer.

INCEPTION NET



INCEPTION NET MODULE



32 filters, $1 \times 1 \times 192$

Andrew Ng

INCEPTION NET MODULE

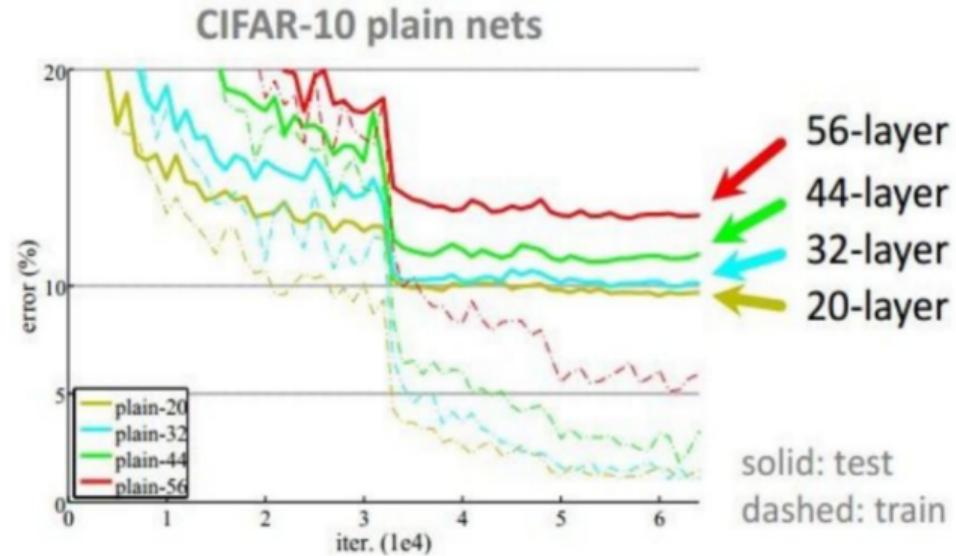
- Apply parallel filter operations on the input from previous layer.
- Multiple receptive field sizes for convolution (1×1 , 3×3 , 5×5) and Pooling operation (3×3).
- Concatenate all filter outputs together channel-wise.
- Very expensive compute
- 1×1 convolutions to reduce feature channel size.
 - ▶ Assume an input of $W \times H \times 64$.
 - ▶ Each filter has size $1 \times 1 \times 64$, and performs a 64-dimensional dot product.
 - ▶ Preserves spatial dimensions, reduces depth.
 - ▶ Projects depth to lower dimension (combination of feature maps).

RESIDUAL NET (RESNET)

- Very deep networks using residual connections
- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15.
- Add residual connections between convolution blocks
- Easy gradient flow for deep models
- Enable very deep models (from 20 layers to 150 and even 1000 layers!)
- No linear layers except classifier.
- Solve the problems of vanishing gradients using Residual connections.

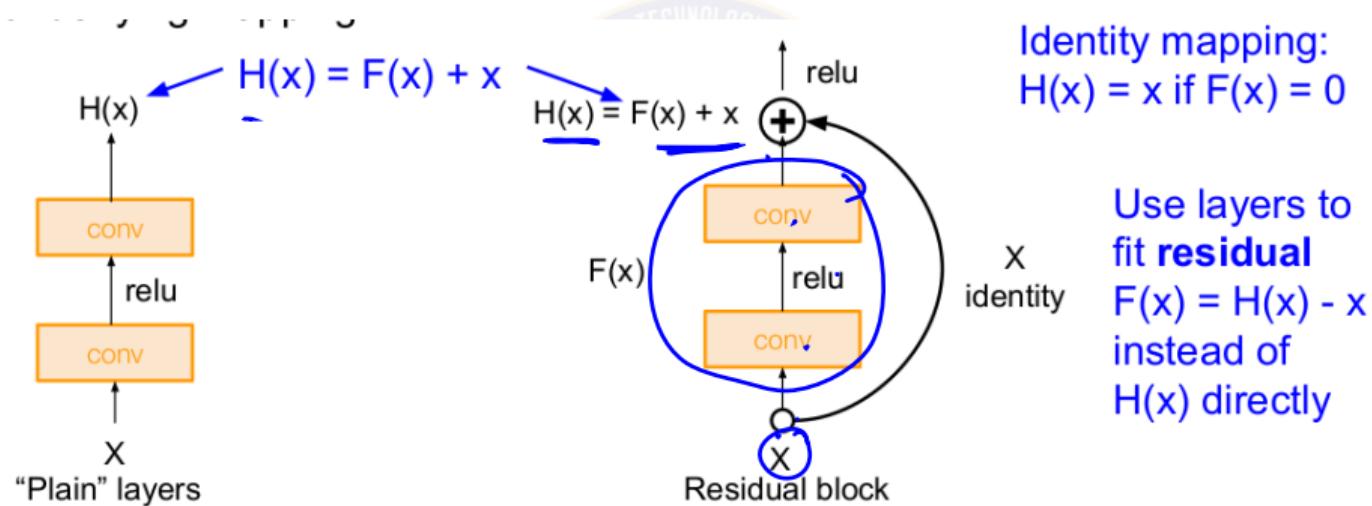
RESNET BLOCK

- What happens when we continue stacking deeper layers on a convolutional neural network?
- The deeper model performs worse, but it's not caused by overfitting.

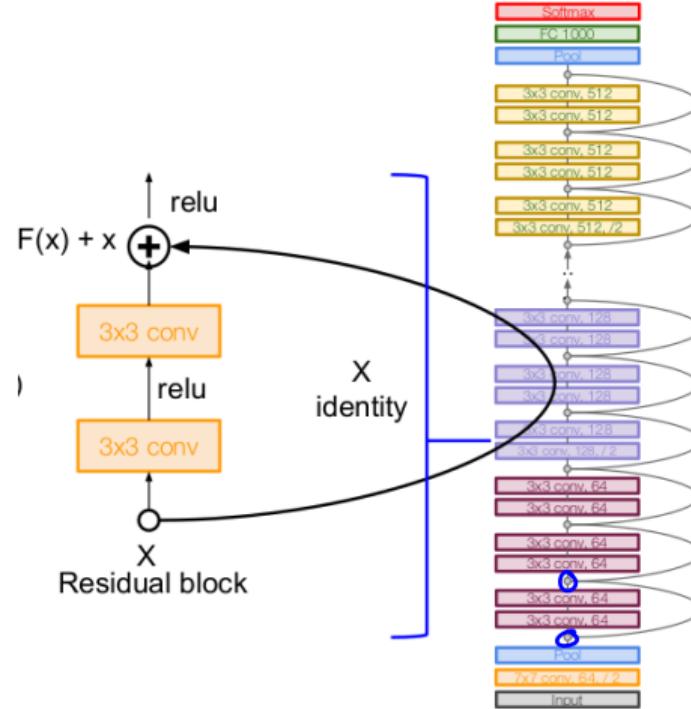


RESNET BLOCK

- Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping.



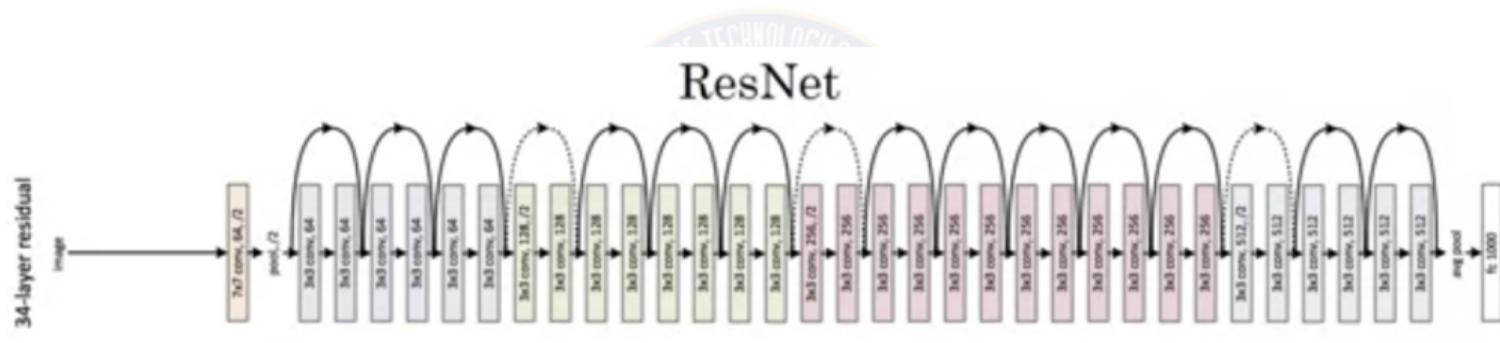
RESNET



RESNET MODULES

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	3×3 max pool, stride 2				
		$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

RESNET



[He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

RESNET ARCHITECTURE

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2.
Reduce the activation volume by half.
- Additional conv layer at the beginning (stem).
- No FC layers at the end (only FC 1000 to output classes)
- Total depths of 18, 34, 50, 101, or 152 layers for ImageNet.
- Training ResNet in practice:
 - ▶ Batch Normalization after every CONV layer
 - ▶ Xavier initialization from He et al.
 - ▶ SGD + Momentum (0.9)
 - ▶ Learning rate: 0.1, divided by 10 when validation error plateaus
 - ▶ Mini-batch size 256
 - ▶ Weight decay of 1e-5
 - ▶ No dropout used

EFFICACY OF RESNET

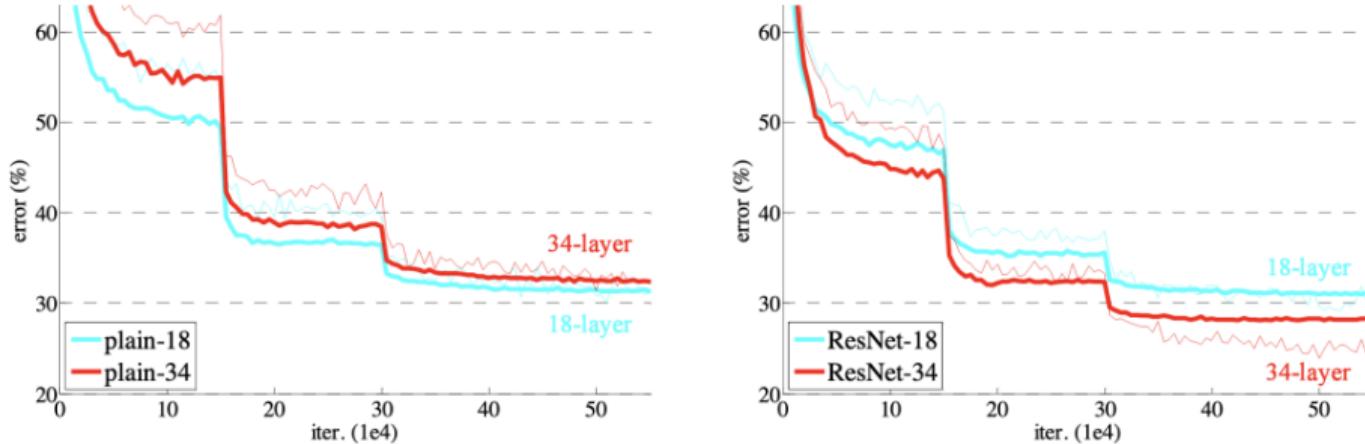
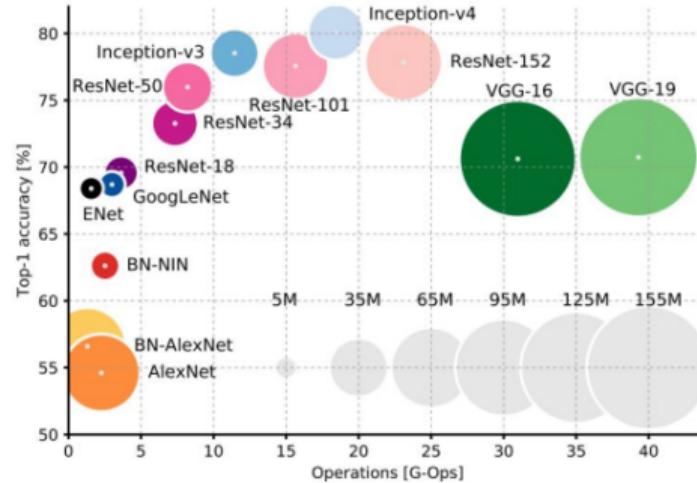
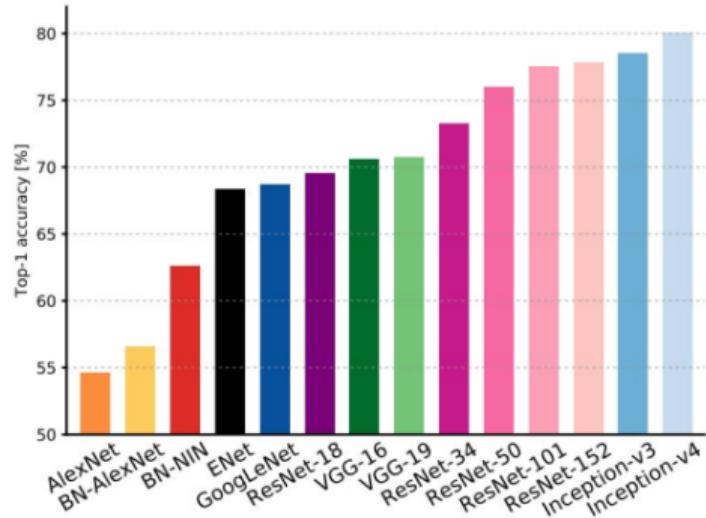


Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

COMPARING COMPLEXITY



COMPARING COMPLEXITY

- Inception-v4: Resnet + Inception – more accuracy
- VGG: most parameters, most operations
- GoogLeNet: most efficient
- AlexNet: Smaller compute, still memory heavy, lower accuracy
- ResNet: Moderate efficiency depending on model, highest accuracy

IN THIS SEGMENT

1 CNN ARCHITECTURES

2 TRANSFER LEARNING

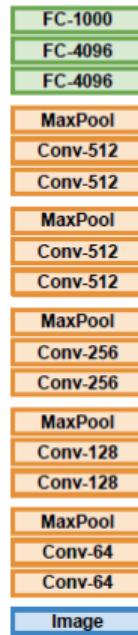
3 FURTHER READING



TRANSFER LEARNING

pre trained w/o .
Transfer Learning with CNNs

1. Train on Imagenet



10⁸

n/k trained on
AlexNet
VGG 16
Yolo v3
Imagenet
weights

h5

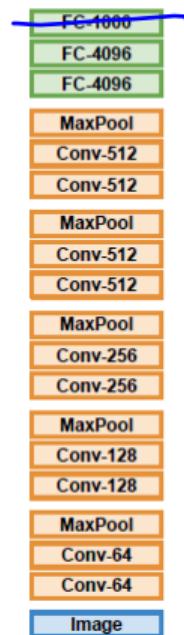
date <--> ImageNet-
Colo.
⋮

- ① My data
<<< ImageNet-
Colo
- ②. c <<< Imagenet-
- ③ Pytek

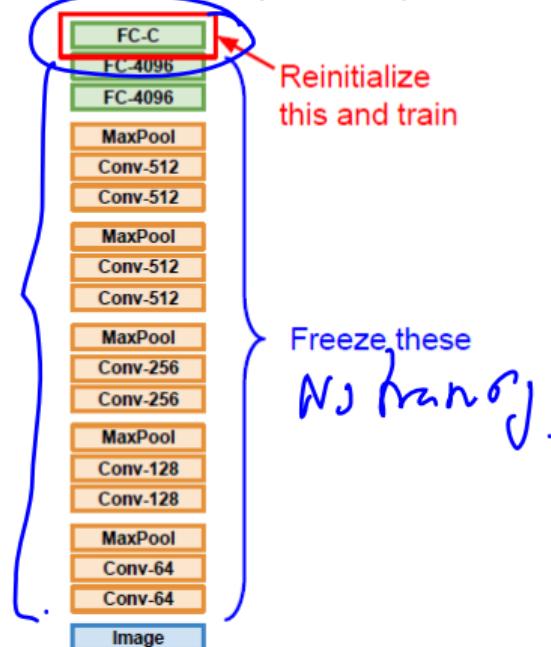
TRANSFER LEARNING

Transfer Learning with CNNs

1. Train on Imagenet



2. Small Dataset (C classes)



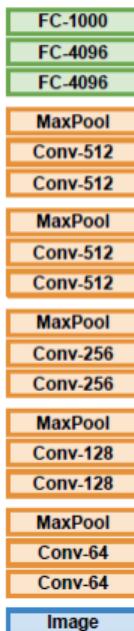
Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

pretrained
↓ .
feature repn
for my data .
↓
FC .
↓

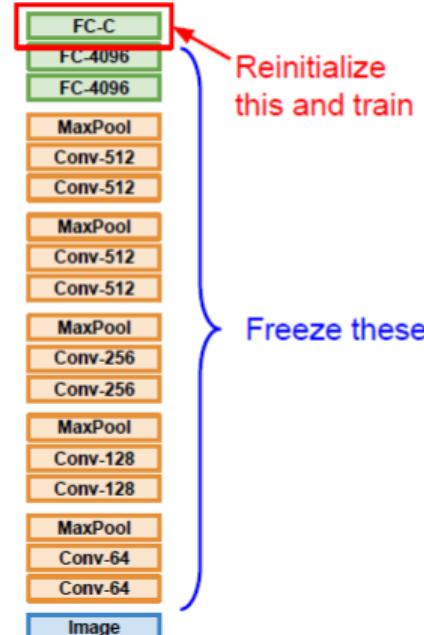
TRANSFER LEARNING

Transfer Learning with CNNs

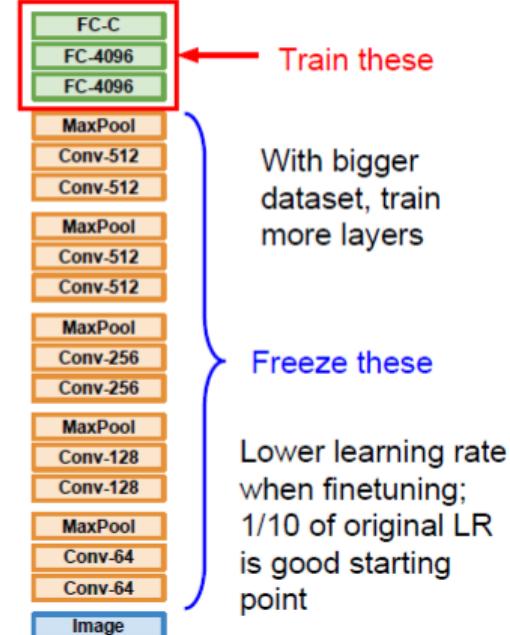
1. Train on Imagenet



2. Small Dataset (C classes)

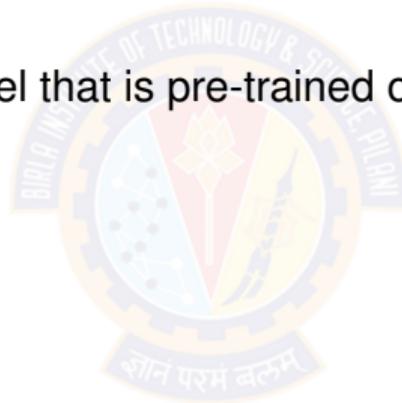


3. Bigger dataset



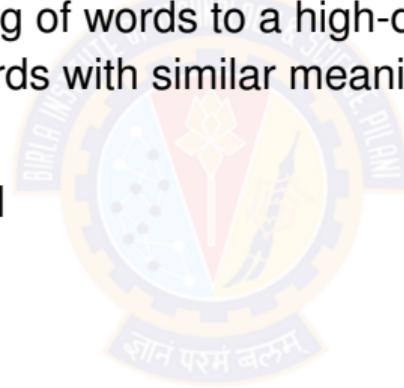
TRANSFER LEARNING FOR IMAGE DATA

- Use a deep learning model that is pre-trained on large dataset like ImageNet or MS Coco.
- Oxford VGG Model
- Google Inception Model
- Microsoft ResNet Model



TRANSFER LEARNING FOR TEXT DATA

- Embedding is the mapping of words to a high-dimensional continuous vector space where different words with similar meanings have similar vector representations.
- Google's word2vec Model
- Stanford's GloVe Model
- FastText
- Gensim

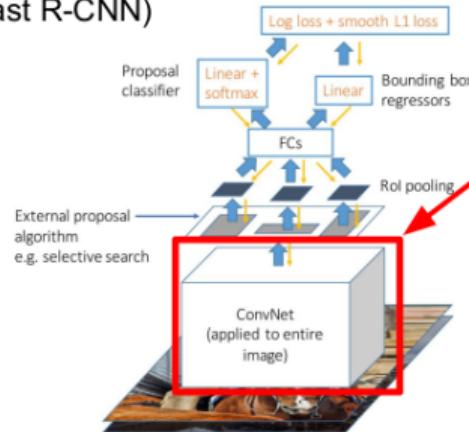


TRANSFER LEARNING - WHEN TO USE?

- You need a lot of data if you want to train/use CNNs / RNNs.
- Task A and Task B have the same type of input. Eg: Input is images for both tasks.
- We have lot of data for training Task A and relatively low data for training in Task B.
- Low level features obtained from Task A could be more helpful for learning Task B.

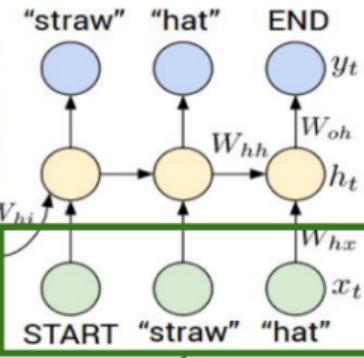
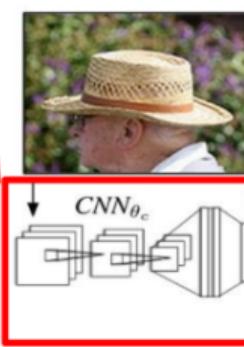
TRANSFER LEARNING EXAMPLE

Object Detection
(Fast R-CNN)



CNN pretrained
on ImageNet

Image Captioning: CNN + RNN



Word vectors pretrained
with word2vec

Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

IN THIS SEGMENT

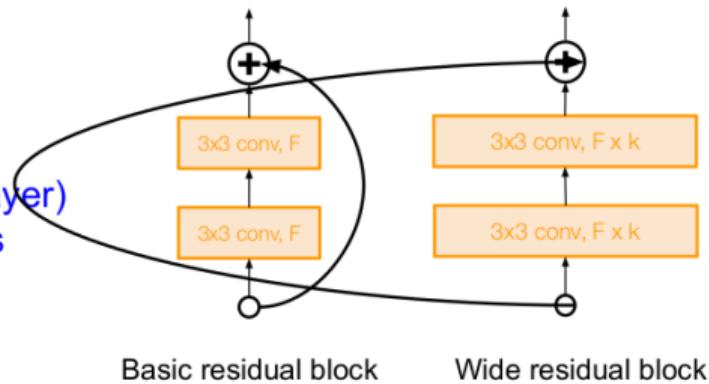
- 1 CNN ARCHITECTURES
- 2 TRANSFER LEARNING
- 3 FURTHER READING



Wide Residual Networks

[Zagoruyko et al. 2016]

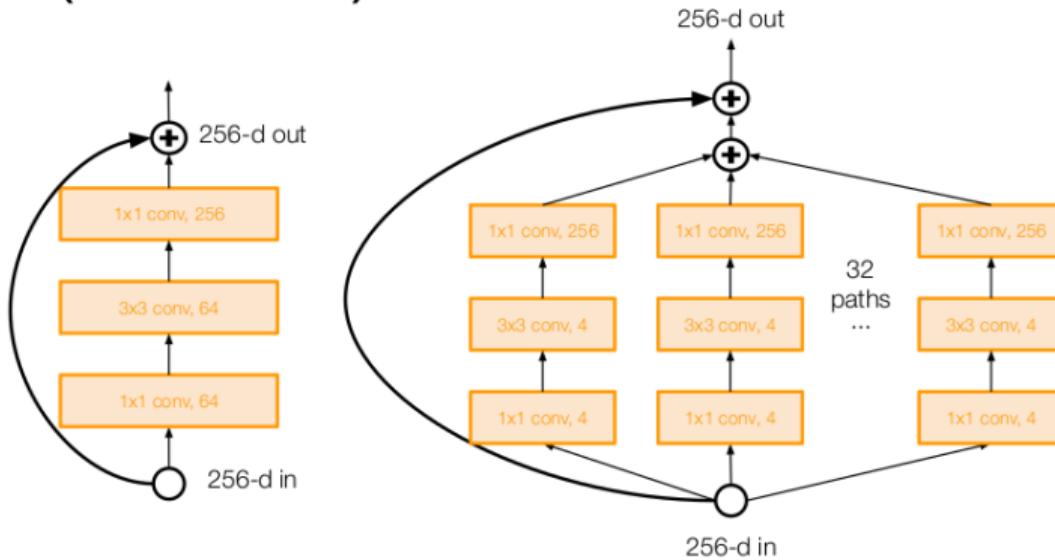
- Argues that residuals are the important factor, not depth
- Use wider residual blocks ($F \times k$ filters instead of F filters in each layer)
- 50-layer wide ResNet outperforms 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)



Aggregated Residual Transformations for Deep Neural Networks (ResNeXt)

[Xie et al. 2016]

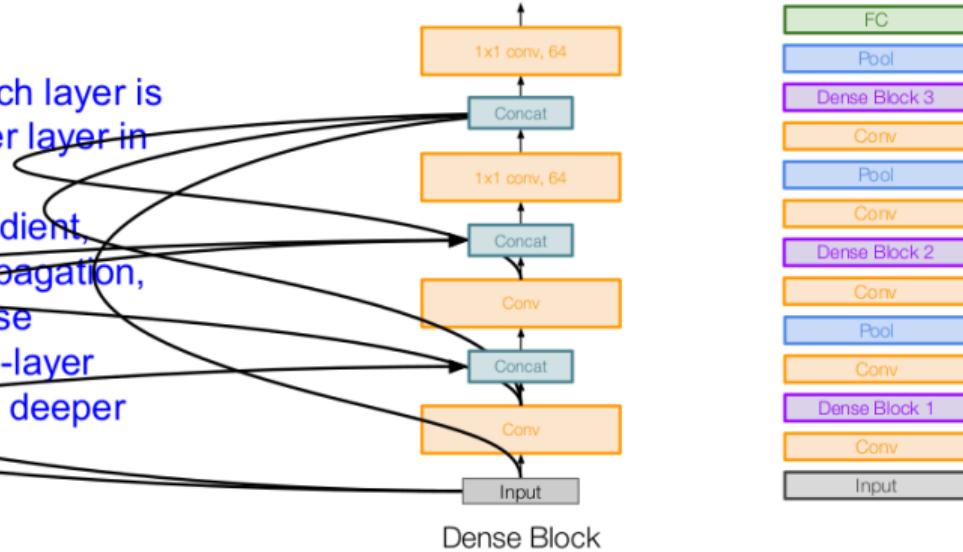
- Also from creators of ResNet
- Increases width of residual block through multiple parallel pathways (“cardinality”)
- Parallel pathways similar in spirit to Inception module



Densely Connected Convolutional Networks (DenseNet)

[Huang et al. 2017]

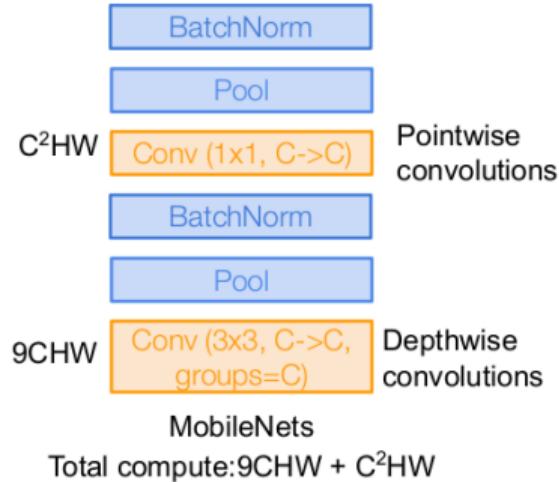
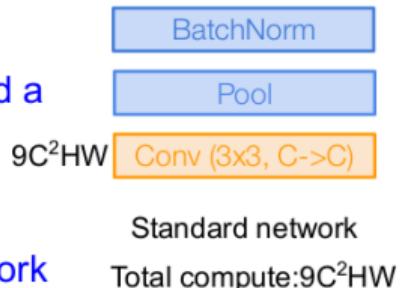
- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse
- Showed that shallow 50-layer network can outperform deeper 152 layer ResNet



MobileNets: Efficient Convolutional Neural Networks for Mobile Applications

[Howard et al. 2017]

- Depthwise separable convolutions replace standard convolutions by factorizing them into a depthwise convolution and a 1×1 convolution
- Much more efficient, with little loss in accuracy
- Follow-up MobileNetV2 work in 2018 (Sandler et al.)
- ShuffleNet: Zhang et al., CVPR 2018



EfficientNet: Smart Compound Scaling

[Tan and Le, 2019]

- Increase network capacity by scaling width, depth, and resolution, while balancing accuracy and efficiency.
- Search for optimal set of compound scaling factors given a compute budget (target memory & flops).
- Scale up using smart heuristic rules

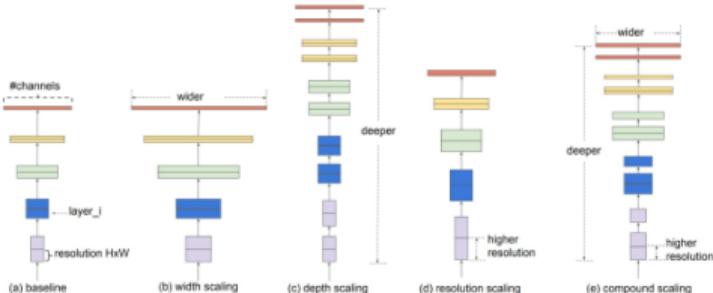
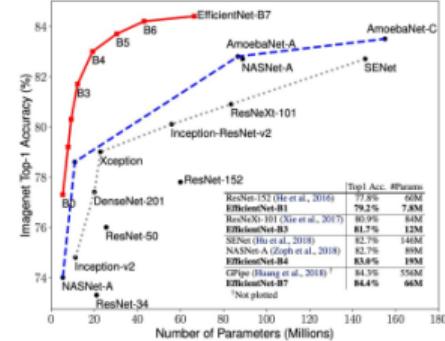
$$\text{depth: } d = \alpha^\phi$$

$$\text{width: } w = \beta^\phi$$

$$\text{resolution: } r = \gamma^\phi$$

$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$



REFERENCES

Refer Chapter 8 of Dive into Deep Learning book.

LENET5 : Yann Lecun and Leon Bottou and Yoshua Bengio and Patrick Haffner,
Gradient-based learning applied to document recognition, IEEE-1998

ALEXNET : Imagenet classification with deep convolutional neural networks,
Krizhevsky, Alex and Sutskever, Ilya and Hinton, Geoffrey E, In: Advances in
neural information processing systems, NIPS 2012

VGG16 : Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional
networks for large-scale image recognition."

INCEPTION v1 : Going deeper with convolutions, Szegedy, Christian et.al CVPR 2015

RESNET : Deep residual learning for image recognition, He et.al CVPR-2016



Thank You!



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

DEEP LEARNING

MODULE # 6 : RECURRENT NEURAL NETWORK [RNN]

Seetha Parameswaran
Asst Prof, BITS Pilani

G ✓
Gl ✓

good ✓

good morning

glass

glad ✓

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.

Gle. ✓

history.

perdition

varying op.

glee

glt

varying
not freed

dependency

Prediction

letter
words
phrases

DNN

CNN

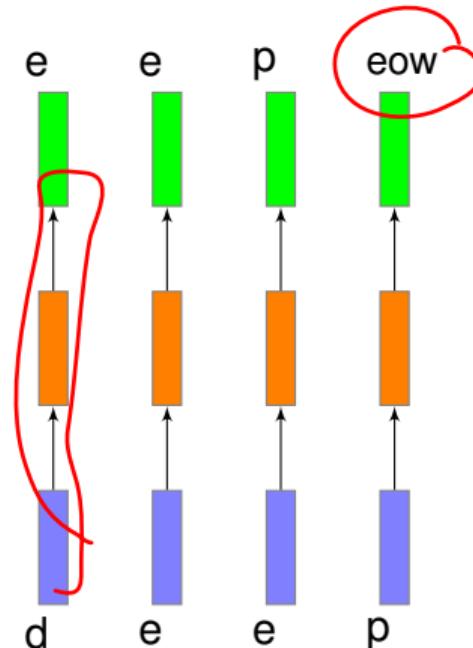
In feedforward and convolutional neural networks

- The size of the input is always fixed.
- Each input to the network is independent of the previous or future inputs.
- The computations, outputs and decisions for two successive inputs / images are completely independent of each other.

This is not true in many applications.

- The size of the input is not always fixed.
- Successive inputs may not be independent of each other.
- Each network (blue - orange - green structure) is performing the same task – input : character output : character.

Example: Auto-completion.



IN THIS SEGMENT

- ① SEQUENCE LEARNING
- ② RECURRENT NEURAL NETWORK (RNN)
- ③ TYPES OF RNN
- ④ LEARNING IN RNN
- ⑤ ISSUES IN RNN
- ⑥ LONG SHORT TERM MEMORY UNIT (LSTM)
- ⑦ GATED RECURRENT UNIT (GRU)
- ⑧ BIDIRECTIONAL RNN



SEQUENCE LEARNING PROBLEMS

To model a sequence we need

- Process an input or sequence of inputs.
- The inputs may have be dependent.
- We may have to maintain the sequence order.
- Each input corresponds to one time step.
- Keep track of long term dependencies.
- Produce an output or sequence of outputs.
- Supervised Learning.
- Share parameters across the sequences.

SEQUENCE MODEL

I/o

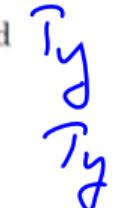
o/p

Speech recognition



text

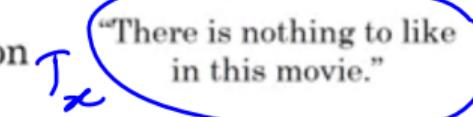
"The quick brown fox jumped over the lazy dog."



Music generation



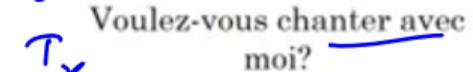
Sentiment classification



DNA sequence analysis



Machine translation



Do you want to sing with
me?



Video activity recognition



Running



Name entity recognition

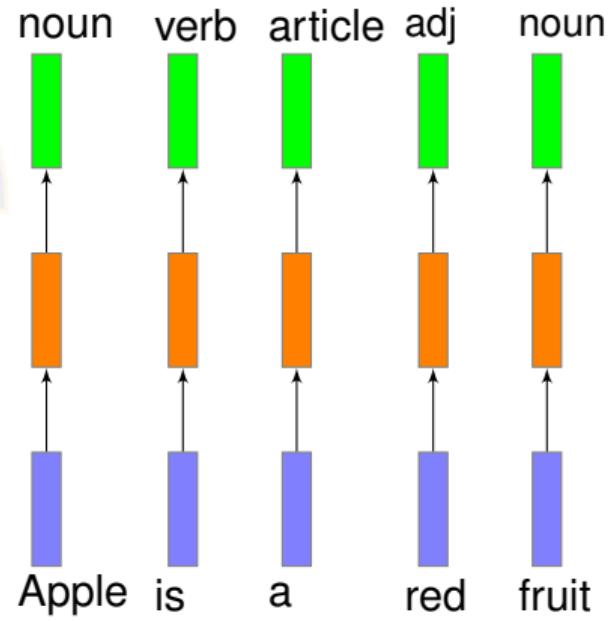


Yesterday, Harry Potter
met $\textcolor{red}{Hermione}$ Granger.



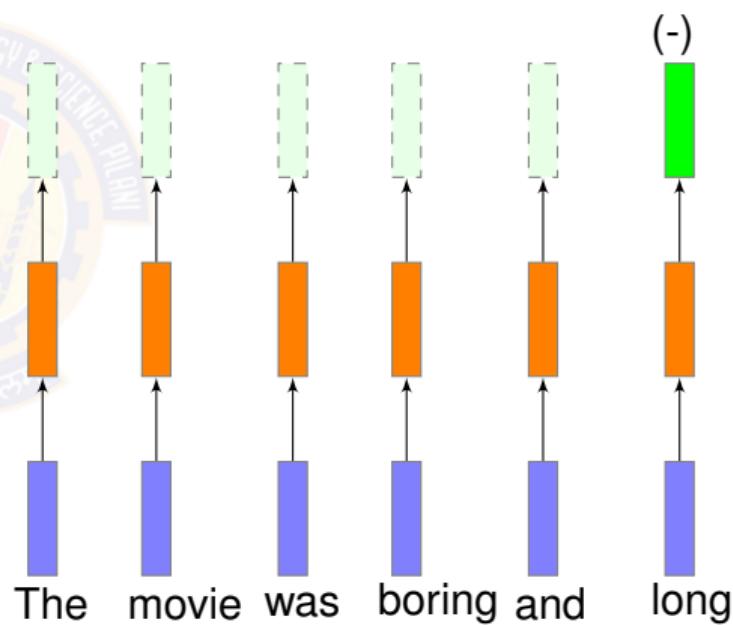
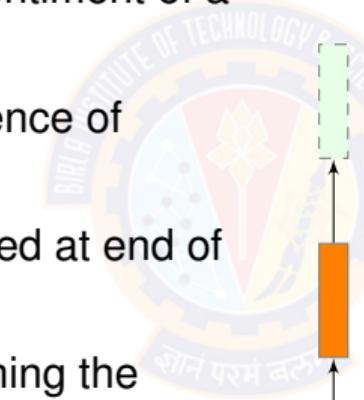
PART OF SPEECH TAGGING

- Task is predicting the part of speech tag (noun, adverb, adjective, verb) of each word in a sentence.
- When we see an adjective we are almost sure, the next word should be a noun.
- The current output depends on the current input as well as the previous input.
- The size of the input is not fixed. Sentences have any number of words.
- An output is produced at end of each time step.
- Each network is performing the same task – input : word, output : tag.



SENTIMENT ANALYSIS

- Task is predicting the sentiment of a whole sentence.
- Input is the entire sequence of inputs.
- An output is **not** produced at end of each time step.
- Each network is performing the same task – input : word, output : polarity +/–.

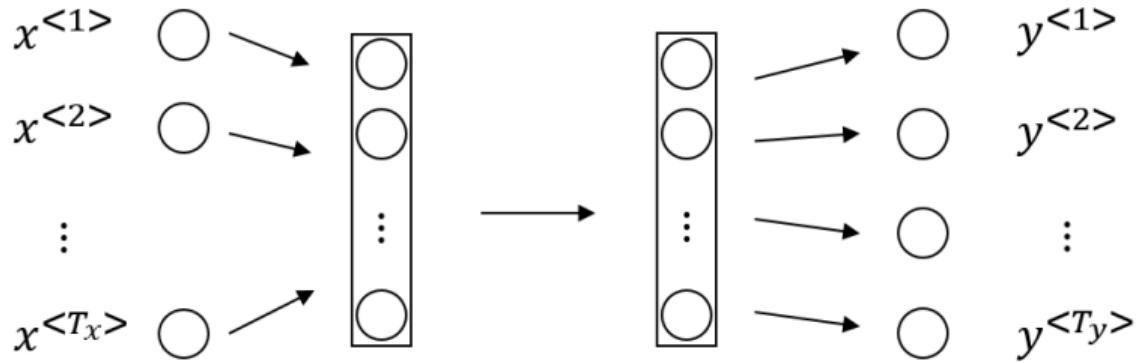


IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN



RECURRENT NEURAL NETWORK (RNN)



Andrew Ng

- Accounts for variable number of inputs.
- Accounts for dependencies between inputs.
- Accounts for variable number of outputs.
- Ensures that the same function executed at each time step.
- The features learned across the inputs at different time step has to be shared.

RNN I

- The function learned at each time step.

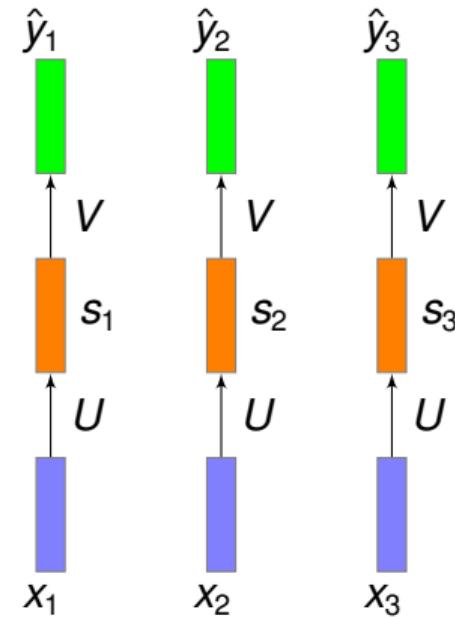
$t = \text{time step}$

$x_t = \text{input at time step } t$

$s_t = \sigma(Ux_t + b)$

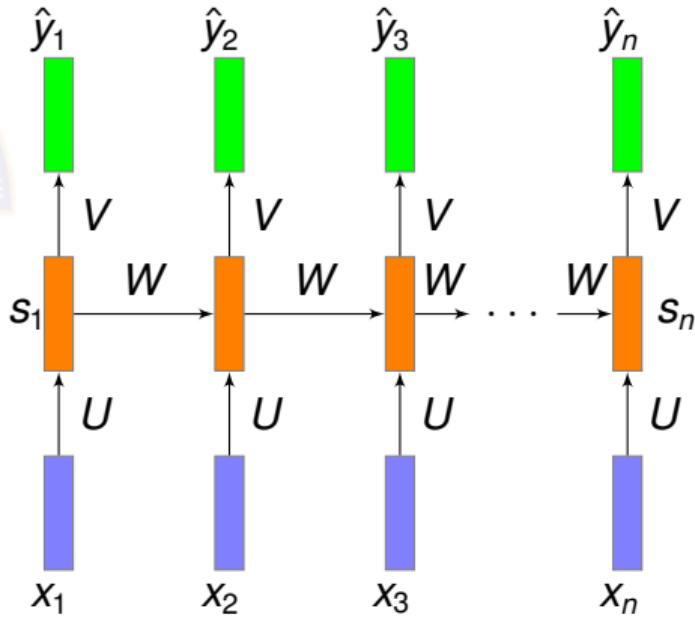
$y_t = g(Vs_t + c)$

- Since the same function has to be executed at each time step we should share the same network i.e., same parameters at each time step.



RNN II

- The parameter sharing ensures that
 - ▶ the network becomes invariant to the length of the input.
 - ▶ the number of time steps doesn't matter.
- Create multiple copies of the network and execute them at each timestep.
 - ▶ i.e. create a loop effect.
 - ▶ i.e. add recurrent connection in the network.



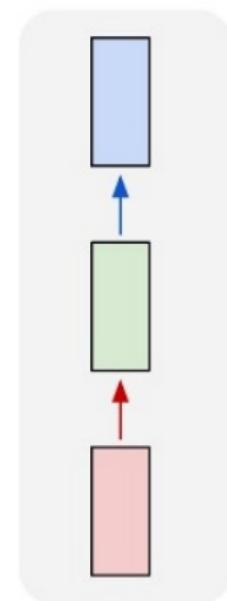
IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN

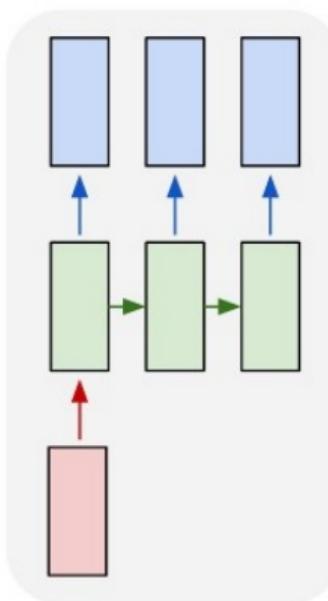


TYPES OF RNN

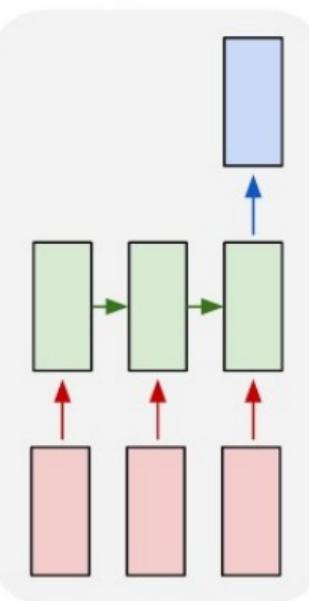
one to one



one to many

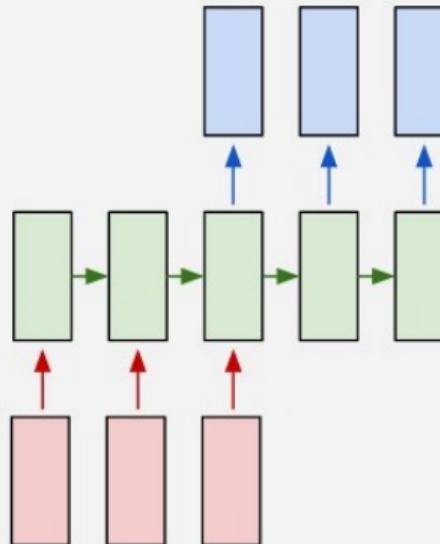


many to one

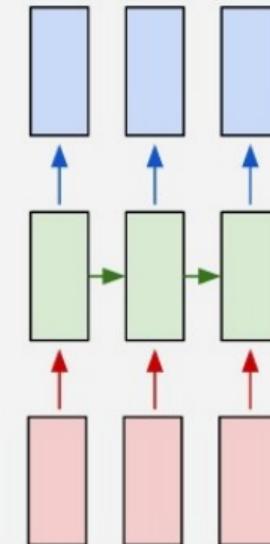


TYPES OF RNN

many to many



many to many



TYPES OF RNN AND APPLICATIONS

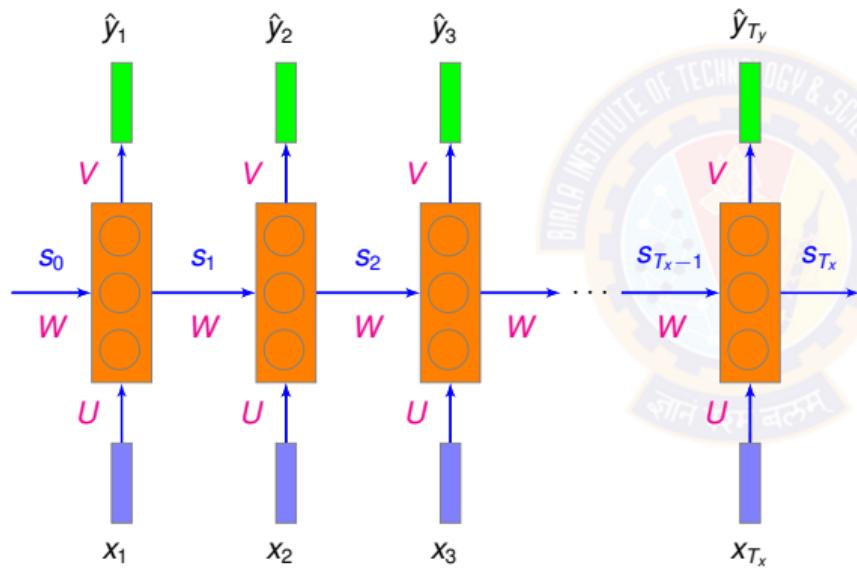
- One to one – Generic neural network, Image classification
- One to many – Music generation, Image Captioning
- Many to one – Movie review or Sentiment Analysis
- Many to many – Machine translation
- Synced Many to many – Video classification

IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN



FORWARD PROPAGATION IN RNN



- s_t is the **state** of the network at time step t .

$$s_0 = 0$$

$$s_t = \sigma(Ux_t + Ws_{t-1} + b)$$

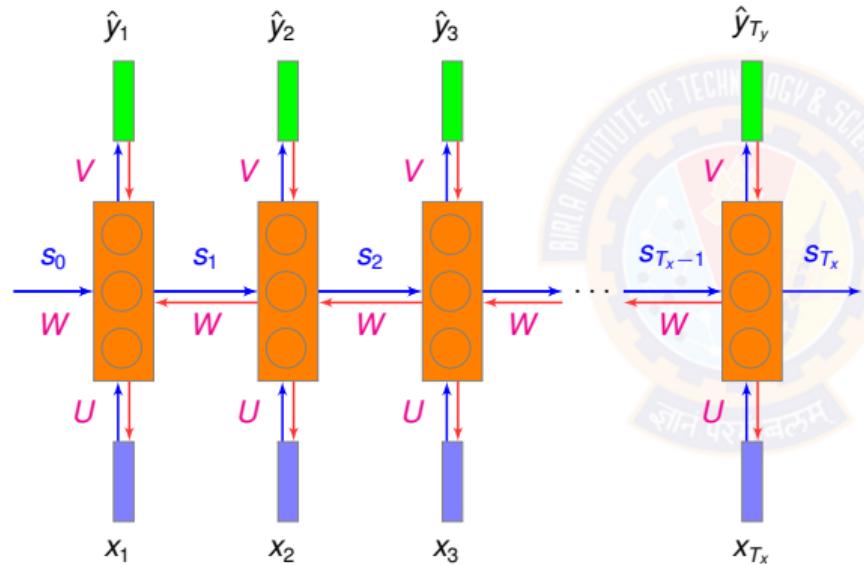
$$\hat{y}_t = g(Vs_t + c)$$

or

$$\hat{y}_t = f(x_t, s_{t-1}, W, U, V, b, c)$$

- The parameters W, U, V, b, c are shared across time steps.

BACK PROPAGATION IN RNN



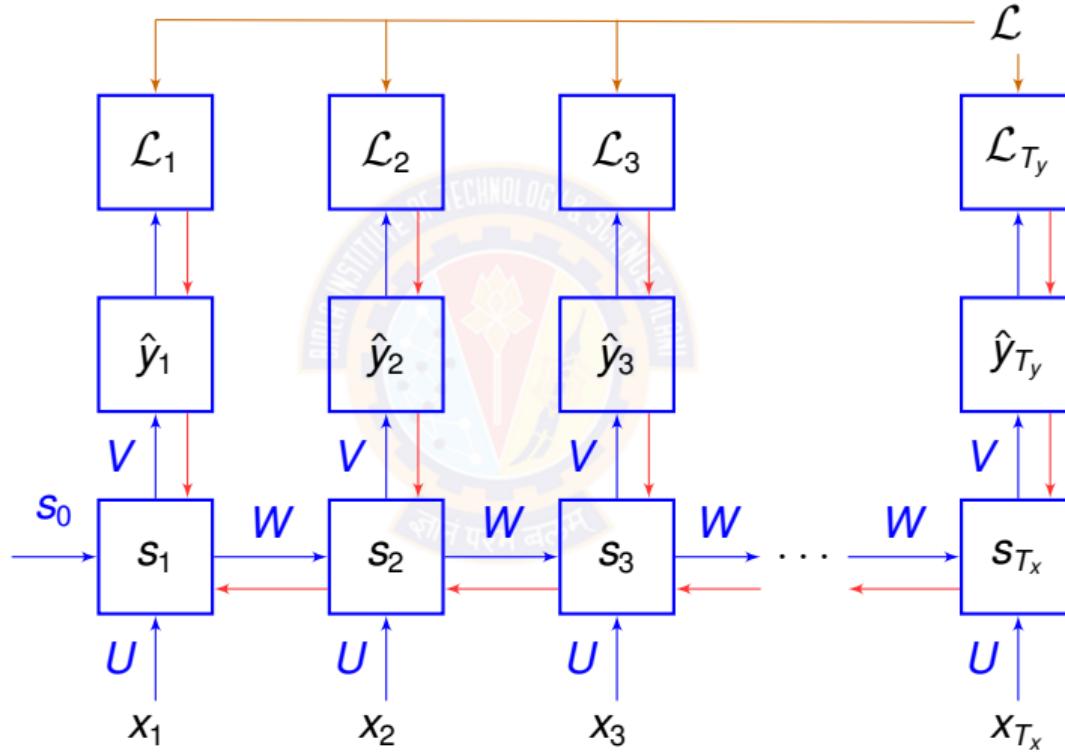
Loss function

$$\mathcal{L}_t(\hat{y}_t, y_t) = \prod_{t=1}^{T_y} P(\hat{y}_t | \hat{y}_{t-1}, \dots, \hat{y}_1)$$

Overall Loss

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}_t(\hat{y}_t, y_t)$$

BACK PROPAGATION IN RNN



Back-propagation through time.

IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN

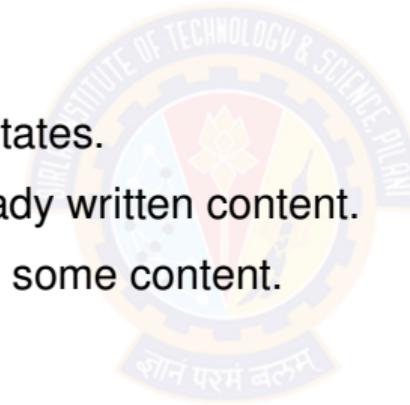


ISSUE OF MAINTAINING STATES

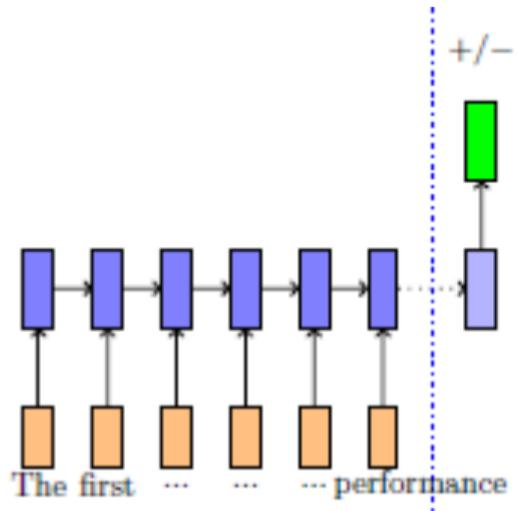
- The old information gets morphed by the current input at each new time step.
- After t steps the information stored at time step $t - k$ (for some $k < t$) gets completely morphed so much that it would be impossible to extract the original information stored at time step $t - k$.
- It is very hard to assign the responsibility of the error caused at time step t to the events that occurred at time step $t - k$.
- Basically depends on the size of memory that is available.

STRATEGY TO MAINTAIN STATES

- Selectively write on the states.
- Selectively read the already written content.
- Selectively forget (erase) some content.



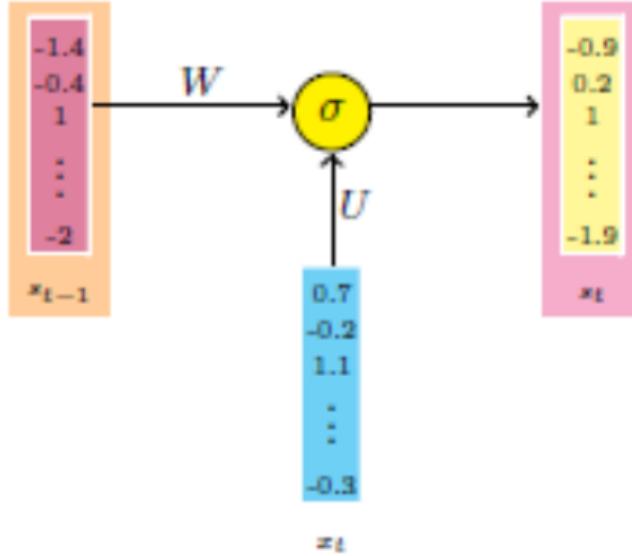
SENTIMENT ANALYSIS



Review: The first half of the movie was dry but the second half really picked up pace. The lead actor delivered an amazing performance

- RNN reads the document from left to right and after every word updates the state.
- By the time we reach the end of the document the information obtained from the first few words is completely lost.
- Ideally we want to
 - ▶ forget the information added by stop words (a, the, etc.).
 - ▶ selectively read the information added by previous sentiment bearing words (awesome, amazing, etc.)
 - ▶ selectively write new information from the current word to the state.

SELECTIVE WRITE



- Recall that in RNNs we use s_{t-1} to compute s_t .

$$s_t = \sigma(Ws_{t-1} + Ux_t + b)$$

SELECTIVE WRITE



- Introduce a vector o_{t-1} which decides what fraction of each element of s_{t-1} should be passed to the next state.
- Each element of o_{t-1} gets multiplied with the corresponding element of s_{t-1} .
- Each element of o_{t-1} is restricted to be between 0 and 1.
- The RNN has to learn o_{t-1} along with the other parameters (W, U, V).

SELECTIVE WRITE

- Compute o_{t-1} and h_{t-1} as

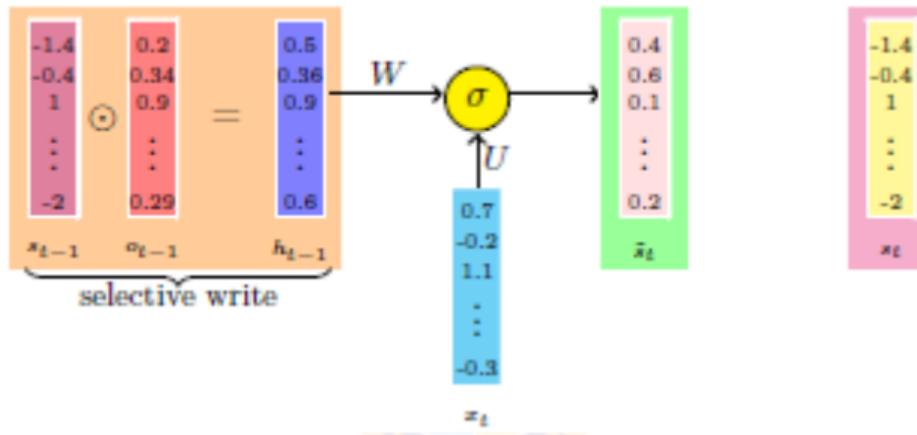
$$o_{t-1} = \sigma(W_o h_{t-2} + U_o x_{t-1} + b_o)$$

$$h_{t-1} = o_{t-1} \odot \sigma(s_{t-1})$$



- The parameters (W_o, U_o, b_o) are learned along with the existing parameters (W, U, V) .
- The sigmoid function ensures that the values are between 0 and 1.
- o_t is called the **output gate** as it decides how much to pass (write) to the next time step.

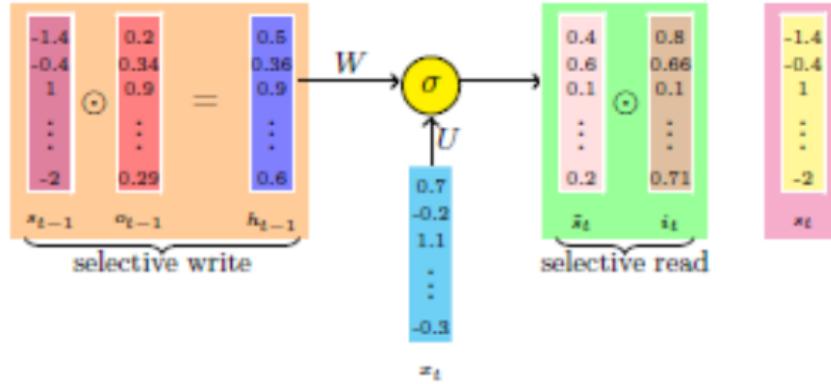
COMPUTE STATE



- h_{t-1} and x_t are used to compute the new state at the next time step.

$$\tilde{s}_t = \sigma(Wh_{t-1} + Ux_t + b)$$

SELECTIVE READ

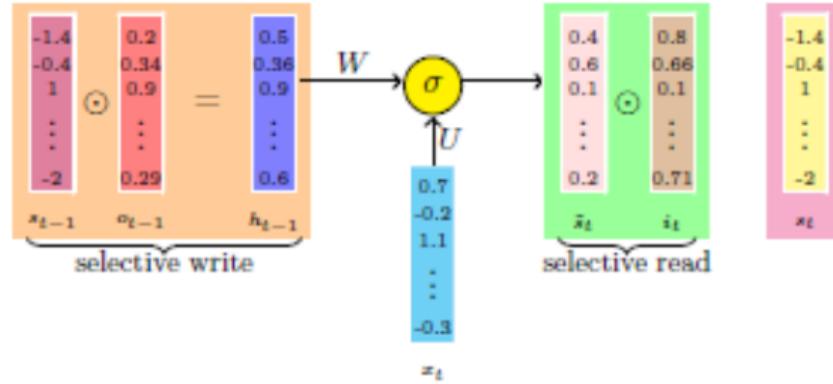


- \tilde{s}_t captures all the information from the previous state h_{t-1} and the current input x_t .
- To do selective read, introduce another gate called the **input gate**.

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$\text{Selectively Read} = i_t \odot \tilde{s}_t$$

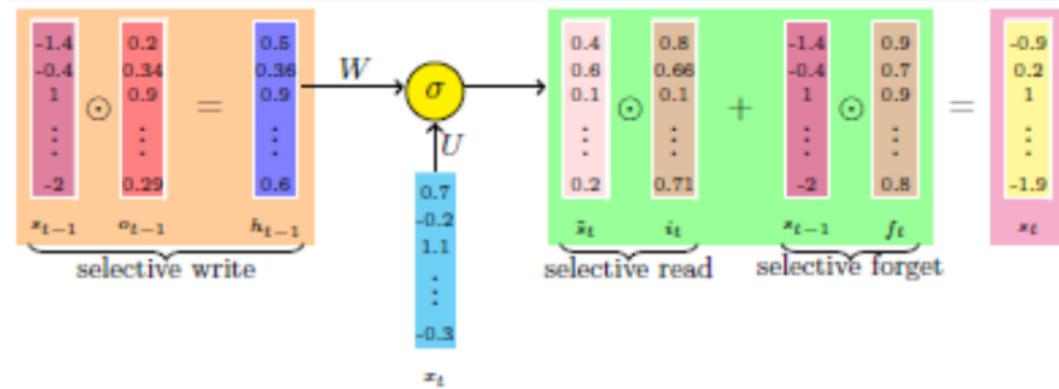
SELECTIVE READ



- \tilde{s}_t captures all the information from the previous state h_{t-1} and the current input x_t .
- To do selective read, introduce another gate called the **input gate**.

$$s_t = s_{t-1} + i_t \odot \tilde{s}_t$$

SELECTIVE FORGET



- To do selective forget, introduce another gate called the **forget gate**.

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

FULL LSTM



- 3 gates

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

- 3 states

$$\tilde{s}_t = \sigma(W h_{t-1} + U x_t + b)$$

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

$$h_t = o_t \odot \sigma(s_t)$$

$$\hat{y}_t = g(V s_t + c)$$

IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN



LONG SHORT TERM MEMORY UNIT (LSTM)

- Another representation
- 3 gates are used – Update gate Γ_u , Forget gate Γ_f and Output gate Γ_o .

$$\tilde{c}^{} = \tanh(W_c[a^{}, x^{}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{}, x^{}] + b_u)$$

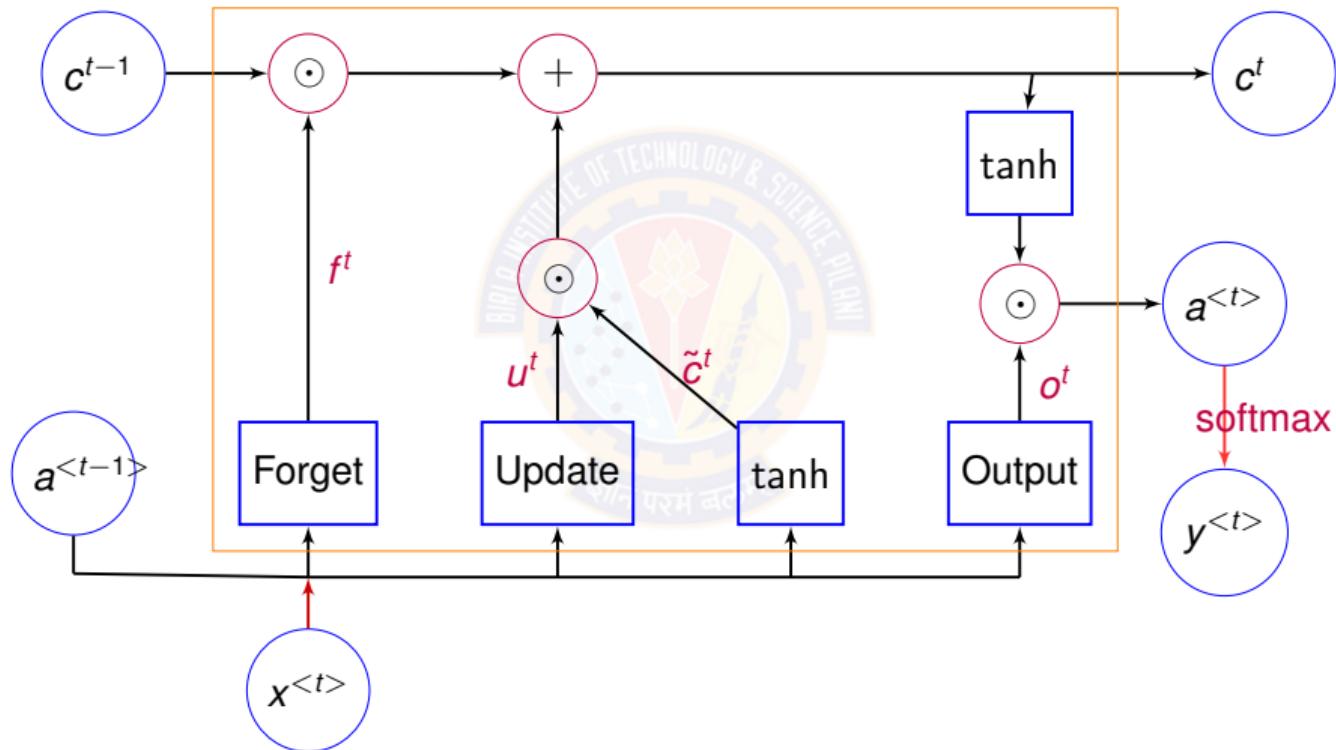
$$\Gamma_f = \sigma(W_f[a^{}, x^{}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{}, x^{}] + b_o)$$

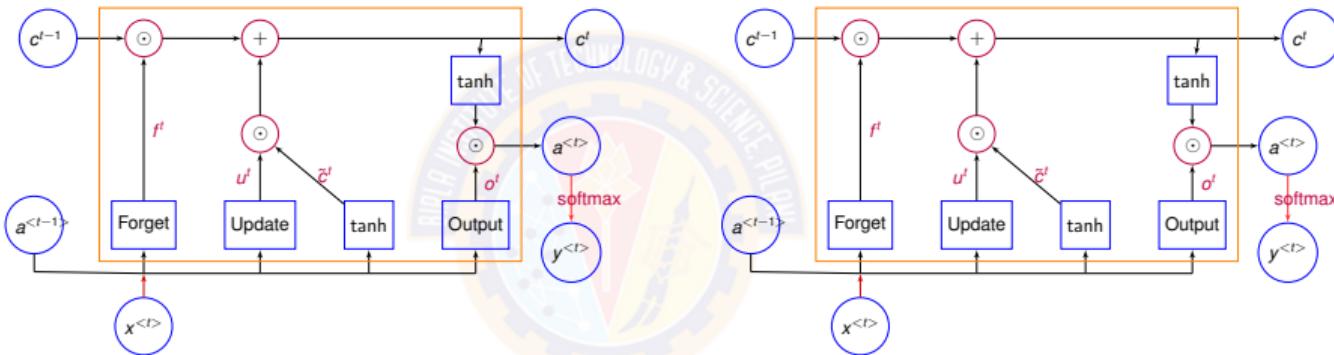
$$c^{} = \Gamma_u * \tilde{c}^{} + \Gamma_f * c^{}$$

$$a^{} = \Gamma_o * \tanh(c^{})$$

LSTM



LSTM



IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN



GATED RECURRENT UNIT (GRU)

- Introduce a memory cell $c^{<t>} = a^{<t>}.$
- Candidate for replacing $c^{<t>}$ is given as $\tilde{c}^{<t>}.$
- The decision whether to update $c^{<t>}$ with $\tilde{c}^{<t>}$ is given by the **update gate** $\Gamma_u.$ Γ_u takes the value of 0 or 1.

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

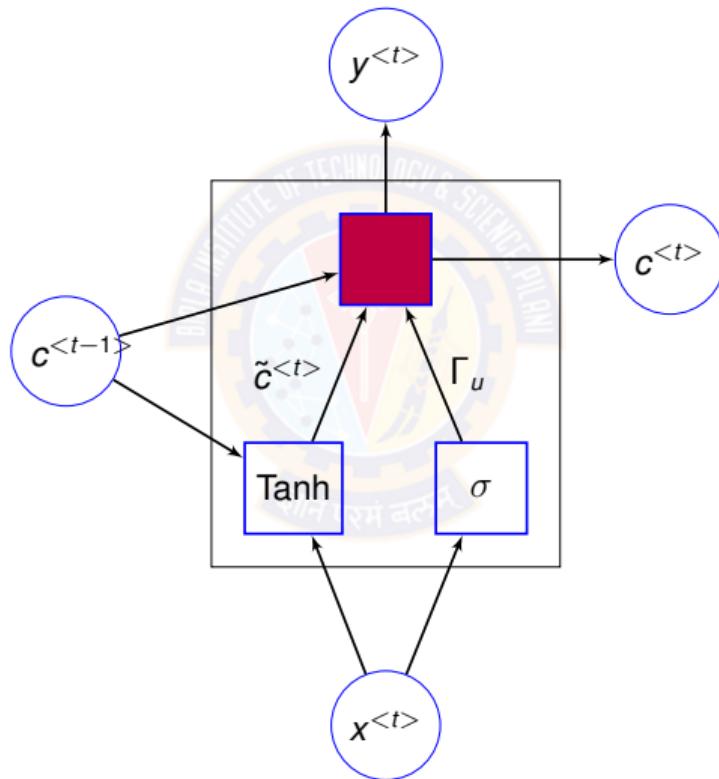
$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>} + b_u])$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>} + b_r])$$

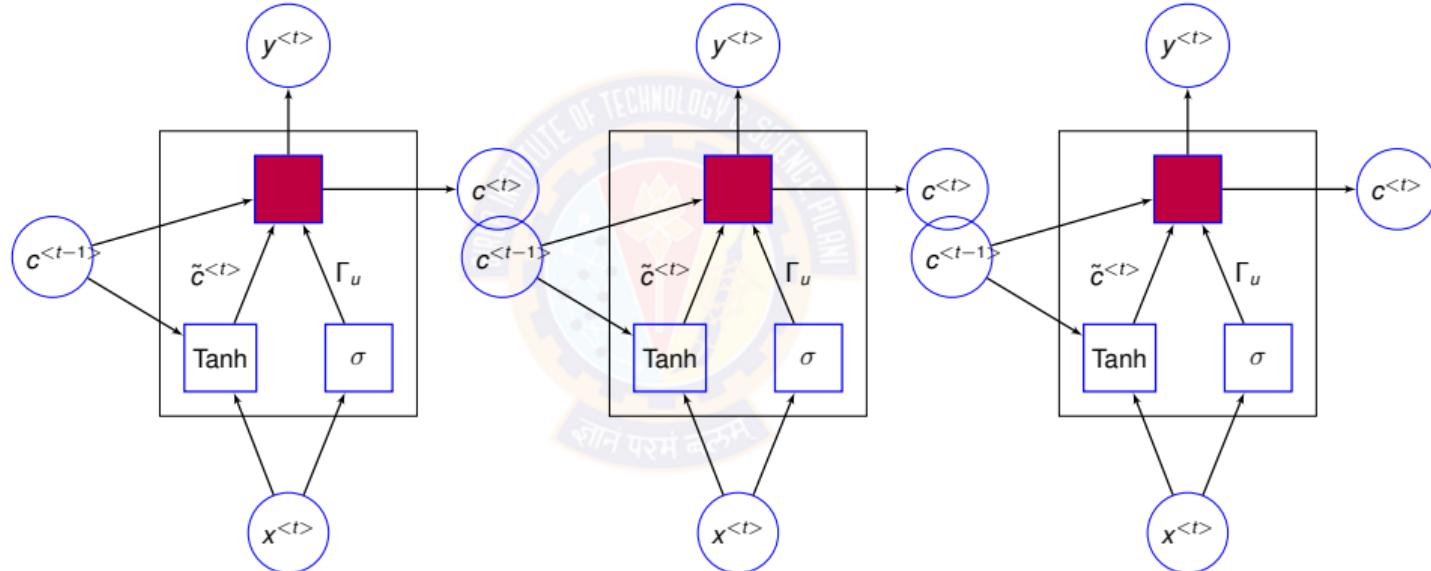
$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

GATED RECURRENT UNIT (GRU)



GATED RECURRENT UNIT (GRU)



IN THIS SEGMENT

- 1 SEQUENCE LEARNING
- 2 RECURRENT NEURAL NETWORK (RNN)
- 3 TYPES OF RNN
- 4 LEARNING IN RNN
- 5 ISSUES IN RNN
- 6 LONG SHORT TERM MEMORY UNIT (LSTM)
- 7 GATED RECURRENT UNIT (GRU)
- 8 BIDIRECTIONAL RNN



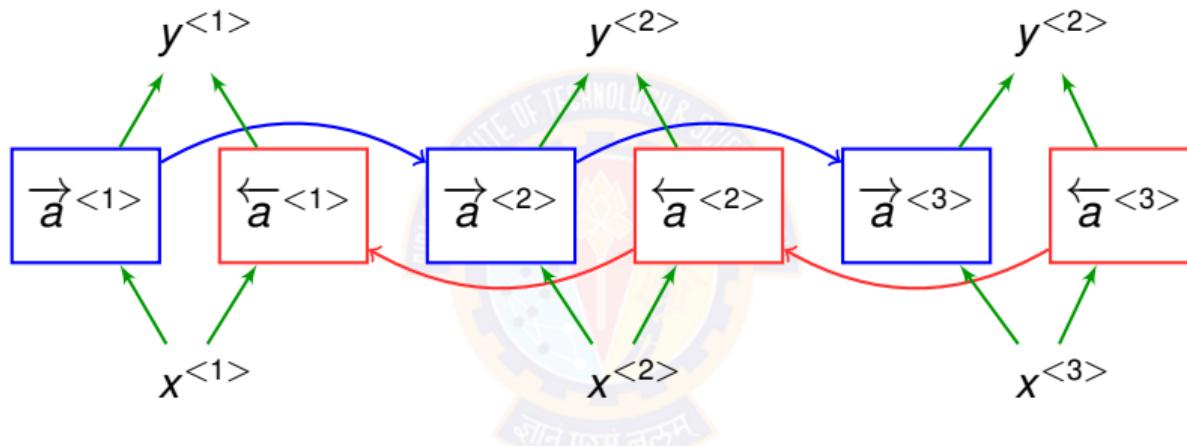
BIDIRECTIONAL RNN (BRNN)

- Forward and backward connections.
- The blocks can be RNN, GRU, LSTM.
- Mostly used in the NLP.
- Acyclic graph
- Example: Name entity recognition
He said “Teddy bear is soft.”
He said “Teddy Roosevelt was a President.”

common noun

proper name

BRNN ARCHITECTURE



$$\hat{y}^{<t>} = g(W_y [\vec{a}^{<t>} \leftarrow \vec{a}^{<t>} x^{<t>}] + b_y)$$

SUMMARY

- Use GRU, when dependency is short. Eg: Weather forecasting
- Use LSTM, when dependency is long. Eg: NLP Translation
- Use BRNN, dependency is in both direction. Eg: Stock prediction

References

- ① Deep Learning by Ian Goodfellow, Yoshua Bengio, Aaron Courville
<https://www.deeplearningbook.org/>
- ② Deep Learning with Python by Francois Chollet.
<https://livebook.manning.com/book/deep-learning-with-python/>

Thank You!



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

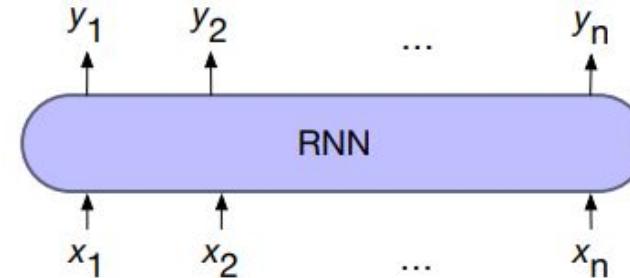
Attention Models & Transformers

S. P. Vimal, CSIS Department (WILP Division)
vimalsp@pilani.bits-pilani.ac.in

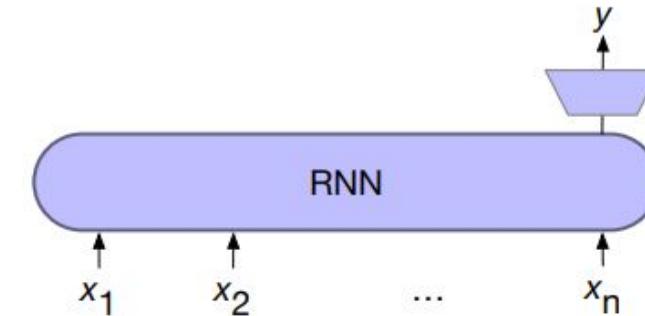
Topics

- ★ Brief review of Sequence Models (RNN, LSTMS)
- ★ Encoder - Decoder Models
- ★ Attention
- ★ Transformers
- ★ Vision Transformers

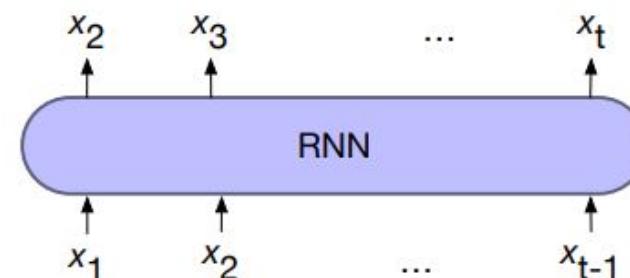
RNN Architectures for NLP Tasks



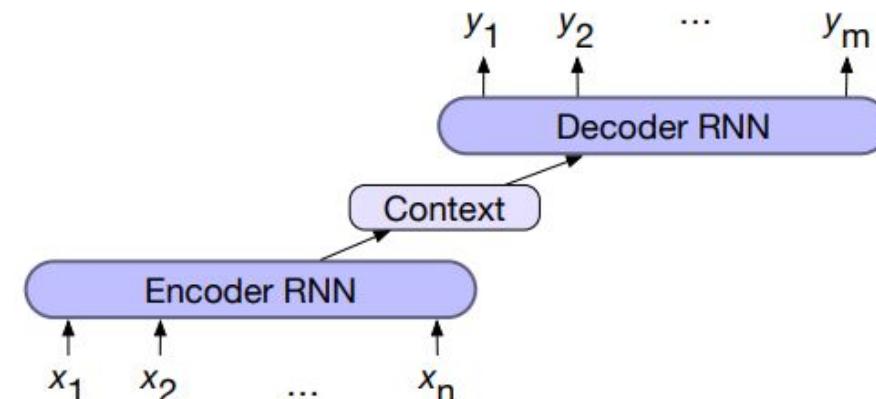
a) sequence labeling
Ex: POS Tagging, Named Entity Tagging



b) sequence classification
Ex: Sentiment Analysis



c) language modeling
Ex: Predict Next Word



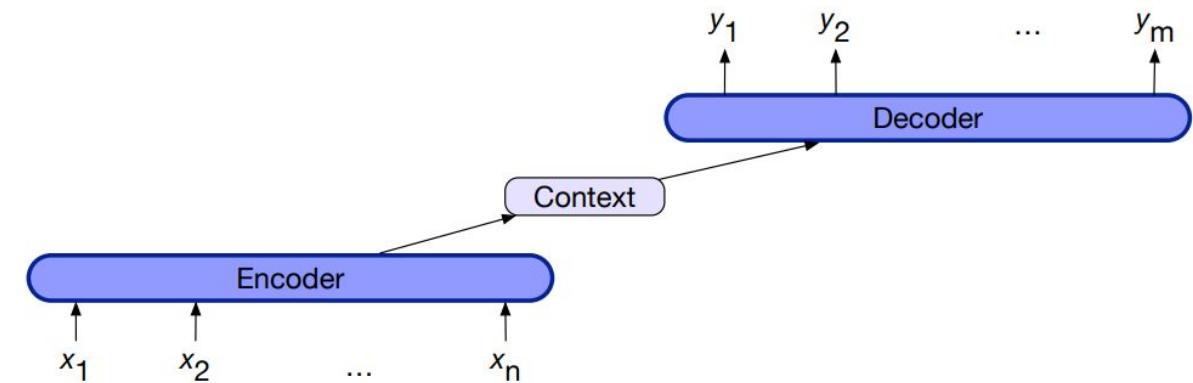
d) encoder-decoder
Ex: Language Translation

Encoder - Decoder

- **Goal:** Develop an architecture capable of generating *contextually appropriate, arbitrary length*, output sequences

- **Applications:**

- Machine translation
- Summarization
- Question answering
- Dialogue modeling.



Encoder - Decoder

Encoder:

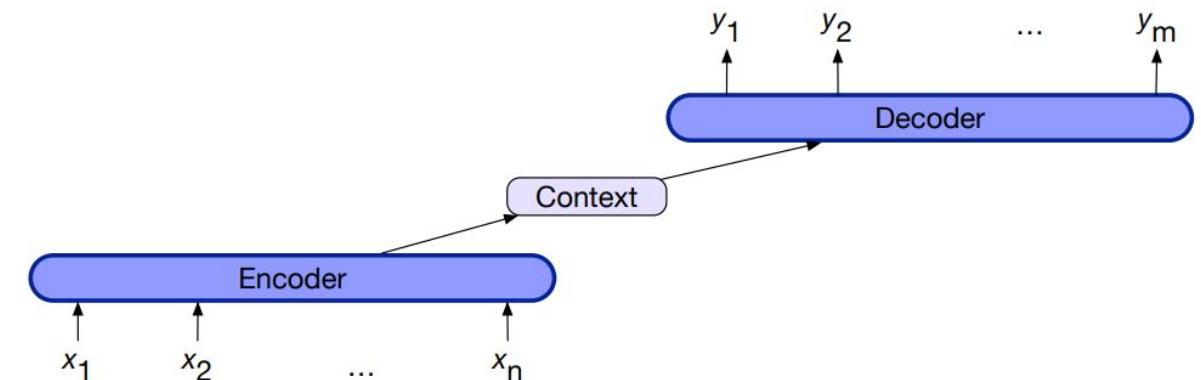
- Input sequence $(x_{1 \dots n}) \rightarrow$ Sequence of contextualized representations, $(h_{1 \dots n})$
- Ex: LSTM, CNN, Transformers etc.

Context:

- c , a function of $(h_{1 \dots n})$

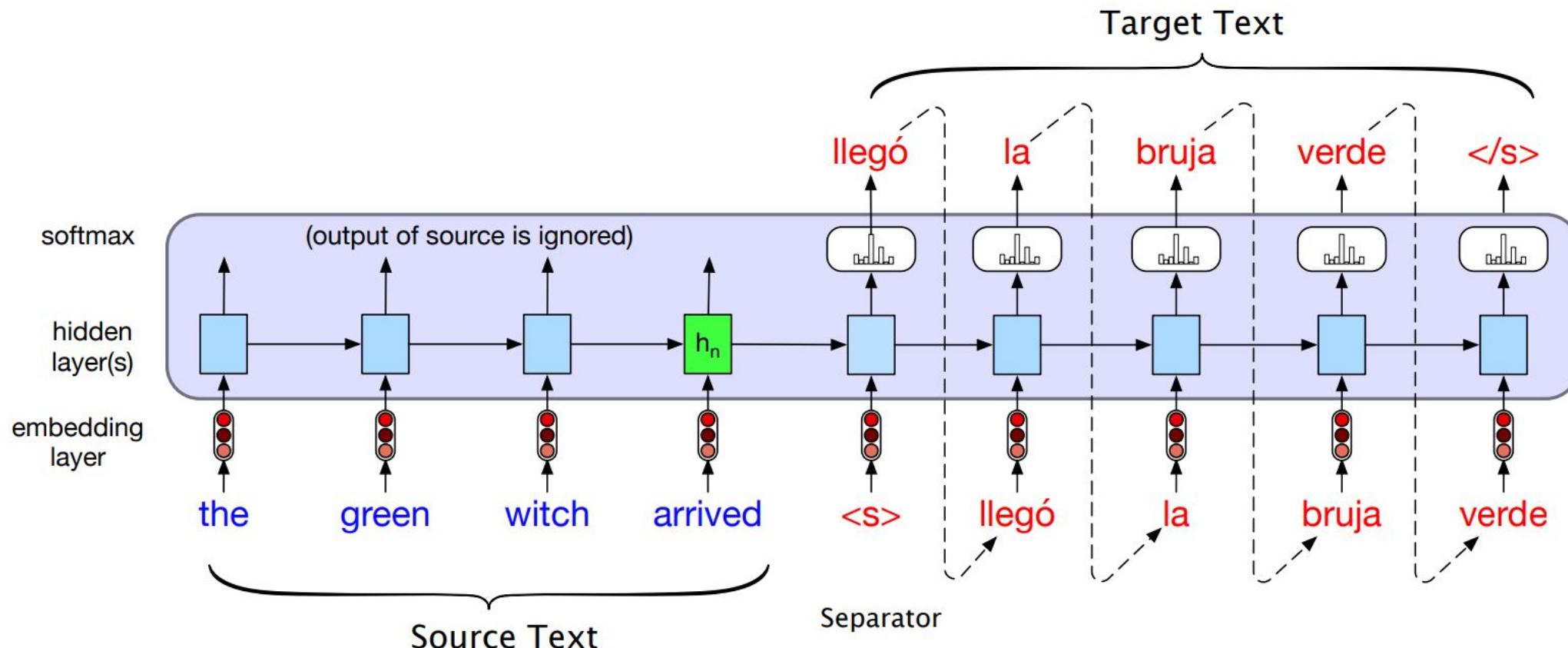
Decoder:

- $c \rightarrow$ arbitrary length sequence of hidden states $(h_{1 \dots m}) \rightarrow$ sequence of output states $(y_{1 \dots m})$



Encoder - Decoder

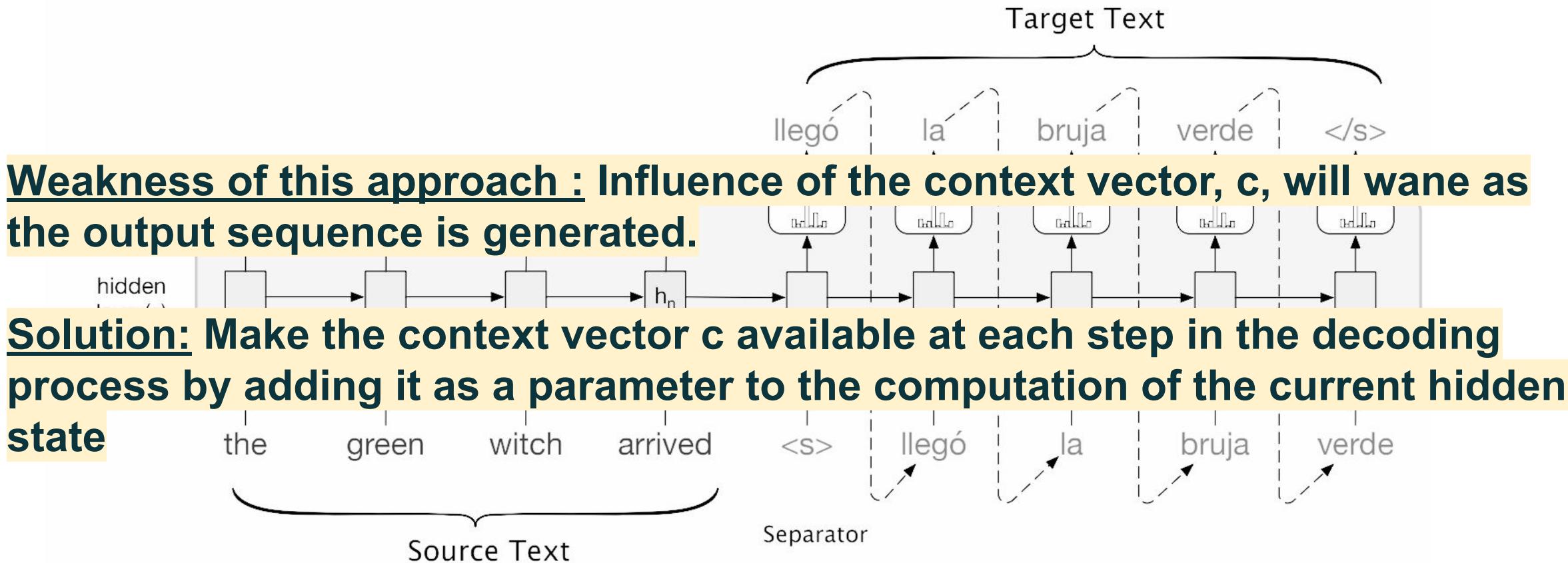
Encoder - Decoder for Language Translation



Source and target sentences are concatenated with a separator token in between, and the decoder uses context information from the encoder's last hidden state

Encoder - Decoder

Encoder - Decoder for Language Translation



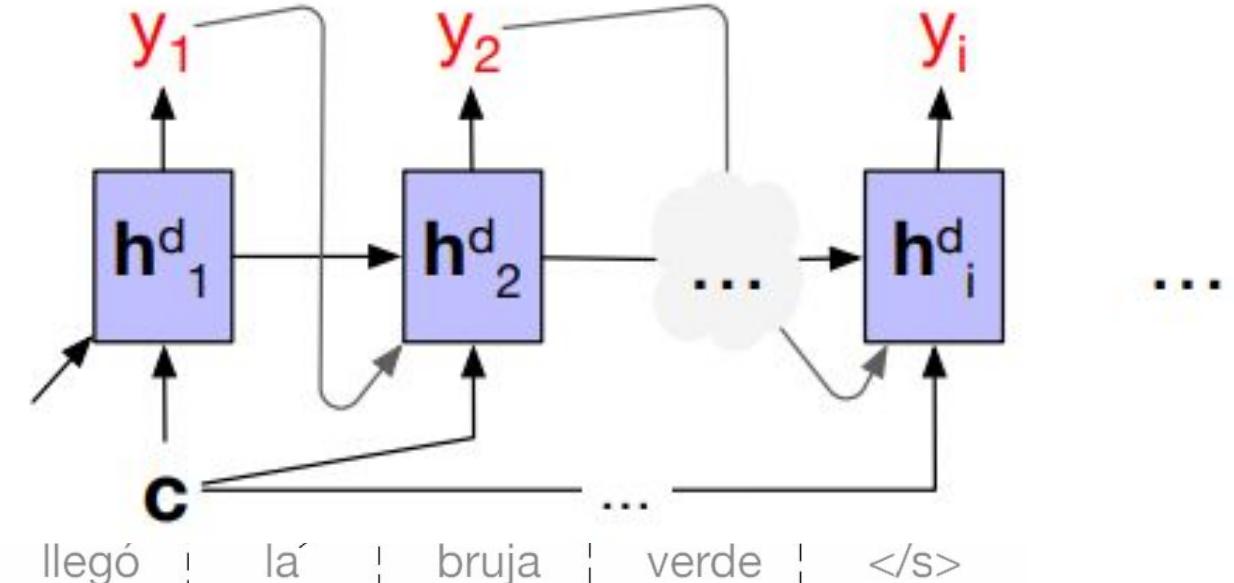
Weakness of this approach : Influence of the context vector, c , will wane as the output sequence is generated.

Solution: Make the context vector c available at each step in the decoding process by adding it as a parameter to the computation of the current hidden state

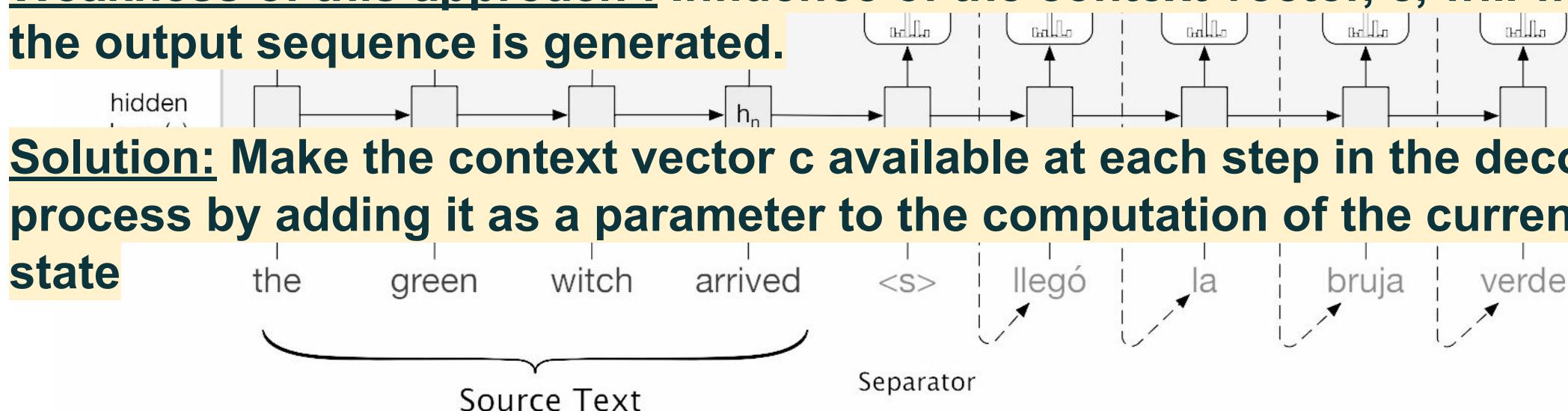
Source and target sentences are concatenated with a separator token in between, and the decoder uses context information from the encoder's last hidden state

Encoder - Decoder

Encoder - Decoder for Language Translation



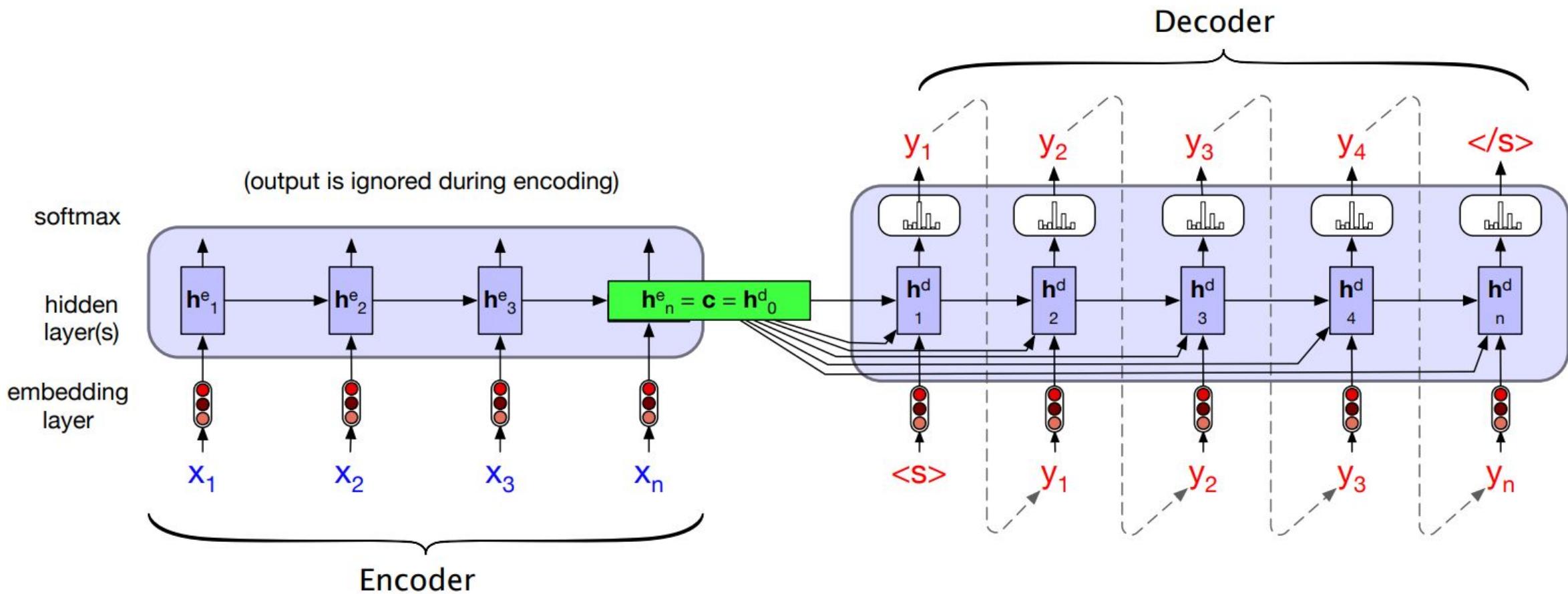
Weakness of this approach : Influence of the context vector, c , will wane as the output sequence is generated.



Source and target sentences are concatenated with a separator token in between, and the decoder uses context information from the encoder's last hidden state

Encoder - Decoder

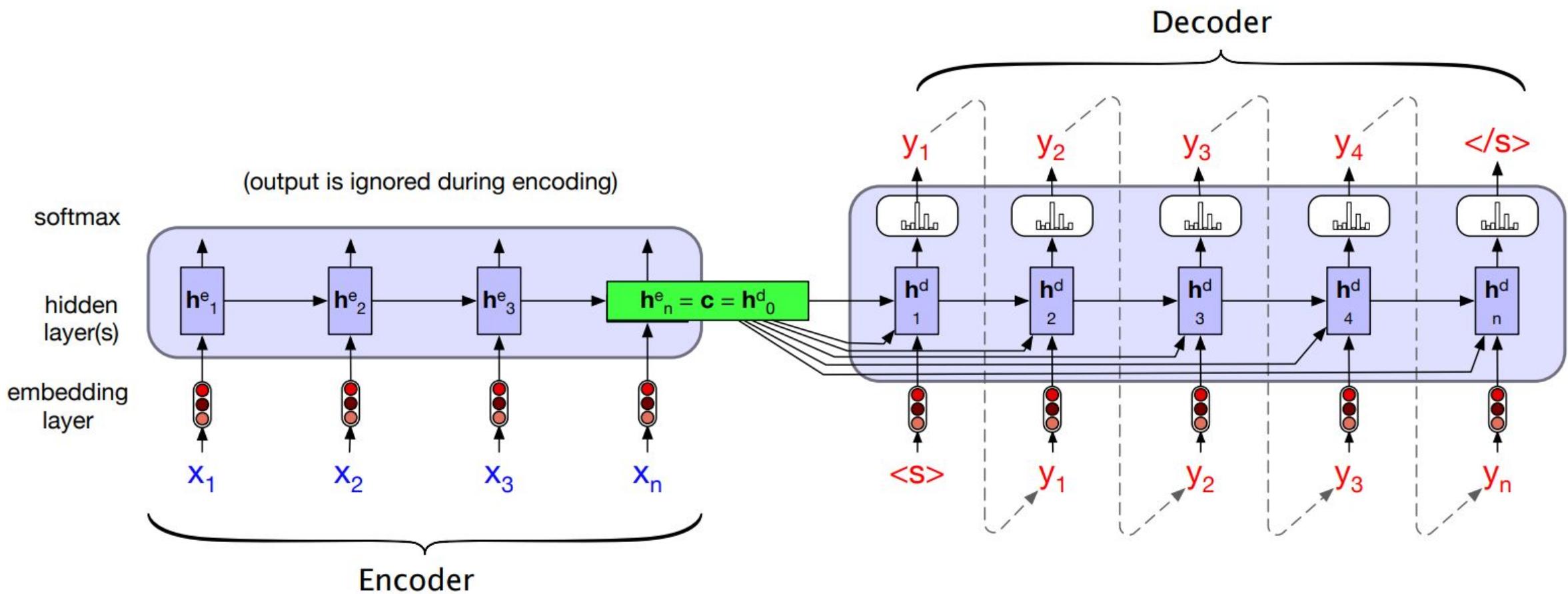
Encoder - Decoder for Language Translation



Encoder - Decoder

Encoder - Decoder for Language Translation

$$\begin{aligned}
 \mathbf{c} &= \mathbf{h}_n^e \\
 \mathbf{h}_0^d &= \mathbf{c} \\
 \mathbf{h}_t^d &= g(\hat{y}_{t-1}, \mathbf{h}_{t-1}^d, \mathbf{c}) \\
 \mathbf{z}_t &= f(\mathbf{h}_t^d) \\
 y_t &= \text{softmax}(\mathbf{z}_t)
 \end{aligned}$$

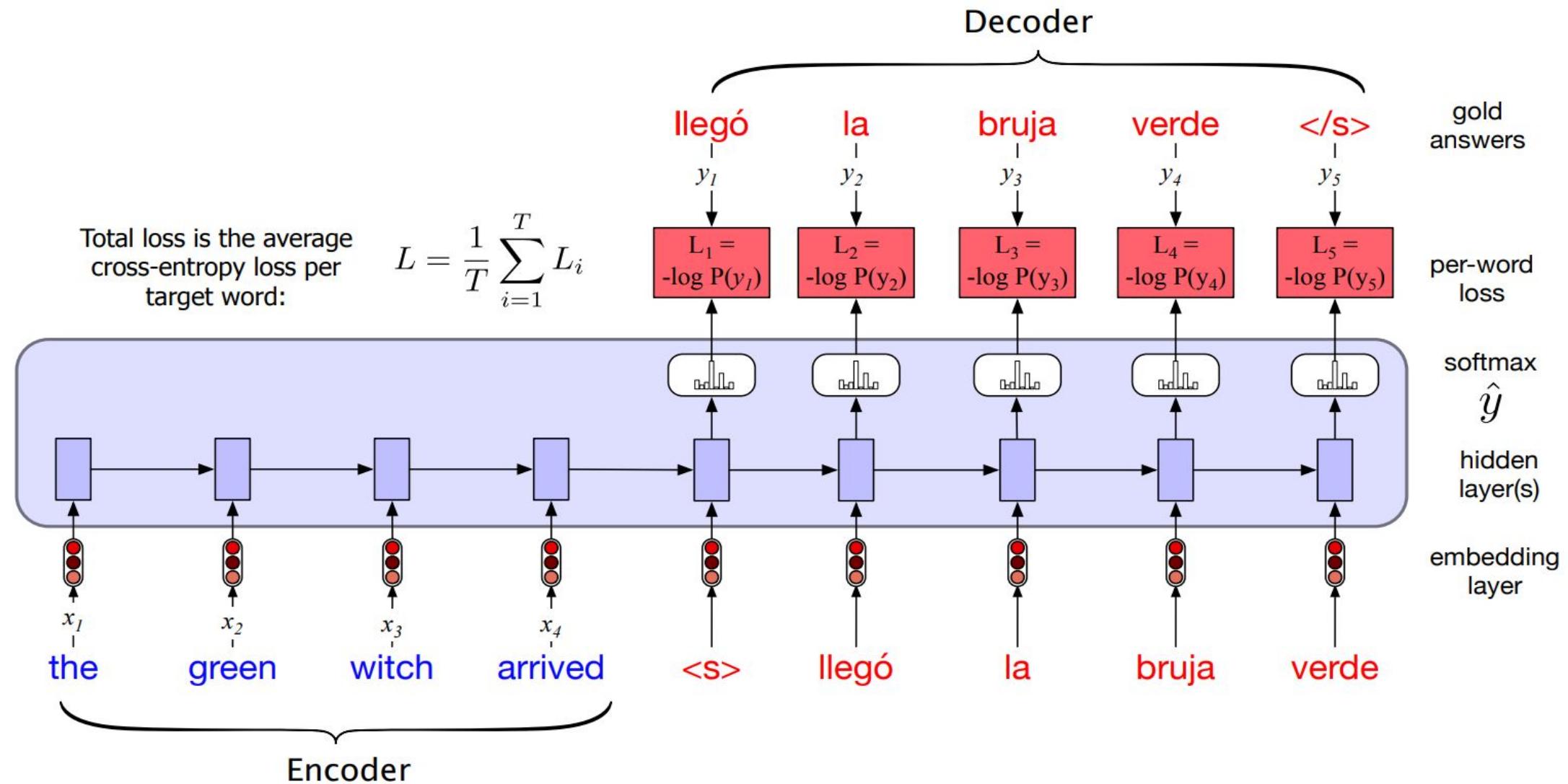


Encoder - Decoder

Training

Total loss is the average cross-entropy loss per target word:

$$L = \frac{1}{T} \sum_{i=1}^T L_i$$



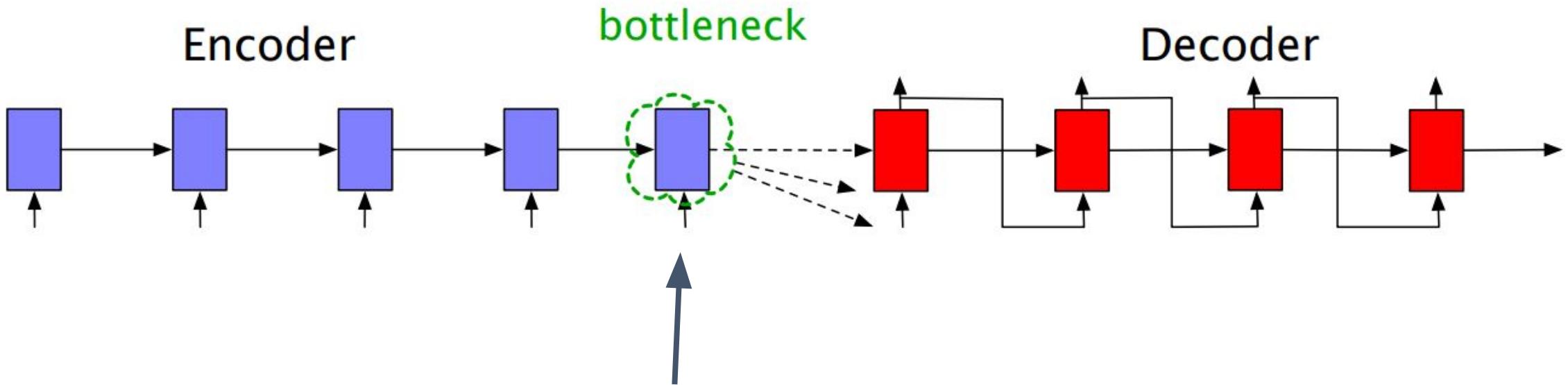


Encoder - Decoder

Teaching Forcing

- Force the system to use the gold target token from training as the next input x_{t+1} , rather than allowing it to rely on the (possibly erroneous) decoder output \hat{y}_t .
- Speeds up training

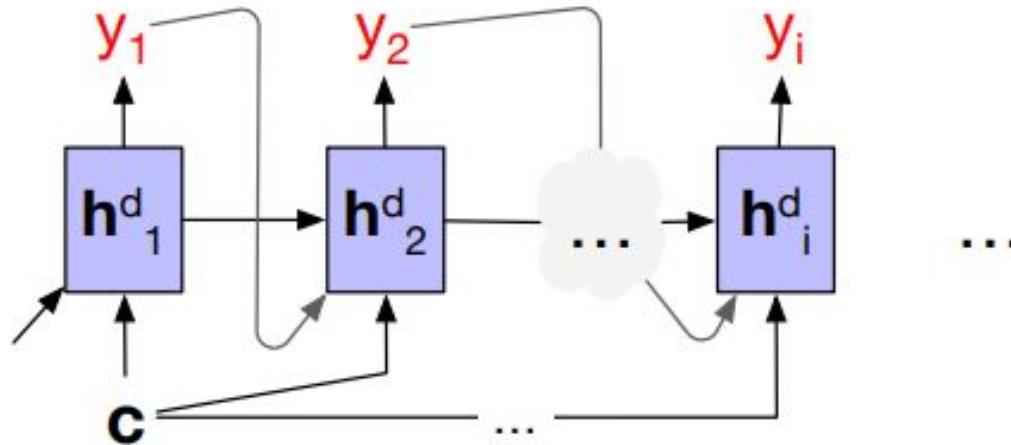
Attention !



In an encoder-decoder arch, the final hidden state acts as a bottleneck:

- It must represent absolutely everything about the meaning of the source text
- The only thing the decoder knows about the source text is what's in this context vector

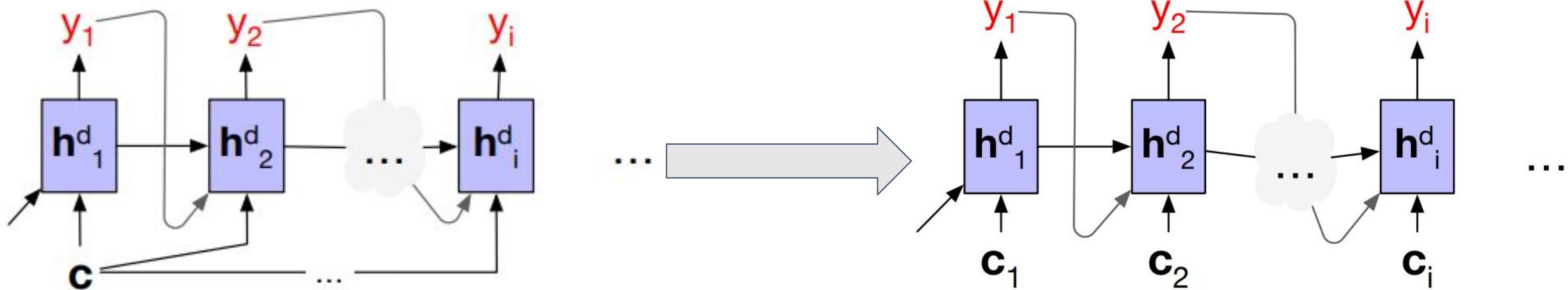
Attention !



Without attention, a decoder sees the same context vector ,
which is a static function of all the encoder hidden states

$$\mathbf{h}_t^d = g(\hat{y}_{t-1}, \mathbf{h}_{t-1}^d, \mathbf{c})$$

Attention !



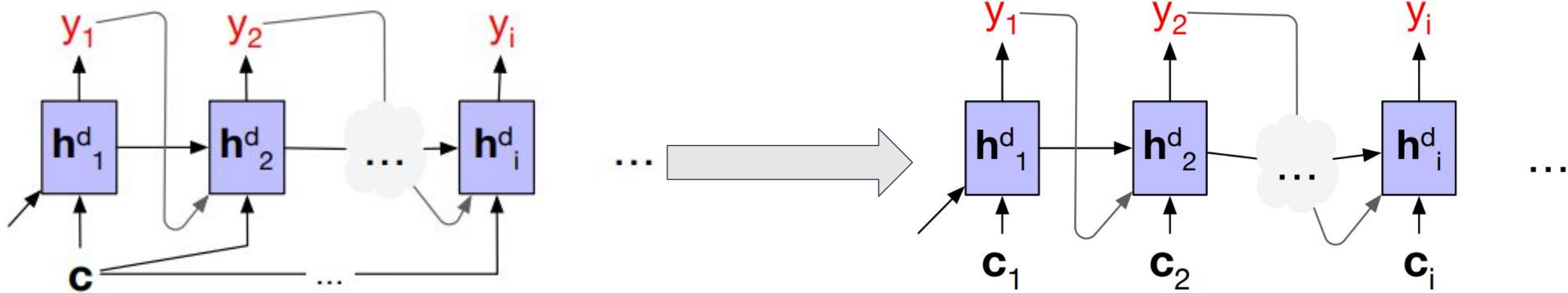
Without attention, a decoder sees the same context vector ,
which is a static function of all the encoder hidden states

$$\mathbf{h}_t^d = g(\hat{y}_{t-1}, \mathbf{h}_{t-1}^d, \mathbf{c})$$

With attention, decoder to sees a different, dynamic,
context, which is a function of all the encoder hidden states

$$\mathbf{h}_i^d = g(\hat{y}_{i-1}, \mathbf{h}_{i-1}^d, \mathbf{c}_i)$$

Attention !



Without attention, a decoder sees the same context vector , which is a static function of all the encoder hidden states

$$\mathbf{h}_t^d = g(\hat{y}_{t-1}, \mathbf{h}_{t-1}^d, \mathbf{c})$$

With attention, decoder gets information from all the hidden states of the encoder, not just the last hidden state of the encoder

Each context vector is obtained by taking a weighted sum of all the encoder hidden states.

The weights focus on ('attend to') a particular part of the source text that is relevant for the token the decoder is currently producing

$$\mathbf{h}_i^d = g(\hat{y}_{i-1}, \mathbf{h}_{i-1}^d, \mathbf{c}_i)$$

Attention !

Step -1 : Find out how relevant each encoder state is to the present decoder state \mathbf{h}_{i-1}^d

Compute a score of similarity between \mathbf{h}_{i-1}^d and all the encoder states : $score(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e)$

Dot Product Attention : $score(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e) = \mathbf{h}_{i-1}^d \cdot \mathbf{h}_j^e$

Step -2 : Normalize all the scores with softmax to create a vector of weights, $\alpha_{i,j}$

$\alpha_{i,j}$ indicates the proportional relevance of each encoder hidden state j to the prior hidden decoder state, \mathbf{h}_{i-1}^d

$$\begin{aligned}\alpha_{ij} &= \text{softmax}(\text{score}(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e) \quad \forall j \in e) \\ &= \frac{\exp(\text{score}(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e))}{\sum_k \exp(\text{score}(\mathbf{h}_{i-1}^d, \mathbf{h}_k^e))}\end{aligned}$$



Attention !

Step -3 : Given the distribution in α , compute a fixed-length context vector for the current decoder state by taking a weighted average over all the encoder hidden states

$$\mathbf{c}_i = \sum_j \alpha_{ij} \mathbf{h}_j^e$$

Attention !

Step -3 : Given the distribution in α , compute a fixed-length context vector for the current decoder state by taking a weighted average over all the encoder hidden states

$$\mathbf{c}_i = \sum_j \alpha_{ij} \mathbf{h}_j^e$$

Plus : In step-1, we can get a more powerful scoring function by parameterizing the score with its own set of weights, \mathbf{W}_s :

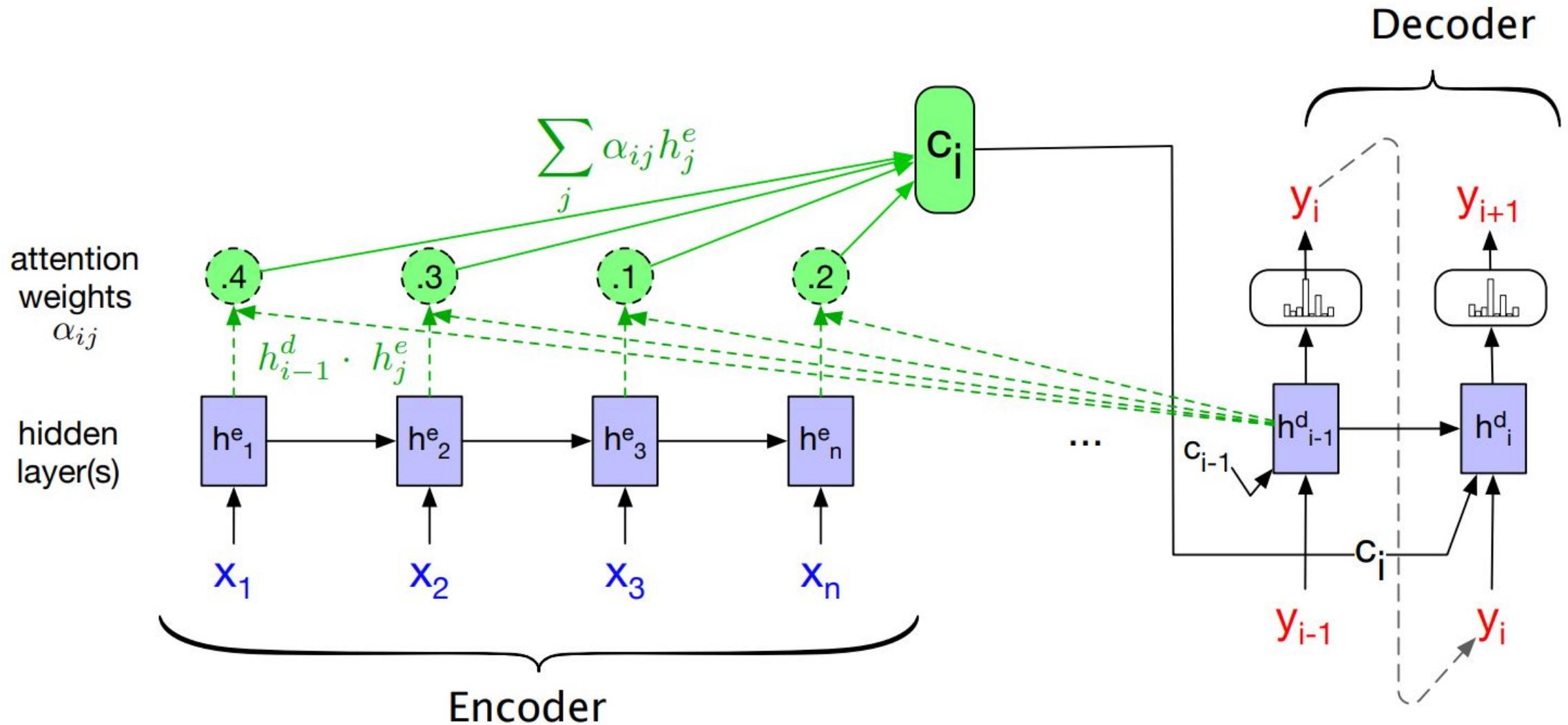
$$\text{score}(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e) = \mathbf{h}_{i-1}^d \mathbf{W}_s \mathbf{h}_j^e$$

\mathbf{W}_s , is trained during normal end-to-end training,

\mathbf{W}_s , gives the network the ability to learn which aspects of similarity between the decoder and encoder states are important to the current application.

Attention !

$$\mathbf{h}_i^d = g(\hat{y}_{i-1}, \mathbf{h}_{i-1}^d, \mathbf{c}_i)$$





Transformers

- 2017, NIPS, Vaswani et. al., **Attention Is All You Need !!!**
- Transformers **map sequences** of input vectors (x_1, \dots, x_n) to sequences of output vectors (y_1, \dots, y_n) **of the same length**
- Made up of **transformer blocks** in which the key component is **self-attention** layers

[Self-attention allows a network to directly extract and use information from arbitrarily large contexts directly !!!]

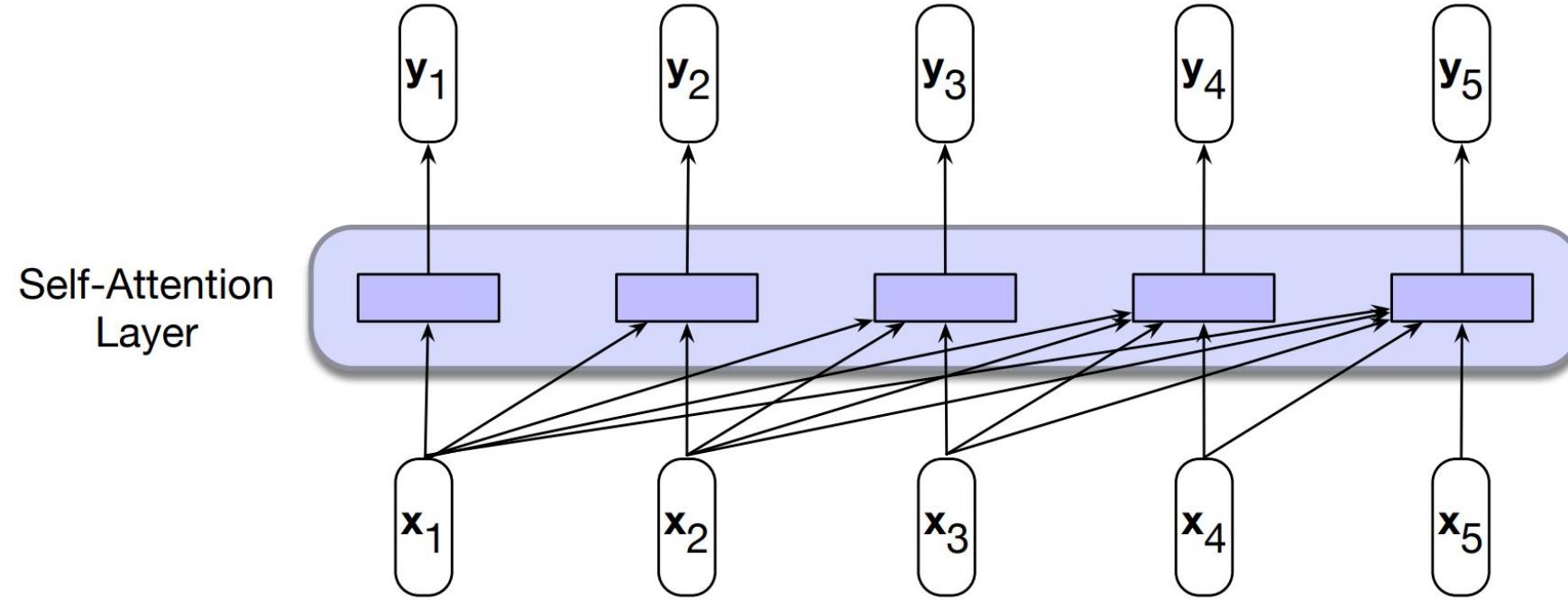
- Transformers are **not based on recurrent connections** \Rightarrow Parallel implementations possible \Rightarrow Efficient to scale (comparing LSTM)



Self-Attention | Transformers

- Attention \Rightarrow Ability to compare an item of interest to a collection of other items in a way that reveals their relevance in the current context.
- Self-attention \Rightarrow
 - > Set of **comparisons** are to other elements **within a given sequence**
 - > Use these comparisons to compute an output for the current input

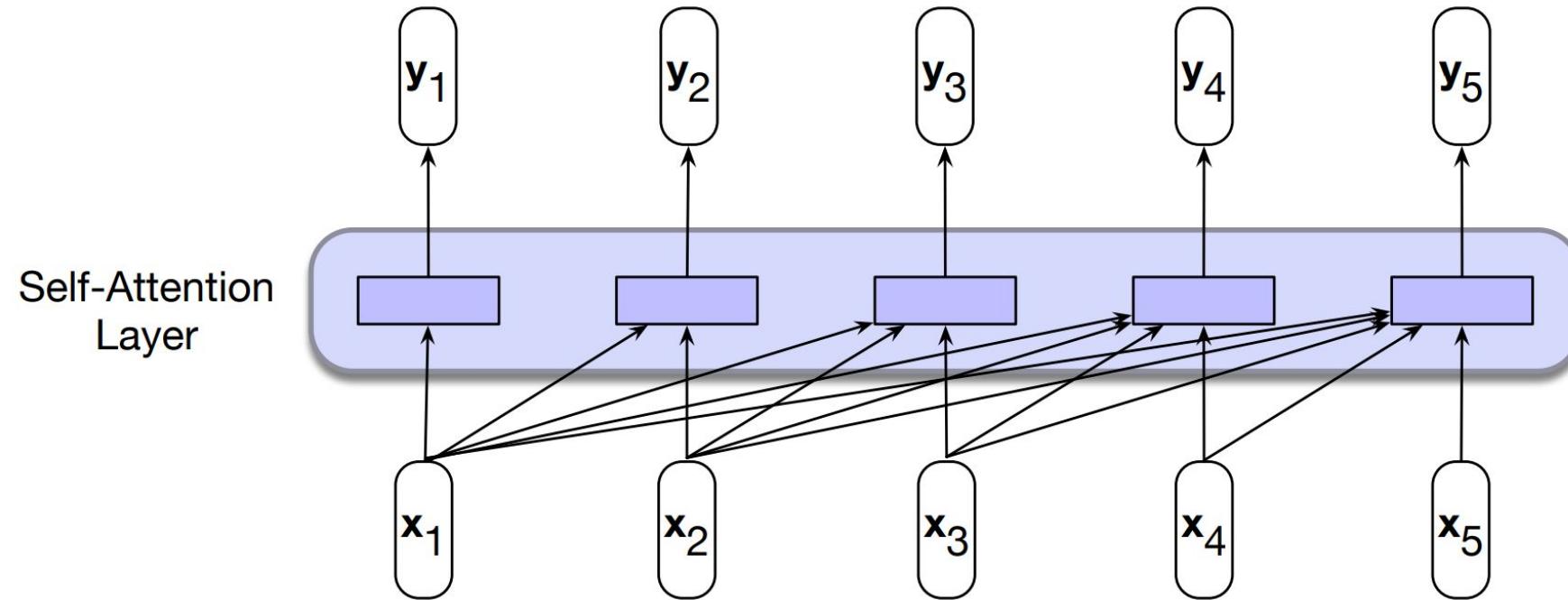
Self-Attention | Transformers



In processing each element of the sequence, the model attends to all the inputs up to, and including, the current one.

Unlike RNNs, the computations at each time step are independent of all the other steps and therefore can be performed in parallel.

Self-Attention | Transformers



$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\begin{aligned} \alpha_{ij} &= \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \quad \forall j \leq i \\ &= \frac{\exp(\text{score}(\mathbf{x}_i, \mathbf{x}_j))}{\sum_{k=1}^i \exp(\text{score}(\mathbf{x}_i, \mathbf{x}_k))} \quad \forall j \leq i \end{aligned}$$

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{x}_j$$



Self-Attention | Transformers

- Let us understand how transformers uses self-attention !

Self-Attention | Transformers

$$\mathbf{W}^V, \mathbf{W}^K, \mathbf{W}^Q \in \mathbb{R}^{d \times d}$$

In Vaswani et al., 2017, d was 1024.

- Let us understand how transformers uses self-attention !

$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i$$

$$\mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i$$

$$\mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$$

Query, Q

As the current focus of attention when being compared to all of the other preceding inputs.

Key, K

In its role as a preceding input being compared to the current focus of attention.

Value, V

As a value used to compute the output for the current focus of attention

Three different roles each \mathbf{x}_i (input embedding) , in the computation of self attention

Self-Attention | Transformers

$$\mathbf{W}^V, \mathbf{W}^K, \mathbf{W}^Q \in \mathbb{R}^{d \times d}$$

In Vaswani et al., 2017, d was 1024.

- Let us understand how transformers uses self-attention !

$$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i \quad \mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i \quad \mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$$

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}$$

The simple dot product can be an arbitrarily large;
scaled dot-product is used in transformers;

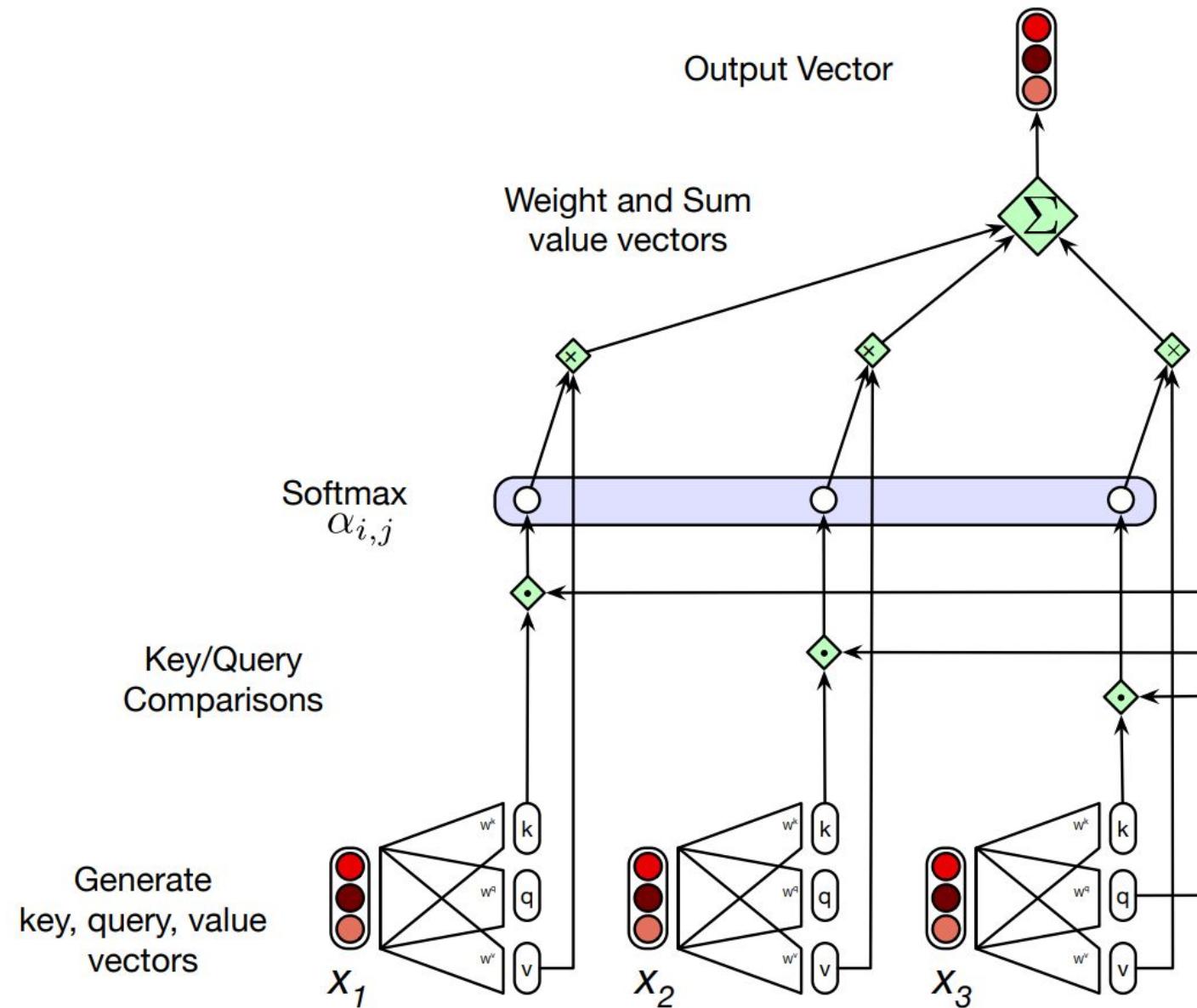
$$\begin{aligned}\alpha_{ij} &= \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \quad \forall j \leq i \\ &= \frac{\exp(\text{score}(\mathbf{x}_i, \mathbf{x}_j))}{\sum_{k=1}^i \exp(\text{score}(\mathbf{x}_i, \mathbf{x}_k))} \quad \forall j \leq i\end{aligned}$$

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

Self-Attention | Transformers

- Each output, y_i , is computed independently
- Entire process can be parallelized

Calculating the value of y_3 , the third element of a sequence using causal (left-to-right) self-attention





Self-Attention | Transformers

- Pack the input embeddings of the N input tokens into a single matrix

$$\mathbf{X} \in \mathbb{R}^{N \times d}$$

> Each row of X is the embedding of one token of the input

- Multiply X by the key, query, and value ($d \times d$) matrices

$$\mathbf{Q} = \mathbf{XW}^{\mathbf{Q}}; \quad \mathbf{K} = \mathbf{XW}^{\mathbf{K}}; \quad \mathbf{V} = \mathbf{XW}^{\mathbf{V}}$$

Self-Attention | Transformers

$$\text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V}$$

QK^T Matrix N

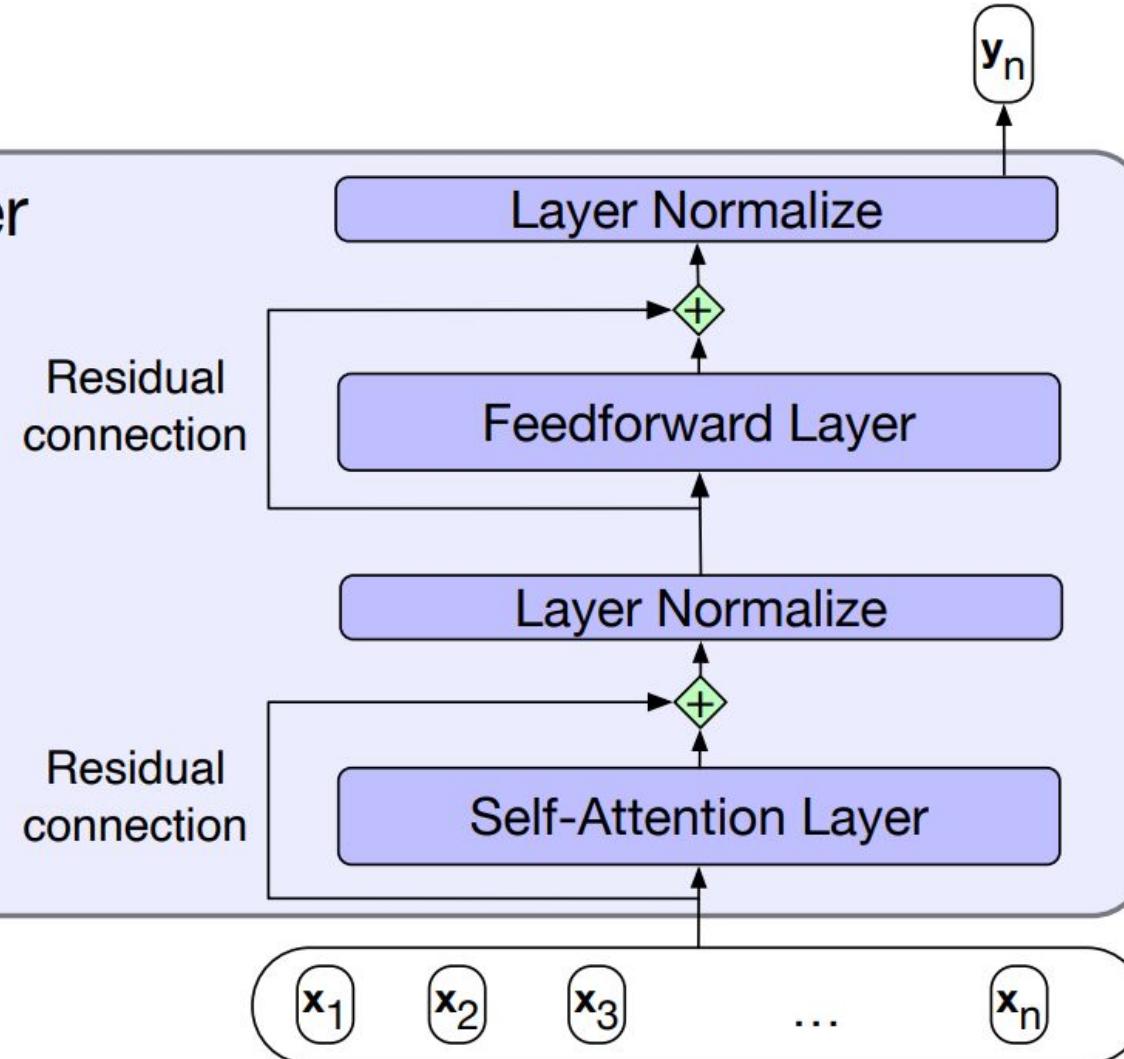
q1•k1	−∞	−∞	−∞	−∞
q2•k1	q2•k2	−∞	−∞	−∞
q3•k1	q3•k2	q3•k3	−∞	−∞
q4•k1	q4•k2	q4•k3	q4•k4	−∞
q5•k1	q5•k2	q5•k3	q5•k4	q5•k5

N

Note: Upper-triangle portion of the comparisons matrix zeroed out (set to $-\infty$, which the softmax will turn to zero)

Transformer Blocks | Transformers

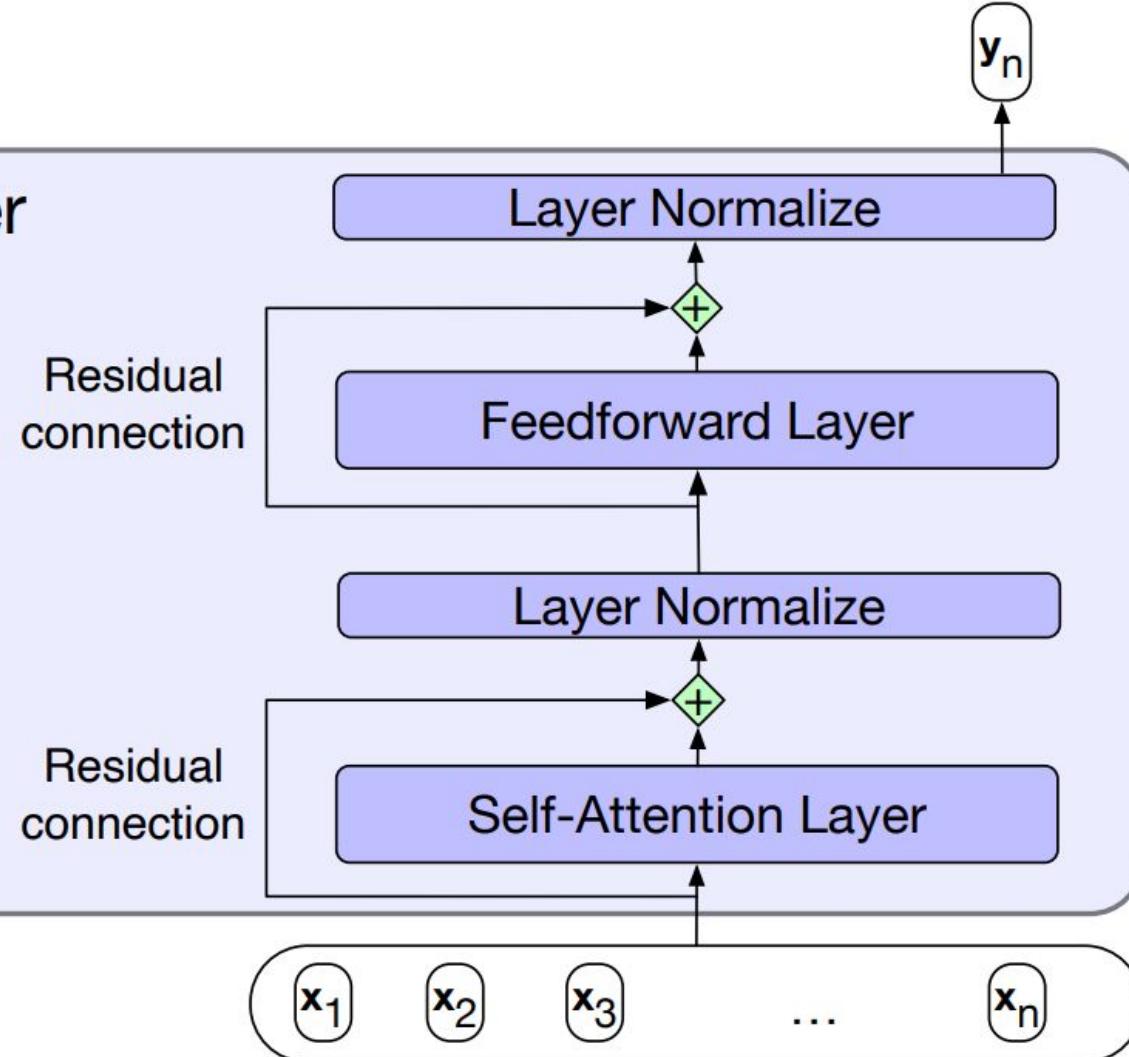
Transformer Block



$$\begin{aligned}\mathbf{z} &= \text{LayerNorm}(\mathbf{x} + \text{SelfAttention}(\mathbf{x})) \\ \mathbf{y} &= \text{LayerNorm}(\mathbf{z} + \text{FFN}(\mathbf{z}))\end{aligned}$$

Transformer Blocks | Transformers

Transformer Block



$$\mathbf{z} = \text{LayerNorm}(\mathbf{x} + \text{SelfAttention}(\mathbf{x}))$$

$$\mathbf{y} = \text{LayerNorm}(\mathbf{z} + \text{FFN}(\mathbf{z}))$$

LayerNorm:

$$\mu = \frac{1}{d_h} \sum_{i=1}^{d_h} x_i$$

$$\sigma = \sqrt{\frac{1}{d_h} \sum_{i=1}^{d_h} (x_i - \mu)^2}$$

$$\hat{\mathbf{x}} = \frac{(\mathbf{x} - \mu)}{\sigma}$$

$$\text{LayerNorm} = \gamma \hat{\mathbf{x}} + \beta$$



Multihead-Attention | Transformers

- Different words in a sentence can relate to each other in many different ways simultaneously
 - >> A single transformer block to learn to capture all of the different kinds of parallel relations among its inputs is inadequate.
- **Multihead** self-attention layers
 - >> **Heads** ⇒ sets of self-attention layers, that reside in parallel layers at the same depth in a model, each with its own set of parameters.
 - >> Each head learn different aspects of the relationships that exist among inputs at the same level of abstraction

Multihead-Attention | Transformers

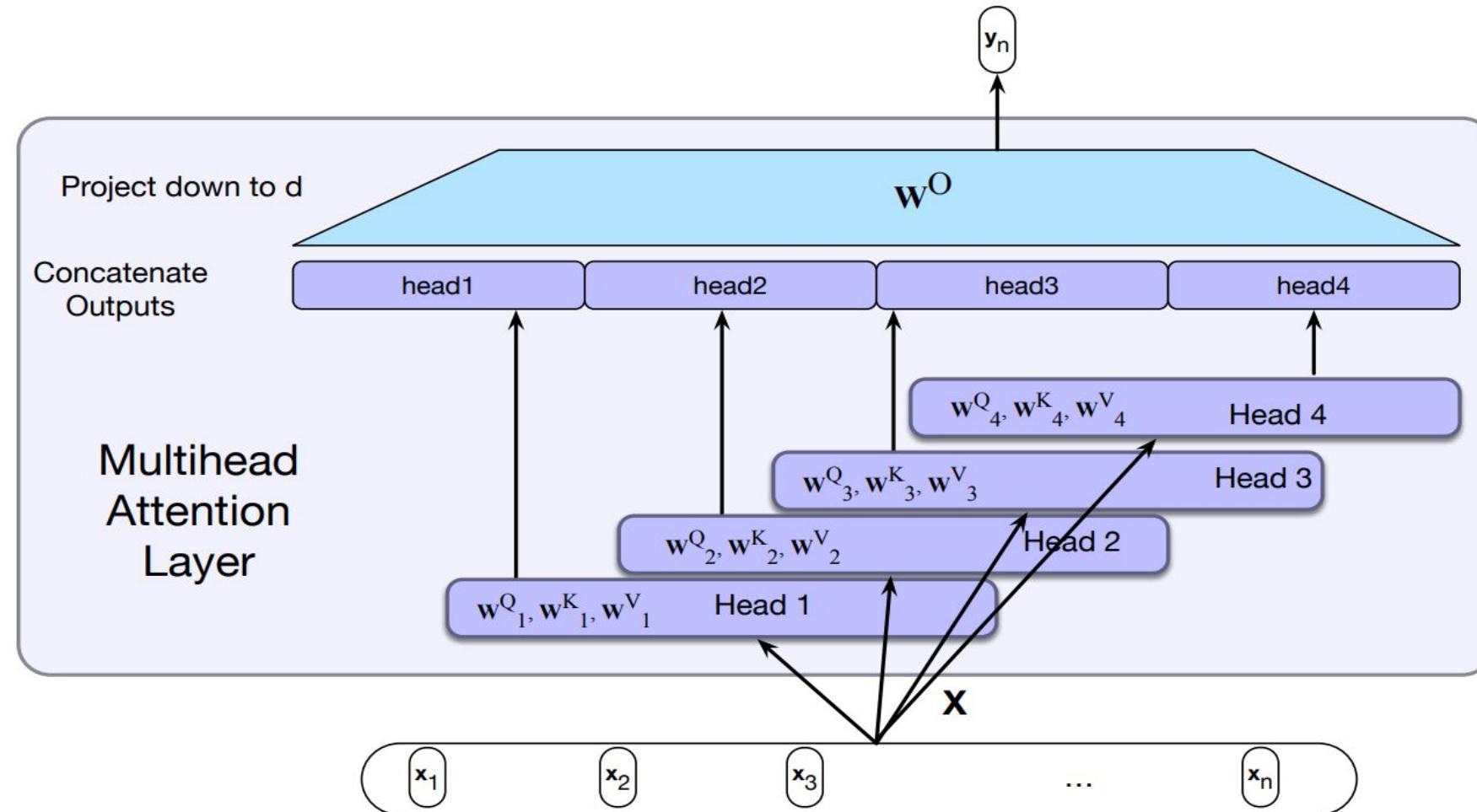
- Different words in a sentence can relate to each other in many different ways simultaneously
 - >> A single transformer block to learn to capture all of the different kinds of parallel relations among its inputs is inadequate.
- **Multihead** self-attention layers
 - >> **Heads** ⇒ sets of self-attention layers, that reside in parallel layers at the same depth in a model, each with its own set of parameters.
 - >> Each head learn different aspects of the relationships that exist among inputs at the same level of abstraction

$$\text{MultiHeadAttention}(\mathbf{X}) = (\mathbf{head}_1 \oplus \mathbf{head}_2 \dots \oplus \mathbf{head}_h) \mathbf{W}^O$$

$$\mathbf{Q} = \mathbf{X} \mathbf{W}_i^Q ; \mathbf{K} = \mathbf{X} \mathbf{W}_i^K ; \mathbf{V} = \mathbf{X} \mathbf{W}_i^V$$

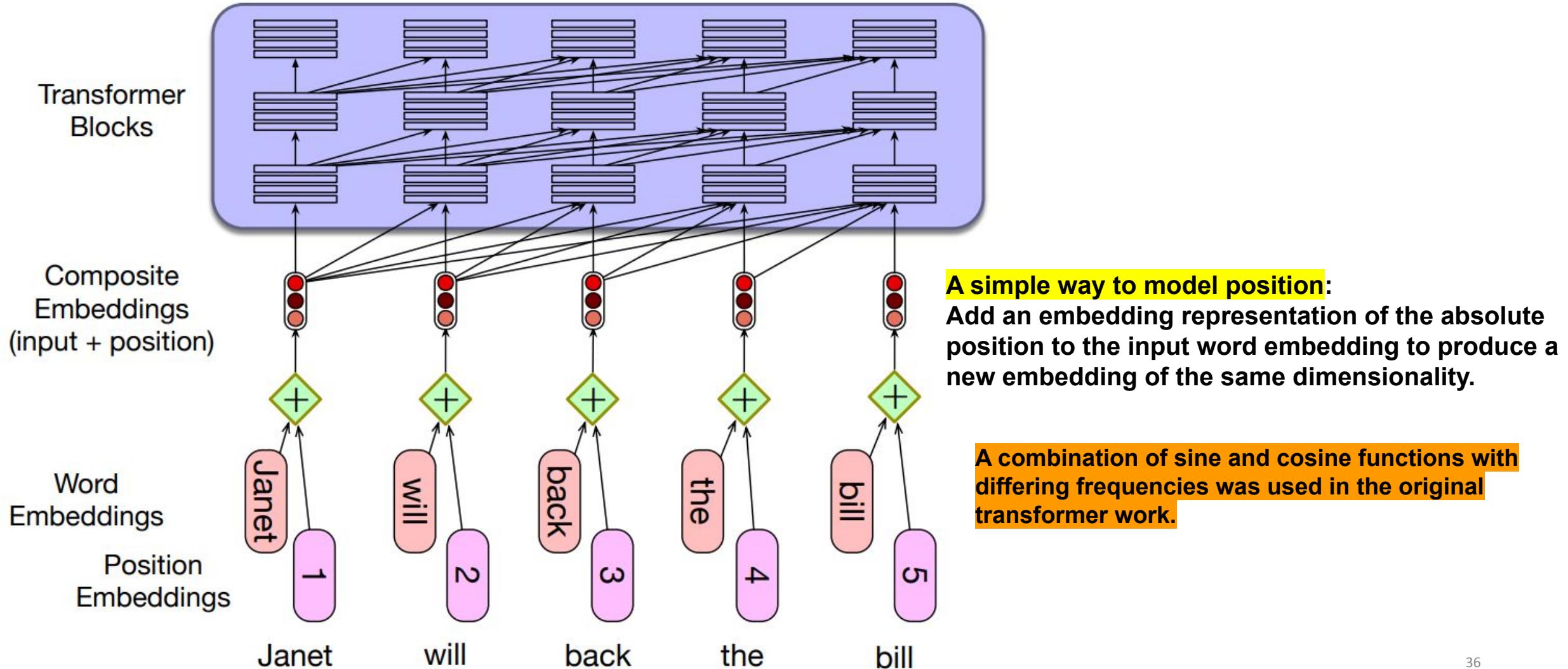
$$\mathbf{head}_i = \text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$$

Multihead-Attention | Transformers

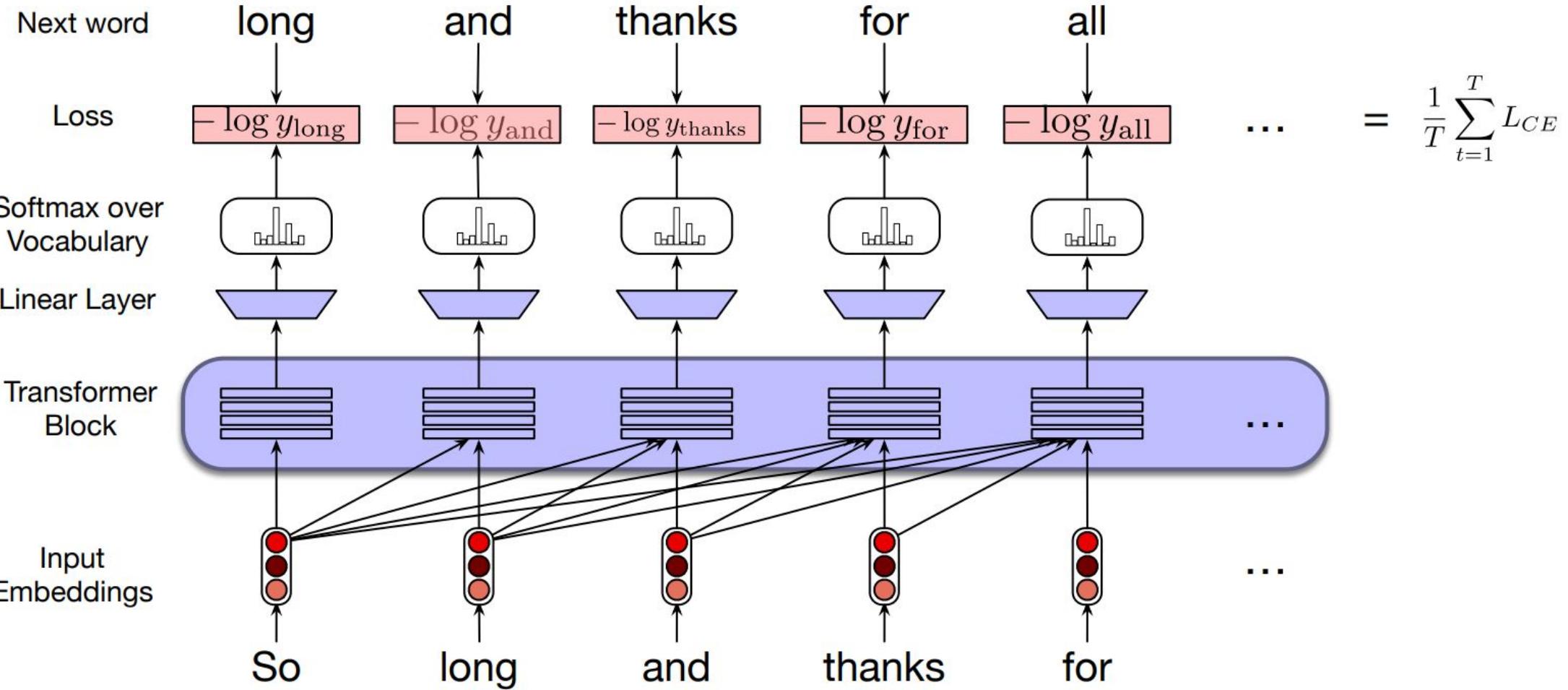


Each of the multihead self-attention layers is provided with its own set of key, query and value weight matrices. The outputs from each of the layers are concatenated and then projected down to d , thus producing an output of the same size as the input so layers can be stacked.

Positional Embeddings | Transformers



Transformers as Language Models





Vision Transformers (ViT)

- Ref: Alexey et. al.(2021) , ICLR, An Image Is Worth 16x16 Words: Transformers For Image Recognition At Scale
-



Readings

1. [Draft Chapter 9](#) and [Draft Chapter 10](#), Speech and Language Processing. Daniel Jurafsky & James H. Martin.
2. Vaswani, Ashish, et al. "Attention is all you need." [Advances in neural information processing systems 30 \(2017\)](#).
3. Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." [arXiv preprint arXiv:2010.11929 \(2020\)](#).



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Thank you



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Neural Network

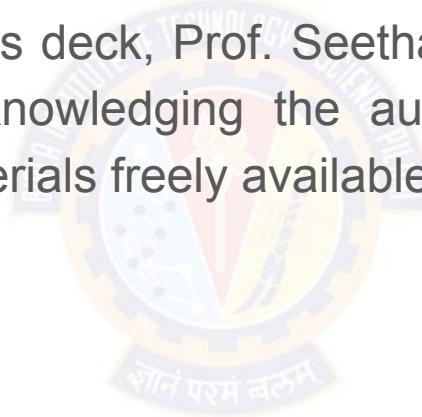
AIML Module 9

Seetha Parameswaran

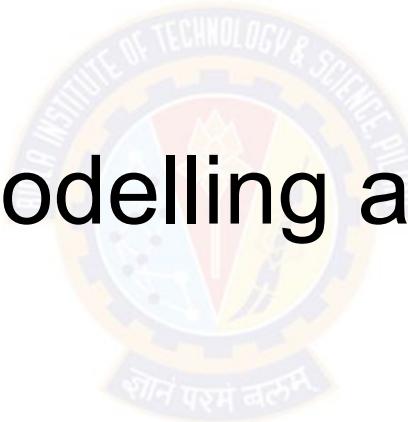
BITS Pilani



The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Time Series Modelling and Forecasting



What we Learn....

9.1 Univariate, Multivariate and Multi-step CNN Models

9.2 Univariate, Multivariate and Multi-step LSTM Models



Time Series Data



Time series Data

- General definition: “**A time series is a collection of observations made sequentially through time, whose dynamics is often characterized by short/long period fluctuations (seasonality and cycles) and/or long period direction (trend)**”.
- Such observations may be denoted by $X_1, X_2, X_3, \dots, X_t, \dots, Y_T$ since data are usually collected at discrete points in time.
 - The interval between observations can be any time interval (seconds, minute, hours, days, weeks, months, quarters, years, etc.) and we assume that these time **periods are equally spaced**.
 - One of the most distinctive characteristics of a time series is the mutual dependence between the observations, generally called **SERIAL CORRELATION OR AUTOCORRELATION**.

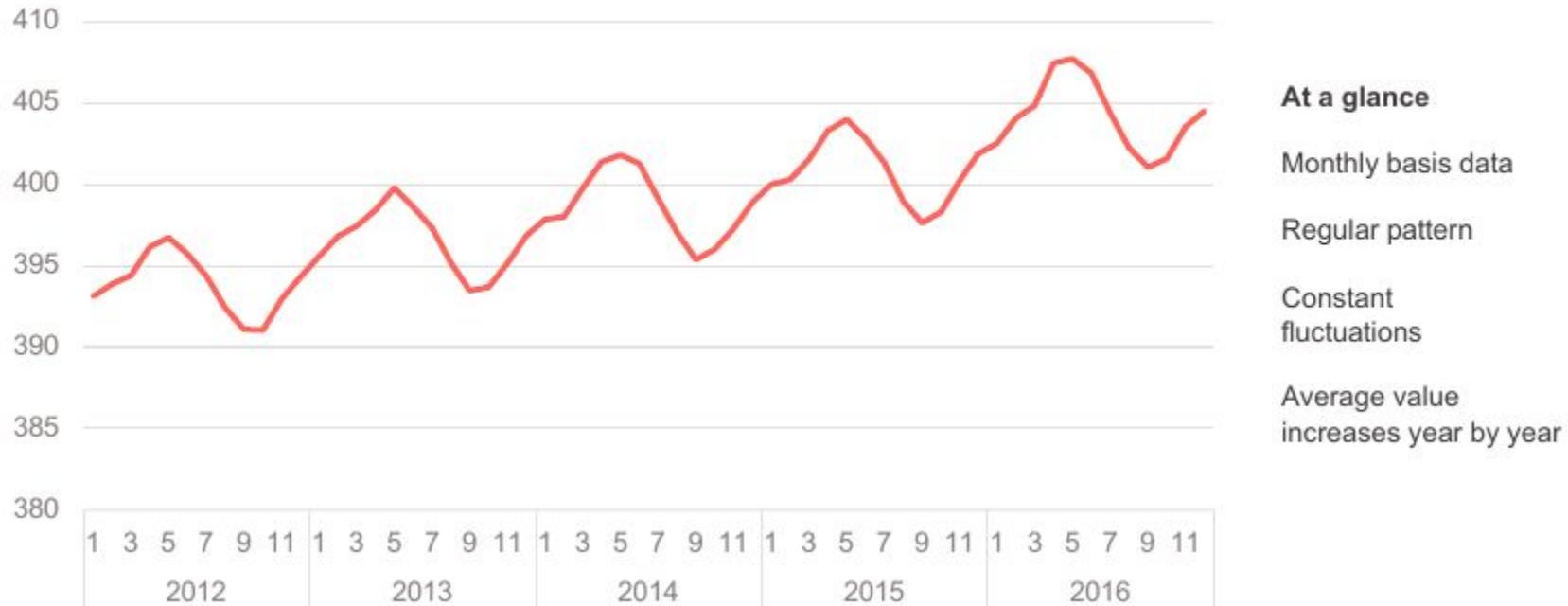
Time series example 1

Numbers of Doctorates Awarded in US, annual data – Engineering Vs. Education



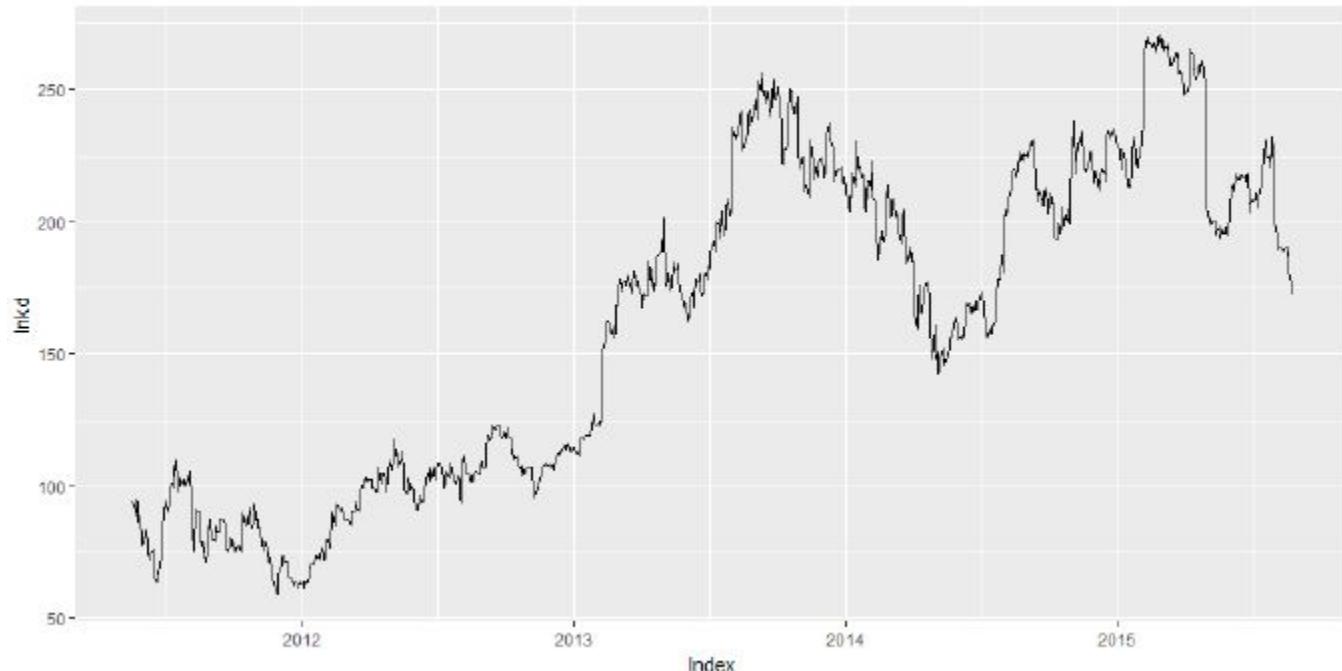
Time series example 2

Monthly carbon dioxide concentration (globally averaged over marine surface sites)



Time series example 3

LinkedIn daily stock market closing price



At a glance

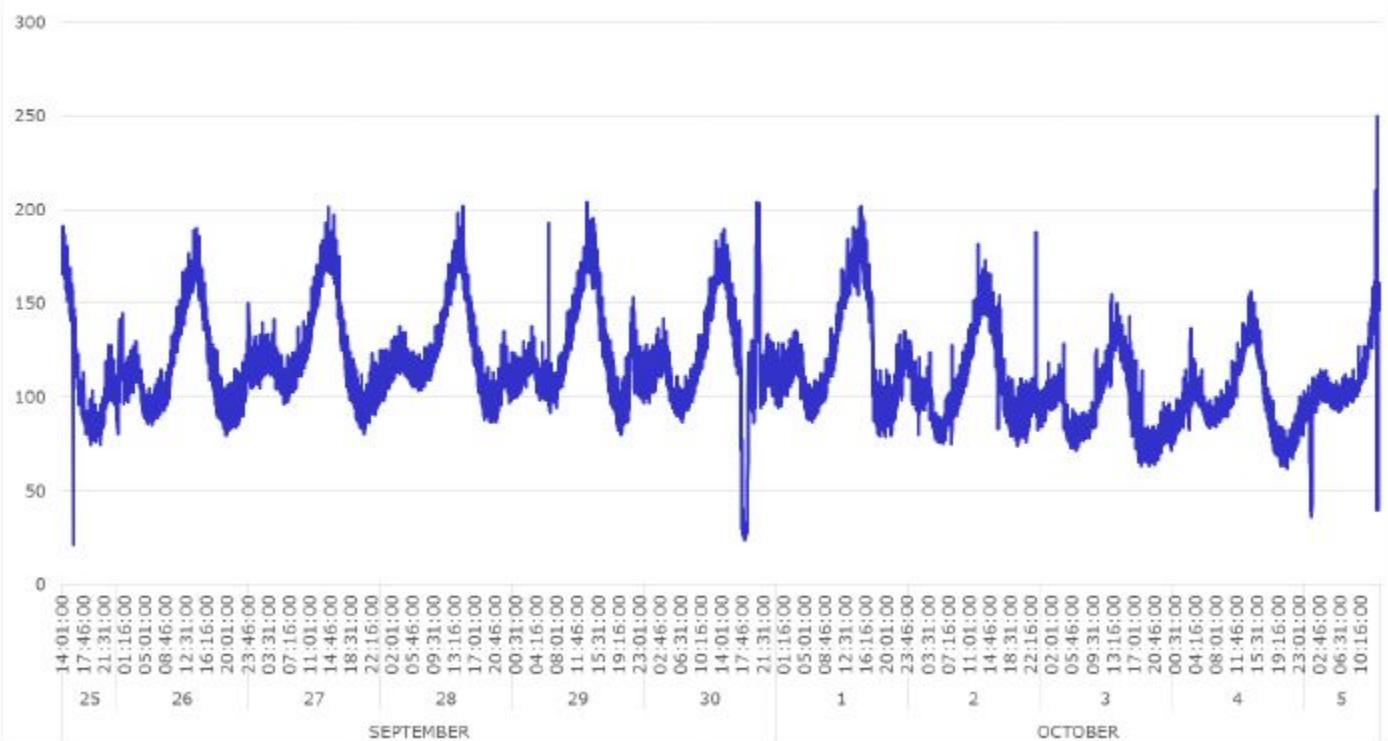
Daily basis data

Very irregular
dynamic

Many sudden
changes

Time series example 4

Number of photos uploaded on the Instagram every minute (regional sub-sample)



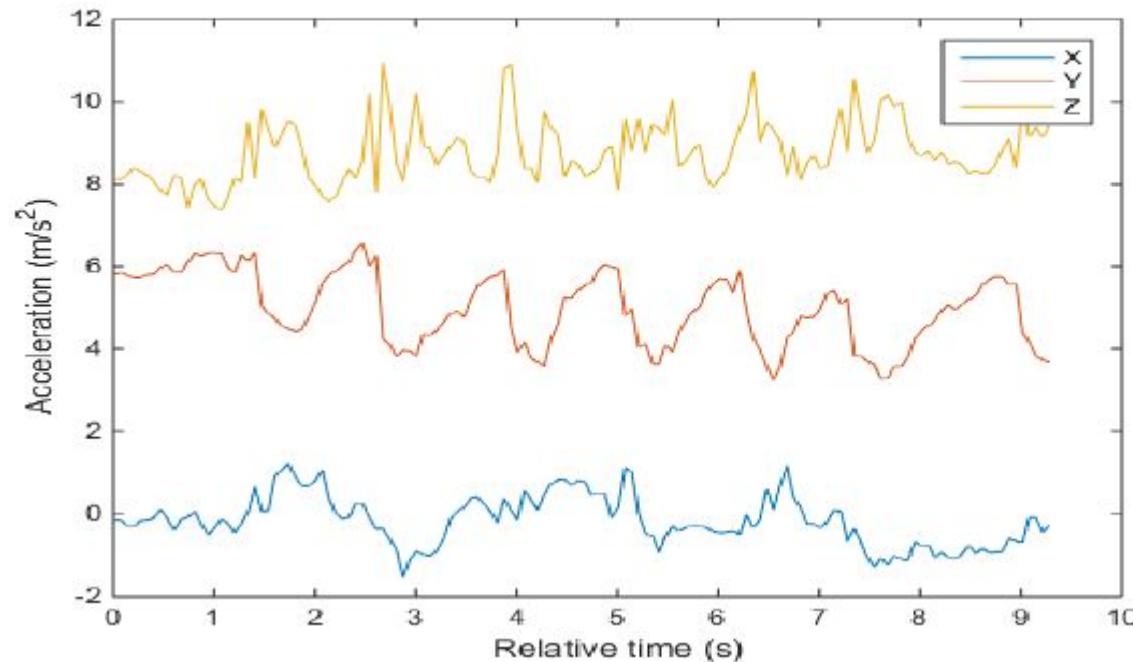
At a glance

Minute basis data

Almost regular daily pattern but with some anomalies and spikes

Time series example 5

Acceleration detected by a smartphone sensors during a workout session (10 seconds)



At a glance

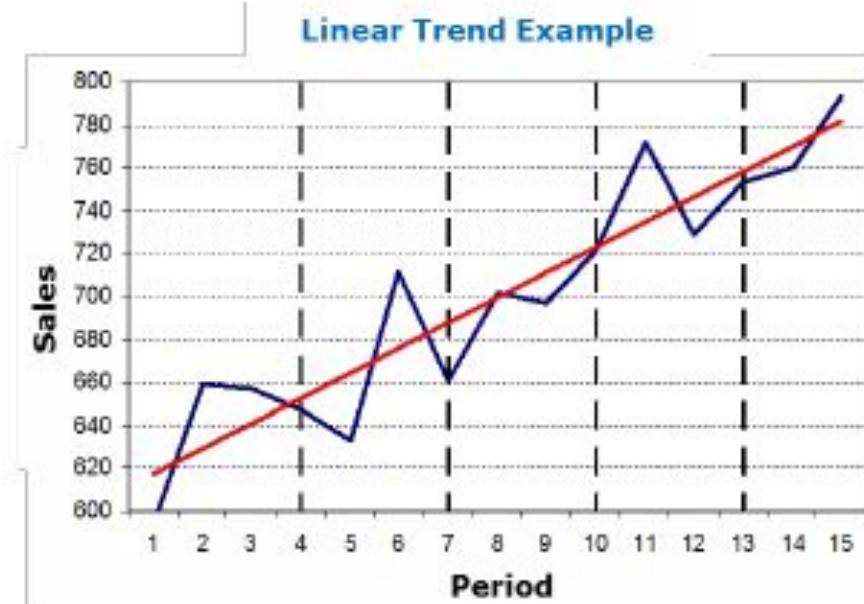
Milliseconds basis
data

Each sensor has its
own dynamics

Time Series Data Properties: Main Elements

1. TREND

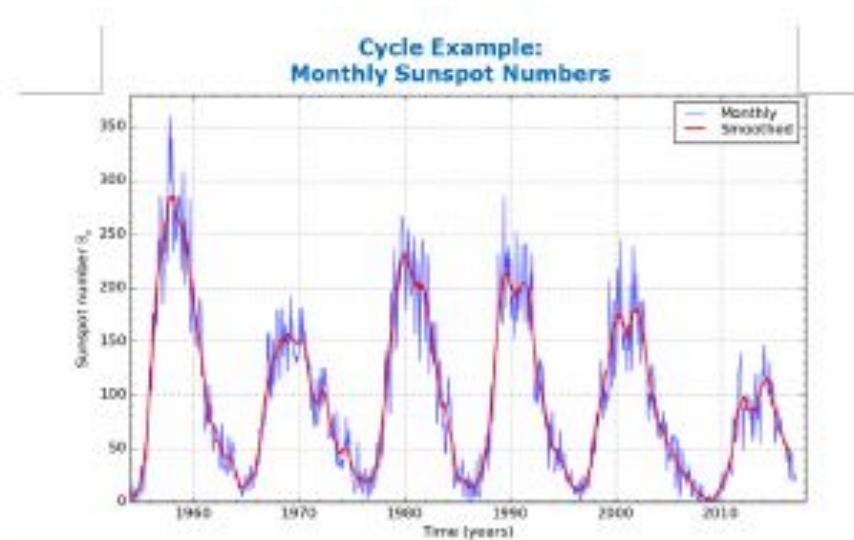
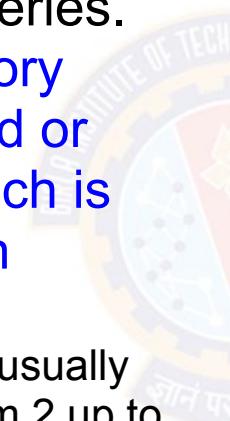
- The general direction in which the series is running during a long period.
- A TREND exists when there is a long-term increase or decrease in the data.
- It does not have to be necessarily linear (could be exponential or others functional form).



Time Series Data Properties: Main Elements

2. CYCLE

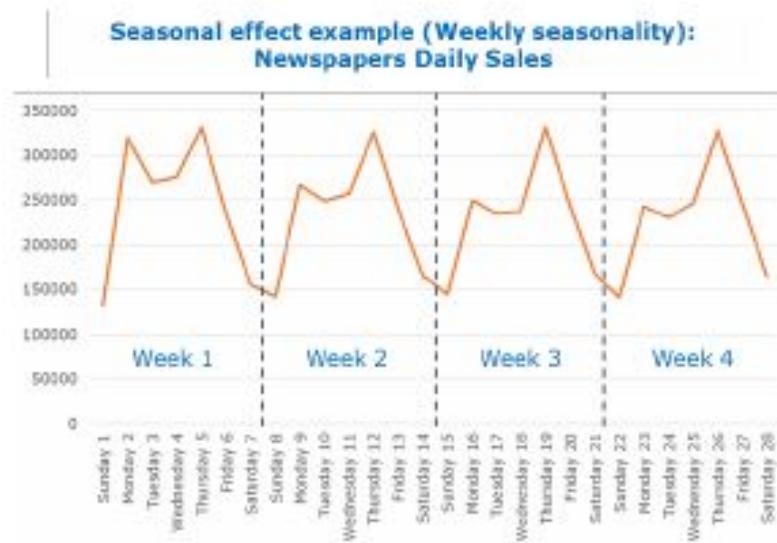
- Long-term fluctuations that occur regularly in the series.
- A CYCLE is an oscillatory component (i.e. Upward or Downward swings) which is repeated after a certain number of years, so:
 - May vary in length and usually lasts several years (from 2 up to 20/30)
 - Difficult to detect, because it is often confused with the trend component.



Time Series Data Properties: Main Elements

3. SEASONAL EFFECTS

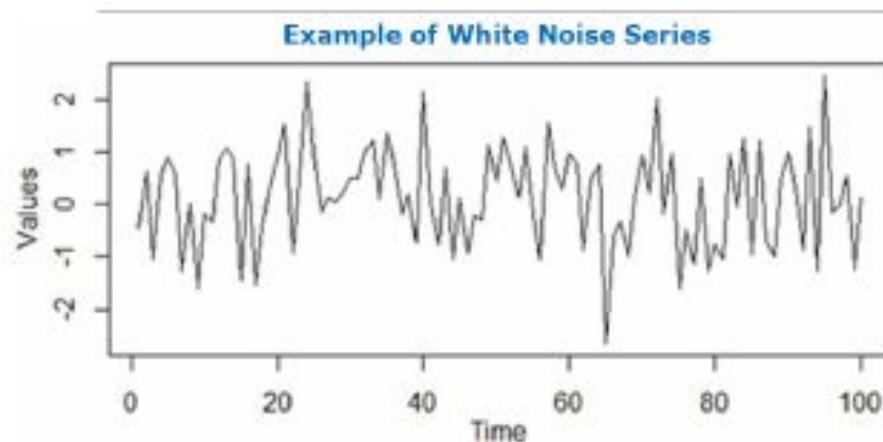
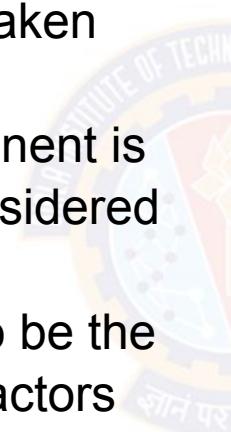
- Short-term fluctuations that occur regularly - often associated with months or quarters.
- A SEASONAL PATTERN exists when a series is influenced by seasonal factors (e.g., the quarter of the year, the month, day of the week).
- Seasonality is always of a fixed and known period.



Time Series Data Properties: Main Elements

4. RESIDUAL

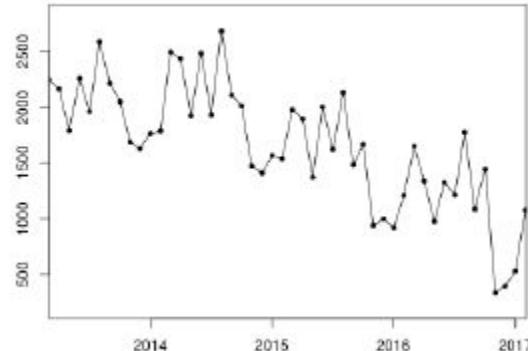
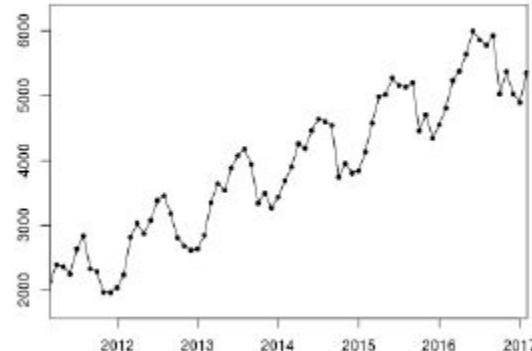
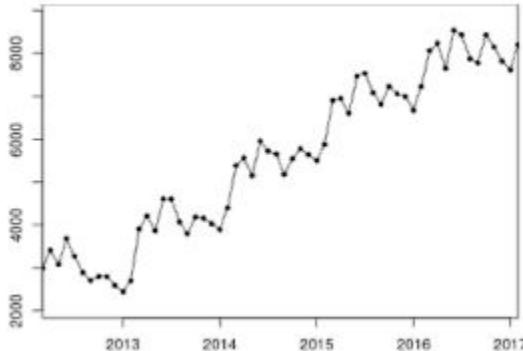
- Whatever remains after the other components have been taken into account.
- The residual/error component is everything that is not considered in previous components.
- Typically, it is assumed to be the sum of a set of random factors (e.g. a white noise series) not relevant for describing the dynamics of the series.



Seasonal effect: Additive Seasonality

- When the seasonality is **Additive**, the dynamics of the components are independent from each other; for instance, an increase in the trend-cycle will not cause an increase in the magnitude of seasonal dips.
- The difference of the trend and the raw data is **roughly constant in similar periods of time** (months, quarters) irrespectively of the tendency of the trend.

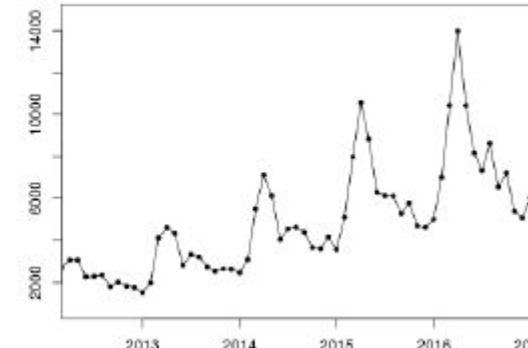
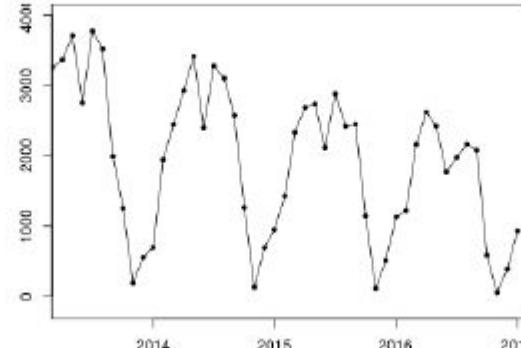
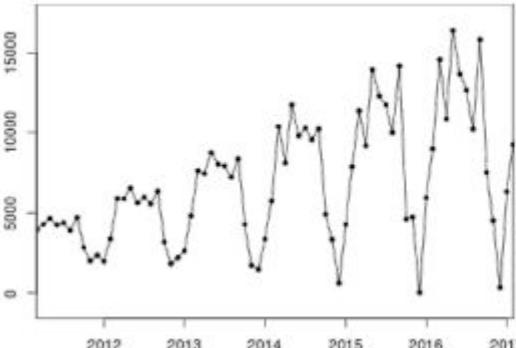
EXAMPLES OF ADDITIVE SEASONALITY



Seasonal effect: Multiplicative Seasonality

- In the multiplicative model the amplitude of the seasonality increase (decrease) with an increasing (decreasing) trend, therefore, on the contrary to the additive case, the **components are not independent from each other**.
- When the variation in the seasonal pattern (or the variation around the trend-cycle) **appears to be proportional** to the level of the time series, then a multiplicative model is more appropriate.

EXAMPLES OF MULTIPLICATIVE SEASONALITY



Seasonal effect: Frequency

- According to the data granularity and to the type of seasonality you want to model, it is important to consider the right seasonal frequency (i.e. how many observations you have for every seasonal cycle)
 - For annual seasonality, the data points are years, quarters, months or weeks (in this case you will face only annual seasonality), but if the frequency of observations is smaller than a week, things get more complicated

Frequency	Cycle type			
	Hour	Day	Week	Year
Data granularity	Annual			1
	Quarterly			4
	Monthly			12
	Weekly		1	52.18
	Daily	1	7	365.25
	Hourly	1	24	168
	Minutes	60	1440	10080
				525960

*Every year, on average, is made up of 365 days and 6 hours → so 365.25 days and $365.25/7=52.18$ weeks

Numerical and Graphical description of Time Series

- The first step in Time Series Analysis is to produce a detailed exploratory analysis of the data to get some insights about the distribution of the series over time.

**Graphical
descriptive
analyses**

- Time plot
- Seasonal plot
- Box plot analysis
- Scatterplots (Lag plots)
- Plotting auto-correlation and cross-correlation functions

**Numerical
descriptive
analyses**

- Sampling period evaluation (start, end, data points features)
- Number of data available
- Missing value and outlier evaluation
- Frequency distribution analysis
- Summary descriptive statistics (overall and by season)

Time Series Data and Analysis

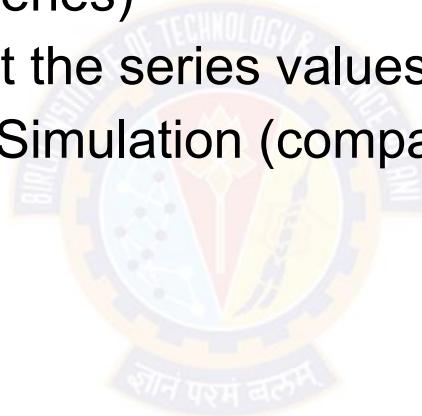


Time Series Analysis

- **Time series analysis** is a specific way of analyzing a sequence of data points collected over an interval of time.
- **Time** is a crucial variable because it shows how the data adjusts over the course of the data points as well as the final results. It provides an additional source of information and a set order of dependencies between the data.
- **Time series analysis** involves developing models to gain an understanding of the data to understand the underlying causes. Analysis can provide the “why” behind the outcomes you are seeing.
- **Forecasting** then takes the next step of what to do with that knowledge and the predictable extrapolations of what might happen in the future.

Main Objectives of Time Series Analysis

- Summary description (graphical and numerical) of data point vs. time
- Interpretation of specific series features (e.g. seasonality, trend, relationship with other series)
- Forecasting (e.g. predict the series values in $(t + 1), (t + 2), \dots, (t + k)$)
- Hypothesis testing and Simulation (comparing different scenarios)

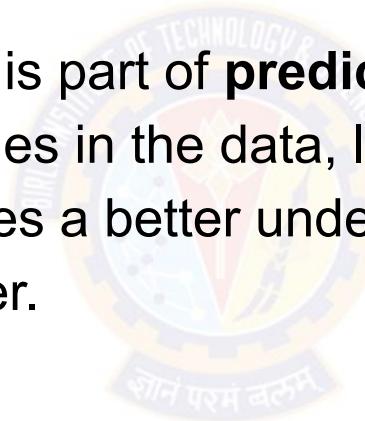


Time Series Forecasting



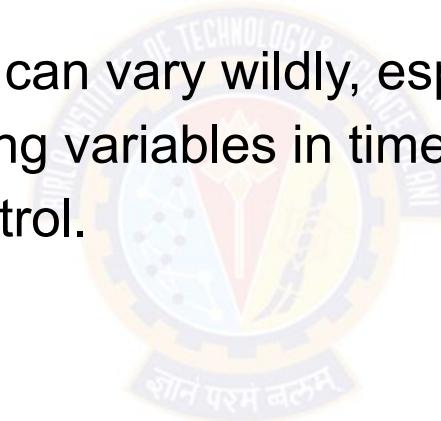
Time Series Forecasting

- **Time series forecasting** is the process of analyzing time series data using statistics and modeling to make predictions and inform strategic decision-making.
- Time series forecasting is part of **predictive** analytics.
- It can show likely changes in the data, like **seasonality or cyclic behavior**, which provides a better understanding of data variables and helps forecast better.



Challenges in Time Series Forecasting

- At the time of the work, the future outcome is completely unavailable and can only be estimated through careful analysis and evidence-based priors.
- Likelihood of forecasts can vary wildly, especially when dealing with the commonly fluctuating variables in time series data as well as factors outside our control.



Applications of Time Series Forecasting

- Weather forecasting
- Climate forecasting
- Economic forecasting
- Healthcare forecasting
- Engineering forecasting
- Finance forecasting
- Retail forecasting
- Business forecasting
- Environmental studies forecasting
- Social studies forecasting



Time series forecasting considerations

1. Time horizons
 - The time frame of your forecast is known as a time horizon.
 - Time horizon is a fixed point in time where the forecast ends.
 - It's much easier to forecast a shorter time horizon with fewer variables than it is a longer time horizon.
2. Dynamic and static states
 - If the forecast is static, it is set in stone once it is made, so make sure your data is adequate for a forecast.
 - Dynamic forecasts can be constantly updated with new information as it comes in. This means you can have less data at the time the forecast is made, and then get more accurate predictions as data is added.

Time series forecasting considerations

3. Data quality

- Typical guidelines for data quality apply.
 - Make sure data is complete,
 - Data is not duplicated or redundant,
 - Data was collected in a timely and consistent manner,
 - Data is in a standard and valid format,
 - Data is accurate for what it is measuring,
 - Data is uniform across sets.
 - Data is collected at consistent intervals over the period of time being tracked. This helps account for trends in the data, cyclic behavior, and seasonality. It also can help identify if an outlier is truly an outlier or if it is part of a larger cycle. Gaps in the data can hide cycles or seasonal variation, skewing the forecast as a result.

Limitations of Classical Methods for TS Forecasting

1. **Focus on complete data:** missing or corrupt data is generally unsupported.
2. **Focus on linear relationships:** assuming a linear relationship excludes more complex joint distributions.
3. **Focus on fixed temporal dependence:** the relationship between observations at different times, and the number of lag observations provided as input, must be diagnosed and specified.
4. **Focus on univariate data:** many real-world problems have multiple input variables.
5. **Focus on one-step forecasts:** many real-world problems require forecasts with a long time horizon.

Strengths of DNN for TS Forecasting

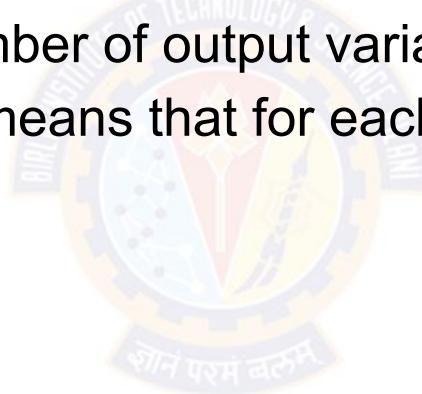
1. **Robust to Noise.** Neural networks are robust to noise in input data and in the mapping function and can support learning and prediction in the presence of missing values.
2. **Nonlinear.** Neural networks do not make strong assumptions about the mapping function and readily learn linear and nonlinear relationships.
3. **Multivariate Inputs.** An arbitrary number of input features can be specified, providing direct support for multivariate forecasting.
4. **Multi-step Forecasts.** An arbitrary number of output values can be specified, providing direct support for multi-step and even multivariate forecasting.

Strengths of DNN for TS Forecasting

5. **Feature Learning in CNN.** Automatic identification, extraction and distillation of salient features from raw input data that pertain directly to the prediction problem that is being modeled.
6. **Native Support for Sequences.** Recurrent neural networks directly add support for input sequence data.
7. **Learned Temporal Dependence.** The most relevant context of input observations to the expected output is learned and can change dynamically.

Limitations of DNN for TS Forecasting

1. The mapping function is fixed or static.
2. Fixed Inputs. The number of lag input variables is fixed, in the same way as traditional time series forecasting methods.
3. Fixed Outputs. The number of output variables is also fixed; although a more subtle issue, it means that for each input pattern, one output must be produced.



Framework for Time Series Forecasting



Framework

Time series forecasting involves **developing and using a predictive model** on data where there is an **ordered relationship between observations**.



Framework

1. What are the inputs and outputs for a forecast?
 - a. Inputs: Historical data provided to the model in order to make a single forecast.
 - b. Outputs: Prediction or forecast for a future time step beyond the data provided as input.
2. What are the endogenous and exogenous variables?
 - a. Endogenous: Input variables that are influenced by other variables in the system and on which the output variable depends.
 - b. Exogenous: Input variables that are not influenced by other variables in the system and on which the output variable depends.
3. Are you working on a regression or classification predictive modeling problem?
 - a. Regression: Forecast a numerical quantity.
 - b. Classification: Classify as one of two or more labels.

Framework

4. Are the time series variables unstructured or structured?
 - a. Unstructured: No obvious systematic time-dependent pattern in a time series variable.
 - b. Structured: Systematic time-dependent patterns in a time series variable (e.g. trend and/or seasonality).
5. Are you working on a univariate or multivariate time series problem?
 - a. Univariate: One variable measured over time.
 - b. Multivariate: Multiple variables measured over time.
 - c. Univariate and Multivariate Inputs: One or multiple input variables measured over time.
 - d. Univariate and Multivariate Outputs: One or multiple output variables to be predicted.

Framework

6. Do you require a single-step or a multi-step forecast?

- a. One-step: Forecast the next time step.
- b. Multi-step: Forecast more than one future time steps.

7. Do you require a static or a dynamically updated model?

- a. Static. A forecast model is fit once and used to make predictions.
- b. Dynamic. A forecast model is fit on newly available data prior to each prediction.

8. Are your observations contiguous or discontiguous?

- a. Contiguous. Observations are made uniform over time.
- b. Discontiguous. Observations are not uniform over time.

Steps to Develop a Forecasting Model

1. Define Problem.
 - a. Inputs vs. Outputs: What are the inputs and outputs for a forecast?
 - b. Endogenous vs. Exogenous: What are the endogenous and exogenous variables?
 - c. Unstructured vs. Structured: Are the time series variables unstructured or structured?
 - d. Regression vs. Classification: Are you working on a regression or classification predictive modeling problem? What are some alternate ways to frame your time series forecasting problem?
 - e. Univariate vs. Multivariate: Are you working on a univariate or multivariate time series problem?
 - f. Single-step vs. Multi-step: Do you require a single-step or a multi-step forecast?
 - g. Static vs. Dynamic: Do you require a static or a dynamically updated model?
 - h. Contiguous vs. Discontiguous: Are your observations contiguous or discontiguous?

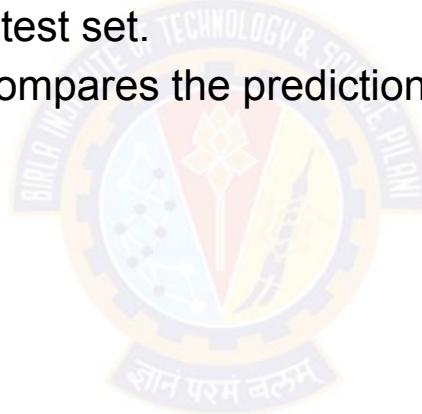
Some useful tools to help get answers include:

- a. Data visualizations (e.g. line plots, etc.).
- b. Statistical analysis (e.g. ACF/PACF plots, etc.).

Steps to Develop a Forecasting Model

2. Design Test Harness

- a. Split the dataset into a train and test set.
- b. Fit a candidate approach on the training dataset.
- c. Make predictions on the test set.
- d. Calculate a metric that compares the predictions to the expected values.



Steps to Develop a Forecasting Model

3. Test Models.: Test many models using your test harness.
 - a. Baseline. Simple forecasting methods such as persistence and averages.
 - b. Autoregression. The Box-Jenkins process and methods such as SARIMA.
 - c. Exponential Smoothing. Single, double and triple exponential smoothing methods.
 - d. Linear Machine Learning. Linear regression methods and variants such as regularization.
 - e. Nonlinear Machine Learning. kNN, decision trees, support vector regression and more.
 - f. Ensemble Machine Learning. Random forest, gradient boosting, stacking and more.
 - g. Deep Learning. MLPs, CNNs, LSTMs, and Hybrid models.

Steps to Develop a Forecasting Model

4. Finalize Model.

- a. List of the top 5 to 10 candidate models that are skillful on the problem.
- b. Pick one or multiple models and finalize them. This involves training a new final model on all available historical data (train and test). The model is ready for use.



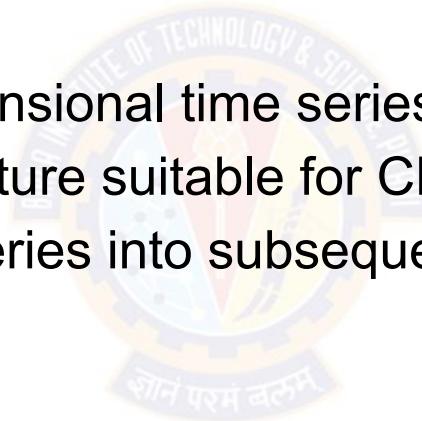
Prepare Time Series Data for DNNs



Prepare Time Series Data for DNNs

Objectives

- Transform a time series dataset into a two-dimensional supervised learning format.
- Transform a two-dimensional time series dataset into a three-dimensional structure suitable for CNNs and LSTMs.
- Split a very long time series into subsequences ready for training a CNN or LSTM model.



Time Series Data to 2D Supervised Learning Data

Univariate Time Series Data

Univariate Time Series Data
for supervised learning

time,	measure
1,	100
2,	110
3,	108
4,	115
5,	120



X,	y
?,	100
100,	110
110,	108
108,	115
115,	120
120,	?

Using sliding window method.

Time Series Data to 2D Supervised Learning Data

Multivariate Time Series Data

time,	measure1,	measure2
1,	0.2,	88
2,	0.5,	89
3,	0.7,	87
4,	0.4,	88
5,	1.0,	90

Multivariate Time Series Data
for supervised learning

X1,	X2,	X3,	y
?,	?,	0.2,	88
0.2,	88,	0.5,	89
0.5,	89,	0.7,	87
0.7,	87,	0.4,	88
0.4,	88,	1.0,	90
1.0,	90,	?,	?



Using sliding window method.

Time Series Data to 2D Supervised Learning Data

Multivariate Time Series Data

time,	measure1,	measure2
1,	0.2,	88
2,	0.5,	89
3,	0.7,	87
4,	0.4,	88
5,	1.0,	90

Multivariate Time Series Data
for sequence prediction
(multi-step)



X1,	X2,	y1,	y2
?,	?,	0.2,	88
0.2,	88,	0.5,	89
0.5,	89,	0.7,	87
0.7,	87,	0.4,	88
0.4,	88,	1.0,	90
1.0,	90,	?,	?

Using sliding window method.

Time Series Data to 3D Supervised Learning Data

- First convert time series data to 2D data for supervised learning.
- Then convert 2D data to 3D data.
 - The input layer for CNN and LSTM models is specified by the input shape argument on the first hidden layer of the network.
 - The input to every CNN and LSTM layer must be three-dimensional.
 - The three dimensions of this input are [samples, timesteps, features]
 - **Samples**. One sequence is one sample.
 - **Time Steps**. One time step is one point of observation in the sample. One sample is comprised of multiple time steps.
 - **Features**. One feature is one observation at a time step. One time step is comprised of one or more features.

Example: Helper function

```
# transform univariate time series to supervised learning problem
from numpy import array

# split a univariate sequence into samples
def split_sequence(sequence, n_steps):
    X, y = list(), list()
    for i in range(len(sequence)):
        # find the end of this pattern
        end_ix = i + n_steps
        # check if we are beyond the sequence
        if end_ix > len(sequence)-1:
            break
        # gather input and output parts of the pattern
        seq_x, seq_y = sequence[i:end_ix], sequence[end_ix]
        X.append(seq_x)
        y.append(seq_y)

    return array(X), array(y)
```

Example: Convert TS data to 2D data

```
# define univariate time series  
series = array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])  
print(series.shape) (10,)  
  
# transform to a supervised learning problem  
X, y = split_sequence(series, 3)  
print(X.shape, y.shape) (7, 3) (7,)  
  
# show each sample  
for i in range(len(X)):  
    print(X[i], y[i])
```

[1 2 3]	4
[2 3 4]	5
[3 4 5]	6
[4 5 6]	7
[5 6 7]	8
[6 7 8]	9
[7 8 9]	10



Example

```
# define univariate time series
series = array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
print(series.shape)
# transform to a supervised learning problem
X, y = split_sequence(series, 3)
print(X.shape, y.shape)
# show each sample
for i in range(len(X)):
    print(X[i], y[i])
```

(10,)
(7, 3) (7,)
(7, 3, 1)

CNNs for Time Series Forecasting



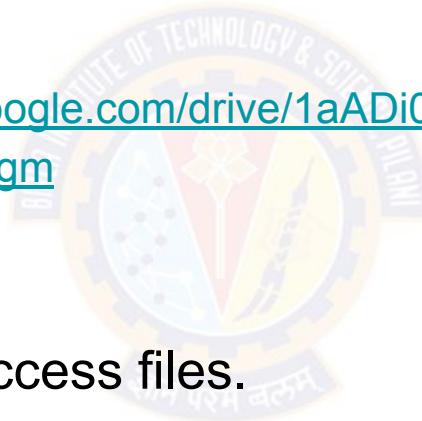
Univariate CNN Models

- Univariate time series are datasets comprised of a single series of observations with a temporal ordering and a model is required to learn from the series of past observations to predict the next value in the sequence.
- The CNN model will learn a function that maps a sequence of past observations as input to an output observation.
- A one-dimensional CNN is a CNN model that has a convolutional hidden layer that operates over a 1D sequence.

Demo Code for CNN

- Univariate CNN
 - <https://colab.research.google.com/drive/1VXCnwxtRdFkDYIxGqhnEuQeXTjQNc0uM#scrollTo=lfms-u8BBrAF>
- Multivariate CNN
 - https://colab.research.google.com/drive/1aADi0xKTExjZSReSFmwTke_Ws4yAwp4k#scrollTo=8kDtsv_uJygm

PS: Use BITS email ID to access files.

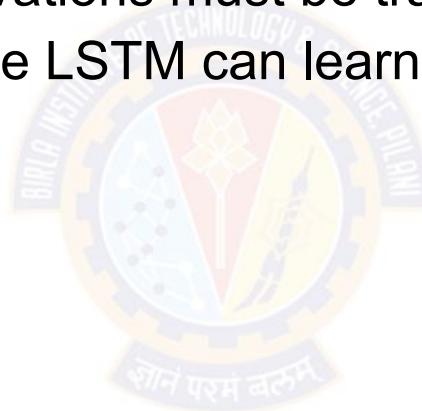


RNNs for Time Series Forecasting



Univariate RNN Models

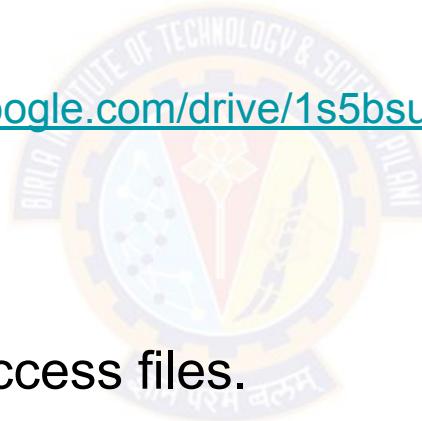
- The LSTM model will learn a function that maps a sequence of past observations as input to an output observation.
- The sequence of observations must be transformed into multiple examples from which the LSTM can learn.



Demo Code for RNN

- Univariate RNN
 - https://colab.research.google.com/drive/1b-ORe5WeYHpTHgcpFKPlowRQcWNLXBq_#scrollTo=vfZ36VoM9Zlh
- Multivariate RNN
 - https://colab.research.google.com/drive/1s5bsuW-wmxyuPkIK_C4mOwc8JPqlMw_wS

PS: Use BITS email ID to access files.



Ref

- Deep Learning for Time Series Forecasting by Jason Brownlee (R4)
- <https://www.tableau.com/learn/articles/time-series-forecasting>

Next Session: CNN





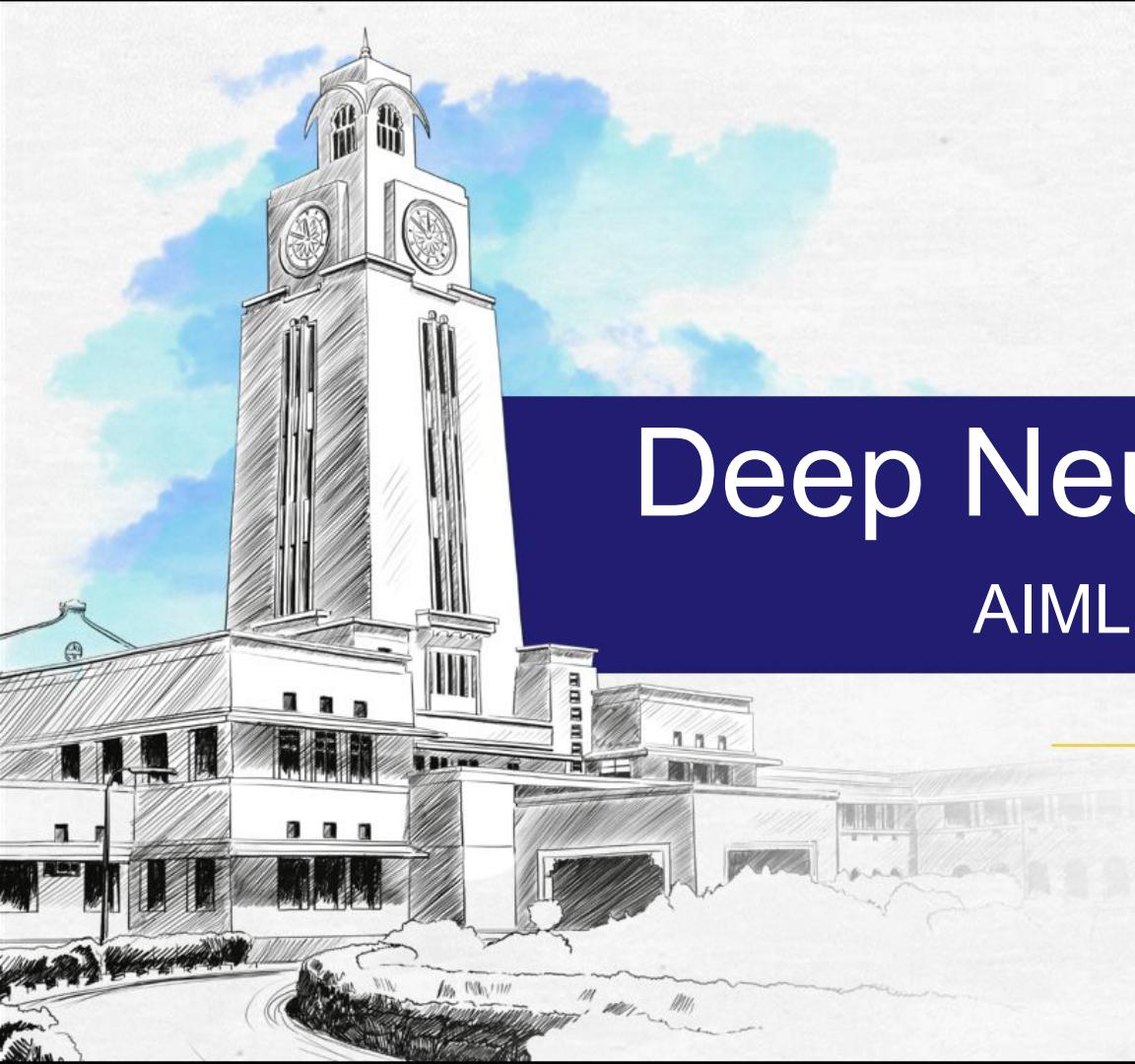
BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Neural Network

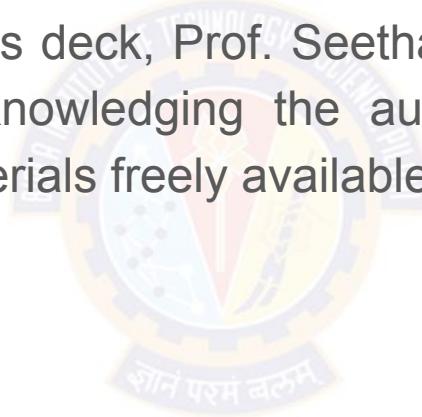
AIML Module 10

Seetha Parameswaran

BITS Pilani



The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



What we Learn....

10.1 Federated learning

10.2 Meta learning

10.3 Online (incremental) learning

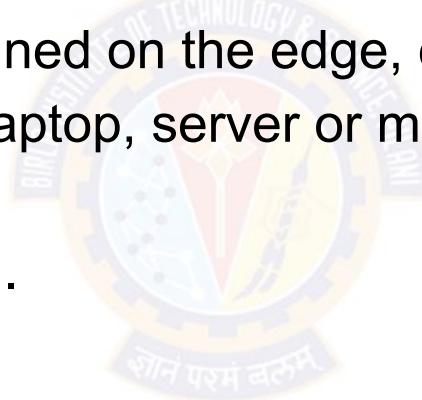


Federated Learning



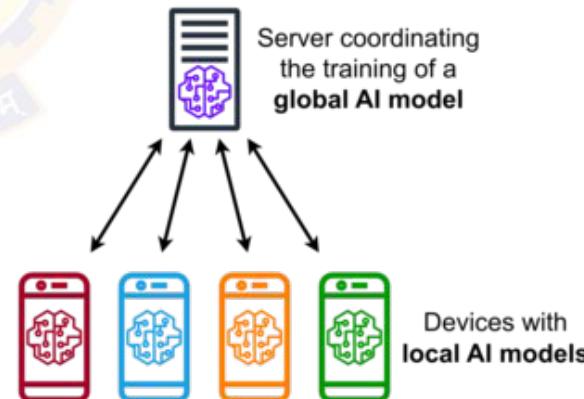
Federated Learning

- Federated learning is a way to train AI models without anyone seeing or touching your data.
- Google introduced the term federated learning in 2016.
- AI models are being trained on the edge, on data that never leave the source of the data like laptop, server or mobile phone or wearable device.
- Decentralized approach.



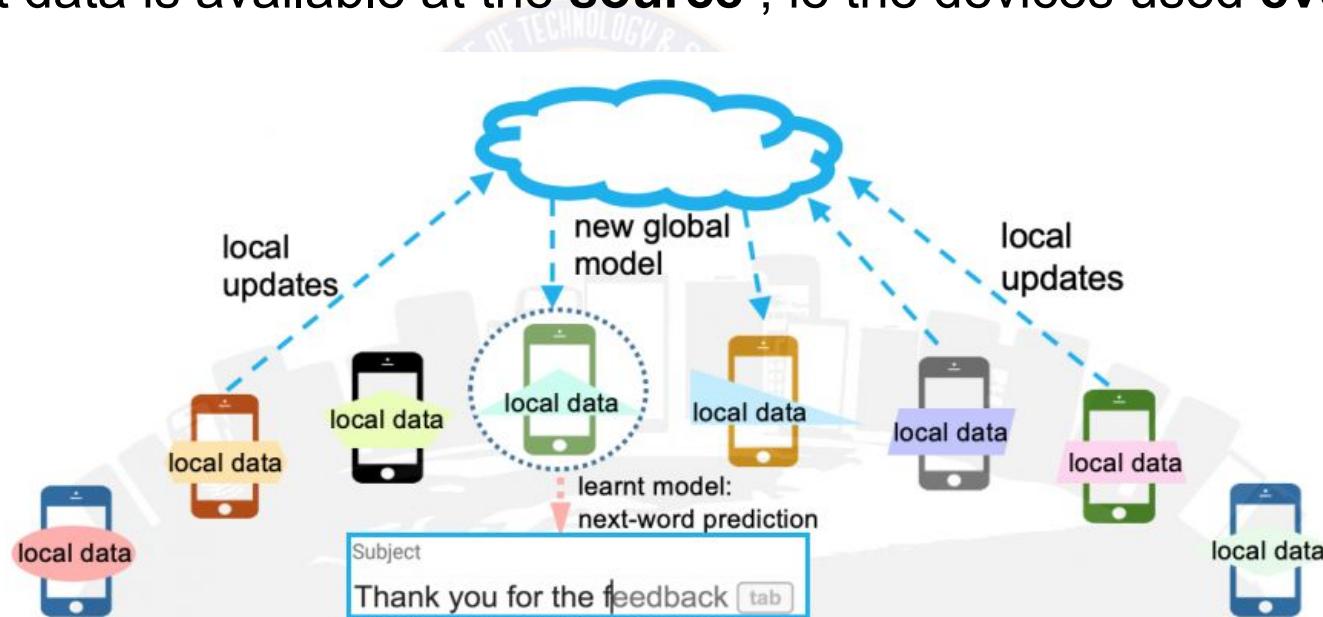
Federated Learning

- Federated Learning is a machine learning approach that allows a model to be trained across multiple decentralized edge devices (such as smartphones, IoT devices, or local servers) holding local data samples, without exchanging them centrally.
- Instead of sending all data to a central server for model training, federated learning keeps data localized, thus addressing privacy, security, and bandwidth concerns.



Federated Learning

- Real world performance of ML model depends on **relevance** of the data used to train it.
- The best data is available at the **source** , ie the devices used **every day**.



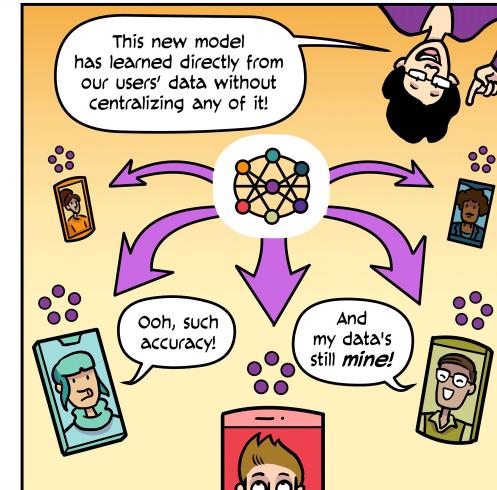
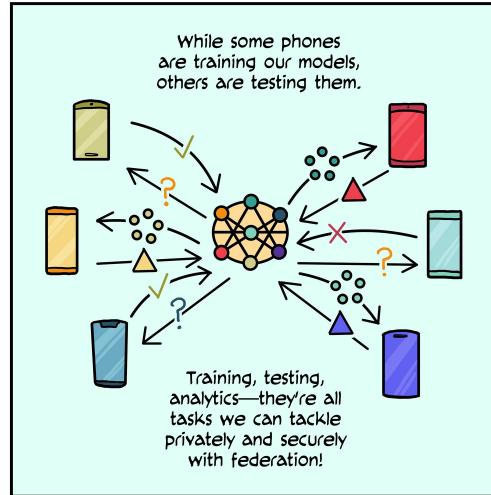
Federated Learning

- Federated Learning = train a centralized model on decentralized data.
 - On-device data is used to train a central model.
Data resides on the device and is not brought to the server on which training is performed.
 - Training is brought to the device. (Some challenges: Battery, wifi, idle, compute)
 - Devices participate when they are eligible. A subset of devices are selected to receive a training model.
The model may be few megs.



Federated Learning

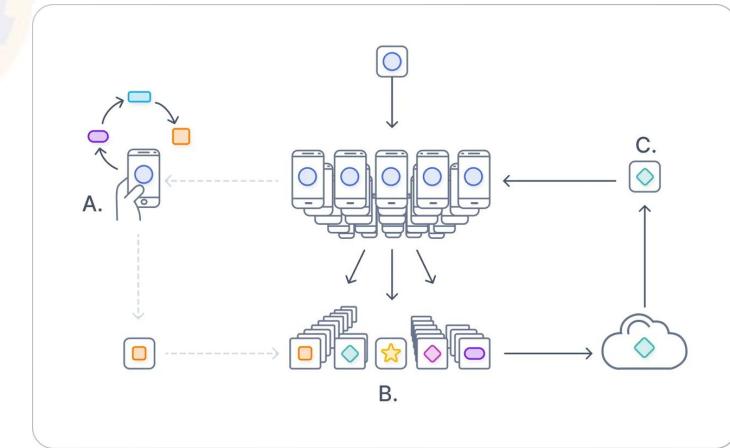
- Federated Learning = train a centralized model on decentralized data.
 - Training on the data on the device in few minutes, sends the training results to the server. The results are encrypted with a key that the server does not have.
 - Uses **Secure aggregation and Differential privacy**
 - ML models are trained on device after multiple rounds of training.
 - Test the ML model on the user devices.
 - New model is **static**. It learned from thousands of users who may have rare unique data.



Types of Federated Learning

1. Centralized federated learning

- Centralized federated learning requires a central server.
- It coordinates the selection of client devices in the beginning and gathers the model updates during training.
- The communication happens only between the central server and individual edge devices.
- Straightforward approach
- Generates accurate models
- Central server poses a bottleneck problem

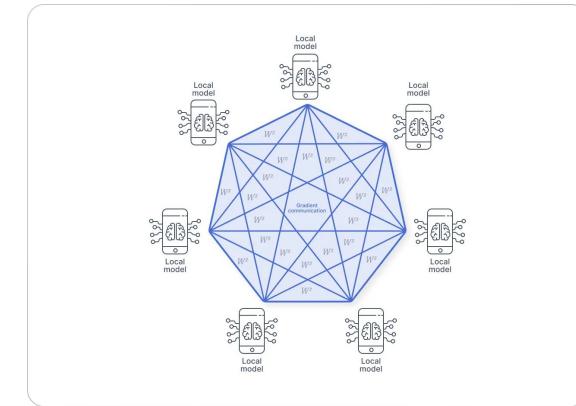


<https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

Types of Federated Learning

2. Decentralized federated learning

- Decentralized federated learning does not require a central server to coordinate the learning.
- The model updates are shared only among the interconnected edge devices.
- The final model is obtained on an edge device by aggregating the local updates of the connected edge devices.
- Prevents the possibility of a single-point failure.
- The model's accuracy is completely dependent on the network topology of the edge devices.

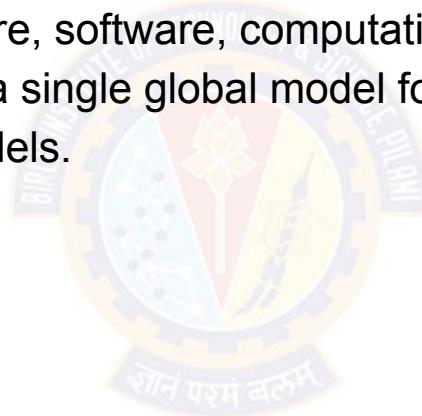


https://www.researchgate.net/figure/Decentralized-Federated-Learning-Environments_fig1_344386575

Types of Federated Learning

3. Heterogeneous federated learning

- Heterogeneous federated learning involves having heterogeneous clients such as mobile phones, computers, or IoT (Internet of Things) devices. These devices may differ in terms of hardware, software, computation capabilities, and data types.
- HeteroFL can generate a single global model for inference after training over multiple varied local models.



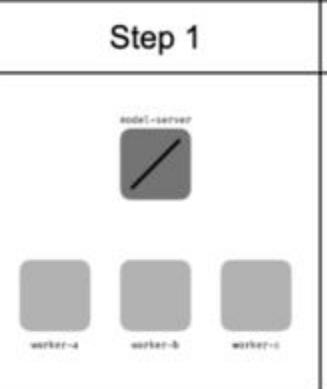
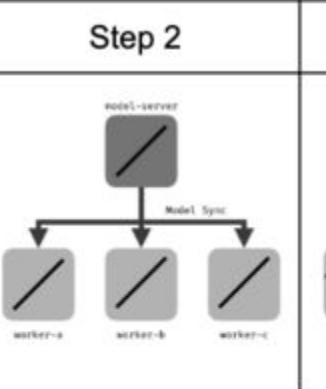
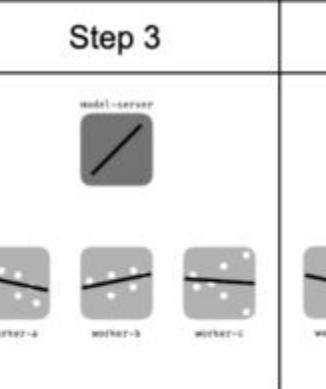
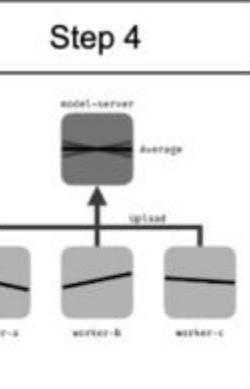
Federated Learning Algorithms

1. Federated stochastic gradient descent (FedSGD)
 - In FedSGD, the central model is distributed to the clients, and each client computes the gradients using local data. These gradients are then passed to the central server, which aggregates the gradients in proportion to the number of samples present on each client to calculate the gradient descent step.
2. Federated averaging (FedAvg)
 - Clients can perform more than one local gradient descent update. Instead of sharing the gradients with the central server, weights tuned on the local model are shared. Finally, the server aggregates the clients' weights (model parameters).
 - If all the clients begin from the same initialization, averaging the gradients is equal to averaging the weights. Therefore, Federated Averaging leaves room for tuning the local weights before sending them to the central server for averaging.
3. Federated learning with dynamic regularization (FedDyn)
 - In federated learning, the global loss must be computed based on local losses generated from heterogeneous devices.
 - Due to the heterogeneity of clients, minimizing global loss is different than minimizing local losses. Therefore, FedDyn method aims to generate the regularization term for local losses by adapting to the data statistics, such as the amount of data or communication cost. This modification of local losses through dynamic regularization enables local losses to converge to the global loss.

Federated Learning

- Secure aggregation enables server to combine the encrypted results and only decrypt the aggregate.
 - On each device, before anything is sent, the secure aggregation protocol adds zero-sum masks to scramble training results.
 - Add up all the training results to cancel out the masks.
- Unique rare data.
 - ML model should not memorize data specific to a device.
 - Add noise to obscure rare data. This is known as differential privacy.
- Differential privacy
 - Deal with the risk of model memorization, where a shared model's parameters may be influenced heavily by a single data.
- Federated Learning = ML model learns from everyone without learning about anyone.

Steps in Federated Learning

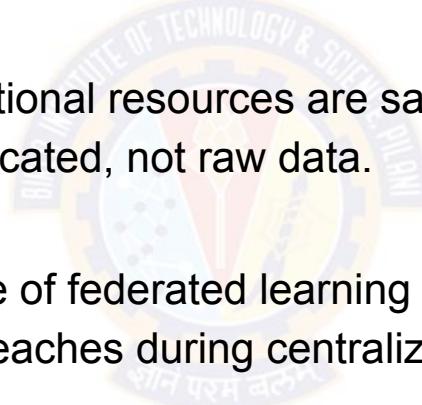
Step 1	Step 2	Step 3	Step 4
			
Central server chooses a statistical model to be trained	Central server transmits the initial model to several nodes	Nodes train the model locally with their own data	Central server pools model results and generate one global mode without accessing any data

Steps in Federated Learning

1. Initialization:
 - a. A global model is initialized on a central server. This model serves as a starting point for training.
 - b. A ML model (e.g., linear regression, neural network, boosting) is chosen to be trained on local nodes and initialized.
2. Client selection
 - a. A fraction of local nodes are selected to start training on local data.
 - b. The selected nodes acquire the current statistical model while the others wait for the next federated round.
3. Local Training:
 - a. Each edge device performs local model updates using its own local data but does not send data samples to the central server.
 - b. The local model is trained on the device using a portion of the data.
 - c. The device computes the gradients of the local model with respect to its data but only shares these gradients with the central server, not the data itself.
4. Aggregation:
 - a. The central server collects gradient updates from all participating devices.
 - b. It aggregates these updates to update the global model.
 - c. Popular aggregation methods include federated averaging, weighted averaging, or other secure and privacy-preserving aggregation techniques.
5. Global Model Update:
 - a. The updated global model is sent back to the participating devices.
6. Local Model Fine-Tuning:
 - a. Devices further fine-tune the global model locally using their data to adapt it to their specific characteristics or changes in data distribution.
7. Repeat:
 - a. Steps 2 to 6 are repeated iteratively for several rounds, allowing the global model to improve based on local data without centralizing it.

Federated Learning - Advantages

- Privacy:
 - User data remains on their devices, reducing privacy concerns associated with centralized data storage and processing.
- Efficiency:
 - Bandwidth and computational resources are saved because only model updates (gradients) are communicated, not raw data.
- Security:
 - The decentralized nature of federated learning makes it more resilient to certain attacks, such as data breaches during centralized data aggregation.
- Personalization:
 - Local fine-tuning allows the global model to adapt to individual device characteristics or data distributions, enabling personalized recommendations or predictions.



Federated Learning

- Use cases for federated learning
 - Transportation: self-driving cars
 - Training self-driving cars on aggregated real-word driver behaviour.
 - Medicine: digital health
 - Help hospitals improve diagnostics while maintaining patient privacy.
 - Industry 4.0: smart manufacturing
 - Improve the efficiency and effectiveness of industrial process while guaranteeing a high level of safety and privacy of sensitive data.
 - Robotics
 - from perception and decision-making to control.
 - Federated Learning is applied to improve multi-robot navigation under limited communication bandwidth scenarios, which is a current challenge in real-world learning-based robotic tasks.
 - Federated Learning is used to learn Vision-based navigation, helping better sim-to-real transfer.

Federated Learning - Challenges

- Dealing with non-IID (non-Independently and Identically Distributed) data
- Communication overhead
- Synchronization issues
- Heterogeneity between the different local datasets: each node may have some bias with respect to the general population, and the size of the datasets may vary significantly
- Temporal heterogeneity: each local dataset's distribution may vary with time
- Interoperability of each node's dataset is a prerequisite
- Each node's dataset may require regular curations
- Hiding training data might allow attackers to inject backdoors into the global model
- Lack of access to global training data makes it harder to identify unwanted biases entering the training e.g. age, gender, sexual orientation
- Partial or total loss of model updates due to node failures affecting the global model
- Lack of annotations or labels on the client side
- Heterogeneity between the processing platform

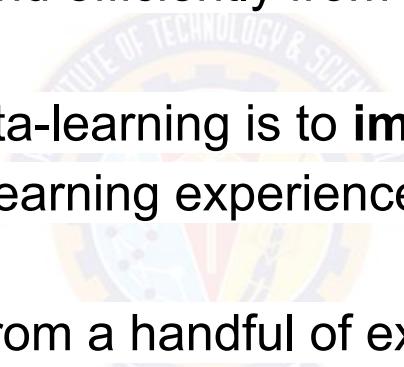
Ref: Federated Learning

- a. <https://federated.withgoogle.com/#:~:text=The%20federated%20learning%20approach%20for,~and%20one%20driven%20by%20AI.>
- b. https://en.wikipedia.org/wiki/Federated_learning
- c. <https://arxiv.org/pdf/1902.04885.pdf>
- d. <https://arxiv.org/pdf/1903.10635.pdf>
- e. Code:
<https://proandroiddev.com/federated-learning-e79e054c33ef>

Meta Learning

Meta-Learning

- **Learning to learn**
- Focus on developing algorithms and techniques that enable models to learn and adapt quickly and efficiently from a limited amount of data or across different tasks.
- The core idea behind meta-learning is to **improve the learning process** itself by leveraging prior learning experiences to tackle new, unseen tasks effectively.
- A model not only learns from a handful of examples, but also learns to classify novel classes during model inference.
- Learning with Limited Labeled Data

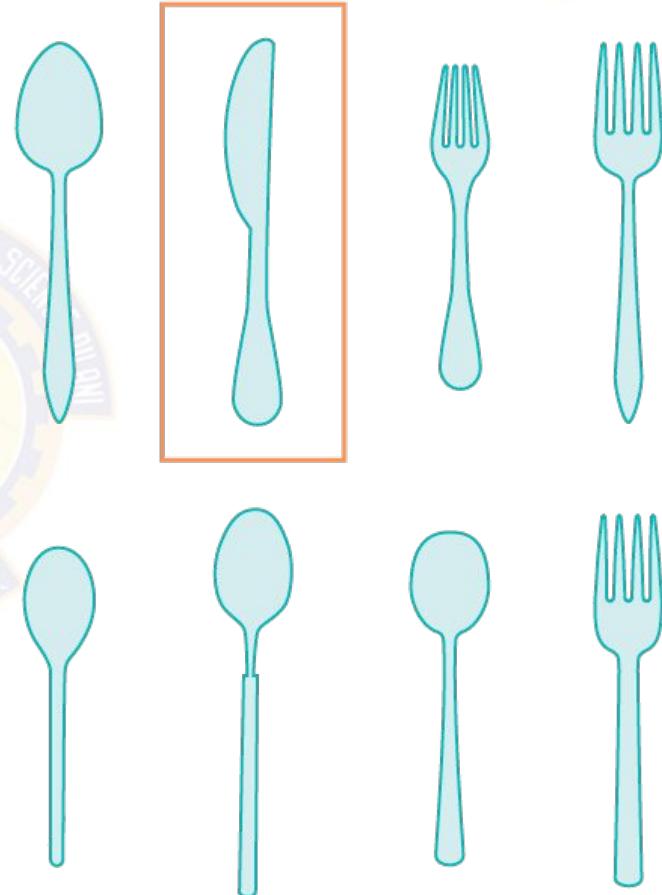
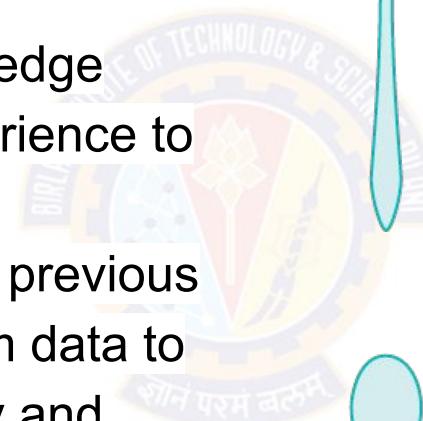


Meta-Learning

	Shallow Learning	Deep Learning	Meta-Learning
Classifier	Learned	Learned	Learned
Feature	Hand-crafted	Learned	Learned
Learner	Hand-crafted	Hand-crafted	Learned

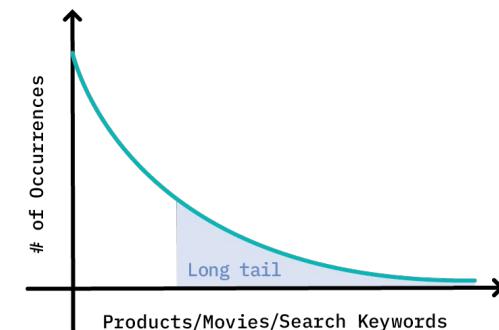
Meta-Learning

- Humans learn quickly.
- Humans adapt to new environments quickly.
- Humans leverage knowledge acquired from prior experience to solve novel tasks.
- Meta-learning leverages previous knowledge acquired from data to solve novel tasks quickly and more efficiently.

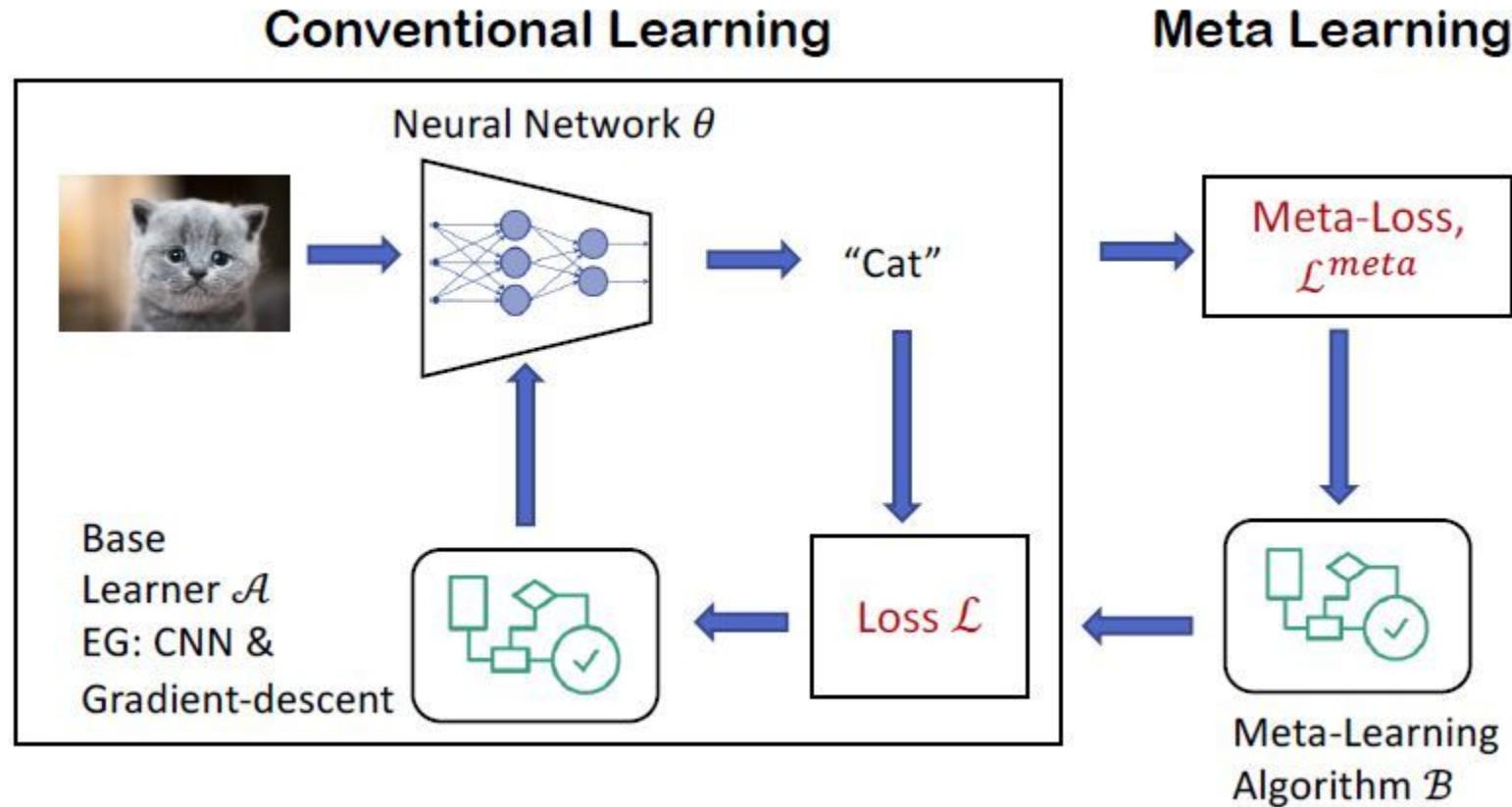


The problem of Generalization

- Supervised learning through deep learning methods requires massive amounts of labeled training data.
- The datasets are expensive to create, especially when one needs to involve a domain expert.
- Pre-training or transfer learning become less effective for domain-specific problems, which still require large amounts of task-specific labeled data to achieve good performance.
- Certain real world problems have long-tailed and imbalanced data distributions.
 - search engines: a few keywords are commonly searched for, whereas a vast majority of keywords are rarely searched for.



Meta-Learning



Meta-Learning

- **Learner** trains a model such as a neural network to solve a task such as visual object recognition.
- **Meta-learner** then wraps the learning process and trains the learner to better solve learning tasks.
- The **learner** typically minimizes a loss function that measures the difference between the true label of an input and the label predicted by a neural network model.
- The **meta-learner** minimizes a meta-loss function such as the generalization error of the trained neural network on novel data.

Why Meta-Learning

1. Ability to learn from a handful of examples
2. Learning or adapting to novel tasks quickly
3. Capability to build more generalizable systems.
4. Applications that require data-efficient approaches
 - a. robots are tasked with learning new skills in the real world, and are often faced with new environments.



Key Aspects of Meta-Learning

1. Learning to Adapt:
 - a. In meta-learning, the goal is to train a model to be highly adaptable, so it can rapidly adapt to new tasks with minimal additional training.
 - b. In traditional machine learning, models are trained for specific tasks and require a substantial amount of labeled data.
2. Few-shot or Zero-shot Learning:
 - a. Meta-learning often deals with scenarios where a model needs to learn from very few examples (few-shot learning) or even no examples (zero-shot learning) for a new task.
 - b. The model leverages knowledge gained from previous tasks to generalize effectively.
3. Task Similarity:
 - a. Meta-learning assumes that there is some degree of similarity or structure among the tasks encountered.
 - b. Models are designed to capture this task similarity and transfer knowledge accordingly.

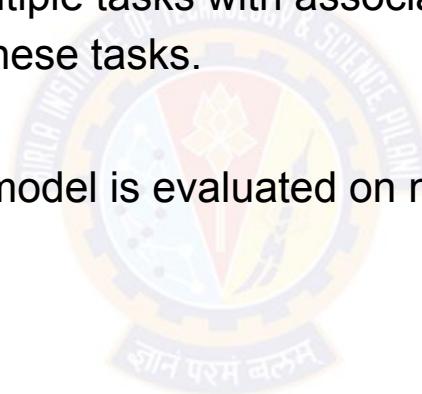
Meta-Learning - Three Approaches

1. Model-based
 - using (cyclic) networks with external or internal memory
 - models updates its parameters rapidly with a few training steps, which can be achieved by its internal architecture or controlled by another meta-learner model.
 - Memory-Augmented Neural Networks - encode new information quickly
 - Meta Networks - learns a meta-level knowledge across tasks
2. Metrics-based
 - learning effective distance metrics
 - Convolutional Siamese Neural Network - twin networks whose output is jointly trained to learn the relationship between input data sample pairs.
 - Matching Networks - learn a network that maps a small labelled support set and an unlabelled example to its label
3. Optimization-based
 - explicitly optimizing model parameters for fast learning
 - LSTM Meta-Learner
 - Model-Agnostic Meta-Learning (MAML) - general optimization algorithm
 - Reptile - meta-learning optimization algorithm

Meta-Data

Two levels of data:

- Meta-Training Data:
 - This data consists of multiple tasks with associated training data.
 - The model learns from these tasks.
- Meta-Test Data:
 - After meta-training, the model is evaluated on new tasks for which it has not seen the training data.

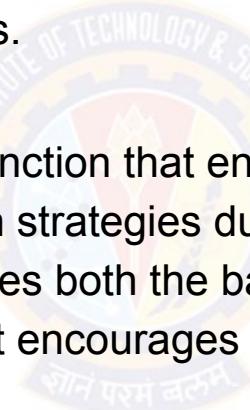


Steps in Meta-Learning

1. Meta-Training:
 - Task Generation:
 - Create a dataset of tasks, each defined by a task description and a few examples.
 - Task Sampling:
 - Randomly select a task from the dataset of tasks.
 - Training:
 - Train a base model (often called the "learner" or "base learner") on the examples provided by the selected task.
2. Meta-Testing:
 - Task Adaptation:
 - Select a new task from the dataset of tasks for which the model was not directly trained.
 - Fine-Tuning:
 - Use the few examples provided by the new task to fine-tune the base model.
 - Inference:
 - Make predictions or perform the task using the adapted model.

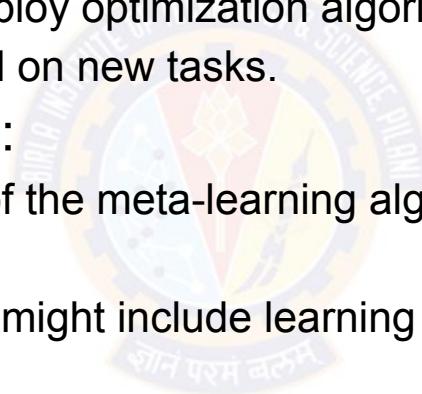
Steps in Meta-Learning

3. Evaluation:
 - Measure the performance of the adapted model on the new task.
 - Repeat the meta-testing process for multiple tasks to assess the model's ability to adapt across a range of tasks.
4. Meta-Learning Objective:
 - Define an objective or loss function that encourages the model to learn useful representations or adaptation strategies during the meta-training process.
 - This objective typically includes both the base model's training loss on each task and a regularization term that encourages the model to adapt well to new tasks during meta-testing.



Steps in Meta-Learning

5. Optimization:
 - Use optimization algorithms, such as gradient descent or reinforcement learning techniques, to update the model's parameters during meta-training.
 - During meta-testing, employ optimization algorithms like gradient descent to fine-tune the base model on new tasks.
6. Hyperparameter Tuning:
 - Tune hyperparameters of the meta-learning algorithm to achieve better adaptation performance.
 - These hyperparameters might include learning rates, regularization strengths, or architectural choices.



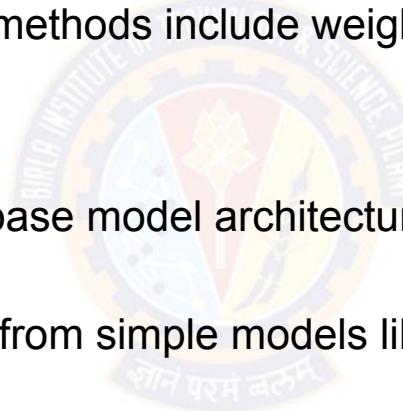
Steps in Meta-Learning

7. Regularization Techniques:

- Apply regularization techniques to prevent overfitting during meta-training and encourage model generalization during meta-testing.
- Common regularization methods include weight decay, dropout, and episodic training.

8. Model Architecture:

- Choose an appropriate base model architecture that can be effectively adapted to various tasks.
- Architectures can range from simple models like linear regressors to complex neural networks.



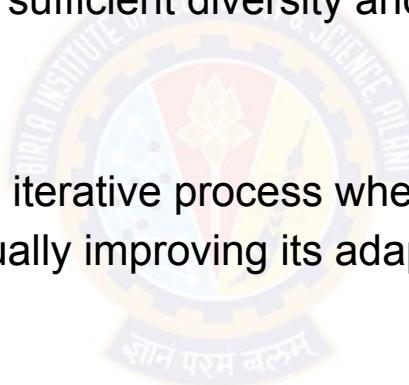
Steps in Meta-Learning

9. Dataset Preparation:

- Curate or generate datasets of tasks and examples that are representative of the problem domain.
- Ensure that tasks exhibit sufficient diversity and complexity to challenge the meta-learning model.

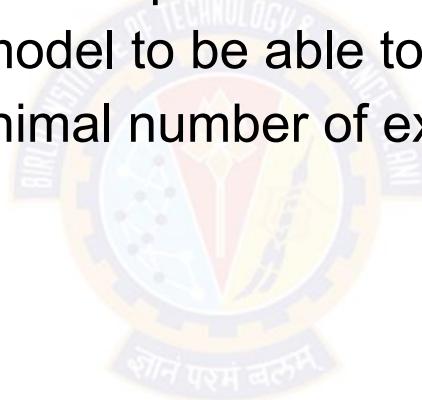
10. Iterative Process:

- Meta-learning is often an iterative process where the model is trained on a sequence of tasks, gradually improving its adaptation capabilities.



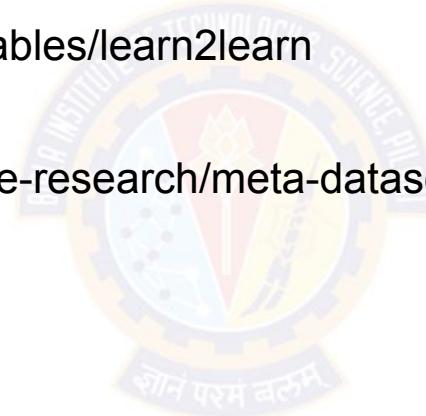
Transfer Learning vs. Meta-Learning

- Both transfer learning and meta-learning involve leveraging prior knowledge
- Transfer learning fine-tunes a pre-trained model on a related task.
- Meta-learning trains a model to be able to adapt to various tasks from scratch, often with a minimal number of examples.



Meta Learning Libraries

- Torch-meta
 - <https://github.com/tristandeleu/pytorch-meta>
- learn2learn
 - <https://github.com/learnables/learn2learn>
- Meta-datasets
 - <https://github.com/google-research/meta-dataset>

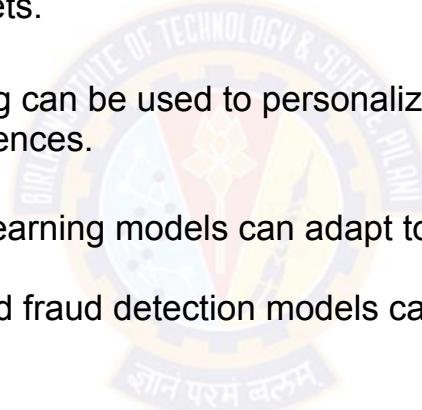


Applications of Meta-Learning

1. Computer Vision:
 - a. Few-Shot Image Classification: Meta-learning can be used to train models that can recognize new object categories with just a few examples.
 - b. Object Detection and Segmentation: It can help in adapting object detection or segmentation models to new classes or domains.
2. Natural Language Processing (NLP):
 - a. Transfer Learning for Text: Meta-learning can be applied to tasks such as sentiment analysis, language modeling, and named entity recognition to adapt models to specific domains or languages.
 - b. Machine Translation: Meta-learning can facilitate adapting translation models to new language pairs.
3. Robotics:
 - a. Robotic Control: Meta-learning can enable robots to adapt their control policies quickly to handle new environments or tasks.
 - b. Object Manipulation: Robots can use meta-learned skills to manipulate objects with different shapes and sizes.
 - c. Autonomous Vehicles:
 - d. Adaptive Driving: Meta-learning can help autonomous vehicles adapt to changing road conditions or new traffic scenarios with limited real-world data.

Applications of Meta-Learning

6. Healthcare:
 - a. Disease Diagnosis: Meta-learning can assist in building diagnostic models that adapt to the characteristics of specific patient populations.
 - b. Drug Discovery: It can facilitate the adaptation of drug discovery models to new chemical compounds or biological targets.
7. Recommendation Systems:
 - a. Personalization: Meta-learning can be used to personalize recommendation algorithms, adapting them to individual user preferences.
8. Finance:
 - a. Financial Forecasting: Meta-learning models can adapt to changing market conditions for improved financial forecasting.
 - b. Fraud Detection: Meta-learned fraud detection models can quickly adapt to new types of fraudulent activities.
9. Education:
 - a. Adaptive Learning: In educational technology, meta-learning can adapt tutoring systems to individual student learning styles and needs.
10. Game Playing:
 - a. Reinforcement Learning: Meta-learned reinforcement learning agents can adapt to new game environments or strategies efficiently.



Applications of Meta-Learning

11. Drug Discovery:
 - a. Molecule Generation: Meta-learning can be applied to generate novel molecules with desired properties by adapting generative models to specific chemical spaces.
12. Meta Reinforcement Learning:
 - a. Robotics Control: In robotics, meta-RL agents can adapt to various tasks or environments by quickly learning new policies.
13. Anomaly Detection:
 - a. Network Security: Meta-learning models can adapt to changing cybersecurity threats and identify anomalies in network traffic.
14. Language Generation:
 - a. Text Generation: Meta-learning can help text generation models adapt to specific writing styles or content domains.
15. Health Monitoring:
 - a. Patient-Specific Monitoring: In healthcare, meta-learning can adapt monitoring models to individual patients' health data.
16. Human-Robot Collaboration:
 - a. Cobotics: Robots can adapt their behavior to collaborate more effectively with humans in various contexts.

Ref: Meta Learning

- a. <https://arxiv.org/pdf/2004.05439.pdf>
- b. <https://meta-learning.fastforwardlabs.com/>
- c. [https://en.wikipedia.org/wiki/Meta-learning_\(computer_science\)](https://en.wikipedia.org/wiki/Meta-learning_(computer_science))
- d. <https://research.samsung.com/blog/Meta-Learning-in-Neural-Networks>

Online / Incremental Learning



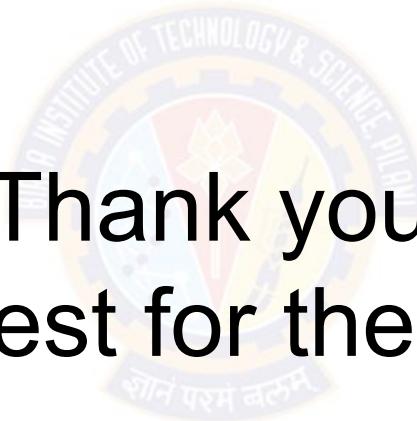
Incremental Learning





Ref: Online Learning

- a. <https://arxiv.org/pdf/1802.07569.pdf>
- b. <https://arxiv.org/pdf/1904.07734v1.pdf>



Thank you
All the best for the future :)



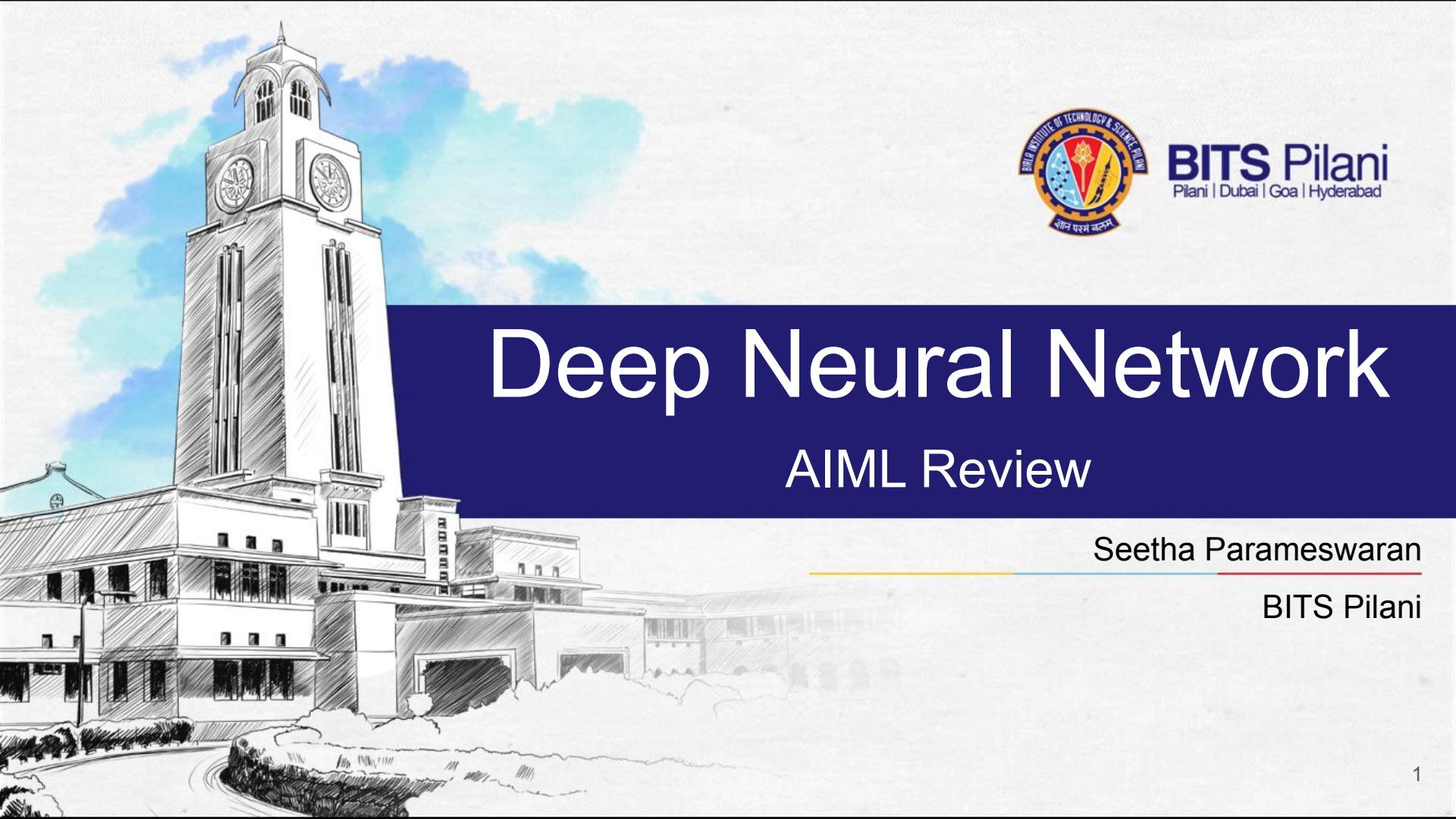
BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Neural Network

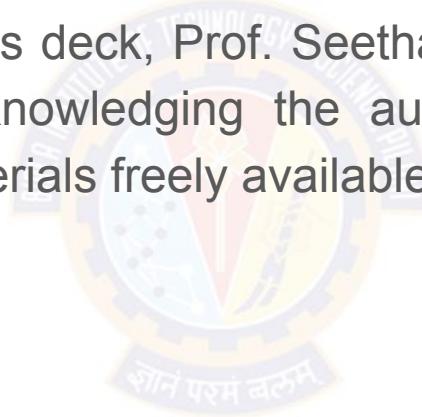
AIML Review

Seetha Parameswaran

BITS Pilani



The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



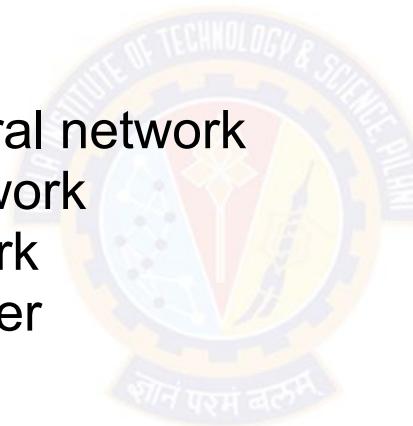
Pattern

Pre Midsem topics - 40%

Post Midsem topics - 60%

Topics

1. Deep Feedforward neural network
2. Convolution neural network
3. Recurrent neural network
4. Attention and transformer
5. Time series forecasting



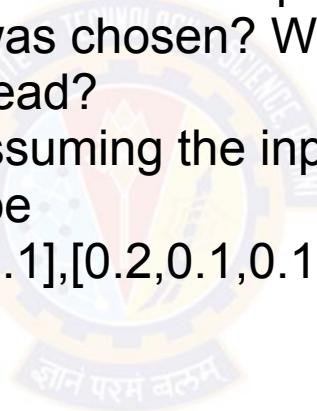
Questions can be conceptual, numerical or descriptive.

**ENSURE THAT YOU ARE HONOURING THE CODE OF CONDUCT
DURING EXAM.**

Question

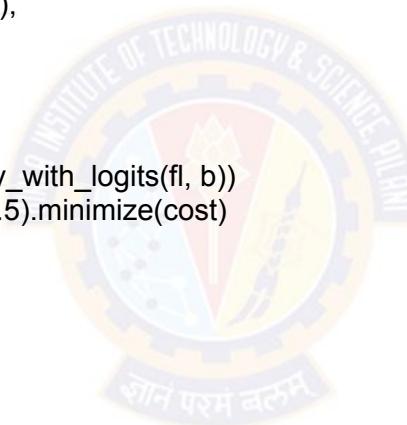
Refer to the partial code of an ANN implementation using Tensorflow and answer following questions:

- a) Which activation function is used in the output layer? Why do you think this particular activation function was chosen? What difference would it make if we used sigmoid function instead?
- b) Calculate the output values assuming the input vector to the output layer to be $[2,0,1,0]$, weight matrix to be $[[0.2,0.3,0.2,0.1],[0.1,0.2,0.5,0.1],[0.2,0.1,0.1,0.1]]$ and bias vector to be $[0.3,0.1,0.1]$



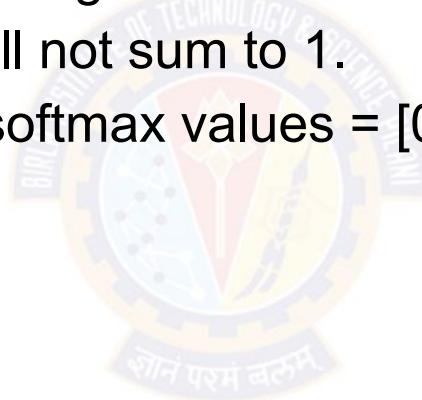
Question

```
import tensorflow as tf
data = pd.read_csv('train.csv')
a = tf.placeholder(tf.float32, shape=[None, 16])
b = tf.placeholder(tf.float32, shape=[None, 3])
phi = { 'alpha': tf.Variable(tf.random_normal([16, 8])),  
'beta': tf.Variable(tf.random_normal([8, 3])) }
omega = { 'alpha': tf.Variable(tf.random_normal([8])),  
'beta': tf.Variable(tf.random_normal([3])) }
tl = tf.add(tf.matmul(x, phi['alpha']), omega['alpha'])
tl = tf.nn.relu(tl)
fl = tf.matmul(tl, phi['beta']) + omega['beta']
cost = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(fl, b))
optimizer = tf.train.AdamOptimizer(learning_rate=0.5).minimize(cost)
init = tf.initialize_all_variables()
```



Answer

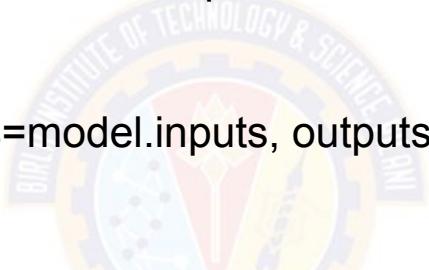
- a) Softmax is used in the output layer. The sum of output values from all the output nodes would be 1 so it is likely that the required output is probability distribution of a given variable. Sigmoid could also be used but the output values will not sum to 1.
- b) $w \cdot x + b = [0.9, 0.8, 0.6]$; softmax values = [0.38, 0.34, 0.28]



Question

Refer to the following code snippet.

```
>>> from tensorflow.keras.applications.vgg16 import VGG16  
>>> from tensorflow.keras.models import Model  
>>> model = VGG16()  
>>> model.layers.pop()  
>>> model = Model(inputs=model.inputs, outputs=model.layers[-2].output)
```

- 
- a) If we add Batch Normalization after every convolution layer in the modified VGG16 model (popping last layer), what will be the total number of additional trainable parameters, beta's and gamma's?
 - b) What will be the total number of non-trainable parameters, i.e., means and variances in the first 3 Batch Normalization layers?

Answer

a) 8448

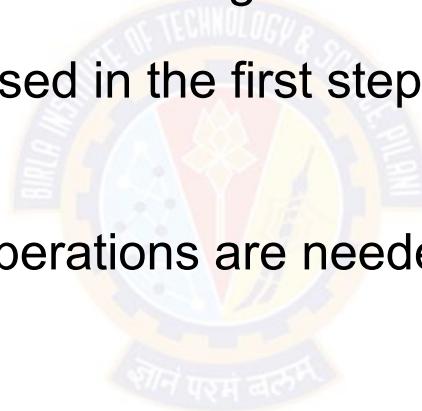
b) $512 (2*64*2 + 128*2)$

Convolution Number	Number of parameters	Number of trainable parameters	Number of parameters per layer	Total number of parameters
64	$2^*64*2 = 256$	$2^*64*2 = 256$	256	512
128	$2^*128*2 = 512$	$2^*128*2 = 512$	512	1,024
256	$3^*256*2 = 1,536$	$3^*256*2 = 1,536$	1,536	3,072
512	$6^*512*2 = 6,144$	$6^*512*2 = 6,144$	6,144	12,288
	8,448	8,448	8,448	16,896

Question

In the following figure, 1×1 operators are used first to process the same 50×50 images of depth 200, and first output 50×50 images of depth 50, and then 5×5 operators are used to output 50×50 images of depth 75.

- a) What is the padding size used in the first step and padding size in the last step?
- b) How many multiplication operations are needed here?



Answer

a) For 1x1 convolution, padding size = 0.

For 5x5 convolution, padding size = 2.

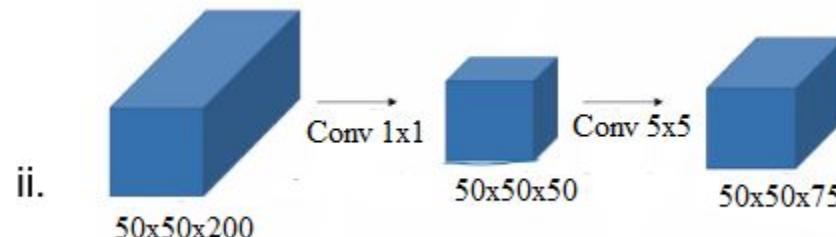
b) For 1x1 convolution,

$$\# \text{ of multiplication operations} = 50 \times 50 \times 50 \times 200 = 2.5 \times 10^7$$

For 5x5 convolution,

$$\# \text{ of multiplication operations} = 50 \times 50 \times 5 \times 5 \times 50 \times 75 = 234375000$$

So, total # of operations = 236,875,000



Question

- A. For what choice of activation functions in hidden and output nodes, autoencoder weights represent the principal component vectors?
- B.. For reconstruction of real inputs, what type of loss function is used in autoencoders?
- C. Gradient ascent is used in training the generative network in GAN. True / False
- D. If reconstructed input are binary (0 or 1), what type of activation function are needed in the output layer of autoencoders?
- E. For reconstruction of binary (0 or 1) inputs, what type of loss function used in autoencoders?
- F. What is the regularization functional used in contractive autoencoder?

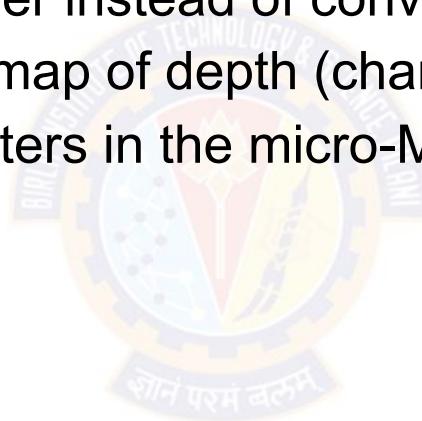
Answer

- A. Linear (i.e., output = input)
- B. Sum of squared error (i.e., SSE)
- C. False
- D. Sigmoid (i.e., $1/(1+\exp(-x))$)
- E. Log-loss / binary cross-entropy
- F. Jacobian



Question

A network-in-network architecture is used to classify gray-scale images of size 64x64 into 100 classes. A micro-MLP with a single hidden layer of 10 nodes is used in the first layer instead of convolution filters of size 5x5 to generate an output feature map of depth (channel) 1. What will be the number of trainable parameters in the micro-MLP based convolution layer?



Answer

Micro-MLP input size: 5*5, Hidden nodes 10, output size 1

Total # of trainable parameters

$$= 5*5*10 \text{ (input-hidden weights)} + 10 \text{ (bias @ hidden)} + 10*1 \\ \text{(hidden-to-output nodes)} + 1 \text{ (bias @ output)}$$

$$= 271$$



Question

Which of the following networks is (are) more suited architecturally for classifying images of varied size (e.g., the input image database has images of size 64x64, 96x96, 128x96, etc.)? No image rescaling is permitted.



Answer

Networks that use global average pooling generate feature vectors whose size is independent of input data, after flattening. Resnet, NiN, Inception all use global average pooling.



Question

Consider an LSTM network with one hidden layer of 20 nodes used for predicting the next word in one hot encoded representation in its output . No bias is used in any of the nodes. The corpus is of length 1000 words and there are 100 unique words. Assume a 10 dimensional word embedding module outside of the LSTM network, whose output is fed to the word predictor LSTM network. What will be the total number of trainable weights in the LSTM network?

Answer

Input size to the network: 10

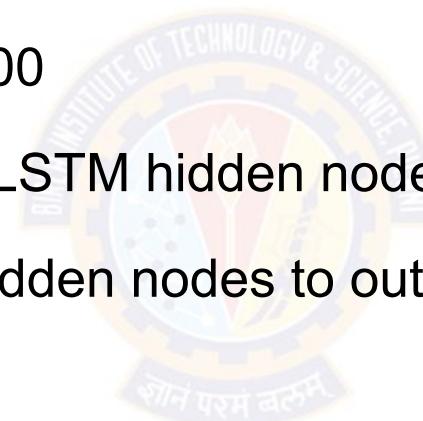
Number of hidden nodes: 20

Number of output nodes: 100

Total weights from input to LSTM hidden nodes = $4 * (20 + 10) * 20 = 2400$

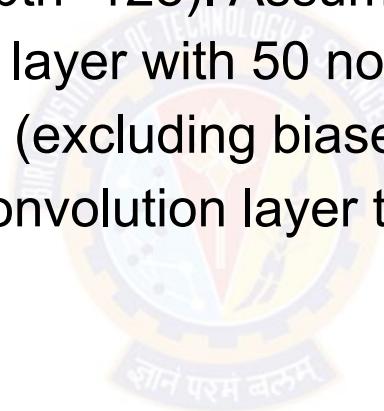
Total weights from LSTM hidden nodes to output nodes = $20 * 100 = 2000$

Total Weights = 4400



Question

In a network in network architecture, one hot output encoding is used to classify 100 classes of objects. The last convolution layer generates $7 \times 7 \times 128$ features maps (depth=128). Assuming the fully connected subnetwork has one hidden layer with 50 nodes, what will be the total number of trainable weights (excluding biases) in the fully connected subnetwork (from the last convolution layer to the output layer)?

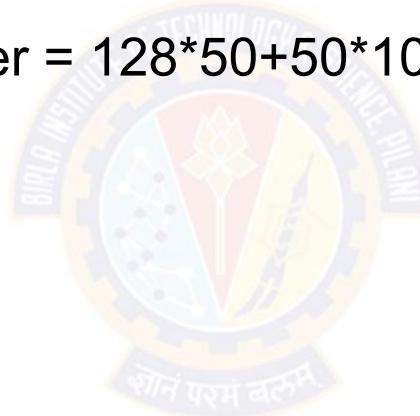


Answer

NiN uses global averaging pooling.

After flattening, feature vector size = 128 hitting the FC layer

Total # of weights in FC layer = $128 \times 50 + 50 \times 100 = 11400$



Consider a vanilla RNN network with one hidden layer of 20 nodes used for predicting the next word in a text corpus. No bias is used in any of the nodes. The corpus is of length 1000 words and there are 50 unique words. Assume a 10 dimensional word embedding module outside of the RNN network, whose output is fed to the word predictor RNN network. One hot encoding is used in the output. What will be the total number of trainable parameters in the RNN network?

Answer

Input-to-hidden weights = $(20+10)*20$

Hidden-to-output weights = $20*50$

Total weights = 1600



Question

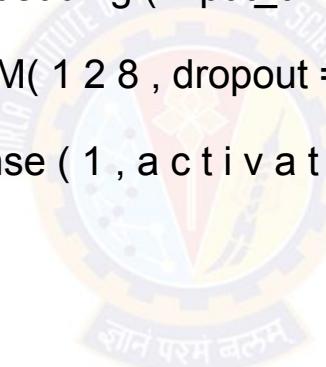
For an image classification problem which classifies images into two categories, 128 by 128 pixel colour images is fed through four convolution layers with 32, 64, 128, 256 kernels of 5 by 5 respectively with max pooling for each layer. Then the tensors are fed through two layers of fully connected neurons with 1024 and 512 neurons.

- (a) Draw a CNN for the above.
- (b) Write Tensorflow-Keras code snippet for the above.
- (c) Compute the number of parameters learned in each convolution and maxpooling layers.

Question

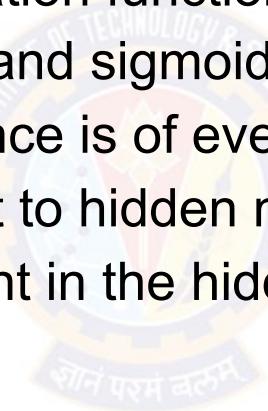
Explain the working of the code snippet given below.

```
model = keras.models.Sequential()  
model.add(keras.layers.Embedding(input_dim = 20000, output_dim = 128))  
model.add(keras.layers.LSTM(128, dropout = 0.2))  
model.add(keras.layers.Dense(1, activation = 'sigmoid'))  
model.summary()
```



Question

Consider a RNN where the hidden state equation is $h(t+1) = W h(t) + x(t+1)$. Calculate and plot the output $y(t)$ of a single hidden node and single input/output RNN, where activation function used in hidden and output nodes are, respectively, linear and sigmoid. Assume that the initial hidden state is 0, and the input sequence is of even length. Assume also that all biases are 0. Weight from input to hidden node is 1 and hidden to output node is 10. The recurrent weight in the hidden node is -1. Input is 1 if $t=0$ or even, and 0 for odd t .



Answer

$$h_0 = h(\text{initial}) + x_0 = 0 + 1 = 1$$

$$y(0) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * 1) = \text{sig}(10) = 1$$

$$h_1 = -h_0 + x_1 = -1 + 0 = -1$$

$$y(1) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * -1) = \text{sig}(-10) = 0$$

$$h_2 = -h_1 + x_2 = -1(-1) + 1 = 2$$

$$y(1) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * 2) = \text{sig}(20) = 1$$

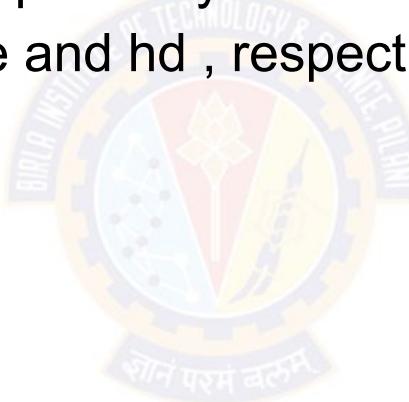
$$h_3 = -h_2 + x_3 = -2 + 0 = 2 = -2$$

$$y(1) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * -2) = \text{sig}(-20) = 0$$

The value of y will be close to 1 when t is even and close to 0 when t is odd. The value will oscillate from 1 to 0 and 0 to 1 for each time stamp starting from $t=0$. The graph will oscillate between 0 and 1.

Question

Consider an RNN based French to Spanish language translation model where attention mechanism is employed. The maximum sentence length in the corpus is I_f and I_s , respectively. The number of hidden nodes in encoder and decoder are h_e and h_d , respectively. How many attention weights will be there?



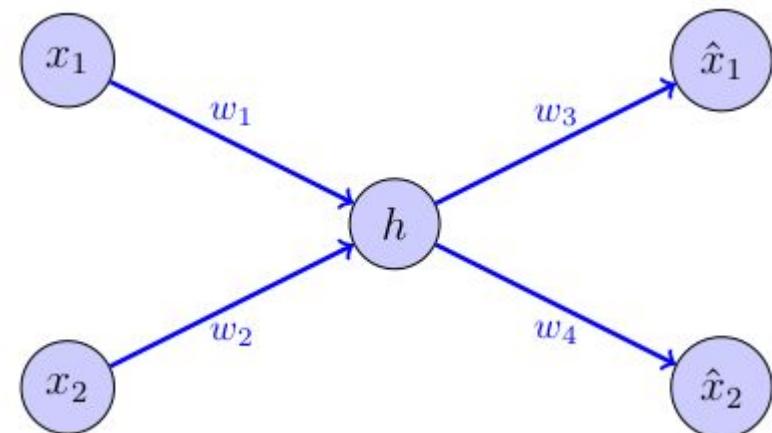
Answer

No of attention weights = Max input seq len * max output seq len = $I_f * I_s$

Question

Consider an autoencoder architecture along with the input vector $x_1 = 0$, $x_2 = 1$ and at (t) , $w_1 = 0$, $w_2 = 1$, $w_3 = 1$ and $w_4 = 0$. Assume learning rate = 0.1 and momentum constant = 0.7. At $(t - 1)$, $w_1 = -0.3$, $w_2 = 0.6$, $w_3 = -0.6$ and $w_4 = 0.3$. Assume derivative of ReLU(z)=0 at $z=0$.

- (a) At iteration t , what will be value of the loss function at iteration t ?
- (b) What will be the weights w_1 and w_3 in iteration $(t + 1)$ assuming gradient descent and gradient descent with momentum?



All the best :)





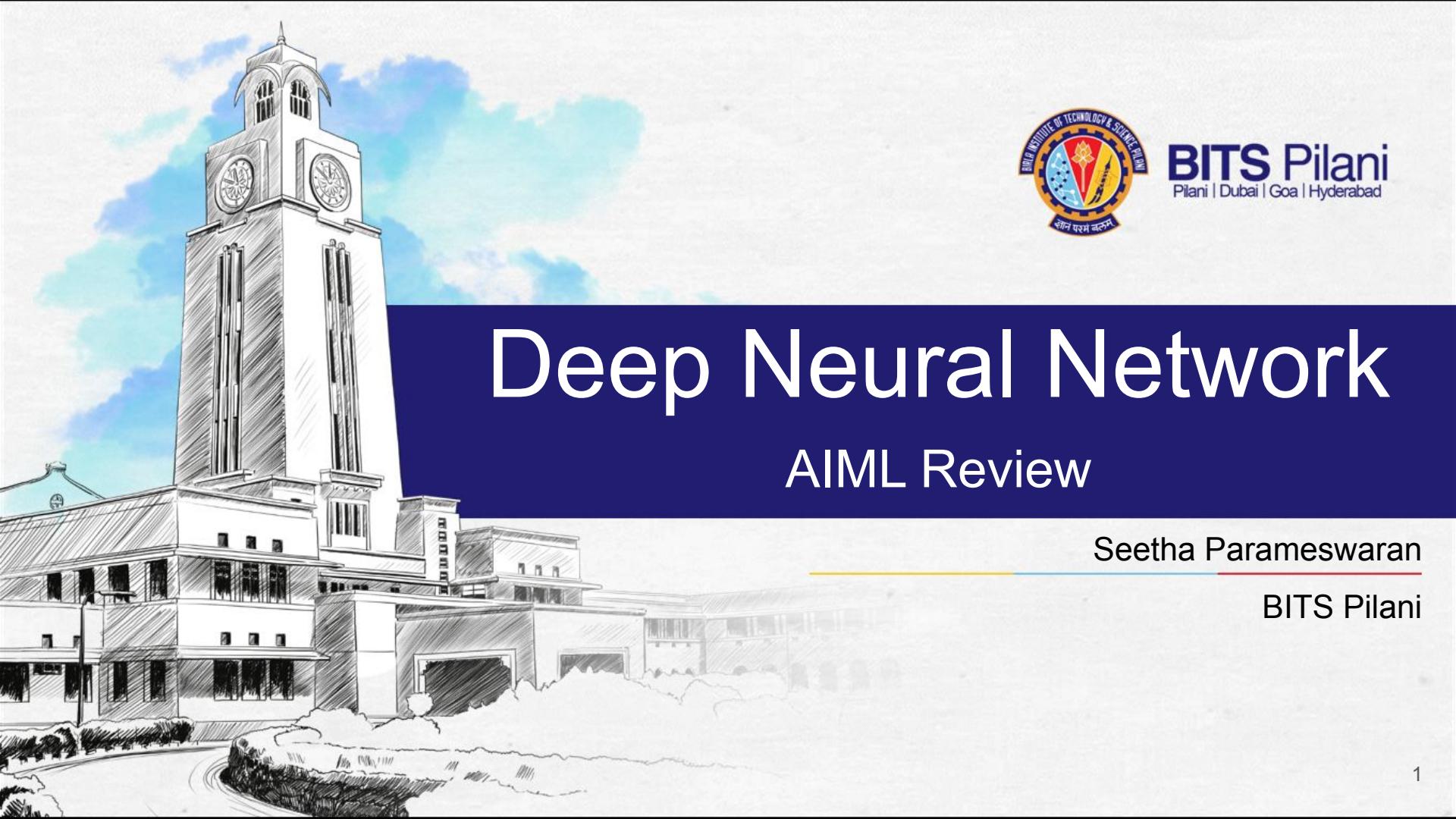
BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Neural Network

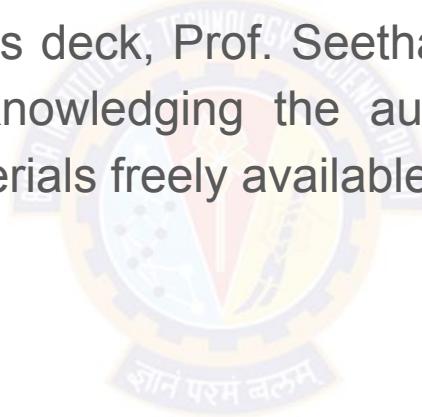
AIML Review

Seetha Parameswaran

BITS Pilani



The author of this deck, Prof. Seetha Parameswaran,
is gratefully acknowledging the authors who made
their course materials freely available online.



Pattern

Pre Midsem topics - 40%

Post Midsem topics - 60%

Topics

1. Deep Feedforward neural network
2. Convolution neural network
3. Recurrent neural network
4. Attention and transformer
5. Time series forecasting



Questions can be conceptual, numerical or descriptive.

**ENSURE THAT YOU ARE HONOURING THE CODE OF CONDUCT
DURING EXAM.**

Question

Refer to the partial code of an ANN implementation using Tensorflow and answer following questions:

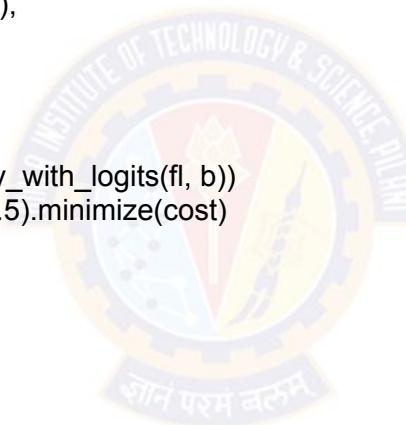
- Which activation function is used in the output layer? Why do you think this particular activation function was chosen? What difference would it make if we used sigmoid function instead?
- Calculate the output values assuming the input vector to the output layer to be $[2,0,1,0]$, weight matrix to be $[[0.2,0.3,0.2,0.1],[0.1,0.2,0.5,0.1],[0.2,0.1,0.1,0.1]]$ and bias vector to be $[0.3,0.1,0.1]$

$$\begin{pmatrix} 0.2 & 0.3 & 0.2 & 0.1 \\ 0.1 & 0.2 & 0.5 & 0.1 \\ 0.2 & 0.1 & 0.1 & 0.1 \end{pmatrix} \begin{pmatrix} 2 \\ 0 \\ 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0.3 \\ 0.1 \\ 0.1 \end{pmatrix} = \begin{pmatrix} 0.9 \\ 0.8 \\ 0.6 \end{pmatrix}$$

sigmoid

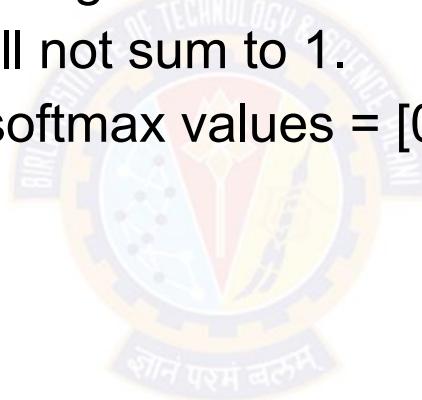
Question

```
import tensorflow as tf
data = pd.read_csv('train.csv')
a = tf.placeholder(tf.float32, shape=[None, 16])
b = tf.placeholder(tf.float32, shape=[None, 3])
phi = { 'alpha': tf.Variable(tf.random_normal([16, 8])),  
'beta': tf.Variable(tf.random_normal([8, 3])) }
omega = { 'alpha': tf.Variable(tf.random_normal([8])),  
'beta': tf.Variable(tf.random_normal([3])) }
tl = tf.add(tf.matmul(x, phi['alpha']), omega['alpha'])
tl = tf.nn.relu(tl)
fl = tf.matmul(tl, phi['beta']) + omega['beta']
cost = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(fl, b))
optimizer = tf.train.AdamOptimizer(learning_rate=0.5).minimize(cost)
init = tf.initialize_all_variables()
```



Answer

- a) Softmax is used in the output layer. The sum of output values from all the output nodes would be 1 so it is likely that the required output is probability distribution of a given variable. Sigmoid could also be used but the output values will not sum to 1.
- b) $w \cdot x + b = [0.9, 0.8, 0.6]$; softmax values = [0.38, 0.34, 0.28]



Question

Refer to the following code snippet.

```
>>> from tensorflow.keras.applications.vgg16 import VGG16  
>>> from tensorflow.keras.models import Model  
>>> model = VGG16()  
>>> model.layers.pop()  
>>> model = Model(inputs=model.inputs, outputs=model.layers[-2].output)
```

- a) If we add Batch Normalization after every convolution layer in the modified VGG16 model (popping last layer), what will be the total number of additional trainable parameters, beta's and gamma's?
- b) What will be the total number of non-trainable parameters, i.e., means and variances in the first 3 Batch Normalization layers?

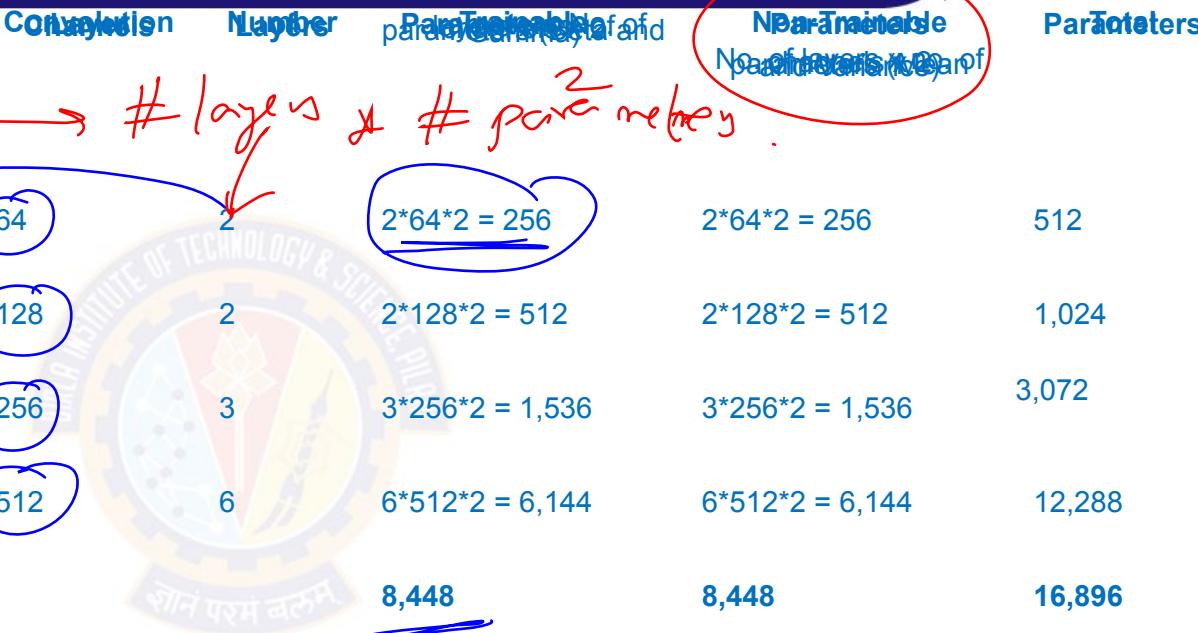
Answer

a) 8448

b) $512(2*64*2 + 128*2)$

Convolution Layer	Number of parameters	Total number of trainable parameters	Total number of parameters
1	$2*64*2 = 256$	256	256
2	$2*128*2 = 512$	512	512
3	$3*256*2 = 1,536$	$1,536$	$1,536$
4	$6*512*2 = 6,144$	$6,144$	$6,144$
	<u><u>8,448</u></u>	<u><u>8,448</u></u>	<u><u>16,896</u></u>

layers *
conv kernels *
parameters = $2(B^2)$

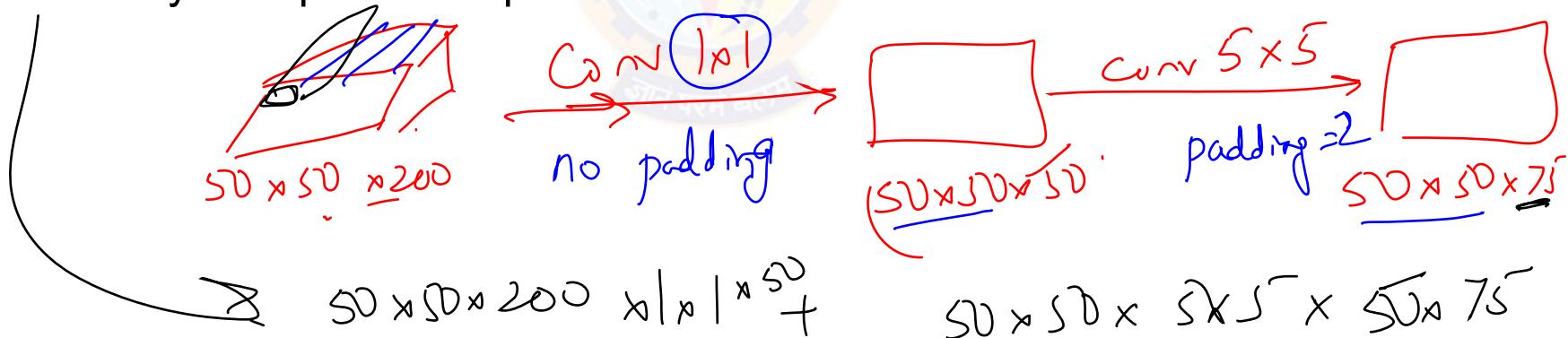


Question

In the following figure, 1×1 operators are used first to process the same 50×50 images of depth 200, and first output 50×50 images of depth 50, and then 5×5 operators are used to output 50×50 images of depth 75.

- a) What is the padding size used in the first step and padding size in the last step?
- b) How many multiplication operations are needed here?

$$P = \frac{k-1}{2}$$



Answer

a) For 1x1 convolution, padding size = 0.

For 5x5 convolution, padding size = 2.

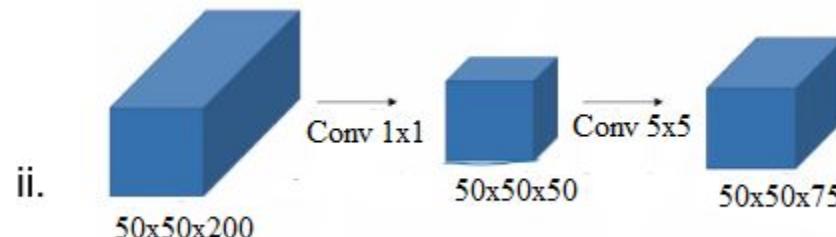
b) For 1x1 convolution,

$$\# \text{ of multiplication operations} = 50 \times 50 \times 50 \times 200 = 2.5 \times 10^7$$

For 5x5 convolution,

$$\# \text{ of multiplication operations} = 50 \times 50 \times 5 \times 5 \times 50 \times 75 = 234375000$$

So, total # of operations = 236,875,000



Question

- A. For what choice of activation functions in hidden and output nodes, autoencoder weights represent the principal component vectors?
- B.. For reconstruction of real inputs, what type of loss function is used in autoencoders?
- C. Gradient ascent is used in training the generative network in GAN. True / False
- D. If reconstructed input are binary (0 or 1), what type of activation function are needed in the output layer of autoencoders?
- E. For reconstruction of binary (0 or 1) inputs, what type of loss function used in autoencoders?
- F. What is the regularization functional used in contractive autoencoder?

Answer

- A. Linear (i.e., output = input)
- B. Sum of squared error (i.e., SSE)
- C. False
- D. Sigmoid (i.e., $1/(1+\exp(-x))$)
- E. Log-loss / binary cross-entropy
- F. Jacobian



Question

A network-in-network architecture is used to classify gray-scale images of size 64x64 into 100 classes. A micro-MLP with a single hidden layer of 10 nodes is used in the first layer instead of convolution filters of size 5x5 to generate an output feature map of depth (channel) 1. What will be the number of trainable parameters in the micro-MLP based convolution layer?

$$\# \text{hidden} = 10$$

$$i/p = 64 \times 64$$

$$\begin{aligned} \text{MLP} \\ i/p - \text{hidden} &= 5 \times 5 \times 10 \end{aligned}$$

$$\text{hidden - o/p} = 10 \times 1$$

$$\text{bias in hidden} = 10$$

$$\text{bias in o/p} = 1$$

Answer

Micro-MLP input size: 5*5, Hidden nodes 10, output size 1

Total # of trainable parameters

$$= 5*5*10 \text{ (input-hidden weights)} + 10 \text{ (bias @ hidden)} + 10*1 \\ \text{(hidden-to-output nodes)} + 1 \text{ (bias @ output)}$$

$$= 271$$



Question

Which of the following networks is (are) more suited architecturally for classifying images of varied size (e.g., the input image database has images of size 64x64, 96x96, 128x96, etc.)? No image rescaling is permitted.

resnet

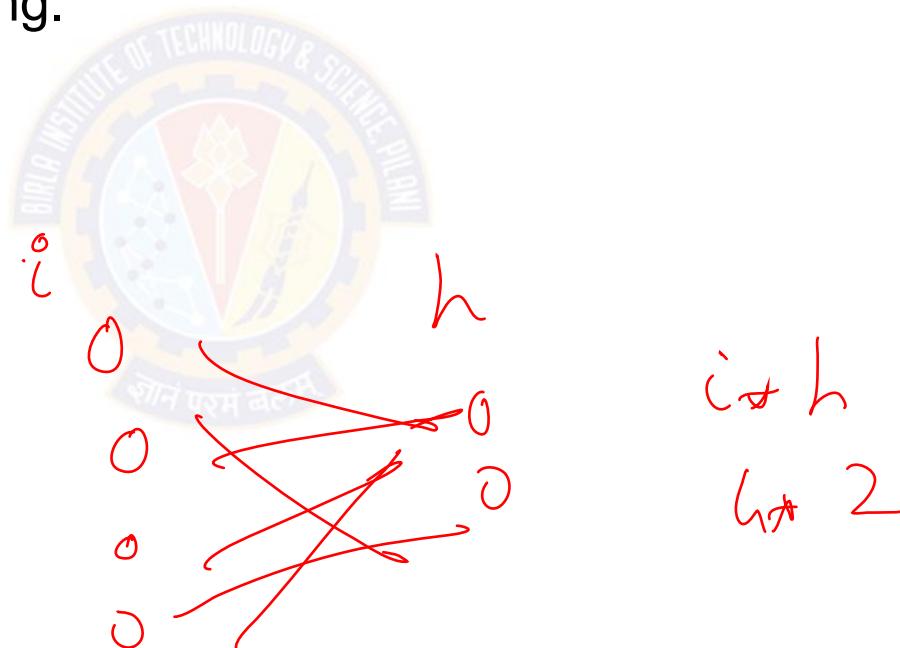
inception net

{ global
average
pooling. }



Answer

Networks that use global average pooling generate feature vectors whose size is independent of input data, after flattening. Resnet, NiN, Inception all use global average pooling.



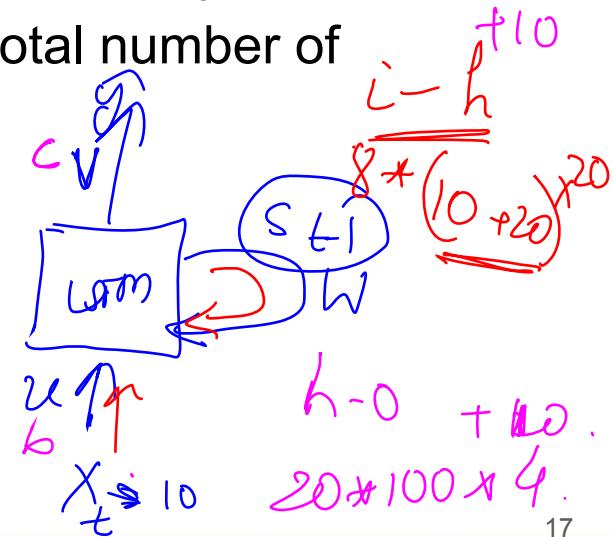
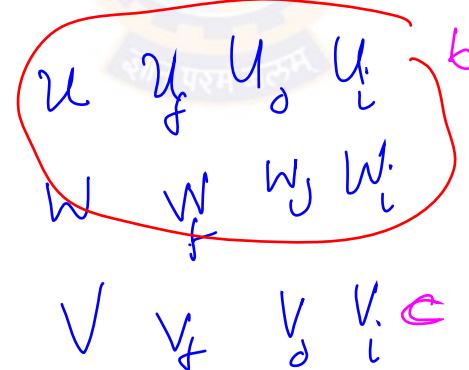
Question

Consider an LSTM network with one hidden layer of 20 nodes used for predicting the next word in one hot encoded representation in its output. No bias is used in any of the nodes. The corpus is of length 1000 words and there are 100 unique words. Assume a 10 dimensional word embedding module outside of the LSTM network, whose output is fed to the word predictor LSTM network. What will be the total number of trainable weights in the LSTM network?

$$i/p = 10$$

$$\text{hidden} = 20$$

$$o/p = 100$$



Answer

Input size to the network: 10

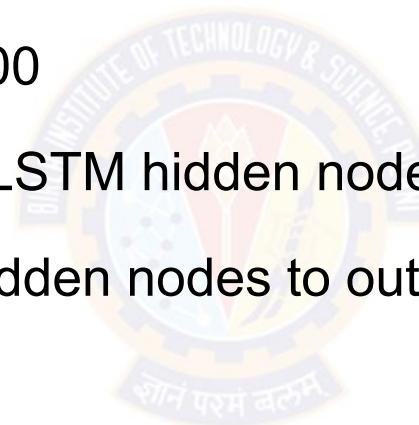
Number of hidden nodes: 20

Number of output nodes: 100

Total weights from input to LSTM hidden nodes = ~~4~~⁸ * (20+10)*20 = 2400

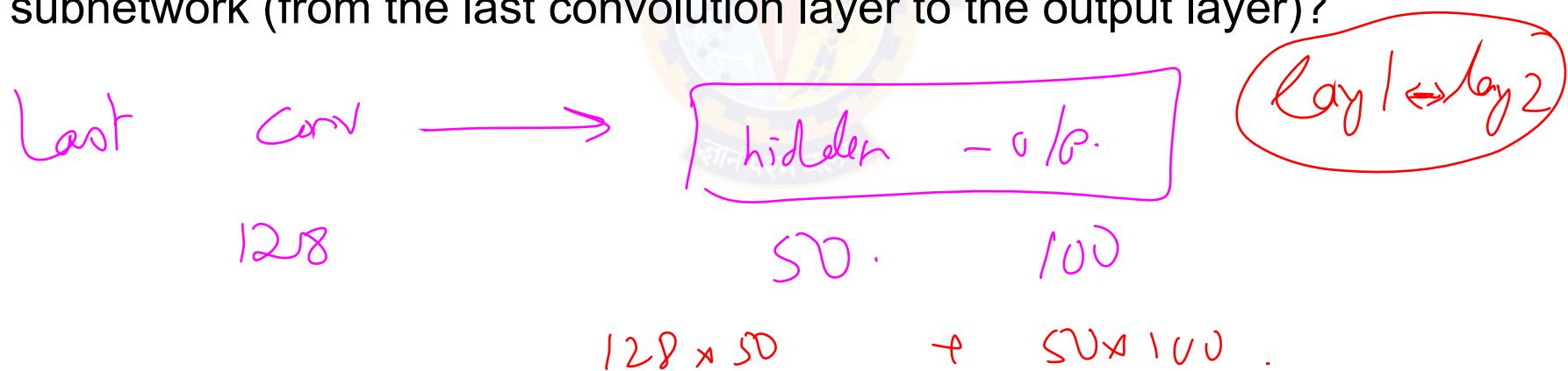
Total weights from LSTM hidden nodes to output nodes = 20*100 = 2000

Total Weights = 4400



Question

In a network in network architecture, one hot output encoding is used to classify 100 classes of objects. The last convolution layer generates 7x7x128 features maps (depth=128). Assuming the fully connected subnetwork has one hidden layer with 50 nodes, what will be the total number of trainable weights (excluding biases) in the fully connected subnetwork (from the last convolution layer to the output layer)?

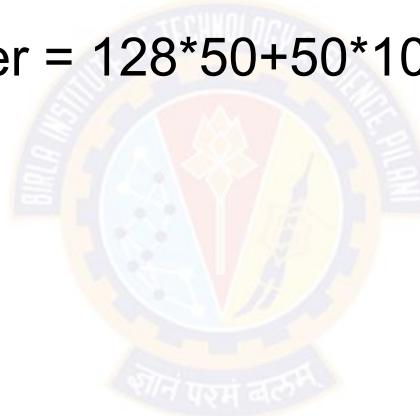


Answer

NiN uses global averaging pooling.

After flattening, feature vector size = 128 hitting the FC layer

Total # of weights in FC layer = $128 \times 50 + 50 \times 100 = 11400$



Consider a vanilla RNN network with one hidden layer of 20 nodes used for predicting the next word in a text corpus. No bias is used in any of the nodes. The corpus is of length 1000 words and there are 50 unique words. Assume a 10 dimensional word embedding module outside of the RNN network, whose output is fed to the word predictor RNN network. One hot encoding is used in the output. What will be the total number of trainable parameters in the RNN network?

Answer

Input-to-hidden weights = $(20+10)*20$

Hidden-to-output weights = $20*50$

Total weights = 1600

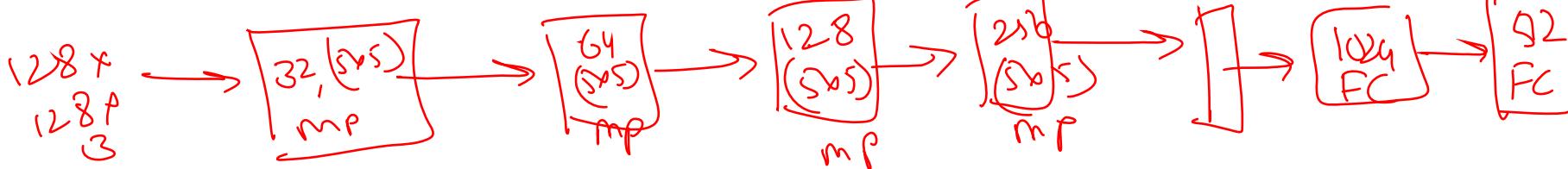


Question

2x2

For an image classification problem which classifies images into two categories, 128 by 128 ³ pixel colour images is fed through four convolution layers with 32, 64, 128, 256 kernels of 5 by 5 respectively with max pooling for each layer. Then the tensors are fed through two layers of fully connected neurons with 1024 and 512 neurons.

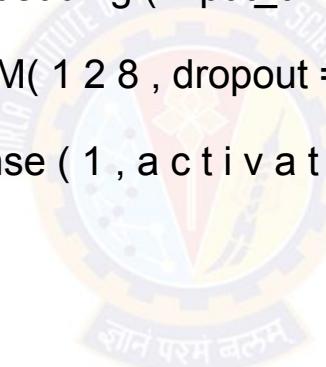
- (a) Draw a CNN for the above.
- (b) Write Tensorflow-Keras code snippet for the above.
- (c) Compute the number of parameters learned in each convolution and maxpooling layers.



Question

Explain the working of the code snippet given below.

```
model = keras.models.Sequential()  
model.add(keras.layers.Embedding(input_dim = 20000, output_dim = 128))  
model.add(keras.layers.LSTM(128, dropout = 0.2))  
model.add(keras.layers.Dense(1, activation = 'sigmoid'))  
model.summary()
```



Question

Consider a RNN where the hidden state equation is $\underline{h(t+1) = W h(t) + x(t+1)}$. Calculate and plot the output $y(t)$ of a single hidden node and single input/output RNN, where activation function used in hidden and output nodes are, respectively, linear and sigmoid. Assume that the initial hidden state is 0, and the input sequence is of even length. Assume also that all biases are 0. (Weight from input to hidden node is 1 and hidden to output node is 10.) The recurrent weight in the hidden node is -1. Input is 1 if $t=0$ or even, and 0 for odd t .

Answer

$$h_0 = h(\text{initial}) + x_0 = 0 + 1 = 1$$

$$y(0) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * 1) = \text{sig}(10) = 1$$

$$h_1 = -h_0 + x_1 = -1 + 0 = -1$$

$$y(1) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * -1) = \text{sig}(-10) = 0$$

$$h_2 = -h_1 + x_2 = -1(-1) + 1 = 2$$

$$y(1) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * 2) = \text{sig}(20) = 1$$

$$h_3 = -h_2 + x_3 = -2 + 0 = 2 = -2$$

$$y(1) = \text{sig}(w_{\text{out}} * h_0) = \text{sig}(10 * -2) = \text{sig}(-20) = 0$$

The value of y will be close to 1 when t is even and close to 0 when t is odd. The value will oscillate from 1 to 0 and 0 to 1 for each time stamp starting from $t=0$. The graph will oscillate between 0 and 1.

Question

Consider an RNN based French to Spanish language translation model where attention mechanism is employed. The maximum sentence length in the corpus is I_f and I_s , respectively. The number of hidden nodes in encoder and decoder are h_e and h_d , respectively. How many attention weights will be there?



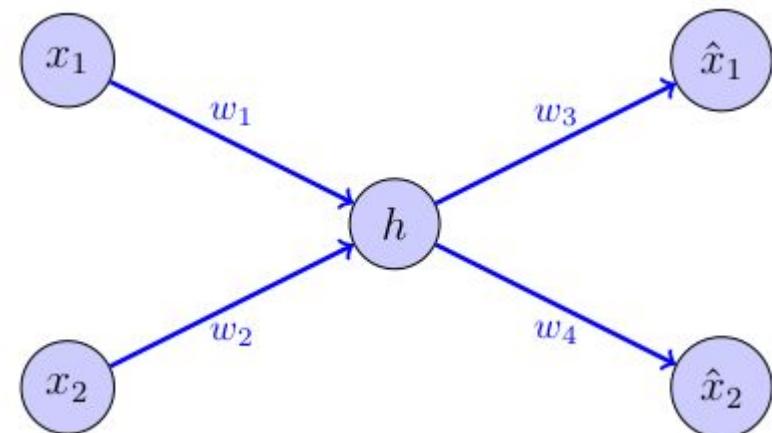
Answer

No of attention weights = Max input seq len * max output seq len = $I_f * I_s$

Question

Consider an autoencoder architecture along with the input vector $x_1 = 0$, $x_2 = 1$ and at (t) , $w_1 = 0$, $w_2 = 1$, $w_3 = 1$ and $w_4 = 0$. Assume learning rate = 0.1 and momentum constant = 0.7. At $(t - 1)$, $w_1 = -0.3$, $w_2 = 0.6$, $w_3 = -0.6$ and $w_4 = 0.3$. Assume derivative of ReLU(z)=0 at $z=0$.

- (a) At iteration t , what will be value of the loss function at iteration t ?
- (b) What will be the weights w_1 and w_3 in iteration $(t + 1)$ assuming gradient descent and gradient descent with momentum?



All the best :)

