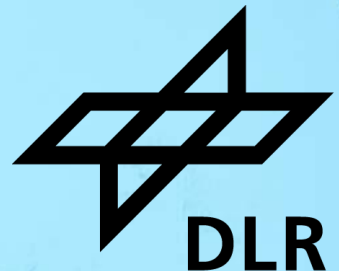# INTRODUCTION TO DEEP LEARNING
## PART I – INTRODUCTION AND BASICS

**Auliya Fitri, Sai Vemuri, Sreerag Naveenachandran**

**Machine Learning Group**
**Institute of Data Science**

DLR

## Auliya Fitri

German Aerospace Center – Institute of Data Science

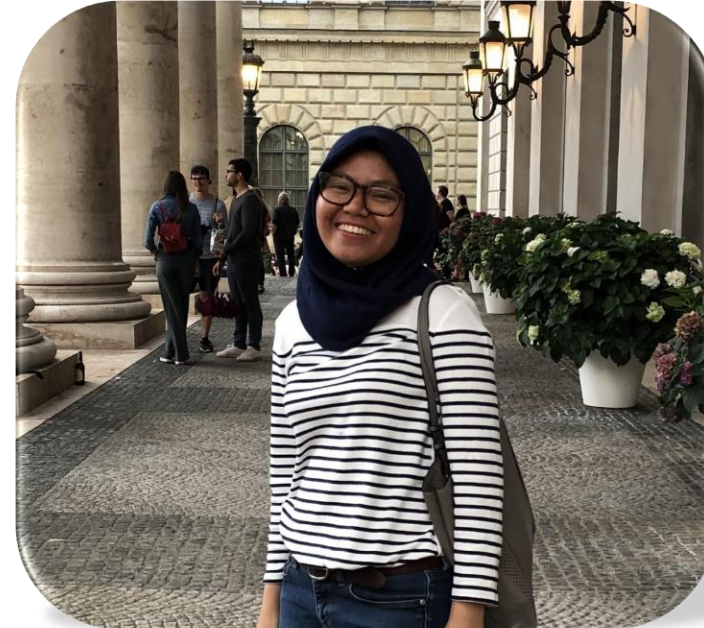Data Analysis and Intelligence

Machine Learning Group

Research Interests:

- Machine Learning

- Explainable Artificial Intelligence

- Uncertainties in Neural Networks

Contact:

Phone: +49 3641 30960 165

Mail: auliya.fitri@dlr.de

## Sai Karthikeya Vemuri

German Aerospace Center – Institute of Data Science

Data Analysis and Intelligence

Machine Learning Group

PhD Student: Computer Vision (FSU Jena)

Research Interests:

- Physics in Machine Learning

- Knowledge Integration

Contact:

Mail: sai.karthikeya.vemuri@uni-jena.de

## Sreerag Vadakkemeppully Naveenachandran

German Aerospace Center – Institute of Data Science

Data Analysis and Intelligence

Machine Learning Group

Research Interests:

- Machine learning for engineering systems
- Anomaly detection in time series
- Diffusion models

Contact:

Phone: +49 3641 30960 206

Mail:  Sreerag.Naveenachandran@dlr.de

# Schedule

| Date | Time | Activity |
|---|---|---|
| 13.11.2025 Day 1 | 09:00 - 10:00 | Introduction and basics |
| | 10:00 - 10:30 | Hands-on I |
| | 10:30 - 10:45 | Coffee break |
| | 10:45 - 11:45 | Advanced concept and Convolutional Neural Network |
| | 11:45 - 12:15 | Hands-on II |
| | 12:15 - 12:30 | Recap Day 1 |
| 14.11.2025 Day 2 | 09:00 - 10:00 | Deep Generative Models |
| | 10:00 - 10:30 | Hands-on III |
| | 10:30 - 10:45 | Coffee break |
| | 10:45 - 11:45 | Transformers, LLM, and other interesting architectures |
| | 11:45 - 12:15 | Hands-on IV |
| | 12:15 - 12:30 | Code and knowledge sources + closing |

# Curriculum

**I. Introduction and basics**

- Application examples

- Machine learning background

- Neural network concepts

- Training procedure
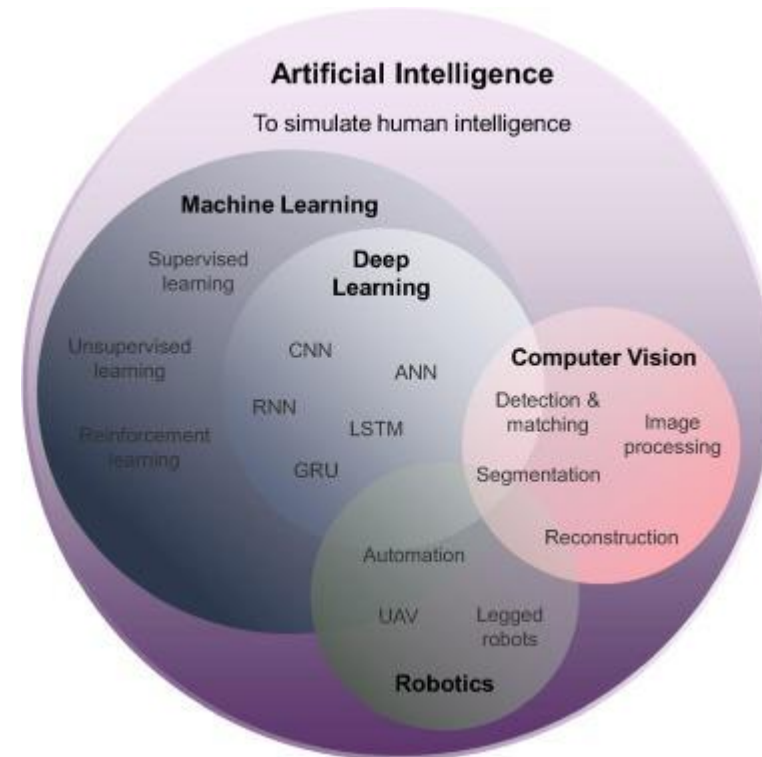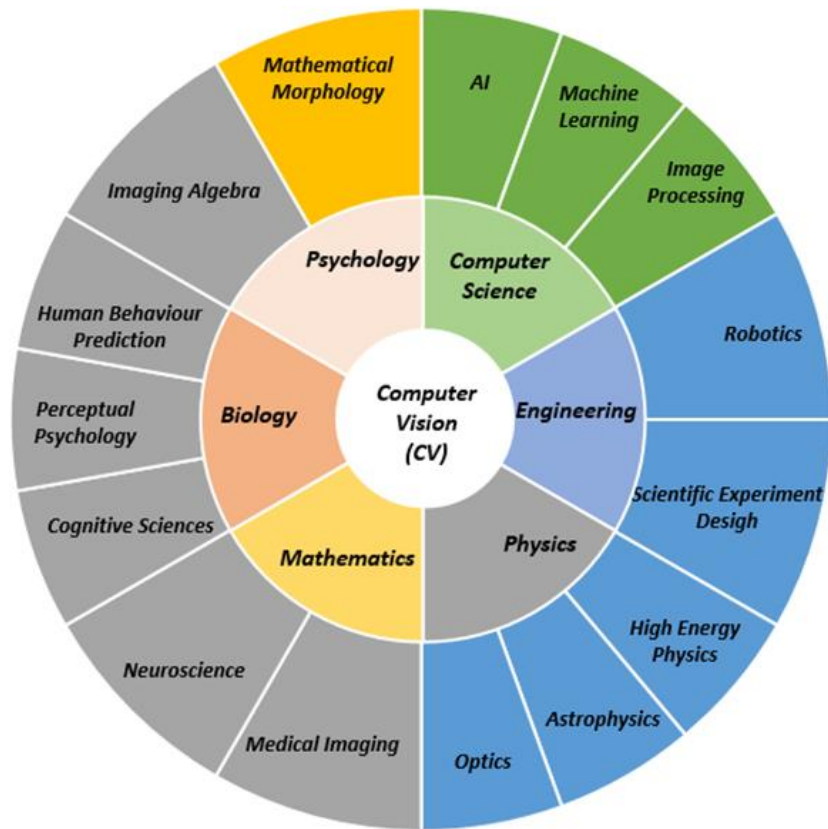
*Inspired by lectures from Paris Saclay and MIT; images taken from these, if not noted otherwise*

# APPLICATION EXAMPLES

Computer Vision is a interdisciplinary field with strong relations to AI and DL

# Application Examples
## Computer Vision - Object Recognition / Image Processing

"Object Recognition" refers to a collection of related tasks for identifying objects in digital photographs."



http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture11.pdf

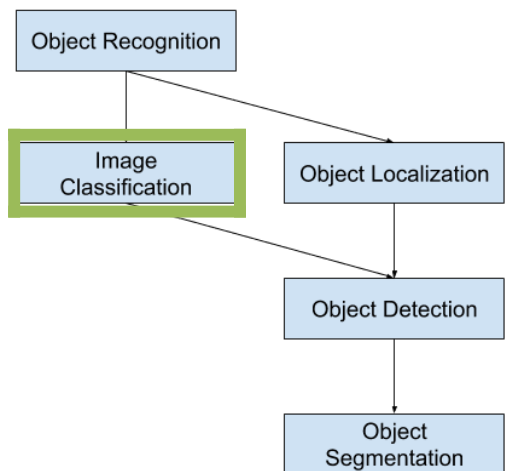A Gentle Introduction to Object Recognition With Deep Learning (machinelearningmastery.com) (2021)

"Image Classification predicts the type or class of an object in an image."





CAT

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

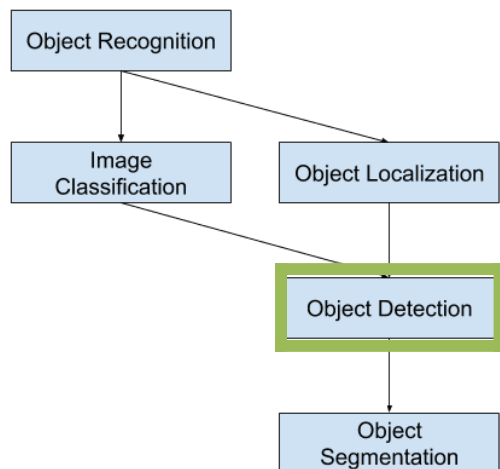"Object Localization locates the presence of objects in an image and indicate their location with a bounding box."

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

"Object Detection locates the presence of objects with a bounding box and types or classes  of the located objects in an image."

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

# Application Examples
## Computer Vision – Semantic and Object Segmentation

Semantic Segmentation
highlights pixels but not objects

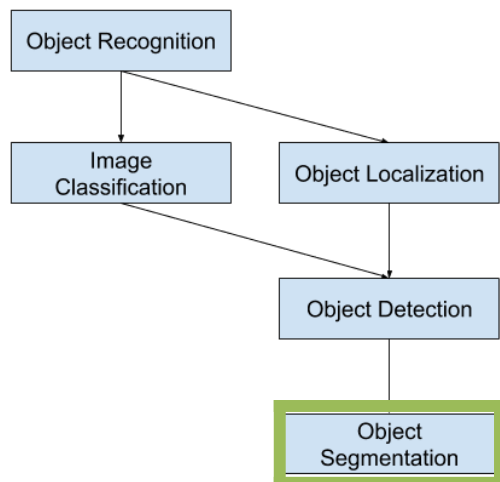"Object Segmentation highlights the specific pixels of the object instead of a coarse bounding box."
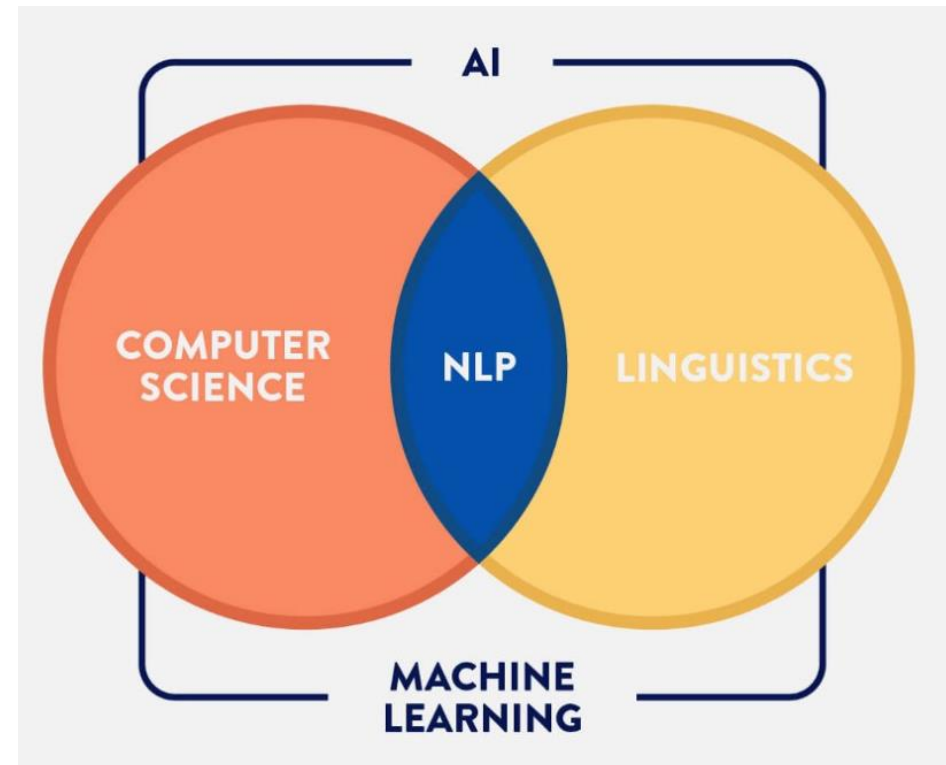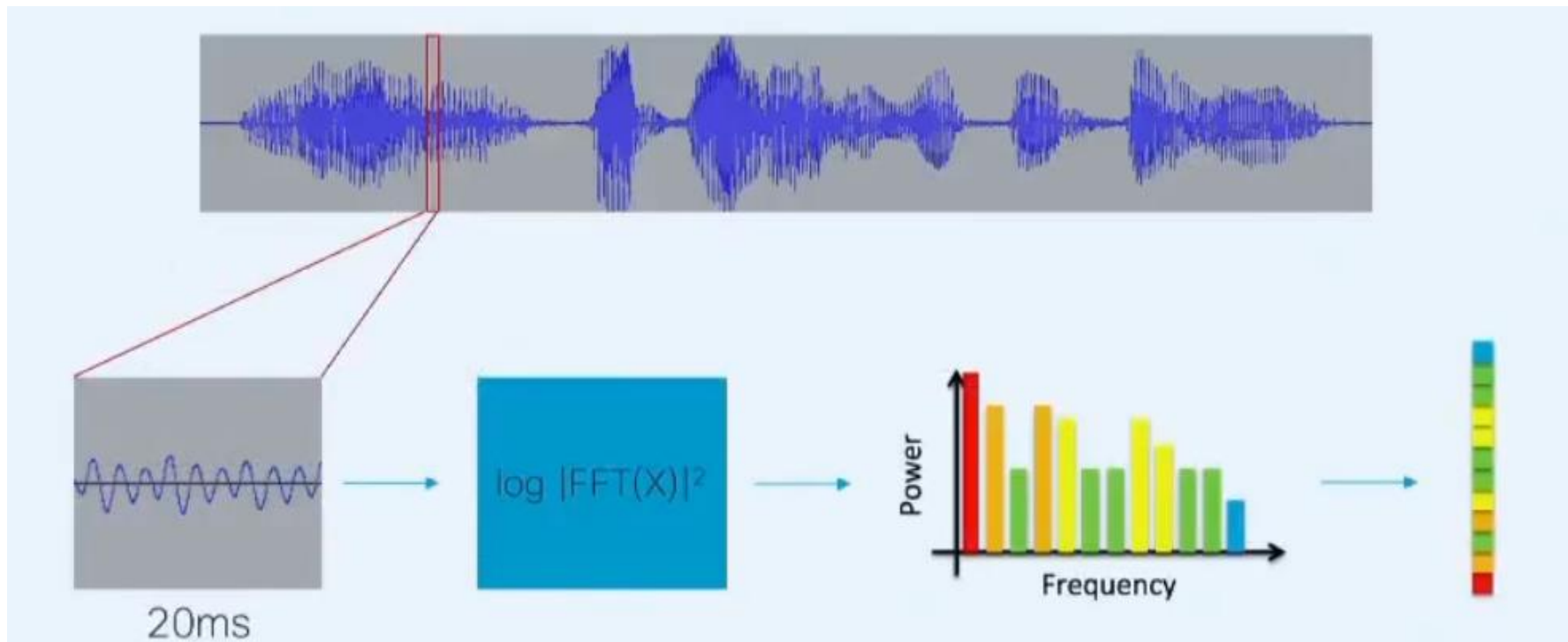
## Natural Language Processing (NLP)

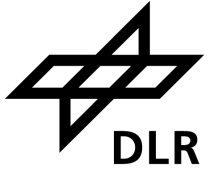NLP is a interdisciplinary field and gives computers the ability to understand human language.

"The process of enabling a computer to identify and respond to the sounds produced in human speech."

"Machine translation is the process of using artificial intelligence (AI) to automatically translate content from one language (the source) to another (the target) without any human input."



A BRIEF HISTORY OF MACHINE TRANSLATION

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

## Google Translate

Machine Translation :: From the Cold War to Deep Learning :: vas3k.com

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

"A virtual assistant is an application that understands voice commands and completes tasks for a user"

A brief history of LLMs (medium.com)

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

AlphaZero (deepmind) is a self-taught computer program for very high complexities. It learns with a complete game information approach solely based on game rules



Go (game) - Wikipedia

Pluribus (facebook) is an AI for multiplayer poker that beats 6 pros at once.
It reaches top performances after 20h training and learns with an incomplete game information approach.

AlphaFold predicts protein structure based on the amino acid sequence of the protein.

# MACHINE LEARNING BACKGROUND

# Machine learning background

Machine Learning

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

# Machine learning background
## Machine Learning Technics

```
Machine Learning  →  Unsupervised
                      Learning

                  →  Supervised
                      Learning
```

## Machine Learning Technics

## Machine Learning Technics



**Machine Learning** → **Unsupervised Learning** → **Clustering**

Clustering:
- **K-Means, K-Medoids**
- **Fuzzy C-Means**
- **Hierarchical**
- **Gaussian Mixture**
- **Hidden Markov Model**
- **Neural Networks**

**Machine Learning** → **Supervised Learning** → **Classification**

Classification:
- **Support Vector Machines**
- **Discriminant Analysis**
- **Naive Bayes**
- **Nearest Neighbor**
- **Neural Networks**

**Supervised Learning** → **Regression**

Regression:
- **Linear Regression, GLM**
- **SVR, GPR**
- **Ensemble Methods**
- **Decision Trees**
- **Neural Networks**

# Machine learning background

## 1. Big Data

+DL superiority with
  increasing big data sets
+Collection & storage
+relevant for
  multiple disciplines

+decreasing error rate

## 2. Hardware

+Parallelization
  (GPU/TPU)
+Cloud-based DL systems
  (e.g. AWS, MS-Azure

+reduce computing Power
  and memory use

## 3. Software

+improved techniques from
  growing research
+optimization SW&HW
+new models and toolboxes

PyTorch   TensorFlow

| Stochastic Gradient Descent | Perceptron -learnable weights | ... | Back-propagation Multi-Layer Perceptron | Deep Con-volutional NN Digit Recogn. | "Long Short-Term memory" (LSTM)- paper | Deep Belief Network-RBM-Layers | GPU Revolution in DL begins | ReLU activation function | Birth of GAN | TPU (Tensor Processing Unit) by Google | Turing Award in Deep learning | DeepSpeed (MS) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1952 | 1958 | ... | 1986 | 1995 | 1997 | 2006 | 2008 | 2011 | 2014 | 2016 | 2018 | 2020 |

# NEURAL NETWORK CONCEPTS

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

# Neural network concepts

**Building the network**
- Input layer
- Hidden layers
- Output layer

**Forward propagate input**
- Propagate input through neurons
- Apply non-linearities

**Evaluate performance**
- Compute the loss of the prediction

**Back propagate**
- Update weights according to the optimization algorithm

**Improving the network**
- Regularization,
- Advanced architectures

**Building the network**


CAT

$$\begin{bmatrix} \chi 1 \\ \chi 2 \\ \chi m \end{bmatrix}$$



$\chi 1$

$\chi 2$

$\chi m$

$z1$

$z2$

$z3$

$zd1$

$ŷ1$ — "CAT"

$ŷ2$ — "DOG"

Inputs          Hidden          Outputs

Calculating weighted sum
Add non-linearity
Make prediction

**Forward Propagate Input**

CAT

$$\begin{bmatrix} \chi1 \\ \chi2 \\ \chi m \end{bmatrix}$$

$\chi1$ $\chi2$ $\chi m$

$z1$ $z2$ $z3$ $zd1$

$\hat{y}1$ "CAT"

$\hat{y}2$ "DOG"

Inputs                    Hidden                    Outputs

**Forward Propagate Input**
A single neuron (perceptron)



Inputs    Weights    Sum    Non-Linearity    Output

Output    Linear combination of inputs

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

Non-linear activation function

Inputs    Weights    Sum     Non-Linearity    Output

$$\hat{y} = g\left(w_0 + \sum_{i=1}^{m} x_i\, w_i\right)$$

Bias

# Forward Propagate Input
A single neuron (perceptron)

$$\hat{y} = g\left(w_0 + \sum_{i=1}^{m} x_i w_i\right)$$



Inputs     Weights     Sum     Non-Linearity     Output

Vector/ Matrix operations

$$\hat{y} = g(w_0 + \boldsymbol{X}^T \boldsymbol{W})$$

$$\boldsymbol{X} = \begin{bmatrix} X_1 \\ \vdots \\ Xm \end{bmatrix}$$

$$\boldsymbol{W} = \begin{bmatrix} W1 \\ \vdots \\ Wm \end{bmatrix}$$

The purpose of activation functions is to **introduce non-linearities** into the network

What if we wanted to build a Neural Network to distinguish green vs red points?

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

The purpose of activation functions is to **introduce non-linearities** into the network



Linear Activation functions produce linear decisions no matter the network size

Non-linearities allow us to approximate arbitrarily complex functions

# Forward Propagate Input
## Non-linear Activation functions

### Sigmoid Function

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

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

### Hyperbolic Tangent

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

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

### Rectified Linear Unit (ReLU)

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

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

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

$$\hat{y} = g \ (w_0 + \boldsymbol{X}^T \ \boldsymbol{W})$$

$$= g \left( 1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix} \right)$$

$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

$\boldsymbol{w_0} = 1$

$\boldsymbol{W} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

This is just a line in 2D

$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

# Forward Propagate Input
## Simplified neuron

$$z = \left( w_0 + \sum_{j=1}^{m} x_j w_j \right)$$



$$\hat{y} = g(z)$$

# Forward Propagate Input
## Multi-output neuron

$$z_i = \left( w_{0,i} + \sum_{j=1}^{m} x_j \, w_{j,i} \right)$$



$\omega_{1,1}$
$\omega_{2,1}$
$\omega_{m,1}$

$y_1 = g(z_1)$

$y_2 = g(z_2)$

$\omega_{1,2}$
$\omega_{2,2}$
$\omega_{m,2}$

# Forward Propagate Input
A single-layer neural network

$$z_2 = \left( w_{0,2}^{(1)} + \sum_{j=1}^{m} x_j \left( w_{j,2}^{(1)} \right) \right)$$

$$= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)}$$



$\mathbf{W}^{(1)}$

Inputs    Hidden    Final Output

# Forward Propagate Input
A single-layer neural network



$$z_i = \left( w^{(1)}_{0,i} + \sum_{j=1}^{m} x_j w^{(1)}_{j,i} \right)$$

$W^{(1)}$

$W^{(2)}$

$g(z1)$

$g(z2)$

$g(z3)$

$g(d_1)$

Inputs

Hidden

$$\hat{y}_i = g\left( w^{(2)}_{0,i} + \sum_{j=1}^{d_1} z_j w^{(2)}_{j,i} \right)$$

# Forward Propagate Input
## A single-layer neural network

# Forward Propagate Input
Deep neural network

$$Z_{k,i} = \left( w^{(k)}_{0,i} + \sum_{j=1}^{d_{k-1}} g\left(z_{k-1,j}\right) w^{(k)}_{j,i} \right)$$



Inputs      Hidden      Final Output

Forward Propagate Input

Calculating weighted sum
Add non-linearity
Make prediction

$$\begin{bmatrix} \chi 1 \\ \chi 2 \\ \chi m \end{bmatrix}$$

CAT

$\chi 1$
$\chi 2$
$\chi m$

$z1$
$z2$
$z3$
$zd1$

$\hat{y}1$ "CAT"
$\hat{y}2$ "DOG"

Inputs          Hidden          Outputs

# Forward Propagate Input
Softmax activation
(= softargmax = normalized exponential function)

$$S(f_{y_i}) = \frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$$

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

$$S(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

Input pixels, $\mathbf{x}$ ⟶ Feedforward output, $\mathbf{y}_i$ ⟶ Softmax output, $\mathbf{S}(\mathbf{y}_i)$

Forward propagation

|  | cat | dog | horse |
|---|---|---|---|
|  | 5 | 4 | 2 |
|  | 4 | 2 | 8 |
|  | 4 | 4 | 1 |

Softmax function

|  | cat | dog | horse |
|---|---|---|---|
|  | 0.71 | 0.26 | 0.04 |
|  | 0.02 | 0.00 | 0.98 |
|  | 0.49 | 0.49 | 0.02 |

https://ljvmiranda921.github.io/notebook/2017/08/13/softmax-and-the-negative-log-likelihood/

Evaluate Prediction

Compute the error and its gradient

$$\begin{bmatrix} \chi 1 \\ \chi 2 \\ \chi m \end{bmatrix}$$

CAT

$\chi 1$

$\chi 2$

$\chi m$

$z1$

$z2$

$z3$

$zd1$

$\hat{y}1$ — "CAT"

$\hat{y}2$ — "DOG"

Inputs          Hidden          Outputs

The **loss** of our network measures the cost incurred from incorrect predictions



$$x^{(1)} = [4, 5]$$

Predicted: **0.1**
Actual: **1**

$$\mathcal{L}(f(x^{(i)}; \boldsymbol{W}), y^{(i)})$$

Predicted    Actual

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

The **empirical loss** measures the total loss over our entire dataset

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$

$z_1$
$z_2$
$z_3$

$\chi_1$
$\chi_2$

$\hat{y}_1$

$f(x)$ $\quad y$

$$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

Also known as:
- Objective function
- Cost function
- Empirical Risk

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

Predicted    Actual

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

The **Cross entropy loss** can be used with models that output a probability between 0 and 1



$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$

$$f(x) \quad y$$

$$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

$$J(W) = \frac{1}{n}\sum_{i=1}^{n} \underset{\text{Actual}}{y^{(i)}}\log(\underset{\text{Predicted}}{f(x^{(i)};W)}) + (1 - \underset{\text{Actual}}{y^{(i)}}) \log(1 - \underset{\text{Predicted}}{f(x^{(i)};W)})$$

## Mean squared error loss

**Mean squared error loss** can be used with regression models that output continuous real numbers

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$\begin{matrix} f(x) & y \end{matrix}$$

$$\begin{bmatrix} 30 \\ 80 \\ 85 \\ \vdots \end{bmatrix} \begin{bmatrix} 90 \\ 20 \\ 95 \\ \vdots \end{bmatrix}$$
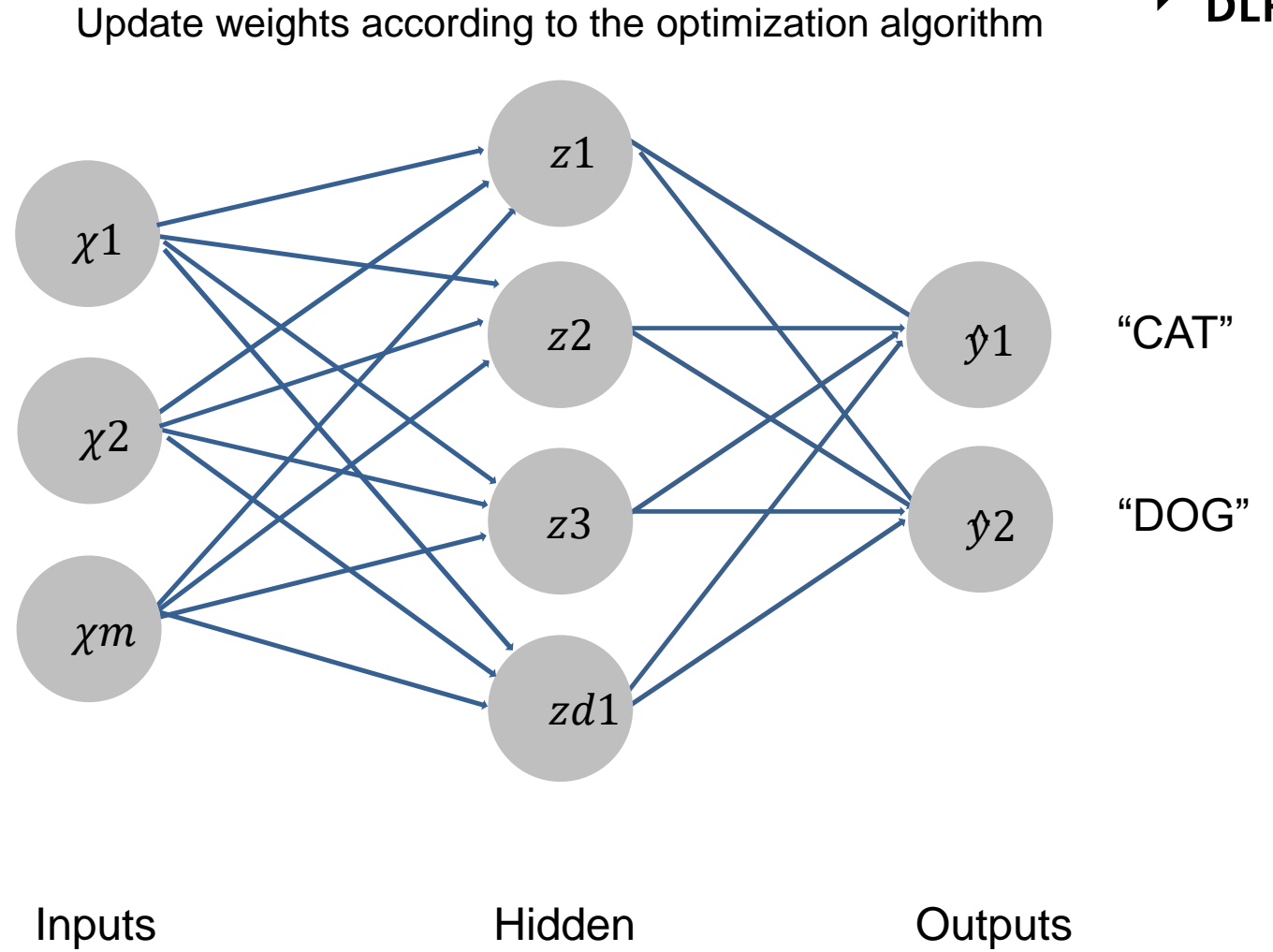
Final Grades
(percentage)

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - f(x^{(i)}; W))^{2}$$

Actual     Predicted

Update weights according to the optimization algorithm

**Back propagate**

CAT

$$\begin{bmatrix} \chi 1 \\ \chi 2 \\ \chi m \end{bmatrix}$$

$\chi 1$

$\chi 2$

$\chi m$

$z1$

$z2$

$z3$

$zd1$

$ŷ1$ "CAT"

$ŷ2$ "DOG"

Inputs          Hidden          Outputs
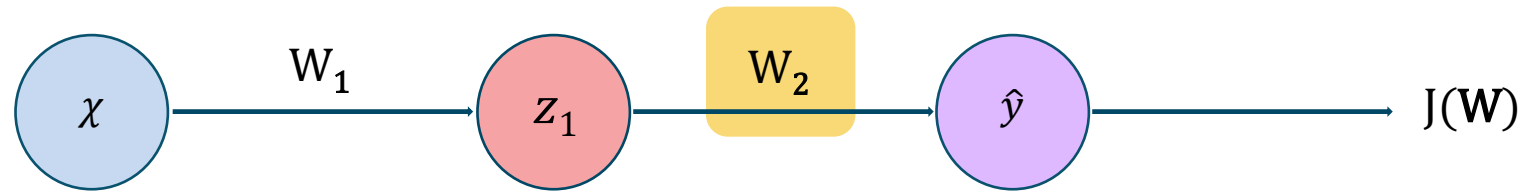
For **optimization** we want to find the network weights that **achieve the lowest loss**

$$W^* = \underset{W}{argmin} \frac{1}{n}\sum_{i=1}^{n} \mathcal{L}(f(x^{(i)}; W), y^{(i)}) \quad w_0 = 1$$

$$W^* = \underset{W}{argmin} J(W) \quad \Big| \quad W = \{W^{(0)}, W^{(1)}, \dots\}$$

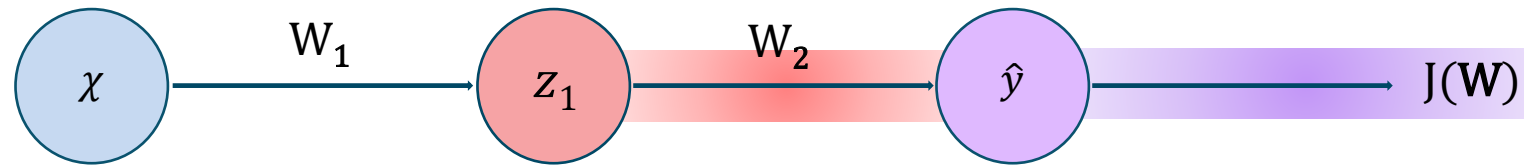How does a small change in one weight (ex. $w_2$) affect the final loss $J(\mathbf{W})$

# Back Propagate



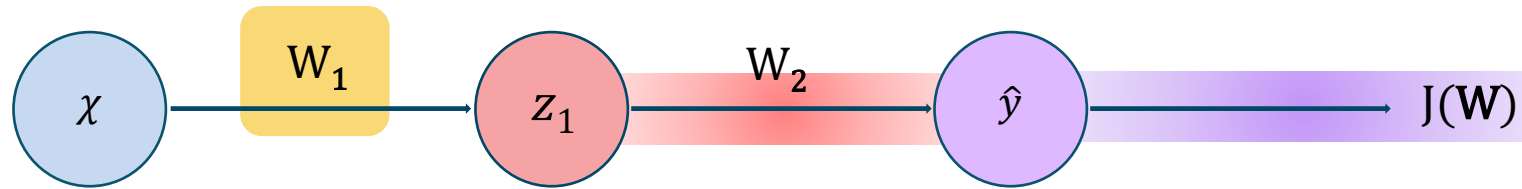$$\frac{\partial J(W)}{\partial w_2} =$$

Let's use the chain rule!

# Back Propagate



$$\frac{\partial J(W)}{\partial w_2} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$
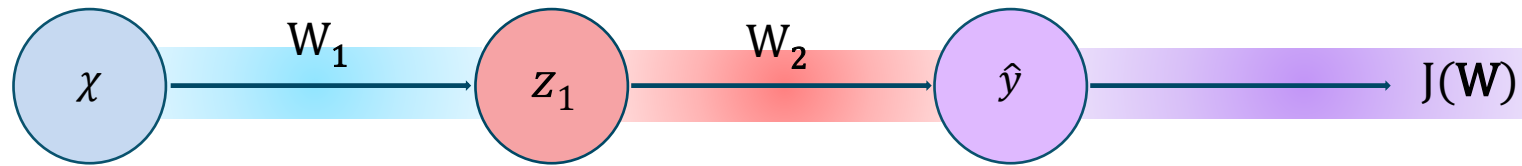
$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

Apply chain rule!        Apply chain rule!

# Back Propagate



$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$
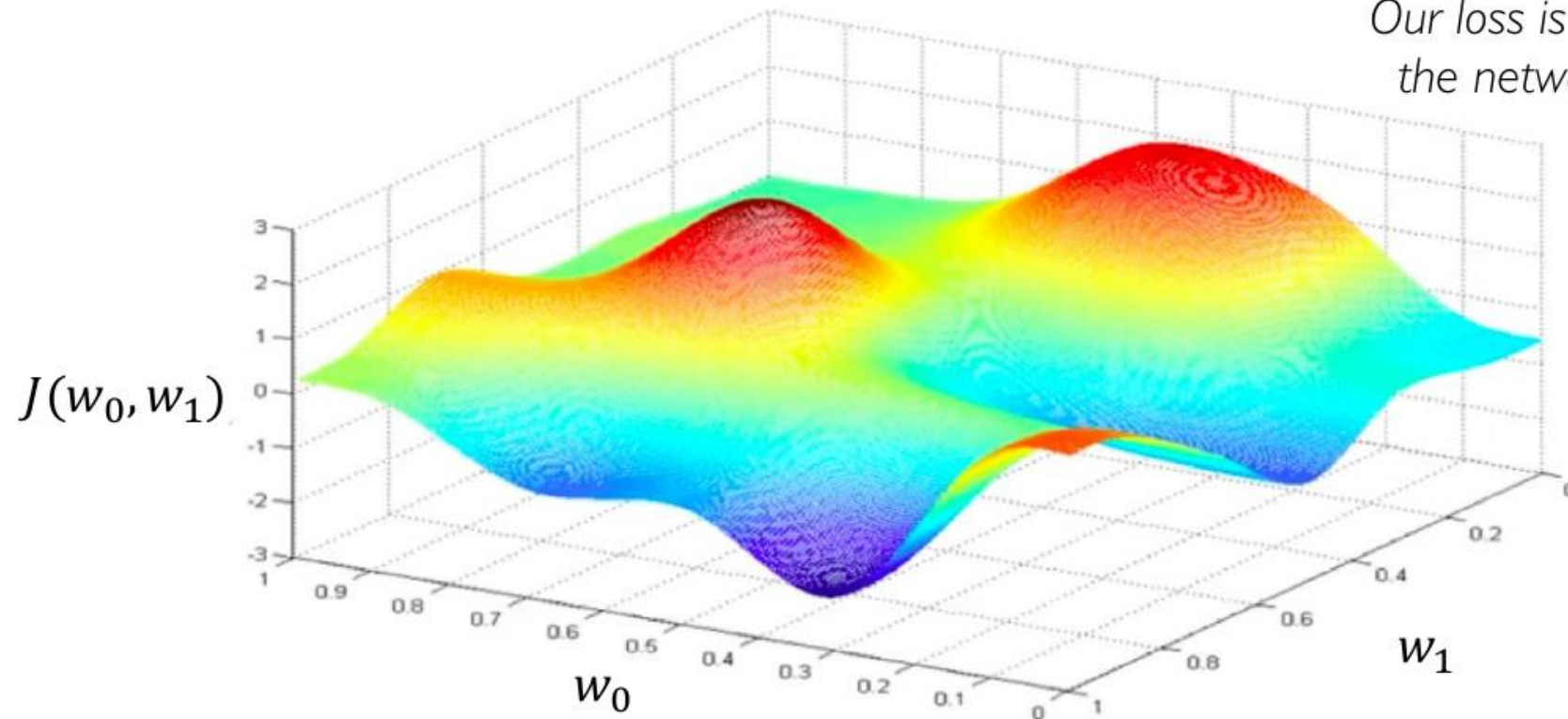
Repeat this for **every weight in the network** using gradients from layers

$$W^* = \underset{W}{\operatorname{argmin}} J(W)$$

Remember:
Our loss is a function of
the network weights!

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

Randomly pick an initial $(w_0, w_1)$

$J(w_0, w_1)$

$w_0$

$w_1$

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

Compute gradient, $\frac{\partial J(W)}{\partial W}$

$J(w_0, w_1)$

$w_0$

$w_1$

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

Take small step in opposite direction of gradient

$J(w_0, w_1)$

$w_0$

$w_1$

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

Repeat until convergence

$J(w_0, w_1)$

$w_0$

$w_1$

A.Fitri, S.Vemuri, S.Naveenachandran, DLR-DW, 13.-14.11.2025

**Algorithm**

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3. Compute gradient, $\dfrac{\partial J(W)}{\partial w}$

4. Update weights, $W \leftarrow W \leftarrow \eta \dfrac{\partial J(W)}{\partial w}$

5. Return weights

# Curriculum

| Building the network | Forward propagate input | Evaluate performance | Back propagate | Improving the network |
|---|---|---|---|---|
| • Input layer<br>• Hidden layers<br>• Output layer | • Propagate input through neurons<br>• Apply non-linearities | • Compute the loss of the prediction | • Update weights according to the optimization algorithm | • Regularization,<br>• Advanced architectures |

**Next:**

▪ Hands-on I

▪ Theory Part II:

  ▪ Advanced concept: Regularization

  ▪ Convolutional Neural Networks

# Imprint

Topic: **Introduction to Deep Learning**
Part I – Introduction and Basics

Date: 2025-11-13

Author: Auliya Fitri, Sai Vemuri, Sreerag Naveenachandran

Institute: Data Science

Image sources: All images "DLR (CC BY-NC-ND 3.0)" unless otherwise stated