UNIVERSITÀ
DEGLI STUDI
DI PADOVA

AMCO
**ARTIFICIAL INTELLIGENCE, MACHINE
LEARNING AND CONTROL RESEARCH GROUP**

# Lecture #23 Introduction to (Artificial) Neural Networks / Deep Learning

**Gian Antonio Susto**

# Before starting: Mock programming exam available on the git page of the course!

https://github.com/amco-unipd/ML_lab_DEI_public

exam_simulation.ipynb

exam_utiles.py

Titanic dataset:

# Before starting: Mock programming exam available on the git page of the course!

[https://github.com/amco-unipd/ML_lab_DEI_public](https://github.com/amco-unipd/ML_lab_DEI_public)

exam_simulation.ipynb

exam_utiles.py

Titanic dataset:

Focus on trees and ensemble

Reference Material (used for this presentation):

- MIT *Introduction to Deep Learning* http://introtodeeplearning.com

- L. Fridman *Deep Learning* https://deeplearning.mit.edu/

- I. Goodfellow, Y. Bengio, A. Courville. *Deep learning*. MIT press, 2016.

- G. Prando. *Seminars on Deep Learning*. UniPD, 2017 (Not available on-line)

- ImageNET http://www.image-net.org/

- C. Szegedy et al. *Going Deeper with Convolutions* CVPR2015

- B. Turovsky *Found in translation: More accurate, fluent sentences in Google Translate*
https://www.blog.google/products/translate/found-translation-more-accurate-fluent-sentences-google-translate/

- A. Ng *Nuts and Bolts of Applying Deep Learning* https://www.youtube.com/watch?v=F1ka6a13S9I

- I. Goodfellow et al. *Explaining and harnessing adversarial examples* ICLR 2015

- A. Siarohin et al. *First Order Motion Model for Image Animation* NeurIPS 2019
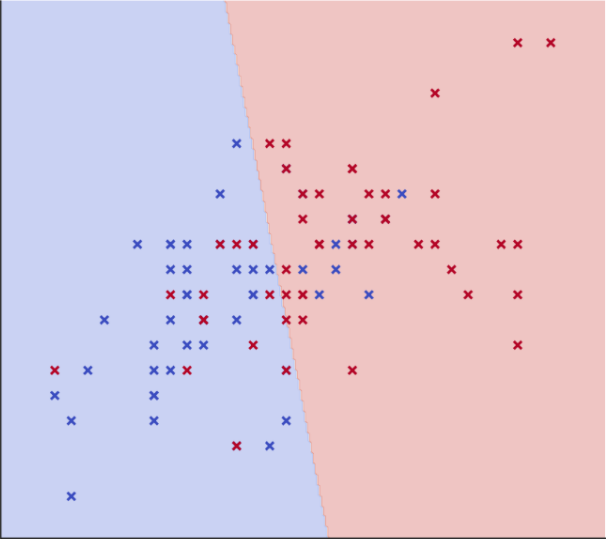

Other references:

- Wired Interview with Geoffry Hinton 13/05/2019

- https://paperswithcode.com/

- The Journal of Open Source Software https://joss.theoj.org/

# Our quest for models with high representational power!

# Machine Learning Approaches



Neural Networks (Deep Learning)

# Neural Networks / Deep Learning



**Conv_1 Convolution** (5 x 5) kernel *valid* padding

**Max-Pooling** (2 x 2)

**Conv_2 Convolution** (5 x 5) kernel *valid* padding

**Max-Pooling** (2 x 2)

**fc_3 Fully-Connected** Neural Network ReLU activation

**fc_4 Fully-Connected** Neural Network

(with dropout)

INPUT (28 x 28 x 1)

n1 channels (24 x 24 x n1)

n1 channels (12 x 12 x n1)

n2 channels (8 x 8 x n2)

n2 channels (4 x 4 x n2)

Flattened

n3 units

OUTPUT

0
1
2
9

Information Retrieval

Sentiment Analysis

Information Extraction

Machine Translation

Text Processing

QuestionAnswering

## A mostly complete chart of Neural Networks

©2016 Fjodor van Veen – asimovinstitute.org

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

Perceptron (P)

Feed Forward (FF)

Radial Basis Network (RBF)

Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)

Long / Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Auto Encoder (AE)

Variational AE (VAE)

Denoising AE (DAE)

Sparse AE (SAE)

Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)
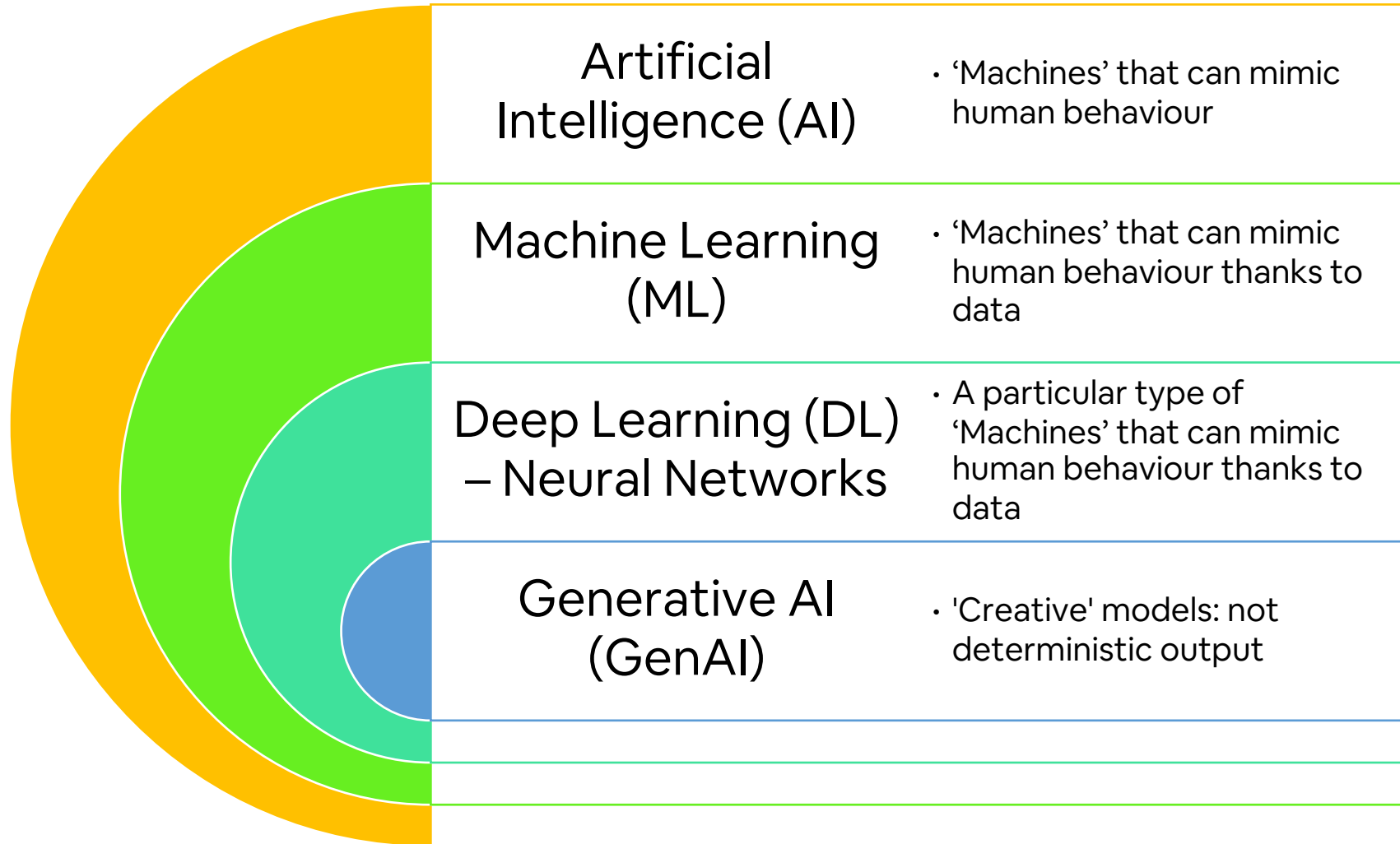
Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Convolutional Network (DCN)

Deconvolutional Network (DN)

Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)

Liquid State Machine (LSM)

Extreme Learning Machine (ELM)

Echo State Network (ESN)

Deep Residual Network (DRN)

Kohonen Network (KN)

Support Vector Machine (SVM)

Neural Turing Machine (NTM)

# The keywords

**Artificial Intelligence (AI)**
- 'Machines' that can mimic human behaviour

**Machine Learning (ML)**
- 'Machines' that can mimic human behaviour thanks to data

**Deep Learning (DL) – Neural Networks**
- A particular type of 'Machines' that can mimic human behaviour thanks to data

**Generative AI (GenAI)**
- 'Creative' models: not deterministic output
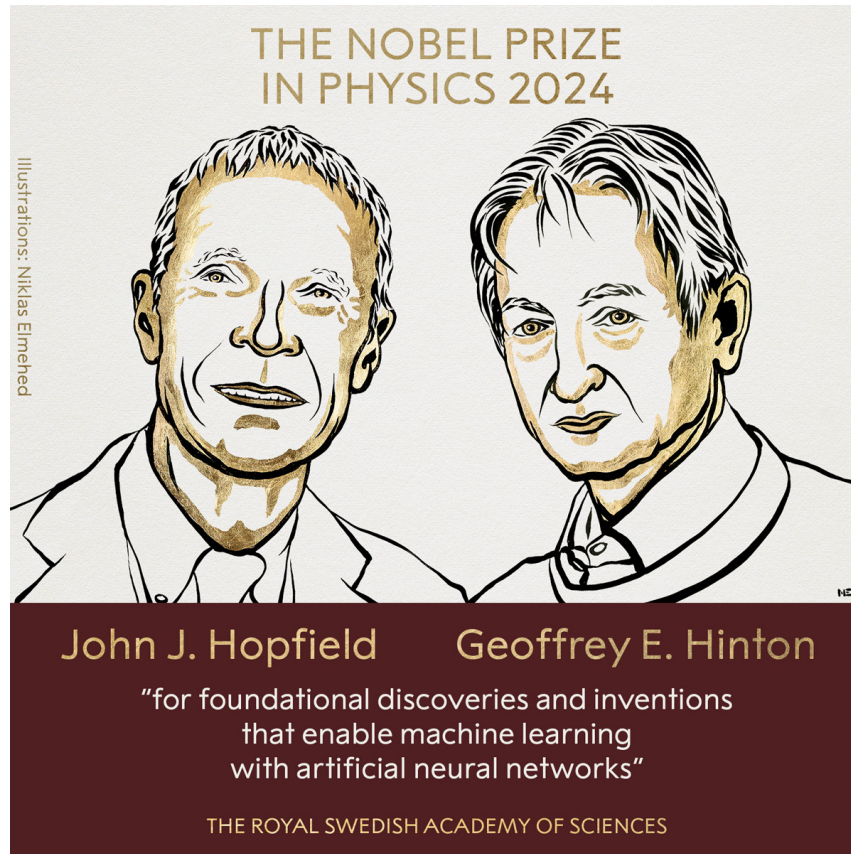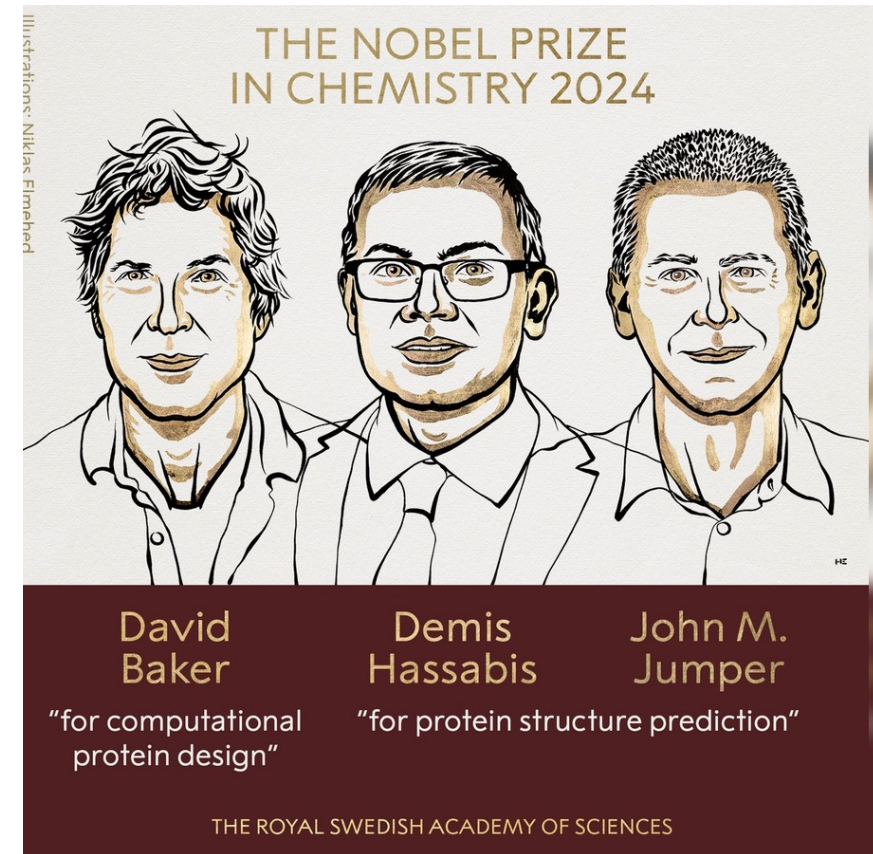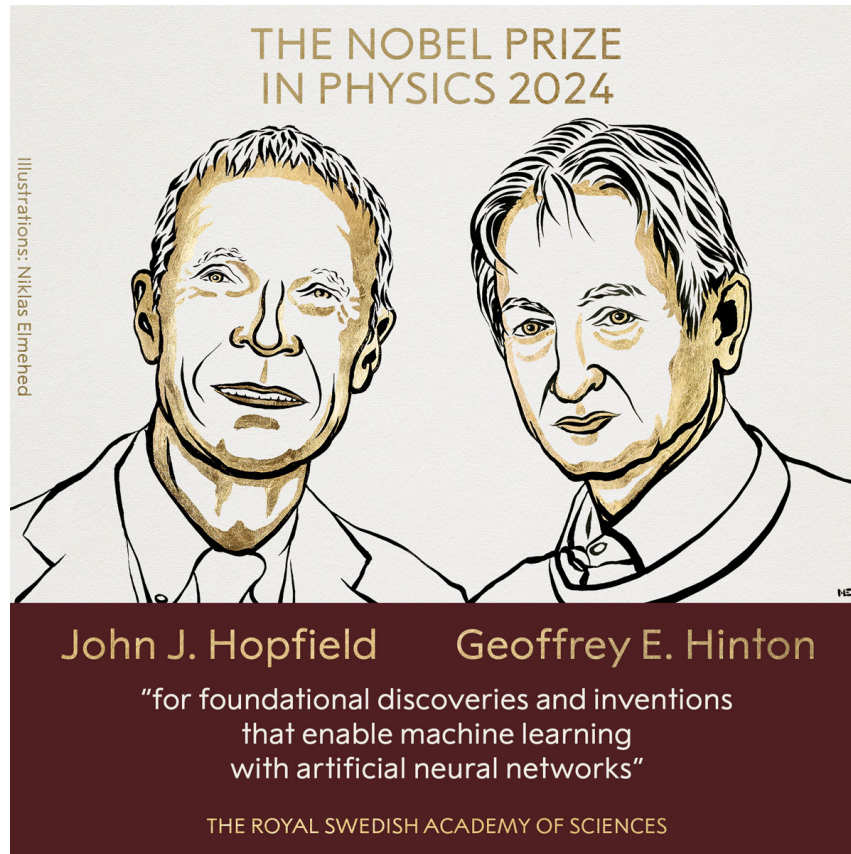
# AI is the hottest topic in science... and mostly because of Deep Learning!

# AI is the hottest topic in science... and mostly because of Deep Learning!





*https://alphafold.ebi.ac.uk/*

Exclusive: the most-cited papers of the twenty-first century

| Rank (median) | Citation | Times cited (range across databases) |
|---|---|---|
| 1 | Deep residual learning for image recognition (2016, preprint 2015) | 103,756–254,074 |
| 2 | Analysis of relative gene expression data using real-time quantitative PCR and the $2^{-\Delta\Delta C_T}$ method (2001) | 149,953–185,480 |
| 3 | Using thematic analysis in psychology (2006) | 100,327–230,391 |
| 4 | Diagnostic and Statistical Manual of Mental Disorders, DSM-5 (2013) | 98,312–367,800 |
| 5 | A short history of SHELX (2007) | 76,523–99,470 |
| 6 | Random forests (2001) | 31,809–146,508 |
| 7 | Attention is all you need (2017) | 56,201–150,832 |
| 8 | ImageNet classification with deep convolutional neural networks (2017) | 46,860–137,997 |
| 9 | Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries (2020) | 75,634–99,390 |
| 10 | Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries (2016) | 66,844–93,433 |

Deep Learning → 1

Machine Learning → 6
Deep Learning → 7
Deep Learning → 8

**Nowadays, when people talk about AI, they typically have deep learning-based technologies in mind...**

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

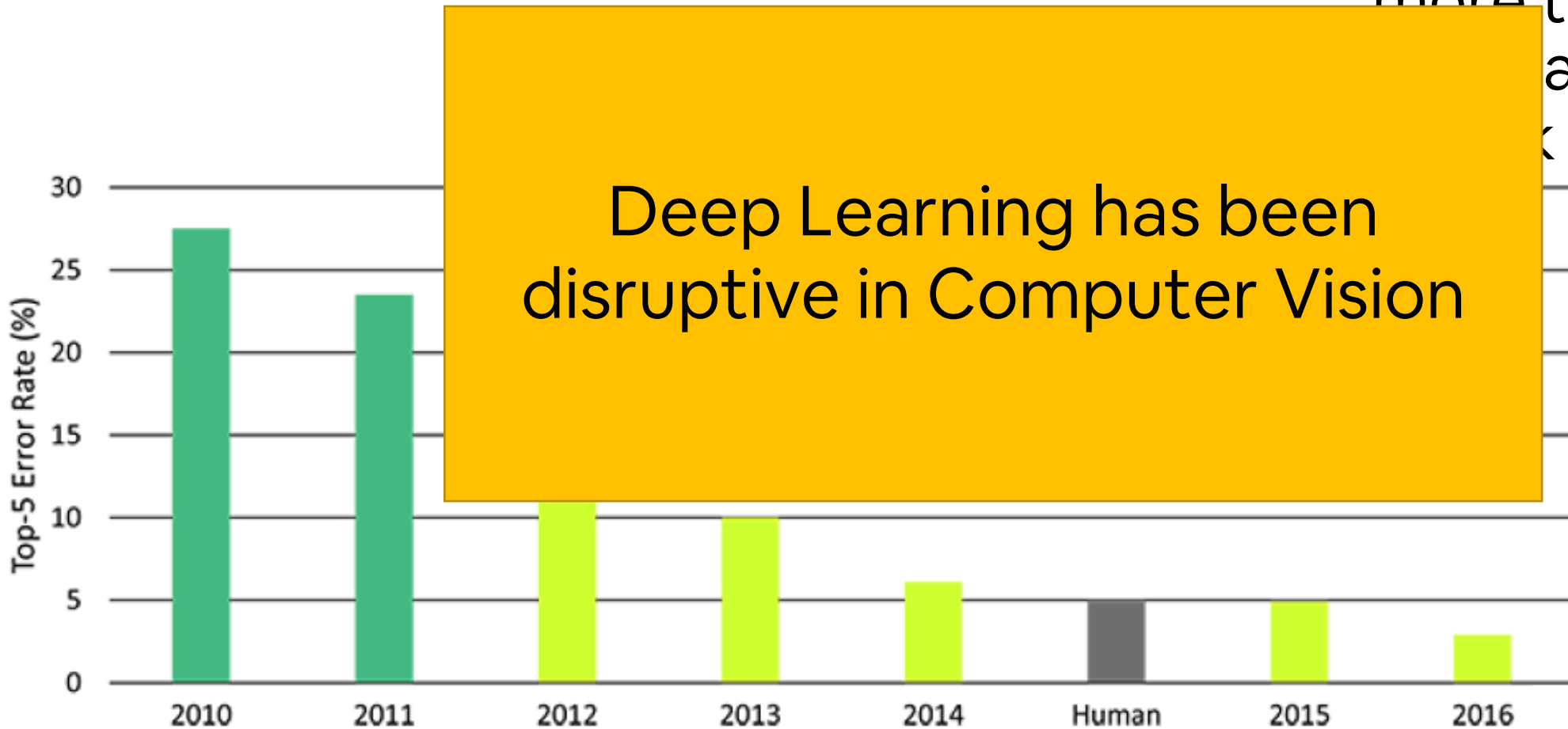A dataset with more than 14 Million annotated images of 20k categories

ILSVRC Top 5 Error on ImageNet

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

A dataset with more than 14 Million ~~ated images~~ ~~categories~~



Deep Learning has been disruptive in Computer Vision

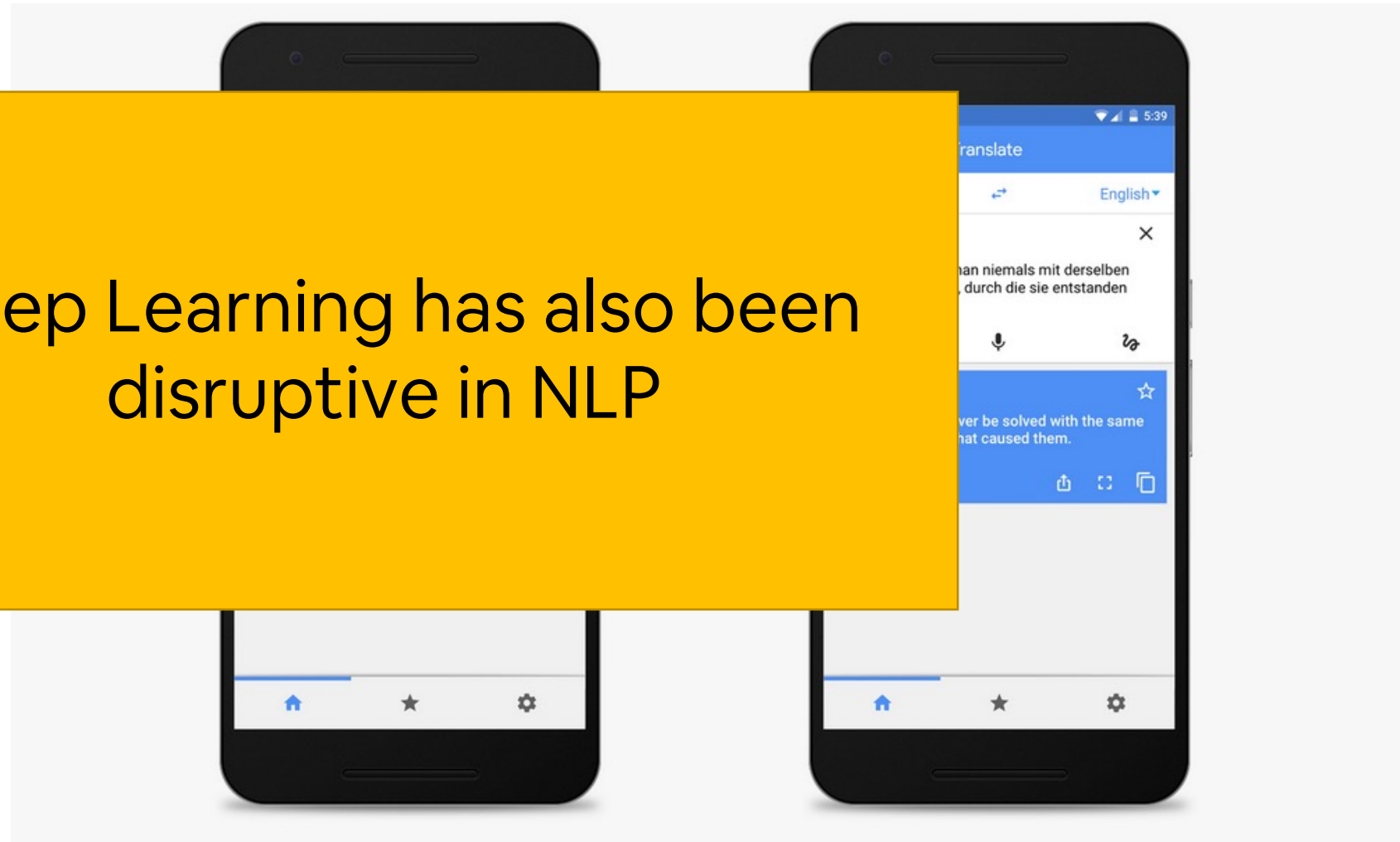G. Prando. *Seminars on Deep Learning*. UniPD, 2017

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

Usage of Neural Machine Translation by Google Translate (15 November 2016)

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

Usage of Neural Machine Translation by Google Translate (15 November 2016)

Deep Learning has also been disruptive in NLP

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

ChatGPT (Generative Pre-trained Transformer): a 'super' chatbot launched by OpenAI in November 2022.

ChatGPT is based built on top of a GPT-X language model

# Everyone talks about ~~AI~~ Deep Learning!
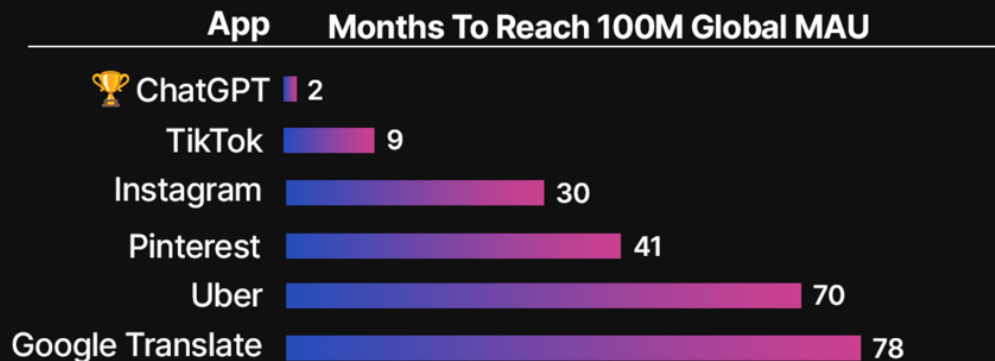


Artificial intelligence
Field of study

Worldwide ▾    2004 – present ▾    All categories ▾    Web Search ▾

Interest over time

OpenAI
ChatGPT

**Time it took to reach 100 million monthly users:**

| App | Months To Reach 100M Global MAU |
|---|---|
| 🏆 ChatGPT | 2 |
| TikTok | 9 |
| Instagram | 30 |
| Pinterest | 41 |
| Uber | 70 |
| Google Translate | 78 |

Source: UBS

@TheRundownAI

Nov 2022
Artificial intelligence       27

1 Apr 2021

**can you write a python code for image preprocessing?**

Sure! Here's some Python code that demonstrates some basic image preprocessing techniques:

Copy code

```python
import cv2
import numpy as np

# Reading image
image = cv2.imread("image.jpg")

# Resizing image
resized_image = cv2.resize(image, (224, 224))

# Converting image to grayscale
gray_image = cv2.cvtColor(resized_image, cv2.COLOR_BGR2GRAY)

# Normalizing image pixels
normalized_image = cv2.normalize(gray_image, None, 0, 255,
cv2.NORM_MINMAX)

# Applying Gaussian blur
blurred_image = cv2.GaussianBlur(normalized_image, (5, 5), 0)
```

In this code, the first thing we do is read the image using the `cv2.imread()` function from the `cv2` library. Then, we use the `cv2.resize()` function to resize the image to a

i want to build a meme related to everybody speaking about chatgpt. any ideas?

Sure, here's one idea:

A picture of a group of people all crowded around a computer, with one person saying "Hey, have you guys heard about ChatGPT? It's the latest and greatest chatbot!" Another person in the group replies, "Oh man, I've been using ChatGPT for months now! It's like having a virtual assistant that never sleeps!" A third person chimes in, "I don't know, I tried ChatGPT and it just kept sending me cat memes. I think I'll stick with my human chatbot."

Caption: When everyone is talking about ChatGPT, but you're still using your trusty human chatbot.

This idea plays on the idea that ChatGPT is a highly advanced chatbot that people are excited to use, but one person in the group is skeptical and prefers a more "human" approach to chatting. The caption adds a humorous twist to the meme.

**i don't want to take out the dog: can you find me an excuse for my wife?**

As an AI, I'm not able to provide you with an excuse, but I can offer you some suggestions on how to communicate with your wife about why you are not able to walk the dog.

Here are a few things you could consider:

- Explain the circumstances and your responsibilities that are causing you to not be able to walk the dog.
- Offer to find a solution together. For example, if you are busy at work, suggest finding a dog walker or asking someone else in the household to take the dog out for a walk.
- Apologize for the inconvenience and assure her that you will make it up to her in some other way.
- And also if possible try to make up for it on weekends or on days when you are free.

It is always important to communicate openly and honestly with your spouse or partner. It will be more likely to find a solution that works for both of you if you approach the conversation with a willingness to work together and understand each other's perspectives.

**i will teach a 12 lecture course on introduction to deep learning. Can you suggest me an outline for the course?**

Sure, here's a basic outline for a 12-lecture course on Introduction to Deep Learning:

Lecture 1: Introduction to Deep Learning

- Overview of the course
- What is deep learning and its applications
- Brief history of deep learning

Lecture 2: Neural Networks

- Basics of artificial neural networks
- Feedforward neural networks
- Activation functions

Lecture 3: Backpropagation

- How neural networks learn
- Gradient descent
- Backpropagation algorithm

Lecture 4: Convolutional Neural Networks

- Introduction to CNNs
- Convolutional layers
- Pooling layers
- Applications of CNNs (image classification, object detection, etc.)

Lecture 5: Recurrent Neural Networks

- Introduction to RNNs
- LSTMs and GRUs
- Applications of RNNs (time series, natural language processing, etc.)

Lecture 6: Generative Models

- Introduction to generative models
- Autoencoders
- Generative Adversarial Networks (GANs)

Lecture 7: Unsupervised Learning

- Clustering
- Dimensionality reduction
- Auto-encoder

Lecture 8: Transfer Learning

- Introduction to transfer learning
- Pretrained models (e.g. VGG, ResNet, etc.)
- Fine-tuning

Lecture 9: Optimization Algorithms

- Stochastic gradient descent (SGD)
- Adaptive gradient algorithms (AdaGrad, Adam, etc.)
- Regularization (L1, L2, Dropout, etc)

## Lecture 10: Deployment and Production

- Introduction to deployment of deep learning models
- Frameworks for deployment (TensorFlow Serving, etc.)
- Tips and tricks for deploying models in production

## Lecture 11: Deep Learning in Computer Vision

- Introduction to computer vision
- Deep learning techniques for image processing
- Applications of deep learning in computer vision

## Lecture 12: Deep Learning in Natural Language Processing

- Introduction to natural language processing
- Deep learning techniques for natural language processing
- Applications of deep learning in natural language processing

Please note that this is just a rough outline, and you may choose to adjust it or add more lectures to fit your specific needs and audience. Additionally, you can also use a software or a framework for deep learning such as Tensorflow, Pytorch, or Keras, to make your students learn the practical implementation of the concepts, if you are planning to use this as a hands-on course.

Also, it is important to make sure students have a strong foundation in math, specifically linear algebra and calculus as it is necessary for understanding many deep learning concepts.

**can you teach the course in my place?**

As an AI, I am able to provide a general outline for a course on deep learning, as well as answer specific questions about various topics related to the field. However, I am not able to physically teach a course in person. Additionally, as an AI, I do not have personal experience or knowledge of any specific course material, curriculum, or teaching methodologies, so it would be best if you taught the class yourself.

That being said, I can provide support to you while teaching the class, answering any questions you may have or explaining certain concepts that might be difficult to understand.

And if your students have any questions, I would be happy to help them as well.

# Multi-modal tasks!

https://openai.com/blog/dall-e-3-is-now-available-in-chatgpt-plus-and-enterprise

I am doing a report on cirrus clouds for my science class. I need photorealistic images that show off how wispy they are. I am going to compare them to photos I took of puffy cumulonimbus clouds at my house yesterday.

Here are the photorealistic images of cirrus clouds showcasing their wispy and delicate nature. You can compare these with the photos of cumulonimbus clouds you took to highlight the differences in their structures and appearances. If you need any further information or assistance, please let me know!

# Multi-modal tasks!

https://openai.com/blog/dall-e-3-is-now-available-in-chatgpt-plus-and-enterprise

My business partner and I are creating a coffee shop inspired by our mutual love for rabbits. We have a logo we like a lot that is minimalist with a rabbit, but anxious to get your take on it.

Here are some minimalist rabbit-inspired logo designs for your coffee shop. Each captures the essence of a rabbit while integrating coffee-related elements. I hope these designs align with the vision you and your business partner have for your coffee shop. If you have any feedback or further preferences, please let me know!

# Multi-modal tasks!

https://openai.com/blog/dall-e-3-is-now-available-in-chatgpt-plus-and-enterprise

I am working with a web designer for a new website to sell bicycle helmets. Could you recommend a few website UI's that are both modern and simple that I can share with her as examples?

Here are some website UI designs tailored for selling bicycle helmets with a modern and simple aesthetic:

# Multi-modal tasks!

# Not only ChatGPT!

a logo for a deep learning course

Bloomberg Professional Services

## Introducing BloombergGPT, Bloomberg's 50-billion parameter large language model, purpose-built from scratch for finance

March 30, 2023

*BloombergGPT outperforms similarly-sized open models on financial NLP tasks by significant margins – without sacrificing performance on general LLM benchmarks*

https://stablecog.com/generate

MISTRAL AI_

## Mixtral of experts

A high quality Sparse Mixture-of-Experts.

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural



L. Fridman MIT Deep Learning https://deeplearning.mit.edu/

I. Goodfellow, Y. Bengio, A. Courville. Deep learning. MIT press, 2016.

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness



https://openai.com/blog/dall-e/



TEXT PROMPT
an armchair in the shape of an avocado […]
AI-GENERATED IMAGES

Edit prompt or view more images↓

TEXT PROMPT
a store front that has the word 'openai' written on it […]
AI-GENERATED IMAGES

Edit prompt or view more images↓

TEXT AND IMAGE PROMPT
the exact same cat on the top as a sketch on the bottom
AI-GENERATED IMAGES

Edit prompt or view more images↓
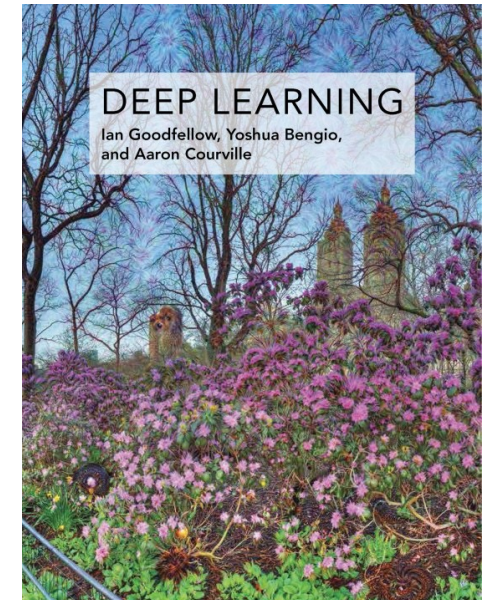
# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness

- Coolness :-)

MIT Introduction to Deep Learning http://introtodeeplearning.com

There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness

- Coolness :-)

There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness

- Coolness :-)



https://thispersondoesnotexist.com/

# There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness

- Coolness :-)

https://medium.com/lessons-from-history/he-used-ai-to-see-todays-looks-of-the-famous-people-from-the-past-3db43021ef39



Maria Milojković, MA
Dec 28, 2022 · 10 min read · Member-only · Listen

## He Used AI to See Today's Looks of The Famous People From the Past

And they look incredibly real

Napoleon Bonaparte | Credit: Hidreley Diao

Ludwig Van Beethoven, a German composer

Ludwig Van Beethoven | Credit: Hidreley Diao

There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness

- Coolness :-)

1. A large community of practitioners and researchers is using deep learning and sharing resources:

- https://paperswithcode.com/

- The Journal of Open Source Software https://joss.theoj.org/

2. A cornerstone book (*Deep Learning* by I. Goodfellow, Y. Bengio, A. Courville) released in late 2016 https://www.deeplearningbook.org/

There are 4 principal reasons for choosing Deep Learning as a modeling tool

- Performances

- Procedural

- Uniqueness

- Coolness :-)

*Dive into Deep Learning* by A. Zhang, Z. Lipton, M. Li, A. Smola
https://d2l.ai/d2l-en.pdf

1. A large community of practitioners and researchers is using deep learning and sharing resources:

-    https://paperswithcode.com/

- The Journal of Open Source Software https://joss.theoj.org/

2. A cornerstone book (*Deep Learning* by I. Goodfellow, Y. Bengio, A. Courville) released in late 2016
https://www.deeplearningbook.org/

**Geoffrey Hinton**

Emeritus Prof. Comp Sci, U.Toronto & Engineering Fellow, Google
Email verificata su cs.toronto.edu - Home page

machine learning    neural networks    artificial intelligence    cognitive

Citazioni per anno



2018 Turing Awards Recipients (+1 Nobel Prize in Physics!)

# Why now?

1952 — Stochastic Gradient Descent

1958 — Perceptron
- Learnable Weights

⋮

1986 — Backpropagation
- Multi-Layer Perceptron

1995 — Deep Convolutional NN
- Digit Recognition

⋮

# Why now?

- Data



| | |
|---|---|
| 1952 | Stochastic Gradient Descent |
| 1958 | Perceptron<br>• Learnable Weights |
| ⋮ | |
| 1986 | Backpropagation<br>• Multi-Layer Perceptron |
| 1995 | Deep Convolutional NN<br>• Digit Recognition |
| ⋮ | |

# Why now?

- Data
- Software



| Year | Event |
|------|-------|
| 1952 | Stochastic Gradient Descent |
| 1958 | Perceptron • Learnable Weights |
| 1986 | Backpropagation • Multi-Layer Perceptron |
| 1995 | Deep Convolutional NN • Digit Recognition |

TensorFlow

Coarse styles $(4^2 - 8^2)$
Middle styles $(16^2 - 32^2)$
Fine styles $(64^2 - 1024^2)$

# Why now?



- Data
- Software
- Hardware

# NVIDIA

BIT: 1NVDA ⋮

**Panoramica** | **Confronto** | **Dati finanziari**

Riepilogo mercato > NVIDIA

## 817,10 EUR

**+ Segui**

+775,35 (1.857,13%) ↑ ultimi 5 anni

11 apr, 15:50 CEST • Limitazione di responsabilità

| 1G | 5G | 1M | 6M | YTD | 1A | **5A** | Max |

732,60 EUR  23 feb 2024



| | | |
|---|---|---|
| Apertura | 813,80 | Capitalizz. | 2,19 Bln USD | Max 52 sett. | 892,30 |
| Massimo | 822,80 | Rapp. P/E | - | Min 52 sett. | 238,50 |
| Minimo | 804,60 | Div./prezzo | - | |

There are 2 principal reasons for **NOT** choosing Deep Learning as a modeling tool

- Complexity (all of the above)



GoogLeNet (incarnation of the Inception architecture)



This CVPR2015 paper is the Open Access version, provided by the Computer Vision Foundation.
The authoritative version of this paper is available in IEEE Xplore.

**Going Deeper with Convolutions**

Christian Szegedy[1], Wei Liu[2], Yangqing Jia[1], Pierre Sermanet[1], Scott Reed[3],
Dragomir Anguelov[1], Dumitru Erhan[1], Vincent Vanhoucke[1], Andrew Rabinovich[4]

[1]Google Inc. [2]University of North Carolina, Chapel Hill
[3]University of Michigan, Ann Arbor [4]Magic Leap Inc.

[1]{szegedy,jiayq,sermanet,dragomir,dumitru,vanhoucke}@google.com

[2]wliu@cs.unc.edu, [3]reedscott@umich.edu, [4]arabinovich@magicleap.com

**References**

[1] Know your meme: We need to go deeper.
http://knowyourmeme.com/memes/we-need-to-go-deeper.
Accessed: 2014-09-15.

C. Szegedy et al. *Going Deeper with Convolutions* CVPR2015

There are 2 principal reasons for **NOT** choosing Deep Learning as a modeling tool

- Complexity (all of the above)

Many desiderata in Machine Learning whrere complexity is typically bad: trust, robustness, interpretability, fairness...



$$x$$
"panda"
57.7% confidence

$$+ .007 \times$$

$$\text{sign}(\nabla_x J(\theta, x, y))$$
"nematode"
8.2% confidence

$$=$$

$$x + \epsilon\,\text{sign}(\nabla_x J(\theta, x, y))$$
"gibbon"
99.3 % confidence

I. Goodfellow *Explaining and harnessing adversarial examples*, ICLR 2015

There are 2 principal reasons for **NOT** choosing Deep Learning as a modeling tool

- Complexity (all of the above)

- Need for huge datasets (but there is *transfer learning/domain adaptation*...)



L. Fridman *Deep Learning* https://deeplearning.mit.edu/

There are 2 principal reasons for **NOT** choosing Deep Learning as a modeling tool

- Complexity (all of the above)

- Need for huge datasets (but there is *transfer learning/domain adaptation*…)



A. Ng *Nuts and Bolts of Applying Deep Learning* https://www.youtube.com/watch?v=F1ka6a13S9I

There are 2 principal reasons for **NOT** choosing Deep Learning as a modeling tool

- Complexity (all of the above)

- Need for huge datasets (but there is *transfer learning/domain adaptation...*)

Why do tree-based models still outperform deep learning on typical tabular data?
Part of Advances in Neural Information Processing Systems 35 (NeurIPS 2022)

Figure 1: **Benchmark on medium-sized datasets**, top only numerical features; bottom: all features. Dotted lines correspond to the score of the default hyperparameters, which is also the first random search iteration. Each value corresponds to the test score of the best model (on the validation set) after a specific number of random search iterations, averaged on 15 shuffles of the random search order. The ribbon corresponds to minimum and maximum scores on these 15 shuffles.

# Basics of Deep Learning

# The building block of Neural Networks (NN): the neuron

- Perceptron (Rosenblatt 1958)

$x_1$

$x_2 \longrightarrow$ output
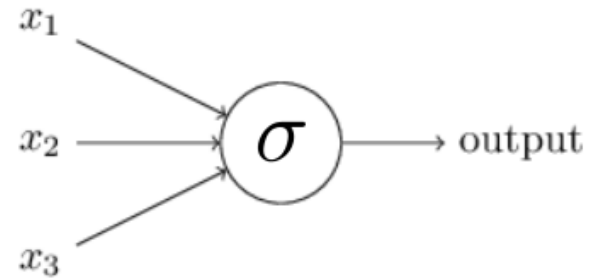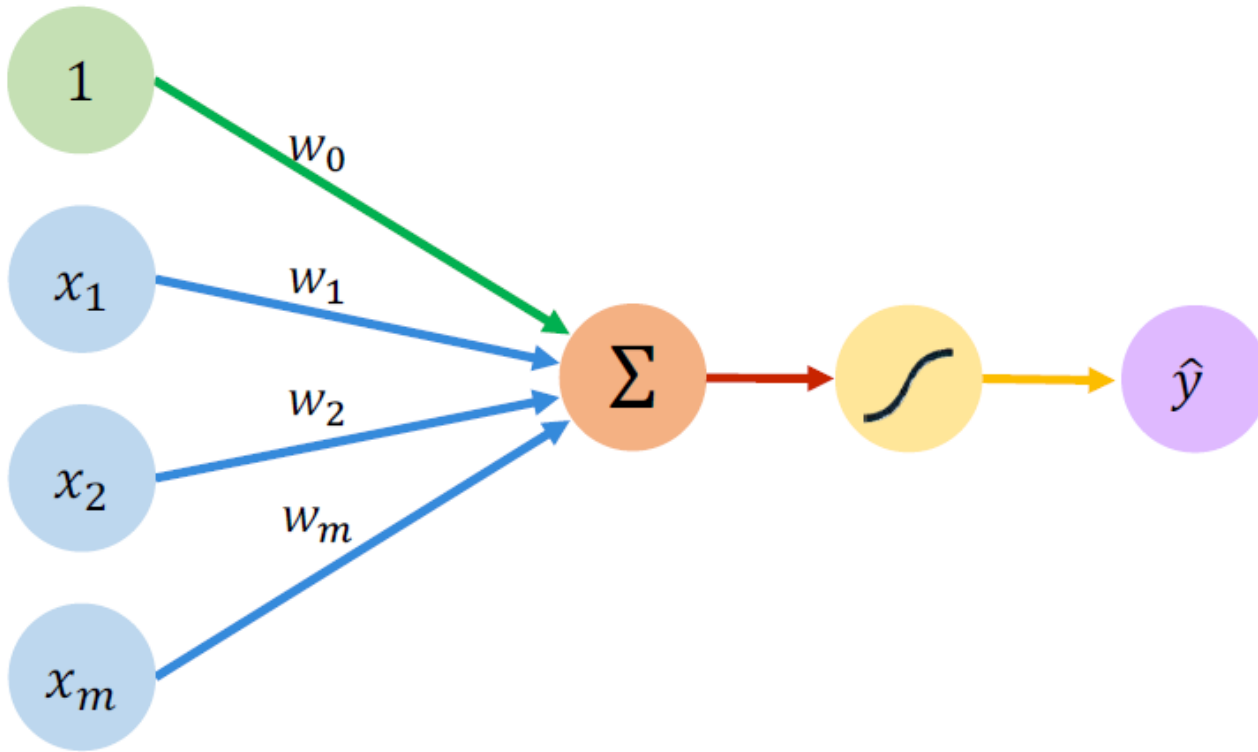
$x_3$

$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j < \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j \geq \text{threshold} \end{cases}$$

impulses carried
toward cell body

branches
of axon

dendrites

nucleus

axon

impulses carried
away from cell body

cell body

axon
terminals

Neuro
block

# The building block of Neural Networks (NN): the neuron

- Perceptron (Rosenblatt 1958)

$x_1$

$x_2 \longrightarrow$ output

$x_3$

$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j < \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j \geq \text{threshold} \end{cases}$$

- Sigmoid unit

$x_1$

$x_2 \longrightarrow \sigma \longrightarrow$ output

$x_3$

$$\text{output} = \sigma\left(\sum_j w_j x_j + b\right)$$

$$= \frac{1}{1 + \exp\left[-\left(\sum_j w_j x_j + b\right)\right]}$$

# The building block of Neural Networks (NN): the neuron



Inputs    Weights    Sum    Non-Linearity    Output

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

$$\text{where: } X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

# The building block of Neural Networks (NN): the neuron



Activation function

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

$$\text{where: } X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

Inputs     Weights     Sum     Non-Linearity    Output

# Activation functions

Linear Activation functions produce linear decisions no matter the network size

Non-linearities allow us to approximate arbitrarily complex functions

# Activation functions

### Sigmoid Function



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

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

### Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$
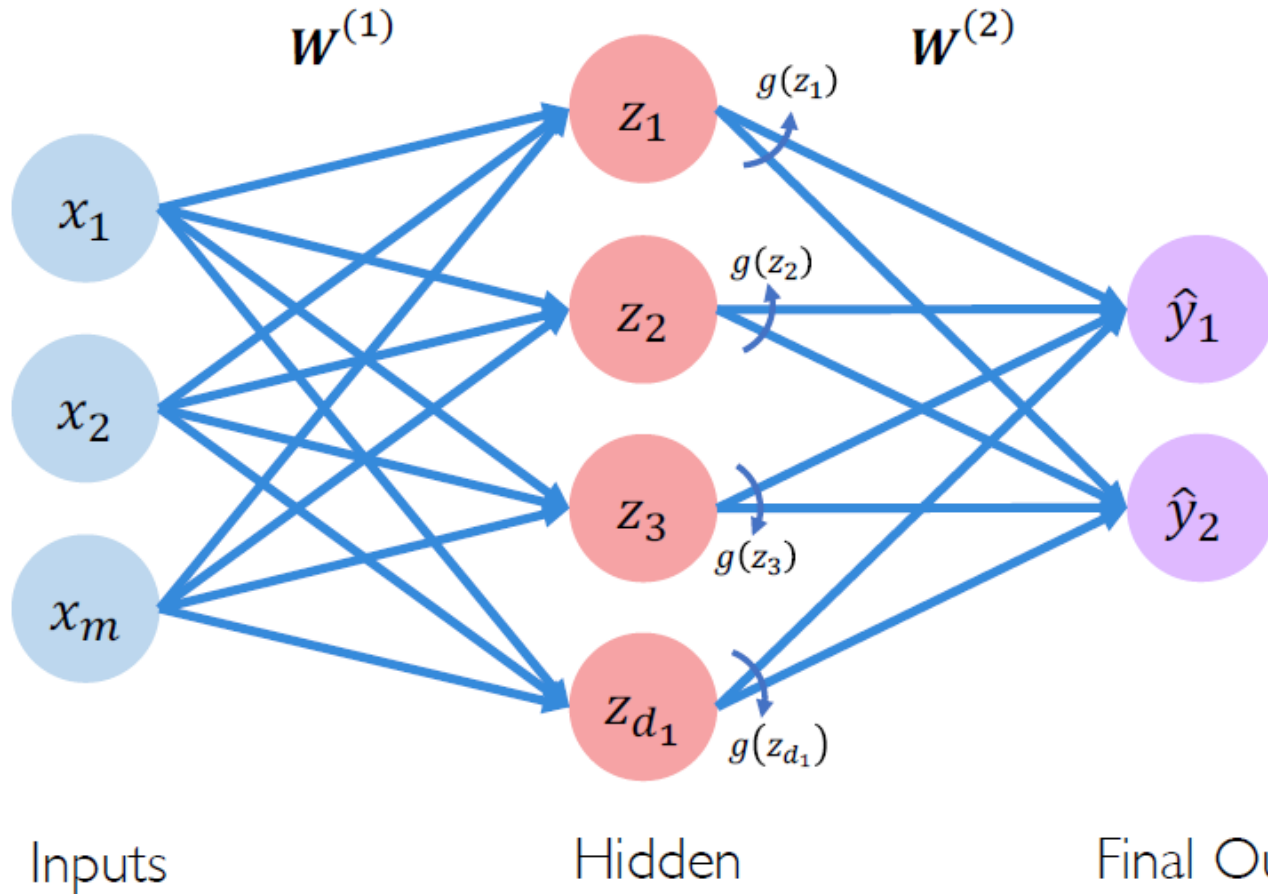
$$g'(z) = 1 - g(z)^2$$

### Rectified Linear Unit (ReLU)



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

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$
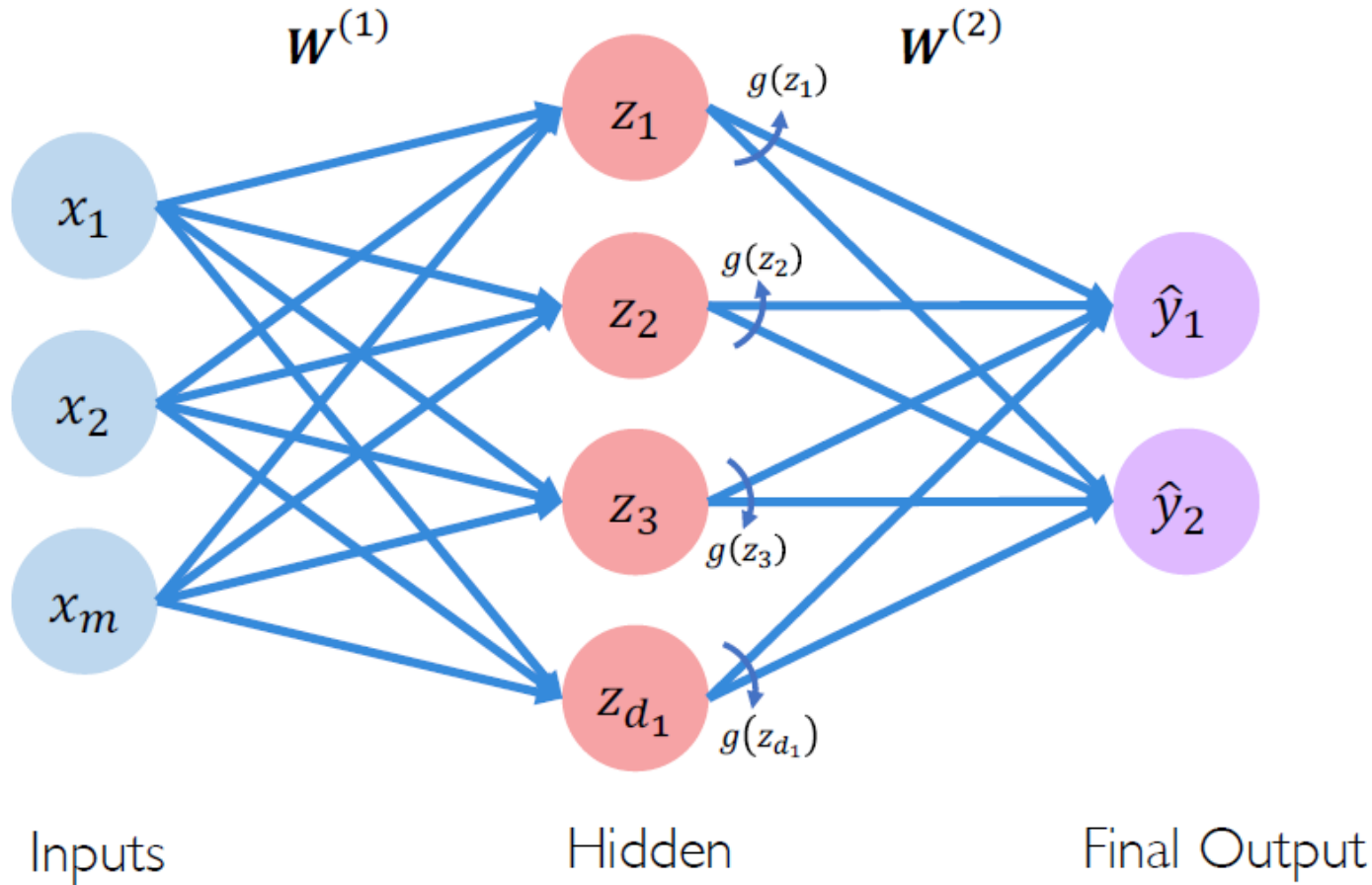
# Feedforward (Single Layer) Neural Network



$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^{m} x_j \, w_{j,i}^{(1)} \qquad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j \, w_{j,i}^{(2)}\right)$$

So-called *Vanilla Neural Network*

T. Hastie, R. Tibshirani, J. Friedman *The Elements of Statistical Learning*
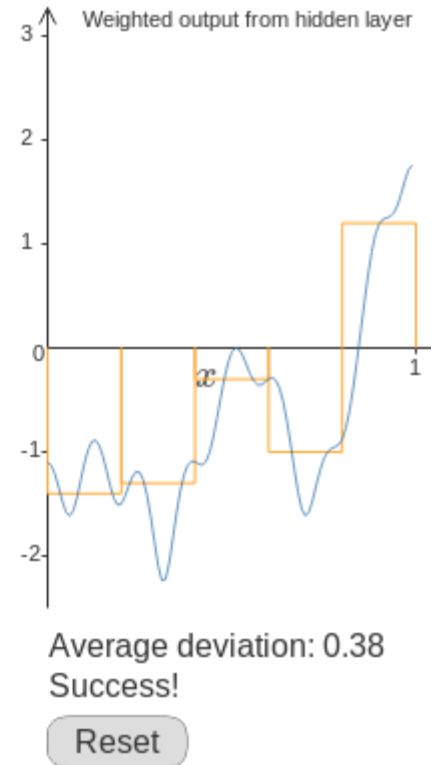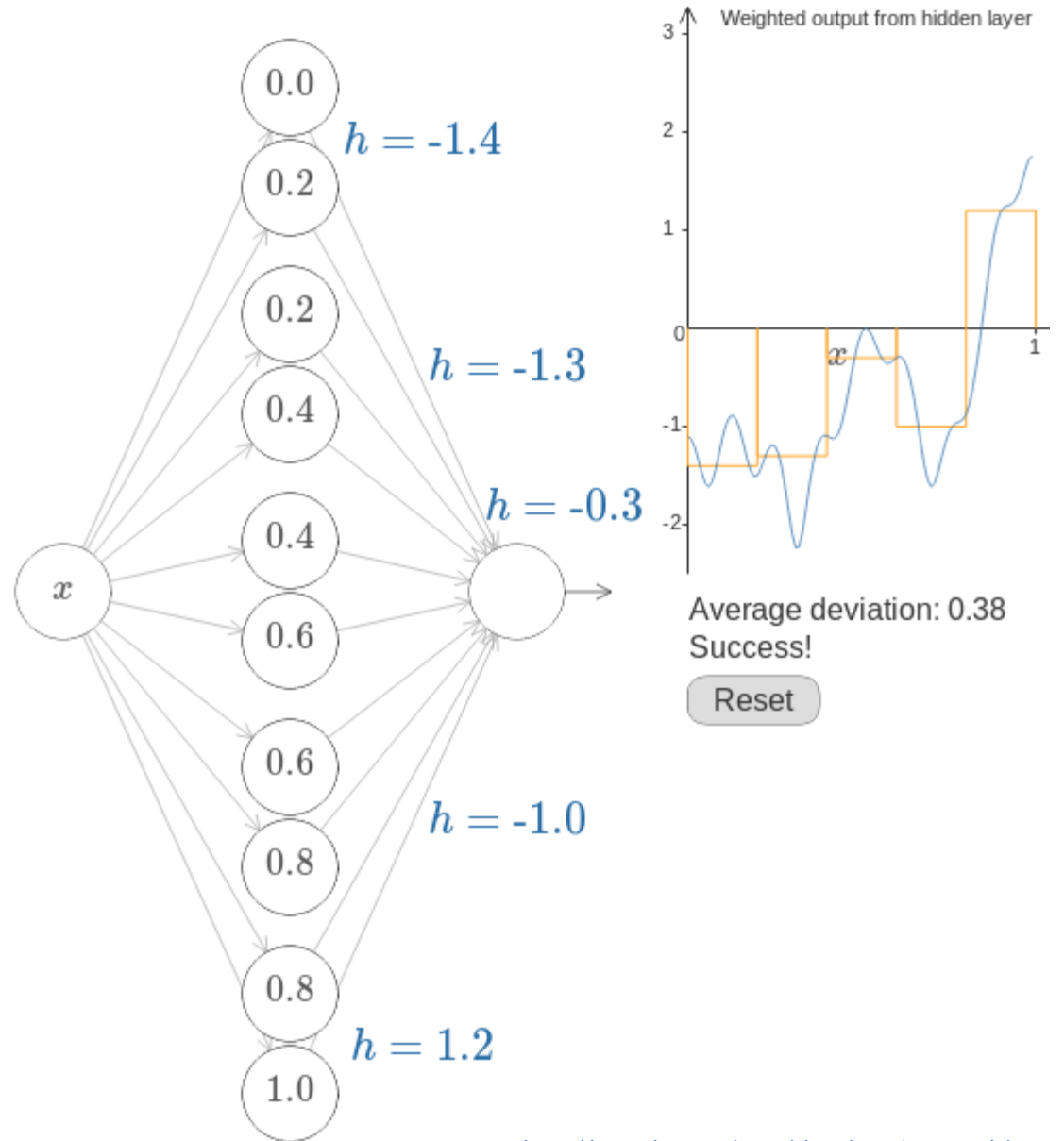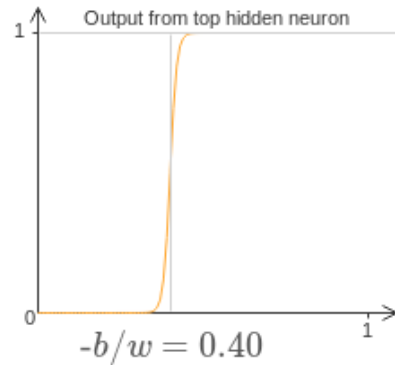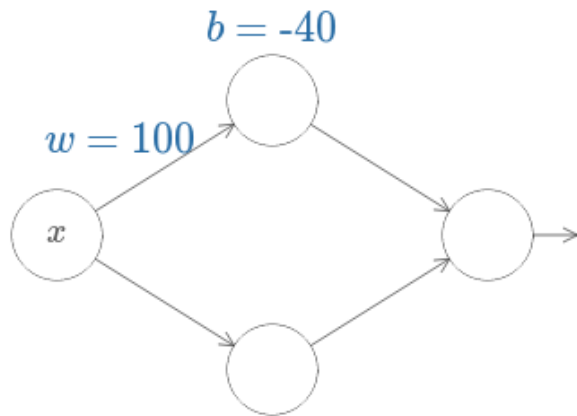
# Feedforward (Single Layer) Neural Network



Concepts:

- Fully connected network

- Architecture design:

• Size of the hidden layer

• Weights

• Initialization
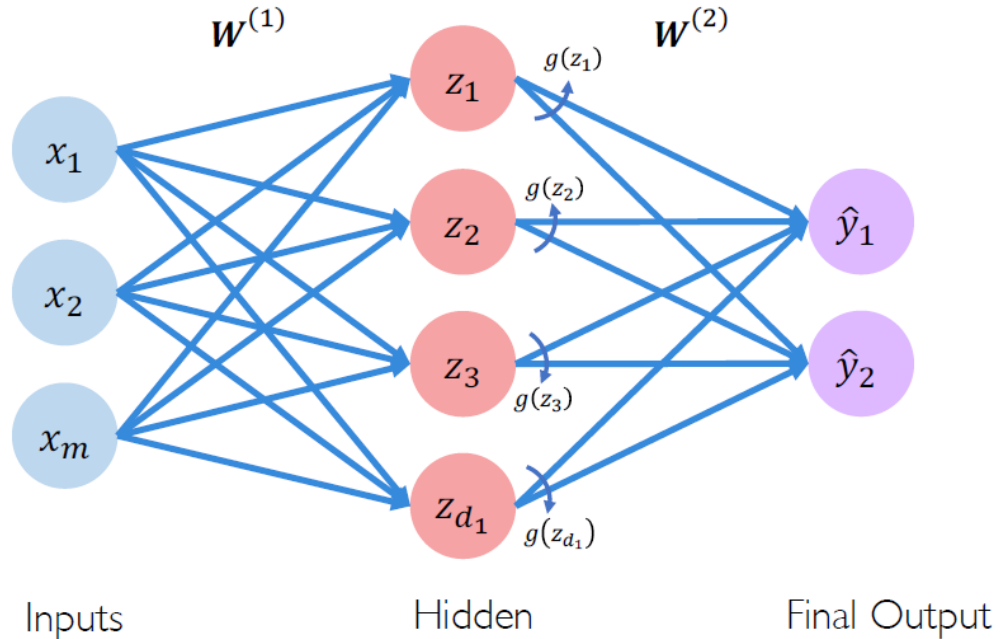
# Universality

Neural nets with a single hidden layer can be used to approximate any continuous function to any desired precision.
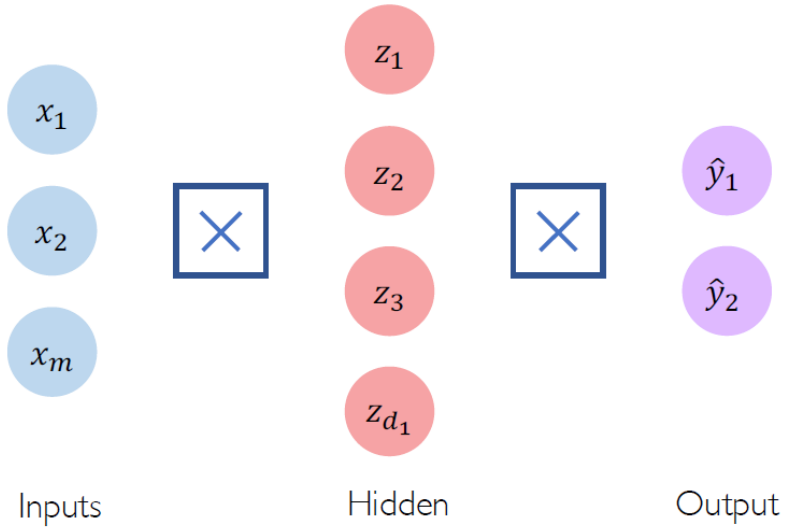


$$b = -40$$

$$w = 100$$

$$x$$

Output from top hidden neuron

$$-b/w = 0.40$$

0.0
$$h = -1.4$$
0.2

0.2
$$h = -1.3$$
0.4

0.4
$$h = -0.3$$

$$x$$

0.6

0.6
$$h = -1.0$$
0.8

0.8
$$h = 1.2$$
1.0

Weighted output from hidden layer

$$x$$

Average deviation: 0.38
Success!

Reset

# From Neural Network to Deep Learning architecture



$W^{(1)}$   $W^{(2)}$

$x_1$
$x_2$
$x_m$

$z_1$ $g(z_1)$
$z_2$ $g(z_2)$
$z_3$ $g(z_3)$
$z_{d_1}$ $g(z_{d_1})$

$\hat{y}_1$
$\hat{y}_2$

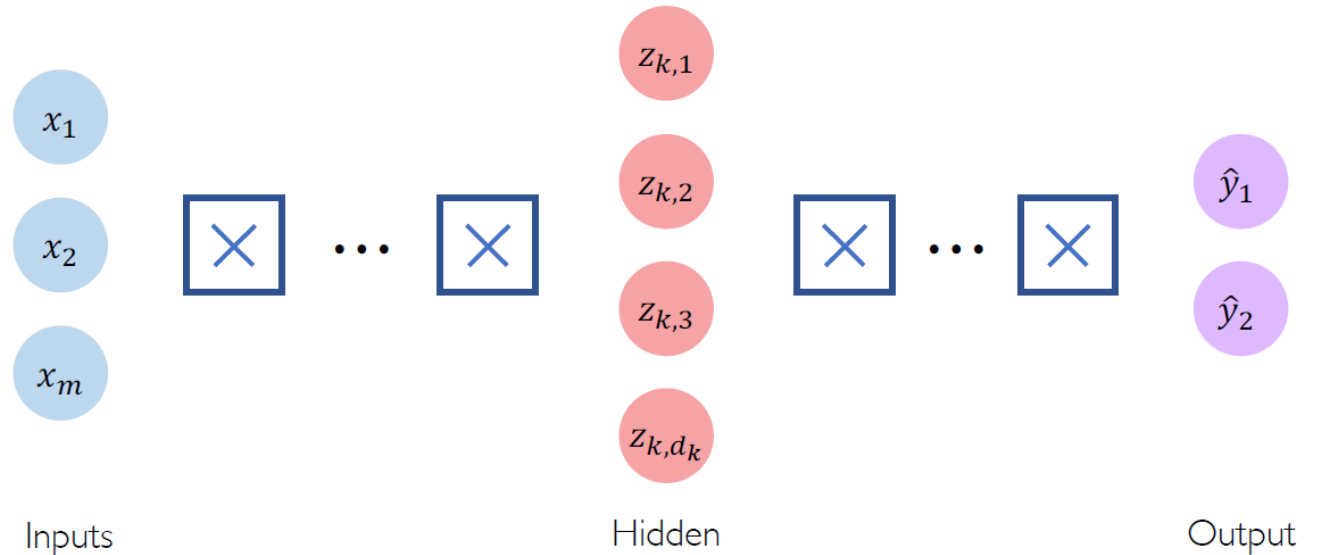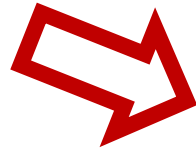Inputs     Hidden     Final Output

MIT *Introduction to Deep Learning* http://introtodeeplearning.com

- A Vanilla Neural Network is a Universal Approximator*: any functions can be described by it

- However, some functions can be compactely represented with k layers, while they may need an exponential number of neurons with 1 laver

• M. Nielsen 'Neural Networks and Deep Learning'Determination Press, 2015 http://neuralnetworksanddeeplearning.com/chap4.html

# From Neural Network to Deep Learning architecture
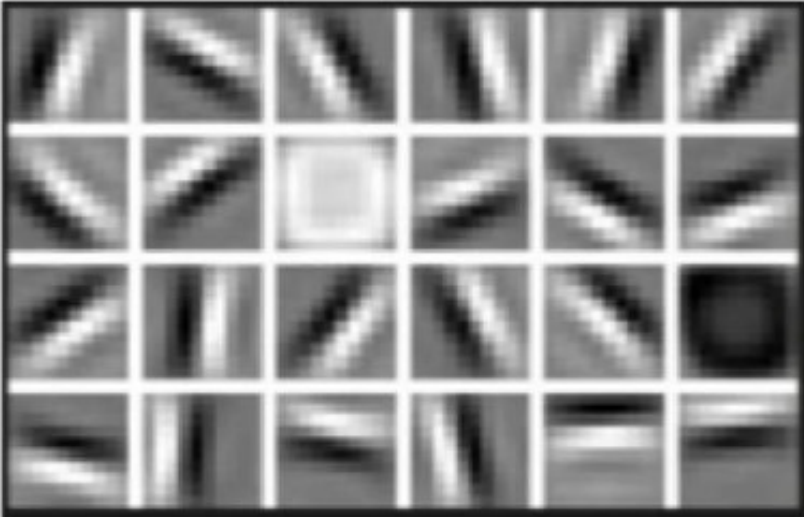


From 1 Layer to 'many' layers

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) \, w_{j,i}^{(k)}$$

# From Neural Network to Deep Learning architecture



Low Level Features — Lines & Edges

Mid Level Features — Eyes & Nose & Ears

High Level Features — Facial Structure

Lee, Honglak, et al. "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations." *Proceedings of the 26th annual international conference on machine learning*. ACM, 2009.

# Training a Neural Network: Loss Function

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

Predicted     Actual

A *Loss* expresses a cost associated with the NN predictions

The *empirical loss* (also known as *objective function, cost function, empirical risk*) measures the total loss over the (training/test) dataset

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

Predicted     Actual

# Training a Neural Network: Loss Function

Depending on the task we can have different losses:

- *Cross entropy loss*: for classification tasks, NNs that have in output a probability

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

$\underset{\text{Actual}}{\underline{\phantom{xx}}}$ $\underset{\text{Predicted}}{\underline{\phantom{xxx}}}$ $\underset{\text{Actual}}{\underline{\phantom{xx}}}$ $\underset{\text{Predicted}}{\underline{\phantom{xxx}}}$

$$f(x) \qquad y$$
$$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

- *Mean squared error loss* for regression tasks

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

$\underset{\text{Actual}}{\underline{\phantom{xx}}}$ $\underset{\text{Predicted}}{\underline{\phantom{xxx}}}$
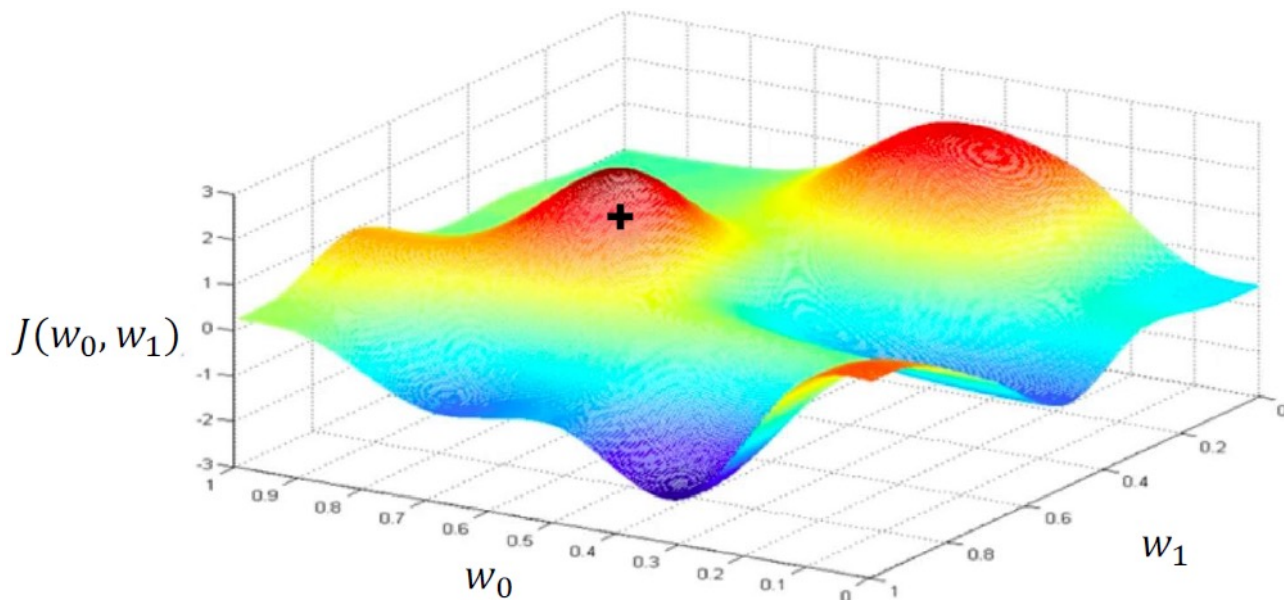
# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

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

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



$J(w_0, w_1)$

$w_0$

$w_1$

# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

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

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



$J(w_0, w_1)$

$w_0$

$w_1$

**Algorithm**

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

# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

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

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



$J(w_0, w_1)$

**Algorithm**

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

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

# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

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

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



**Algorithm**

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

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

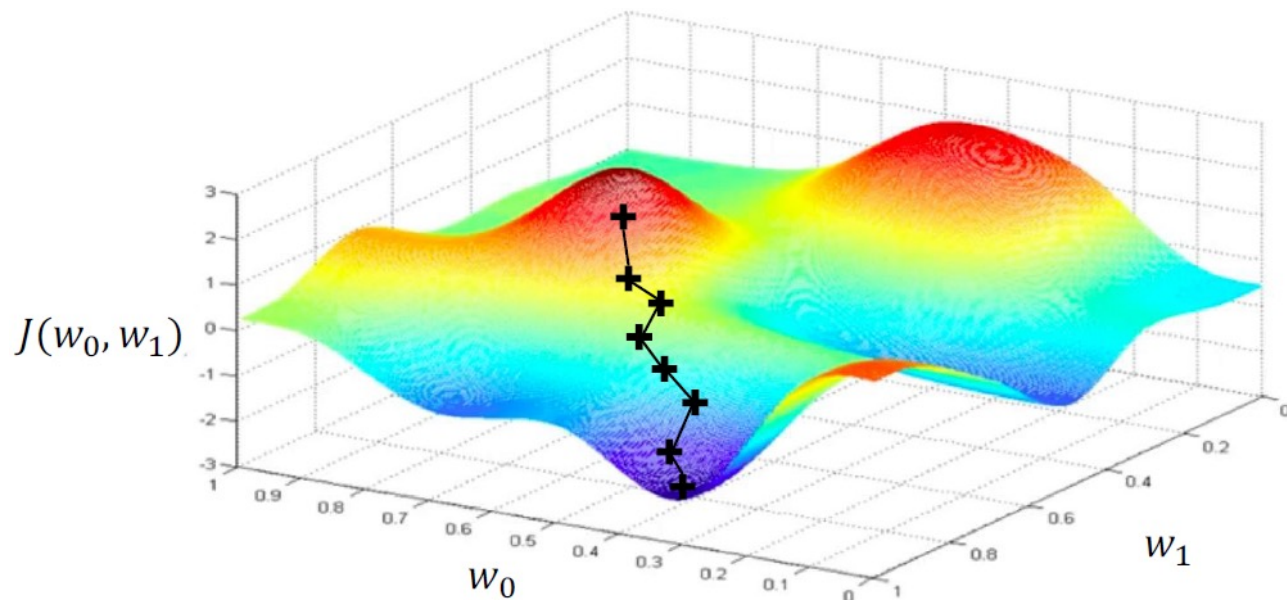   Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

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

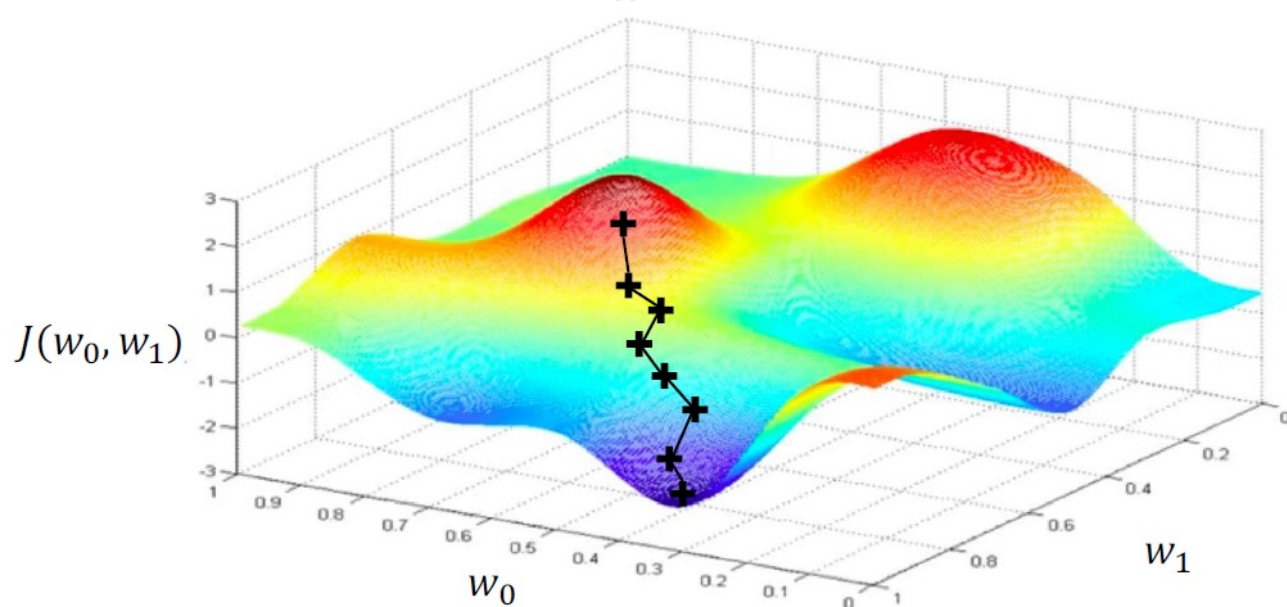$$W^* = \underset{W}{\mathrm{argmin}}\, J(W)$$



**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, $W \leftarrow W - \eta\,\dfrac{\partial J(W)}{\partial W}$

5. Return weights

# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

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

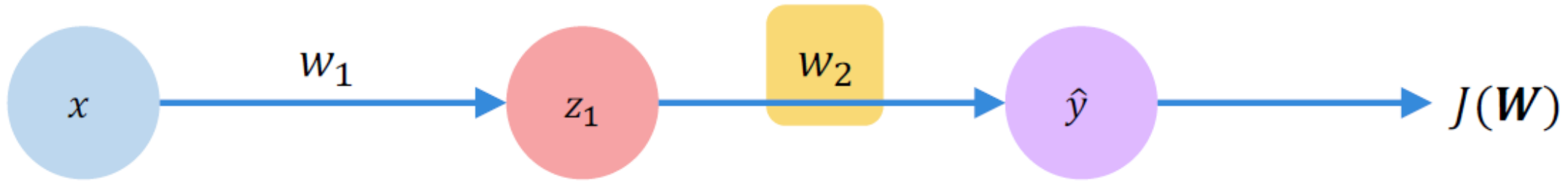$$W^* = \underset{W}{\text{argmin}} \, J(W)$$



**Algorithm** (Gradient Descent)

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, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

                         Learning Rate

5. Return weights
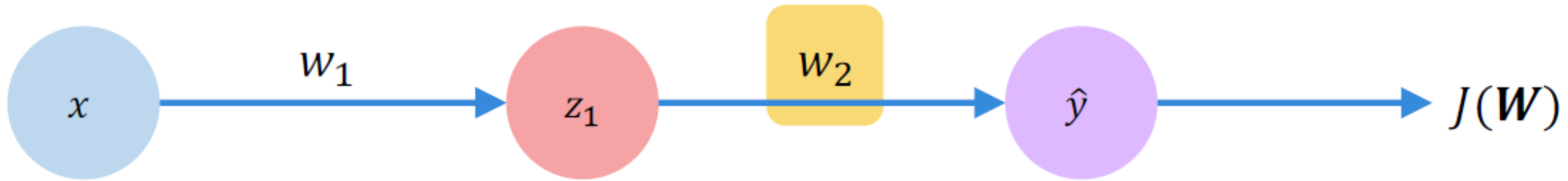
# How to compute the gradient: backpropagation

Backpropagation is about understanding how changing the weights and biases in a network changes the cost function.



Let's consider a simple NN with one node: how the final loss is affected by changes in $w_2$ ?

# How to compute the gradient: backpropagation

Let's consider a simple NN with one node: how the final loss is affected by changes in $w_2$
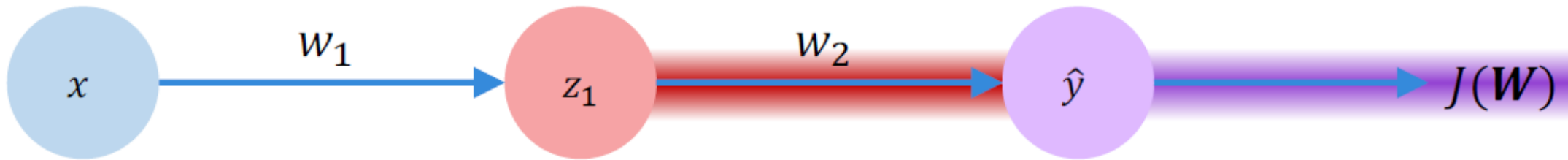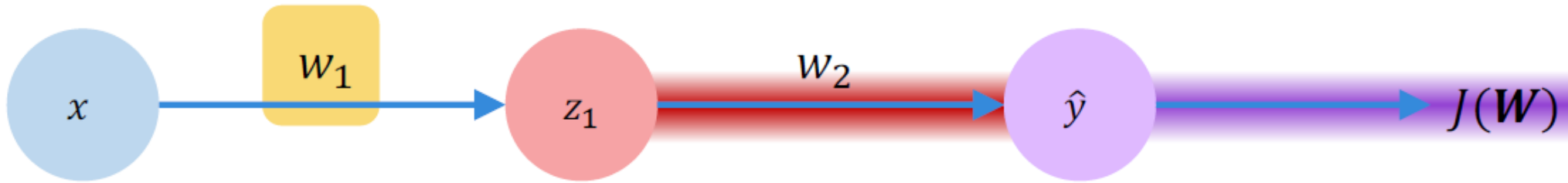


We apply the chain rule!

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

# How to compute the gradient: backpropagation

Let's consider a simple NN with one node: how the final loss is affected by changes in $w_2$



We apply the chain rule!

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

# How to compute the gradient: backpropagation

This simple network is characterized also by the weigth $w_1$



$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$
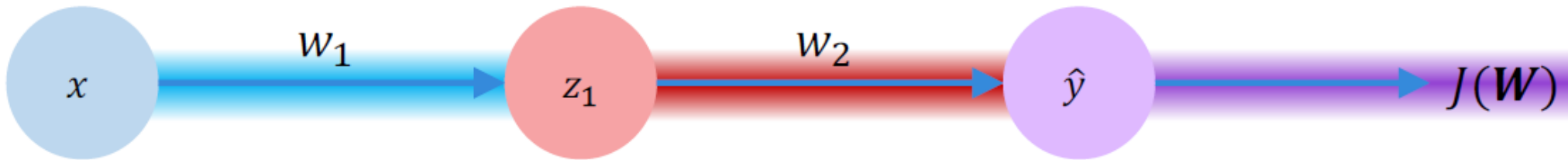
Apply chain rule!          Apply chain rule!

# How to compute the gradient: backpropagation

This simple network is characterized also by the weigth $w_1$



$$\frac{\partial J(\boldsymbol{W})}{\partial w_1} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$
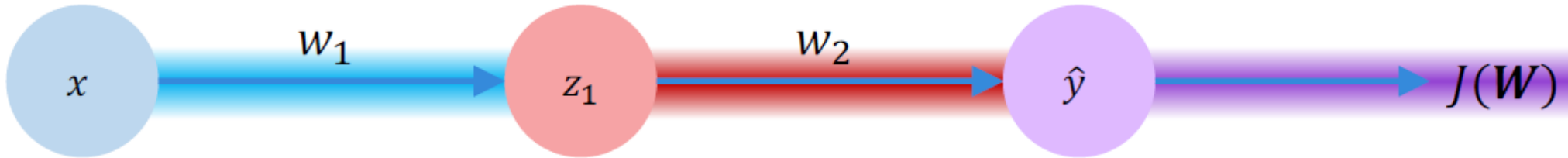
This procedure has to be done for each parameter of the newtwork!

# How to compute the gradient: backpropagation

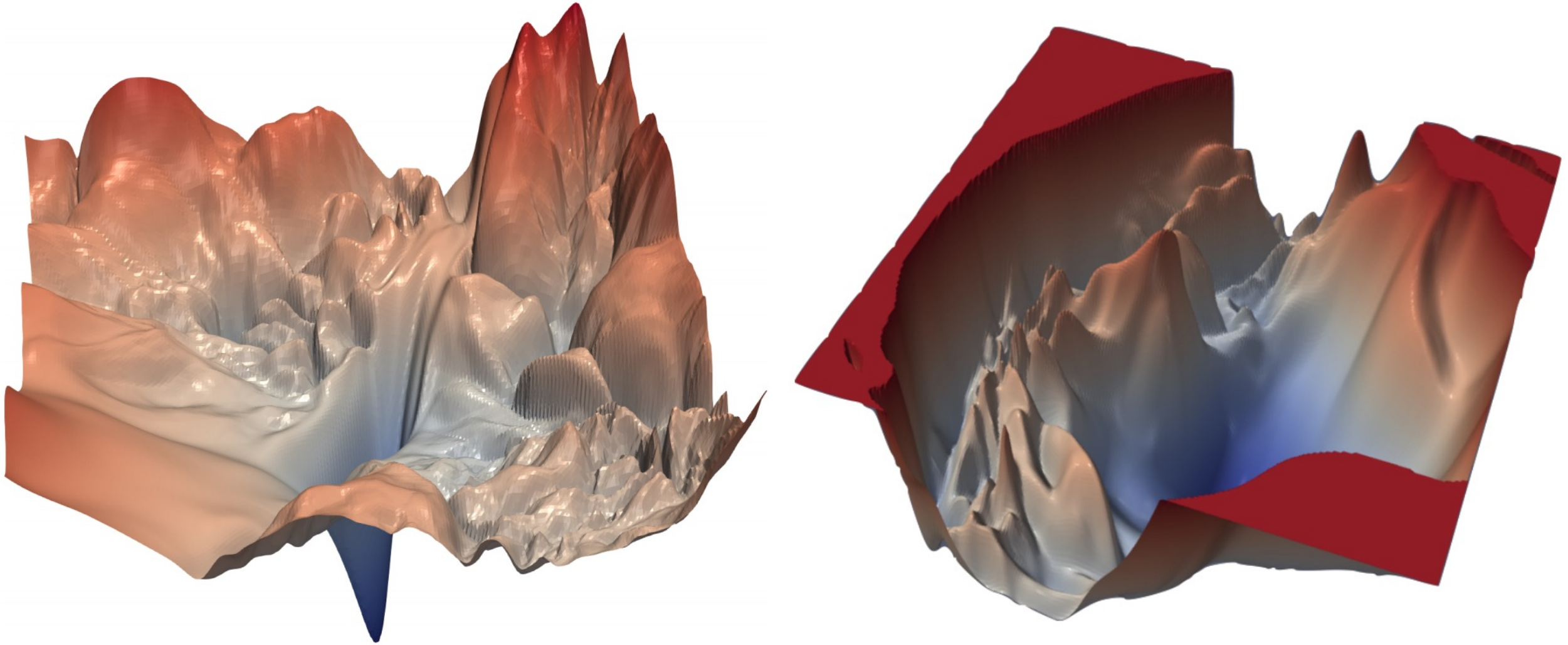This simple network is characterized also by the weigth $w_1$



$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

You can have a look here for the math: http://neuralnetworksanddeeplearning.com/chap2.html

# Loss functions may be quite complex...



H. Li et al. *Visualizing the Loss Landscape of Neural Nets* NIPS 2018

# To cope with complex loss functions and to optimize training the following algorithm is generally sophisticated
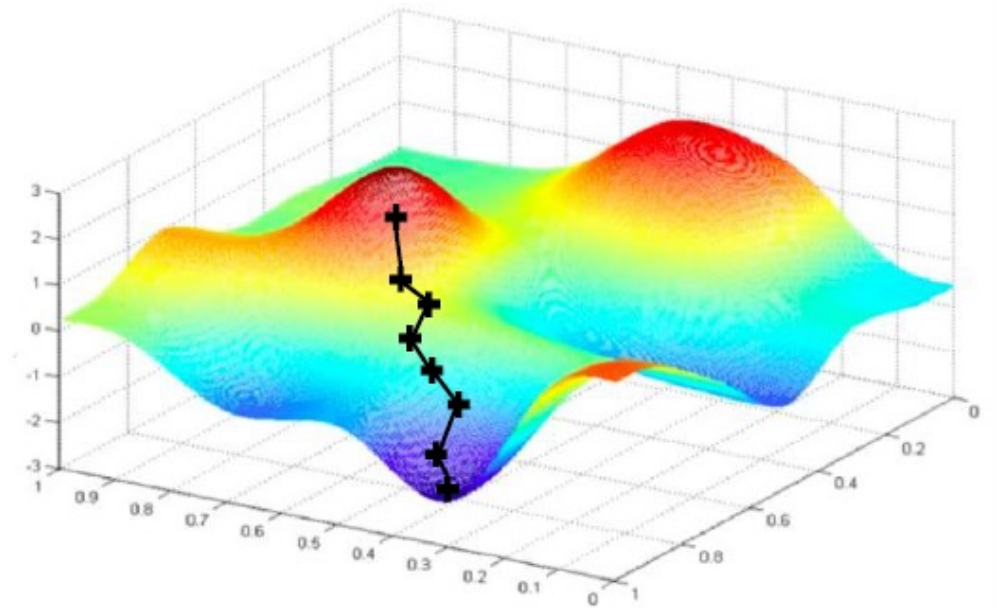
**Algorithm** (Gradient Descent)

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, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

   Learning Rate

5. Return weights

# Stochastic Gradient Descent

**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, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$
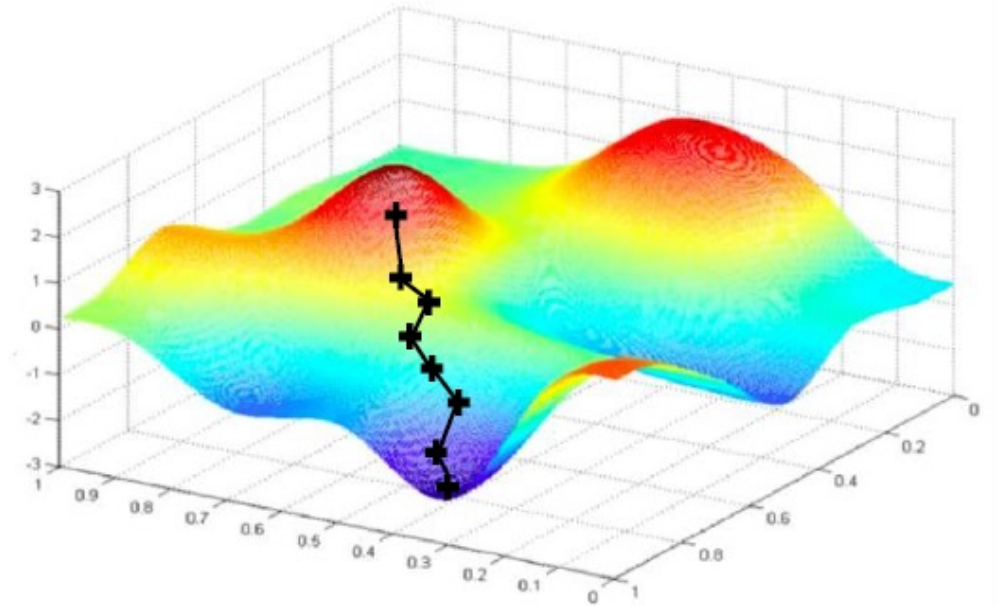
5. Return weights

Can be very burdensome
to compute...
Average over all samples
in the dataset!

# Stochastic Gradient Descent

## Algorithm

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

2. Loop until convergence:

3.   Pick single data point $i$

4.   Compute gradient, $\dfrac{\partial J_i(\boldsymbol{W})}{\partial \boldsymbol{W}}$

5.   Update weights, $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$
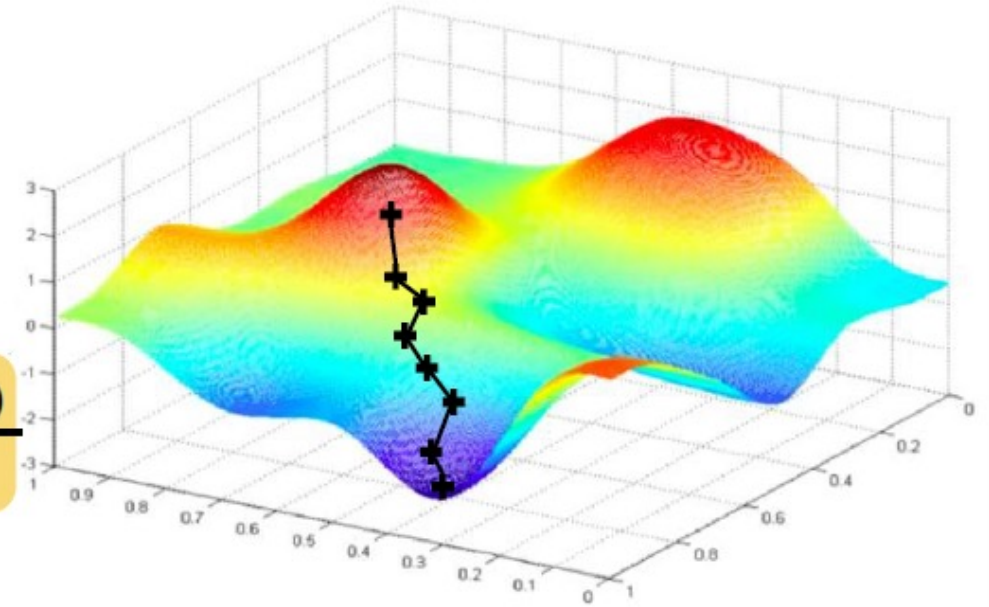
6. Return weights

Easy to compute but **very noisy** (stochastic)!

# Stochastic Gradient Descent

## Algorithm

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

2. Loop until convergence:

3.      Pick batch of $B$ data points

4.      Compute gradient, $\dfrac{\partial J(W)}{\partial W} = \dfrac{1}{B} \sum_{k=1}^{B} \dfrac{\partial J_k(W)}{\partial W}$

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

6. Return weights

Fast to compute and a much better estimate of the true gradient!

# Stochastic Gradient Descent: Mini-batches

The set of $B$ data points is called *mini-batch*

Mini-batches allow to accurated estimation of gradient, smoother convergence, larger learning rates

Mini-batches lead to fast training: computation can be parallelized and significant speed increases can be obtained on GPUs

# Thank you!

# Gian Antonio Susto