

DEEP LEARNING

CSCS – USI SUMMER SCHOOL 2020

Simon Scheidegger – simon.scheidegger@unil.ch

University of Lausanne, Department of Finance

July 22nd – 23rd, 2020

THE RISE OF NEURAL NETWORKS

TOM SIMONITE BUSINESS 01.25.19 01:05 PM

DEEPMIND BEATS PROS AT STARCRAFT IN ANOTHER TRIUMPH FOR BOTS



A pro gamer and an AI bot duke it out in the strategy game StarCraft, which has become a benchmark for research on artificial intelligence. © STARCRAFT

IN LONDON LAST month, a team from Alphabet's UK-based artificial intelligence research unit DeepMind quietly placed a new marker in the contest between humans and computers. On Thursday it revealed the achievement in a three-hour-long YouTube stream, in which aliens and robots fought to the death.

TOM SIMONITE BUSINESS 10.10.17 01:00 PM

THIS MORE POWERFUL VERSION OF ALPHAGO LEARNS ON ITS OWN



NOAH SHELDON FOR WIRED

AT ONE POINT during his historic defeat to AlphaGo last year, world champion Go player Lee Sedol abruptly left the room. The bot had played a move that confounded established theories of the board, a moment that came to epitomize the mysteriousness of AlphaGo.

NEWS BIOLOGY 21 DECEMBER 2017

AI beats docs in cancer spotting

A new study provides a fresh example of machine learning as an important diagnostic tool. Paul Biegler reports.



Deep Learning Software Speeds Up Drug Discovery

Wed, 01/16/2019 - 8:00am 1 Comment by Kenny Walter , Science Reporter - @RandDMagazine



The long, arduous process of narrowing down millions of chemical compounds to just a select few that can be further developed into mature drugs, may soon be shortened, thanks to new artificial intelligence (AI) software.

MUSIC GENERATED BY AI

- <https://openai.com/blog/musenet/>



MuseNet

We've created MuseNet, a deep neural network that can generate 4-minute musical compositions with 10 different instruments, and can combine styles from country to Mozart to the Beatles. MuseNet was not explicitly programmed with our understanding of music, but instead discovered patterns of harmony, rhythm, and style by learning to predict the next token in hundreds of thousands of MIDI files. MuseNet uses the same general-purpose unsupervised technology as GPT-2, a large-scale transformer model trained to predict the next token in a sequence, whether audio or text.

APRIL 26, 2018
A MINUTE READ, 16 MINUTE LISTEN

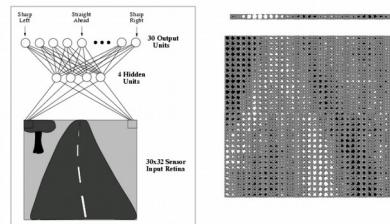
STYLE TRANSFER



Figs. from https://www.tensorflow.org/tutorials/generative/style_transfer

SELF-DRIVING CARS

- Carnegie Mellon University – 1990ies
 - ALVINN: Autonomous Land Vehicle In a Neural Network



- Today (e.g., Waymo)
 - <https://www.youtube.com/watch?v=LSX3qdy0dFg>

APP: COLORING OLD MOVIES

- <https://deepsense.ai/ai-movie-restoration-scarlett-ohara-hd/>



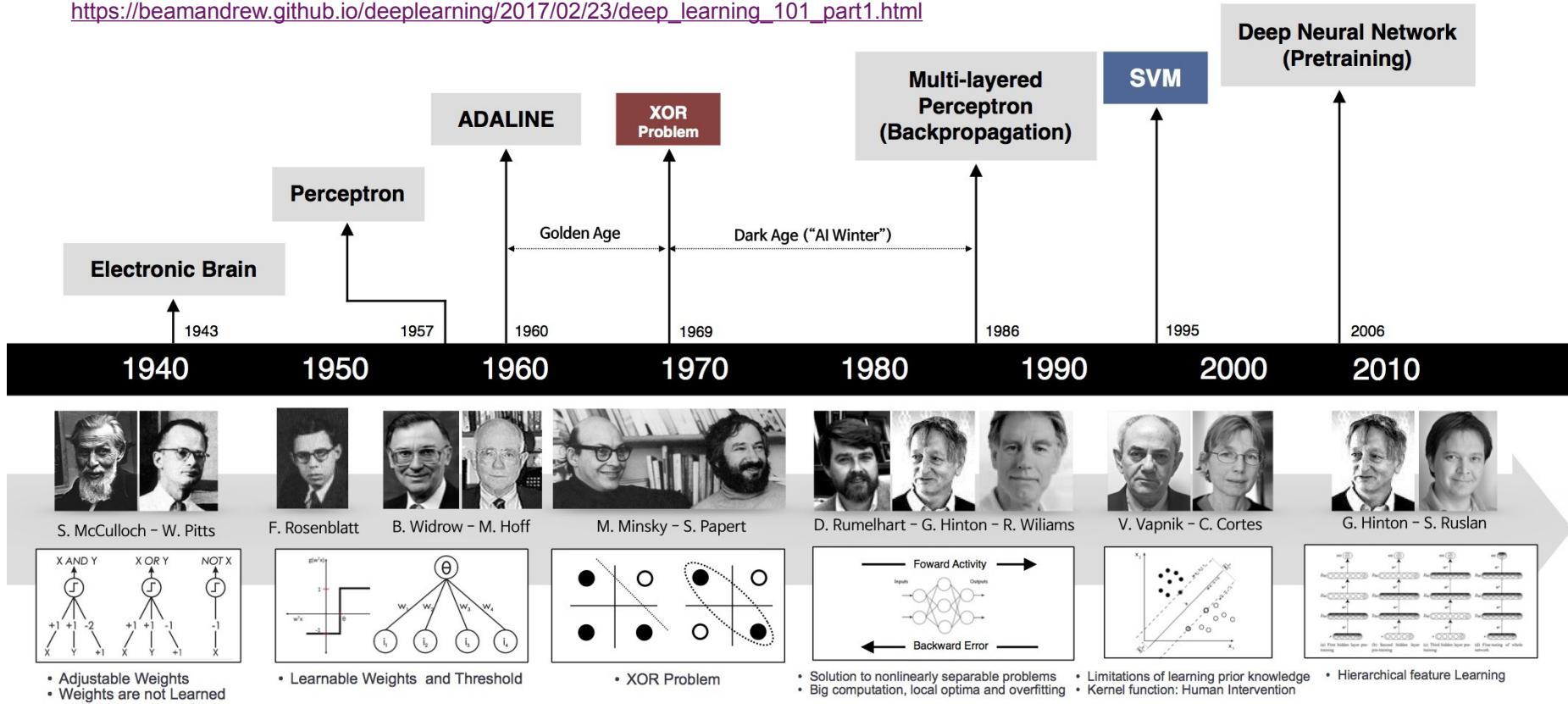
TWO-LEGGED ROBOTS

- <https://www.youtube.com/watch?v=hSjKoEva5bg&feature=youtu.be>



A TIMELINE OF DEEP LEARNING

https://beamandrew.github.io/deeplearning/2017/02/23/deep_learning_101_part1.html



WHY NOW?

- Neural Networks date back decades, so why the resurgence?
(Stochastic Gradient Descent: 1952, Perceptron: 1958, Back-propagation: 1986, Deep Convolutional NN: 1995)
- **Big Data**
 - Large Datasets
 - Easier Collection and Storage
- **Hardware**
 - GPUs, TPUs,...
- **Software**
 - Improved Techniques
 - Toolboxes



 TensorFlow  PyTorch

ROAD-MAP – WED (9.00-12.00, 1.30-3.30)

- Lecture (5 x 45-50 min)
 - A brief recap on Machine Learning Basics
 - Deep Learning Basics
 - White-box examples:
 - The multi-layer perceptron
 - Feed-forward networks
 - Network training – SGD
 - Error back-propagation
 - Some notes on over-fitting
- Throughout lectures – hands-on:
 - Perceptron
 - Gradient descent
 - Artificial neural networks: a simple multi-layer perceptron implementation & several examples

ROAD-MAP – THU (9.00-12.00, 1.30-3.30)

- Lecture (5 x 45-50 min):
 - Deep Learning cont'd
 - More advanced topics:
 - Recurrent neural networks and beyond.
 - Reinforcement learning.
 - Detour – think out of the box: "Deep Equilibrium Nets"
- Throughout lectures – hands-on:
 - Basics on Tensorflow & Keras
 - Examples related to the day's topics in Tensorflow

EXPECTATIONS MANAGEMENT

What are those lectures on deep learning about:

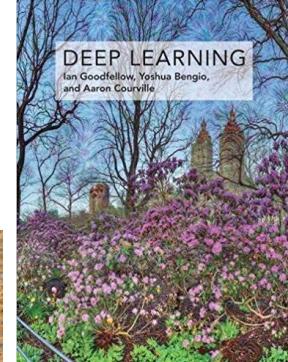
- This is a 2-day introduction to Artificial Neural Networks, Deep Learning, and the related topics.
- We will have only time to cover the basics.
- After those lectures, you will need to practice, practice, practice to become a master.

SOME USEFUL MATERIALS

Deep Learning

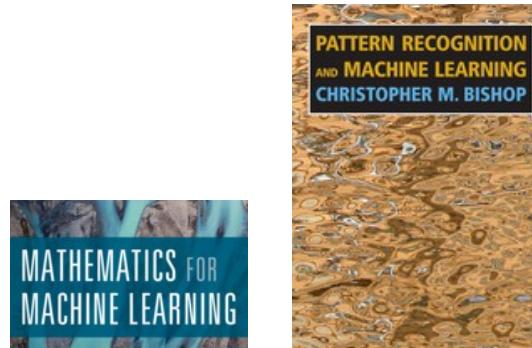
Ian Goodfellow and Yoshua Bengio and Aaron Courville
MIT Press 2016

<http://www.deeplearningbook.org/>



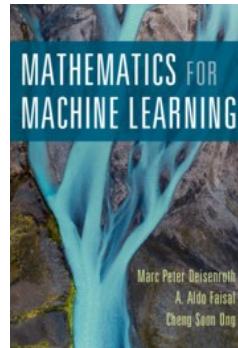
Pattern Recognition and Machine Learning

C. M Bishop, Springer 2006
(pdf freely available)



Mathematics for Machine Learning

Deisenroth, A. Aldo Faisal, and Cheng Soon Ong.
Cambridge University Press 2020



→ *There is a great community out there (use your browser and Google around...)*

FEW WORDS ABOUT MYSELF

- Prof. of advanced data analytics at the University of Lausanne.
- Ph.D. in theoretical physics (Core-collapse supernovae simulations).
- Research interest in computational economics and finance, HPC, and ML applied to those fields.

WEDNESDAY, JULY 22ND, 2020



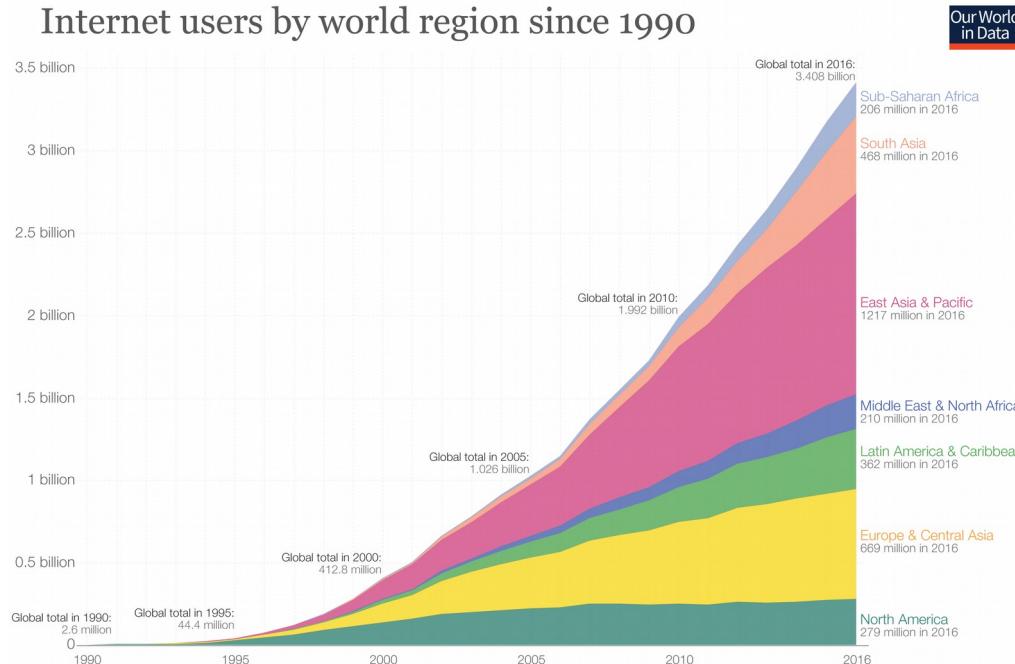
RECAP ON MACHINE LEARNING



BIG DATA AND ITS AVAILABILITY

<https://ourworldindata.org/internet>

Internet users by world region since 1990



Data source: Based on data from the World Bank and data from the International Telecommunications Union. Internet users are people with access to the worldwide network.
The interactive data visualization is available at OurWorldInData.org. There you find the raw data and more visualizations on this topic.

Licensed under CC-BY-SA by the author Max Roser.

OTHER SOURCES OF BIG DATA

■ Scientific experiments

- Cern (e.g., LHC) generates ~ 25 petabytes per year (2012).
- LIGO.
- generates ~ 1 Petabyte per year.



<https://home.cern/>

■ Numerical computations

■ ...



<https://www.ligo.caltech.edu/>

Unil

UNIL | Université de Lausanne

BIG DATA AND ITS AVAILABILITY

- Size of the internet as we speak: **~17.567.134 Petabyte**
 - <http://www.live-counter.com/how-big-is-the-internet>

1 Gigabyte ~ 1000 MB

1 Terabyte ~ 1000 GB

1 Petabyte ~ 1000 TB

1 Exabyte ~ 1000 PB

1 Zettabyte ~ 1000 EB

- **1 Gigabyte:** It takes an author 50 years to write every week a book with about 190 pages, more specifically, with 383,561 characters (with spaces and sentence included). This would be a billion letters or bytes.
- **1 Exabyte:** 212 million DVDs weighing 3,404 tons.

THE NEED FOR DATA ANALYTICS

- Data analysis is a process of inspecting, cleansing, transforming, and modeling data with the goal of*
 - **discovering useful information**
 - **informing conclusions**
 - **supporting decision-making**
- Data analysis has multiple facets and approaches, encompassing diverse techniques under a variety of names, while being used in different business, science, and social science domains.
- In today's business, data analysis is playing a role in
 - making decisions more scientific
 - helping the business achieve effective operation



DATA MINING

- Data mining is the process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.
 - Retail: Market basket analysis, customer relationship management
 - Finance: Credit scoring, fraud detection, trading.
 - Manufacturing: Control, robotics, troubleshooting.
 - Medicine: Medical diagnosis.
 - Telecommunications: Spam filters, intrusion detection.
 - Bio-informatics: Motifs, alignment.
 - Web mining: Search engines.

WHY MACHINE LEARNING

- Machine learning aims at gaining insights from data and making predictions based on it.
- Build a model that is a good and useful approximation to the data.
- Machine learning methods have been investigated for more than 60 years, but became mainstream only recently due to more data being available and advances in computing power (“Moore’s Law”).
- There is no need to “learn” to calculate for example the payroll.
- Learning is used when:
 - Human expertise does not exist (navigating on Mars).
 - Humans are unable to explain their expertise (speech recognition).
 - Solution changes in time (routing on a computer network).
- You’re relying on machine learning every day, maybe without being aware of it!
 - You certainly use a smartphone?

DATA SCIENTIST

David Donoho (2015). 50 years of Data Science

Data Scientist (n.) :

Person who is better at statistics than any software engineer and better at software engineering than any statistician.



DATA SCIENCE

David Donoho (2015). 50 years of Data Science

The activities of Greater Data Science are classified into 6 divisions:

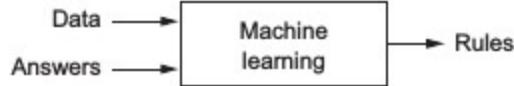
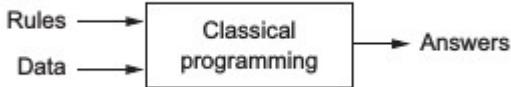
1. Data Exploration and Preparation
2. Data Representation and Transformation
3. Computing with Data
4. Data Modeling
5. Data Visualization and Presentation
- *6. Science about Data Science

SOME TERMINOLOGY

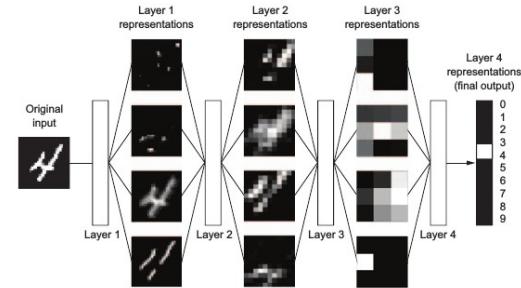
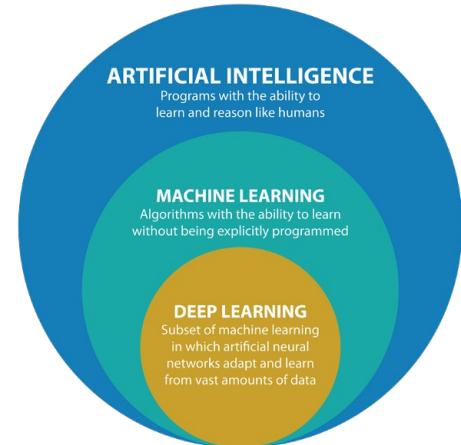
■ Artificial intelligence (AI)

- Can computers be made to “think”—a question whose ramifications we’re still exploring today.
- A concise definition of the field would be as follows: the effort to automate intellectual tasks normally performed by humans.

■ Machine learning (e.g., supervised ML)



■ Deep Learning as a particular example of an ML technique



TYPES OF MACHINE LEARNING

- **Supervised Learning (focus of this series of lectures)**
 - assume that training data is available from which they can learn to predict a target feature based on other features (e.g., monthly rent based on area).
 - Classification
 - Regression
- **Unsupervised Learning**
 - take a given data-set and aim at gaining insights by identifying patterns, e.g., by grouping similar data points.
- **Reinforcement Learning**

SUPERVISED REGRESSION

- Regression aims at predicting a numerical target feature based on one or multiple other (numerical) features.
- Example: Price of a used car.
- x : car attributes
- y : price
- $y = h(x | \theta)$
- $h(\cdot)$: model
- θ : parameters

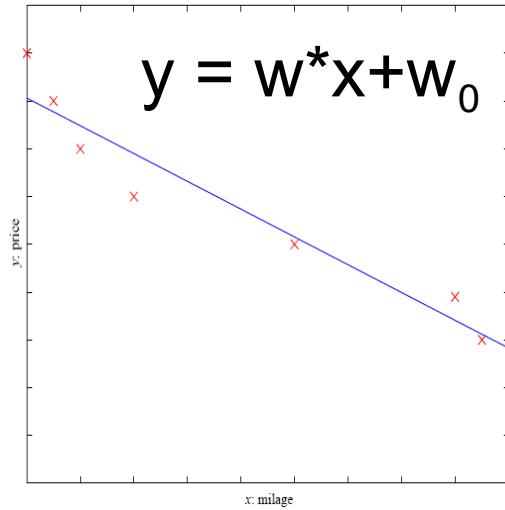


Fig. from Alpaydin (2014)

SUPERVISED CLASSIFICATION

■ Example 1: Spam Classification

- Decide which emails are Spam and which are not.
- Goal: Use emails seen so far to produce a good prediction rule for **future** data.



■ Example 2: Credit Scoring

- Differentiating between low-risk and high-risk customers from their income and savings.

- **Discriminant:** IF $\text{income} > \theta_2$ AND $\text{savings} > \theta_2$
THEN low-risk ELSE high-risk

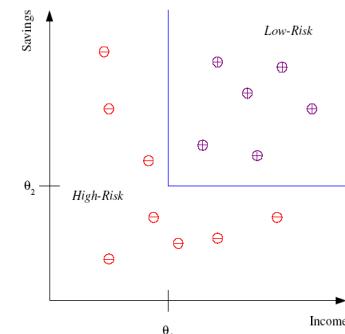


Fig. from Alpaydin (2014)

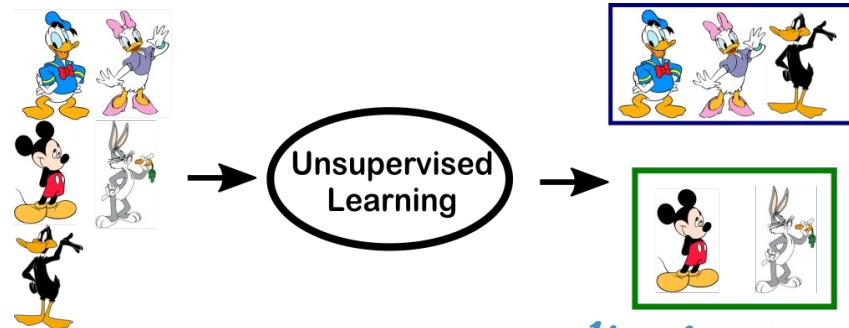
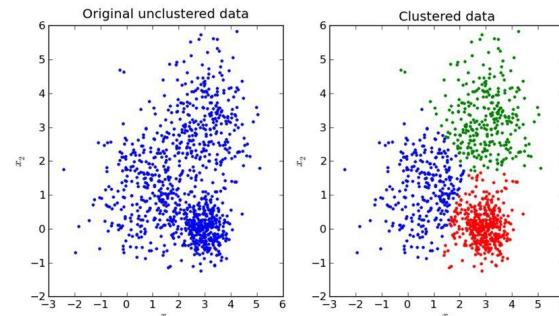
HANDWRITTEN DIGIT CLASSIFICATION

- Early 1990ies
- We have come a long way since then...
 - Handwritten Digit Classification (LeNet) – by Yann Lecun
- See, e.g., <https://www.youtube.com/watch?v=yxuRnBEczUU>

UNSUPERVISED ML

- No output
- Clustering: Grouping similar instances
- Example applications:
 - Customer segmentation
 - Image compression
 - Bio-informatics: Learning motifs
 - ...

Unsupervised Learning



REINFORCEMENT LEARNING

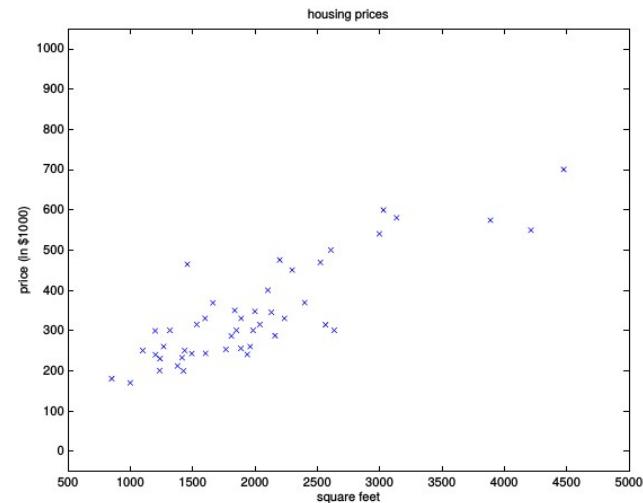
- Learning a policy: A sequence of outputs
- No supervised output but delayed reward
 - Game playing
 - Robot in a maze
 - ...
- See, e.g., <https://www.youtube.com/watch?v=V1eYniJ0Rnk&vl=en>

BUILDING AN ML ALGORITHM

- Optimize a performance criterion using example data or past experience.
- Role of statistics: Inference from a sample.
- Role of computer science: Efficient algorithms to
 - Solve the optimization problem.
 - Representing and evaluating the model for inference.

BUILDING AN ML ALGORITHM (II)

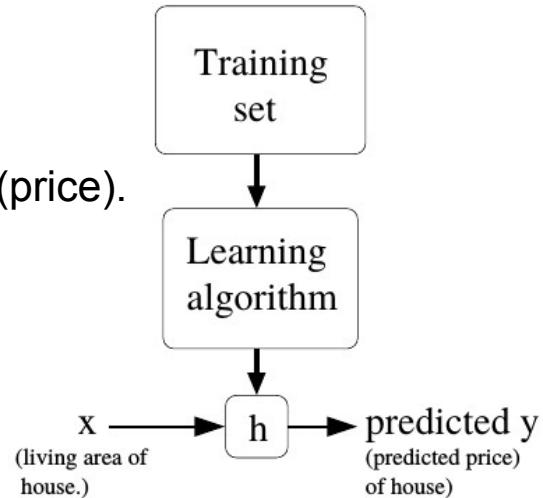
Living area (feet ²)	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
:	:



- Given data like this, how can we learn to predict the prices of other houses as a function of the size of their living areas?

BUILDING AN ML ALGORITHM (III)

- **$x(i)$:** “input” variables (living area in this example), also called **input features**
- **$y(i)$:** “output” / **target variable** that we are trying to predict (price).
- **Training example:** a pair $(x(i) , y(i))$.
- **Training set:** a list of m training examples
 $\{(x(i), y (i)); i = 1, \dots, m\}$
 - To perform supervised learning, we must decide how we’re going to represent **functions/hypotheses h** in a computer.

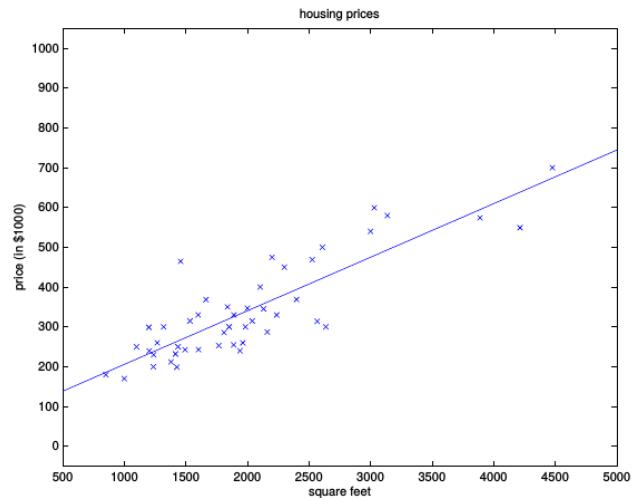


BUILDING AN ML ALGORITHM (IV)

- Model / Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x_1$
 - θ 's: parameters

- Cost Function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$



- Minimize $J(\theta)$ in order to obtain the coefficients θ .

BUILDING AN ML ALGORITHM (V)

- In General: Machine learning in 3 steps:
 - Choose a **model** $h(x|\theta)$.
 - Define a **cost function** $J(\theta|x)$.
 - **Optimization procedure** to find θ^* that minimizes $J(\theta)$.

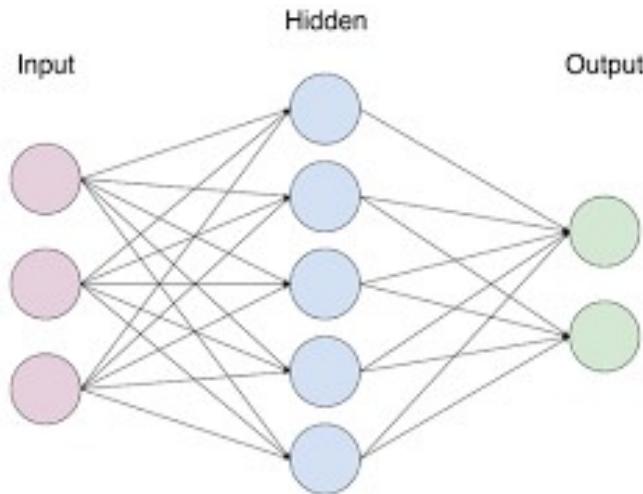
- Computationally, we need:
 - data, linear algebra, statistics tools, and optimization routines.

DON'T RE-INVENT THE WHEEL

■ Plenty of Frameworks out there

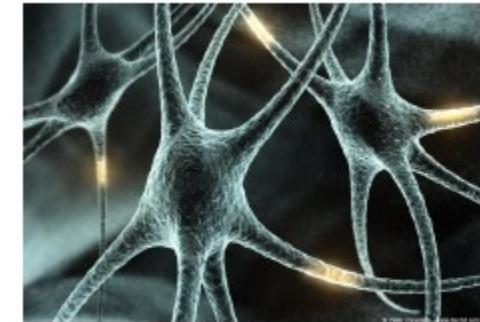
- Tensorflow
- Pytorch
- Caffee
- Scikit-learn
- ...

ARTIFICIAL NEURAL NETWORKS



THE BRAIN AND THE NEURON

- **Biological systems** built of very complex webs of interconnected neurons.
- Highly connected to other neurons, and perform (parallel) computations by combining signals from other neurons.
- Outputs of these computations may be transmitted to one or more other neurons.
- **Artificial Neural Networks** (ANN) built out of a densely interconnected set of simple units (e.g., sigmoid units).
- Each unit takes real-valued inputs (possibly the outputs of other units) and produces a real-valued output (which may become input to many other units).



CONNECTIONIST MODEL

- Consider humans
 - Neuron switching time ~ 0.001 second
 - Number of neurons $\sim 10^{10}$
 - Connections per neuron $\sim 10^{4-5}$
 - Scene recognition time ~ 0.1 second
- 100 inference steps doesn't seem like enough
→ a lot of parallel computation
- Properties of artificial neural nets (ANN's):
 - Many neuron-like threshold switching units
 - Many weighted interconnections among units
 - Highly parallel, distributed process

HEBB'S RULE

- Hebb's rule says that the changes in the strength of synaptic connections are proportional to the correlation in the firing of the two connecting neurons.
- So if **two neurons consistently fire simultaneously**, then any connection between them will change in strength, becoming stronger.
- However, if the two neurons never fire simultaneously, the connection between them will die away.
- The idea is that if two neurons both respond to something, then they should be connected.

HEBB'S RULE – INTUITION

- Suppose that you have a neuron somewhere that recognizes your grandmother (this will probably get input from lots of visual processing neurons, but don't worry about that).
- Now if your grandmother always gives you a chocolate bar when she comes to visit, then some neurons, which are happy because you like the taste of chocolate, will also be stimulated.
- Since these neurons fire at the same time, they will be connected together, and the connection will get stronger over time.
- **So eventually, the sight of your grandmother, even in a photo, will be enough to make you think of chocolate.** Sound familiar?

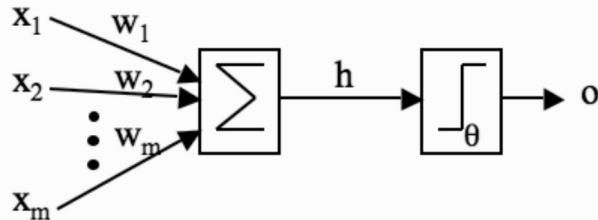


ARTIFICIAL NEURAL NETWORKS

- Artificial Neural networks arise from attempts to model human/animal brains
 - Many models, many claims of biological plausibility.
- We will focus on multi-layer perceptron
 - Mathematical properties rather than plausibility.



MODEL OF A NEURON (1943)



- A picture of McCulloch and Pitts' (1943) mathematical model of a neuron.
 1. a set of weighted inputs w_i that correspond to the synapses.
 2. an adder that sums the input signals (equivalent to the membrane of the cell that collects electrical charge).
 3. an activation function (initially a threshold function) that decides whether the neuron fires ('spikes') for the current inputs.
- The inputs x_i are multiplied by the weights w_i , and the neurons sum their values.
- If this sum is greater than the threshold θ then the neuron fires; otherwise it does not.

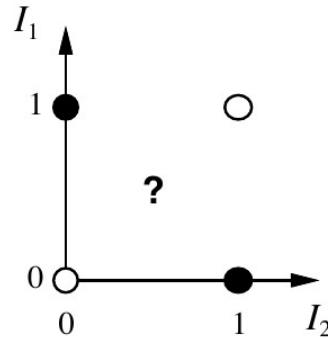
FEED-FORWARD NETWORKS

- In feed-forward networks (a.k.a. multi-layer perceptrons) we let each basis function be another non-linear function of linear combination of the inputs:

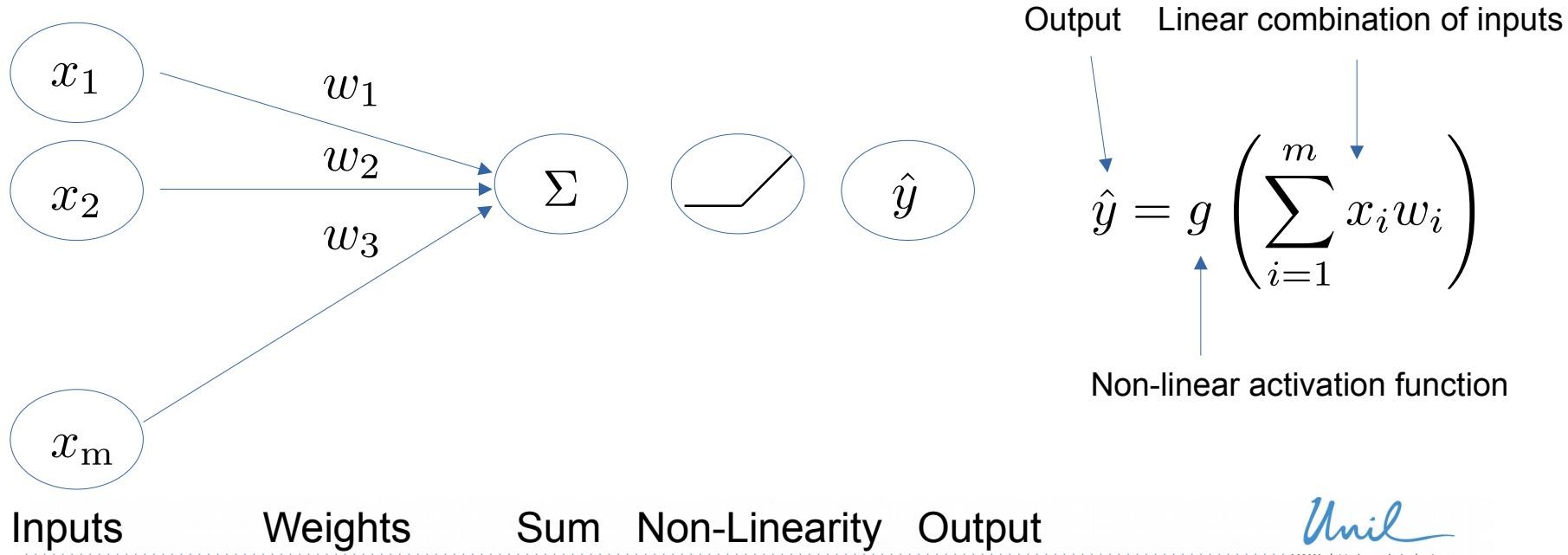
$$\phi_j(\mathbf{x}) = f \left(\sum_{j=1}^M \dots \right)$$

LIMITATIONS OF PERCEPTRONS

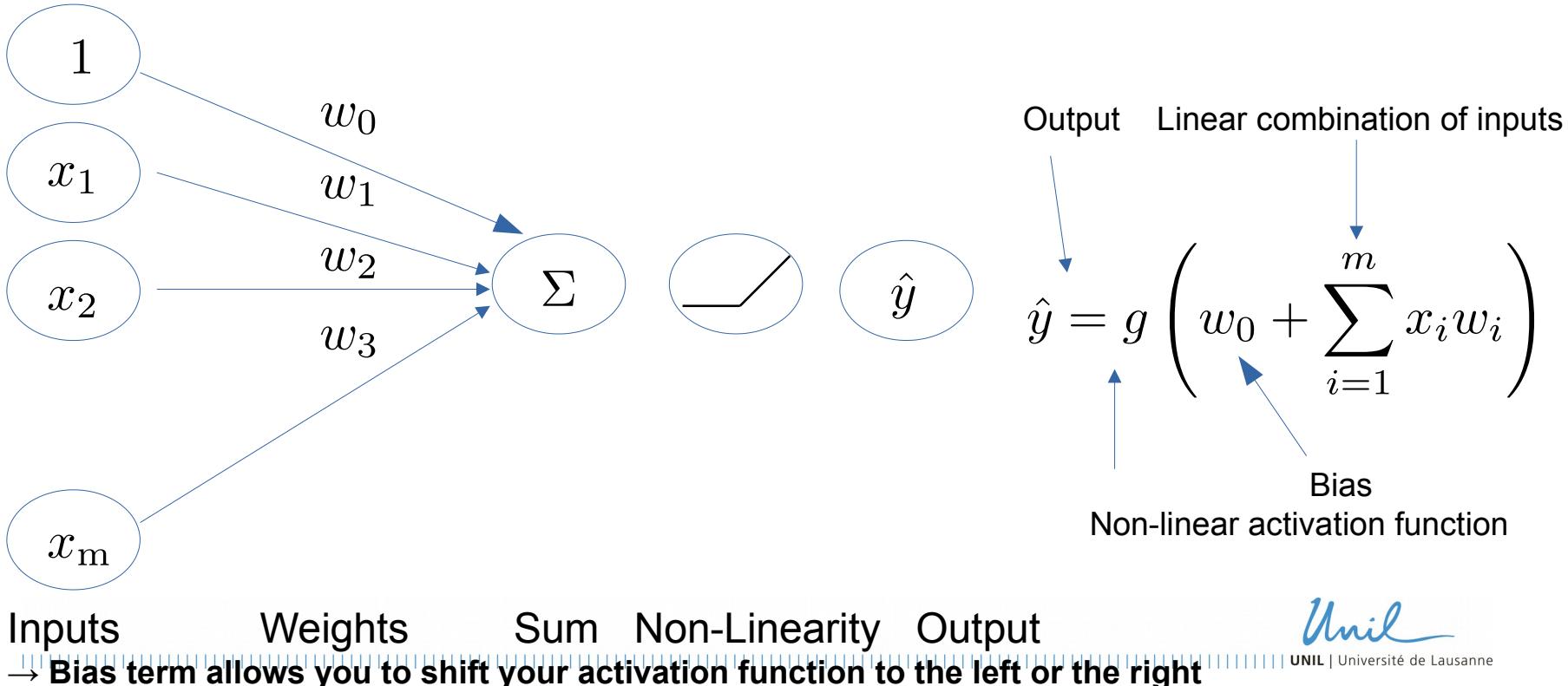
- Perceptrons can only solve linearly separable problems in feature space.
- A canonical example of non-separable problem is X-OR.
- Real data sets can look like this too.



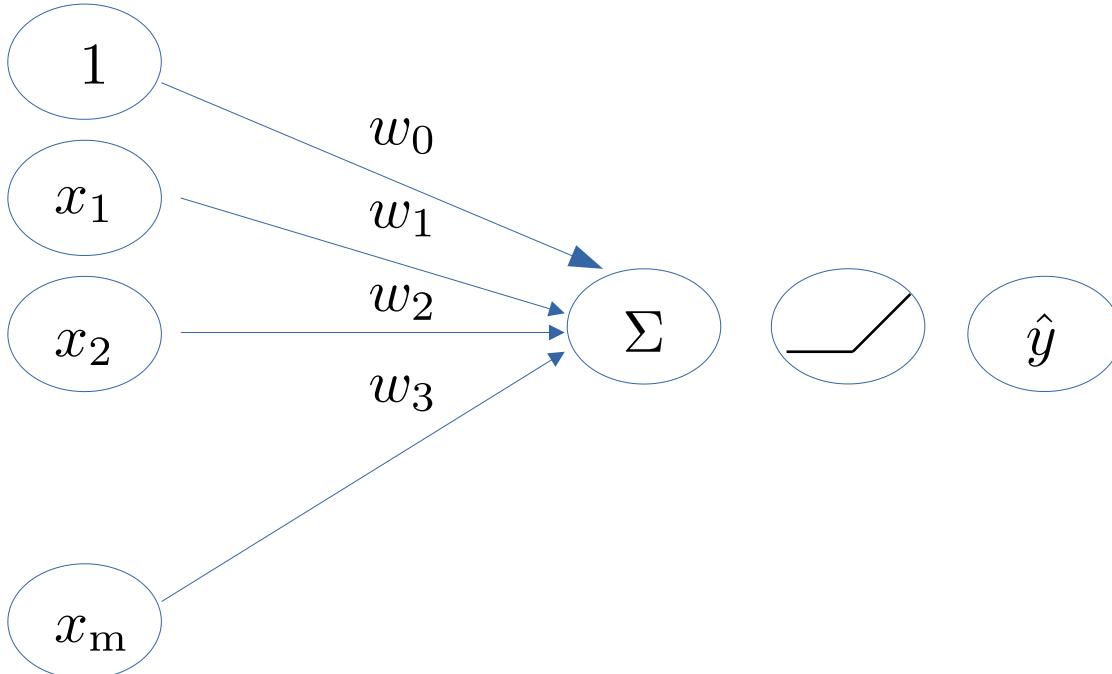
THE PERCEPTRON: FORWARD PROPAGATION



THE PERCEPTRON: FORWARD PROPAGATION



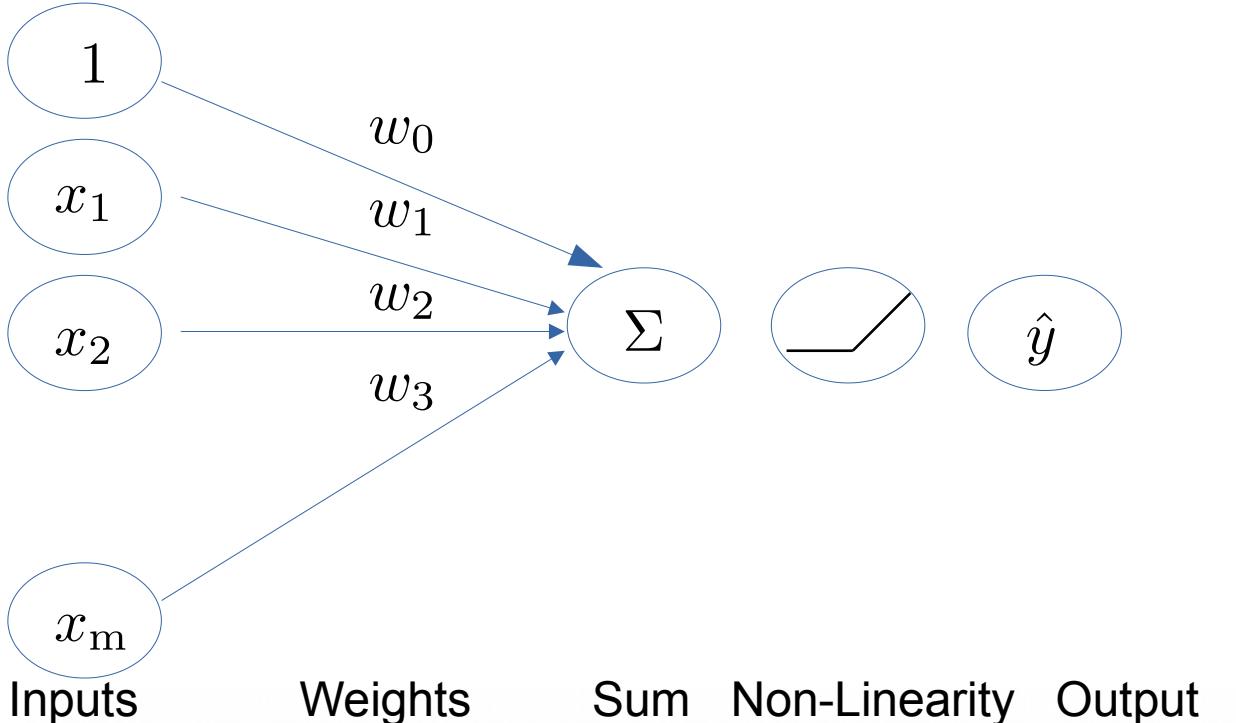
THE PERCEPTRON: FORWARD PROPAGATION



→ Bias term allows you to shift your activation function to the left or the right

$$\begin{aligned}\hat{y} &= g(w_0 + \sum_{i=1}^m x_i w_i) \\ \hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{W}) \\ \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } \mathbf{W} &= \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}\end{aligned}$$

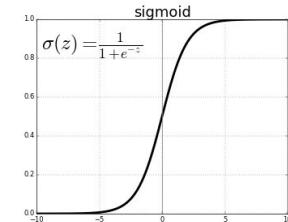
THE PERCEPTRON: FORWARD PROPAGATION



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

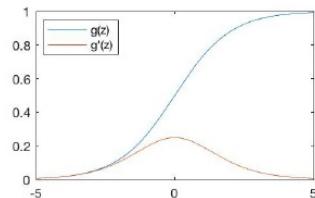
Activation Functions
e.g. sigmoid function

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



FEW ACTIVATION FUNCTIONS

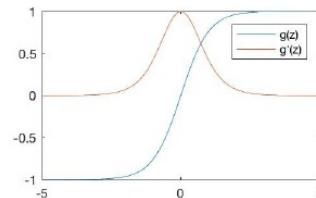
Sigmoid Function



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

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

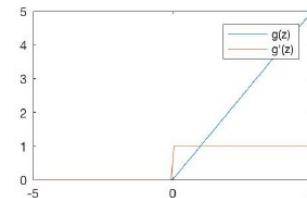
Hyperbolic Tangent



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

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

Rectified Linear Unit (ReLU)



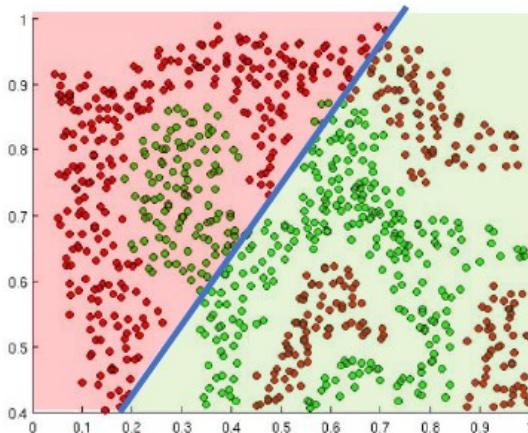
$$g(z) = \max(0, z)$$

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

- Needs to be differentiable for gradient-based learning (later)
 - Very useful in practice.
 - Sigmoid function, e.g., useful for classification (Probability).

IMPORTANCE OF ACTIVATION FCT.

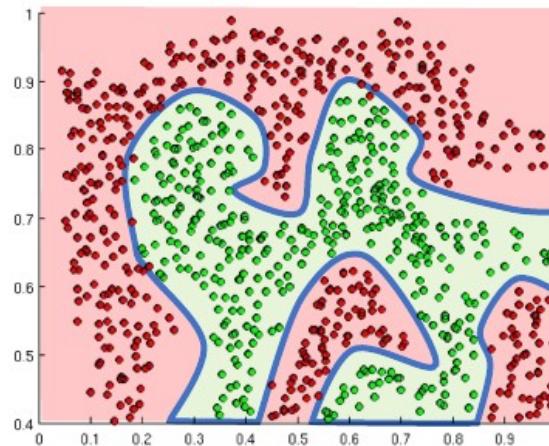
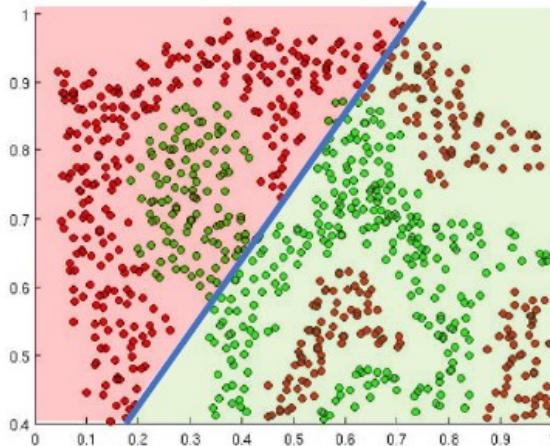
- The purpose of activation functions is to introduce non-linearities into the network.



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

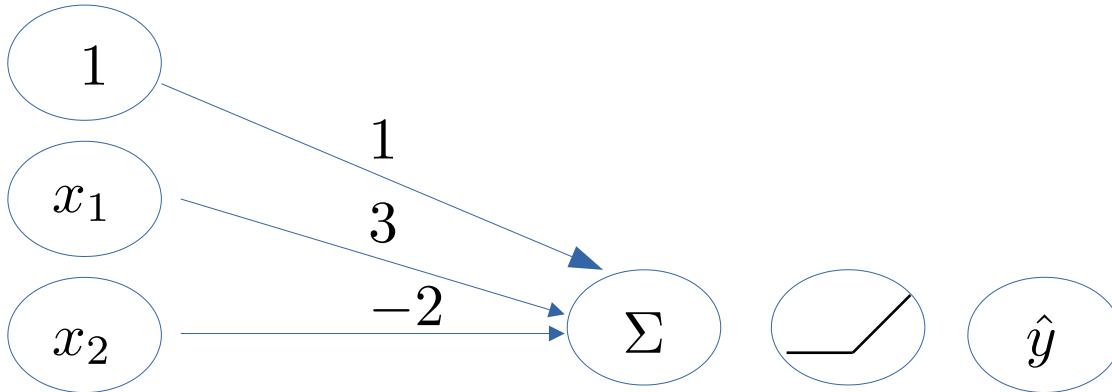
IMPORTANCE OF ACTIVATION FCT.

- 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.

PERCEPTRON – AN EXAMPLE



We have: $w_0 = 1$ and $W = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

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

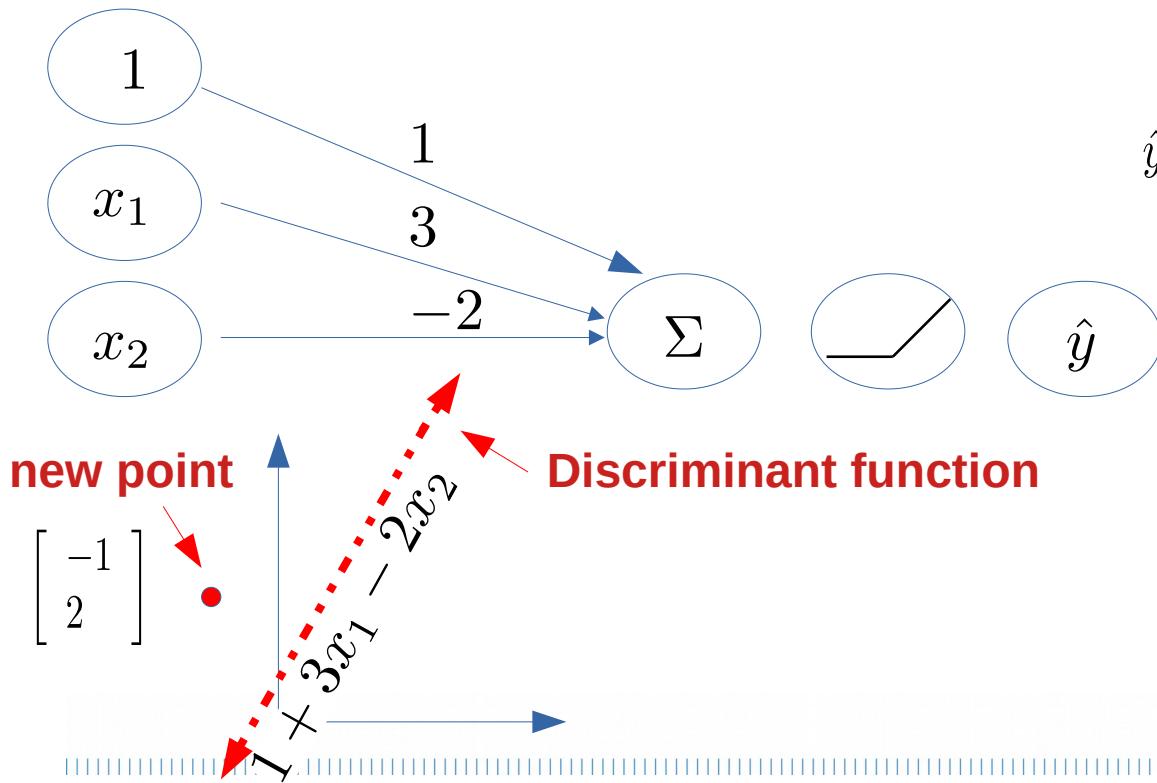
$$= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right)$$

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

This is just a line in 2D

Imagine we have a trained network with weights given.
→ how do we compute the output?

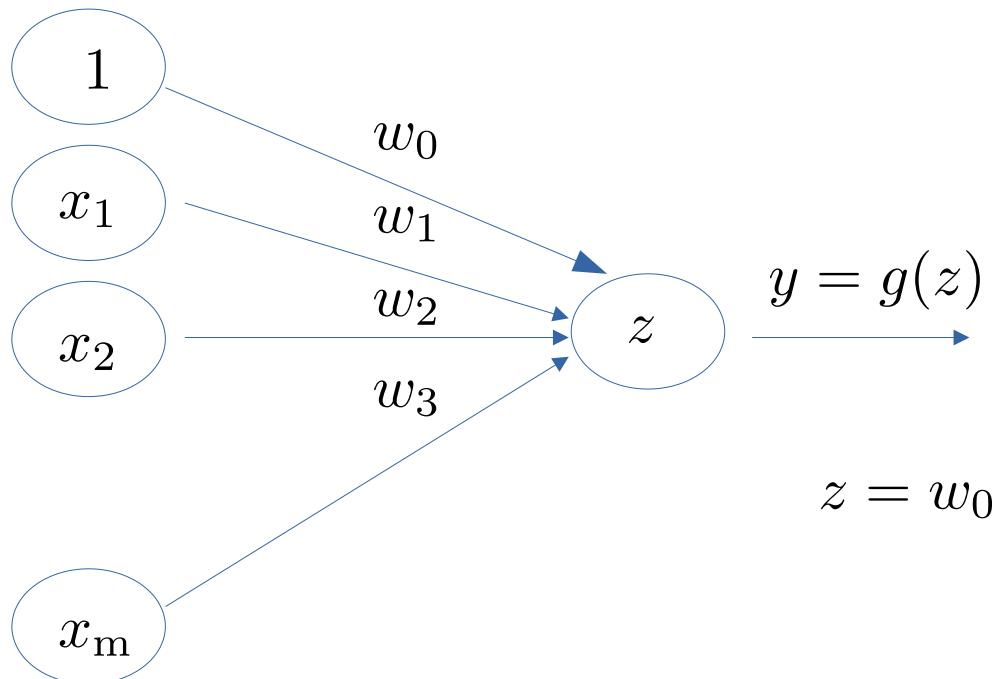
PERCEPTRON – AN EXAMPLE



Assume we have input: $X = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

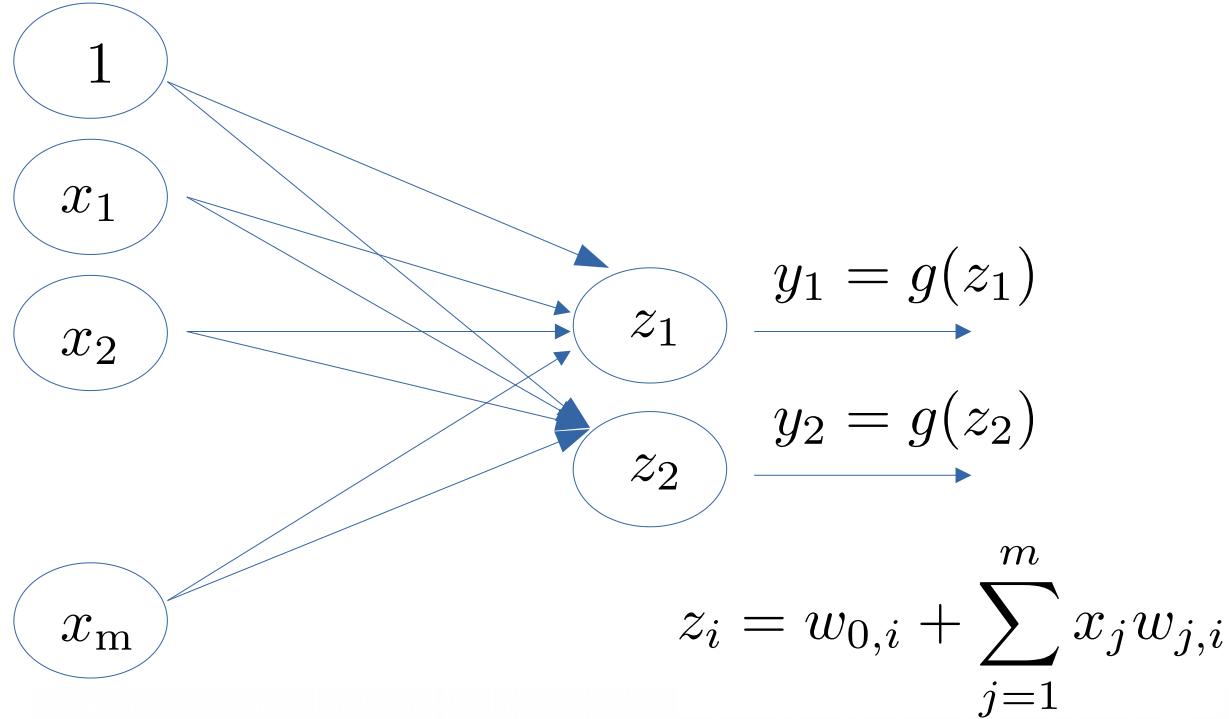
$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$

A PERCEPTRON – SIMPLIFIED

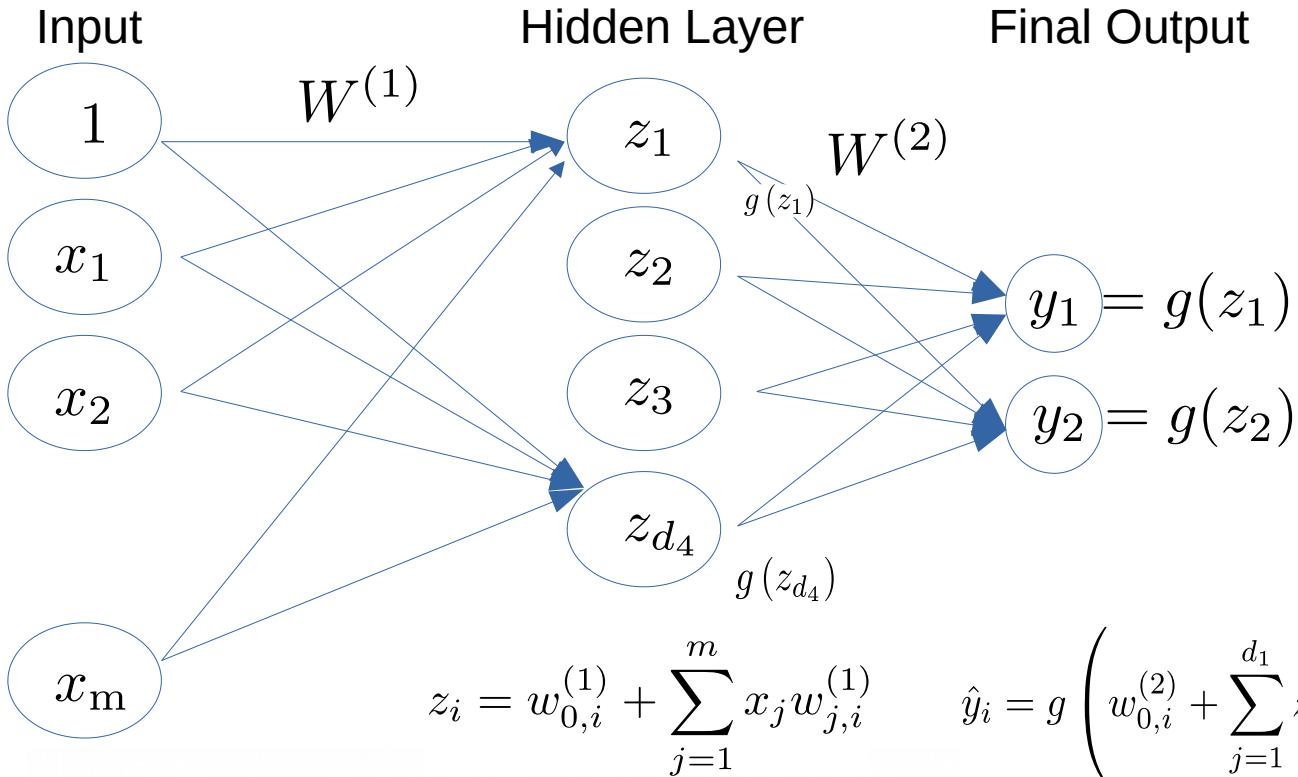


$$z = w_0 + \sum_{j=1}^m x_j w_j$$

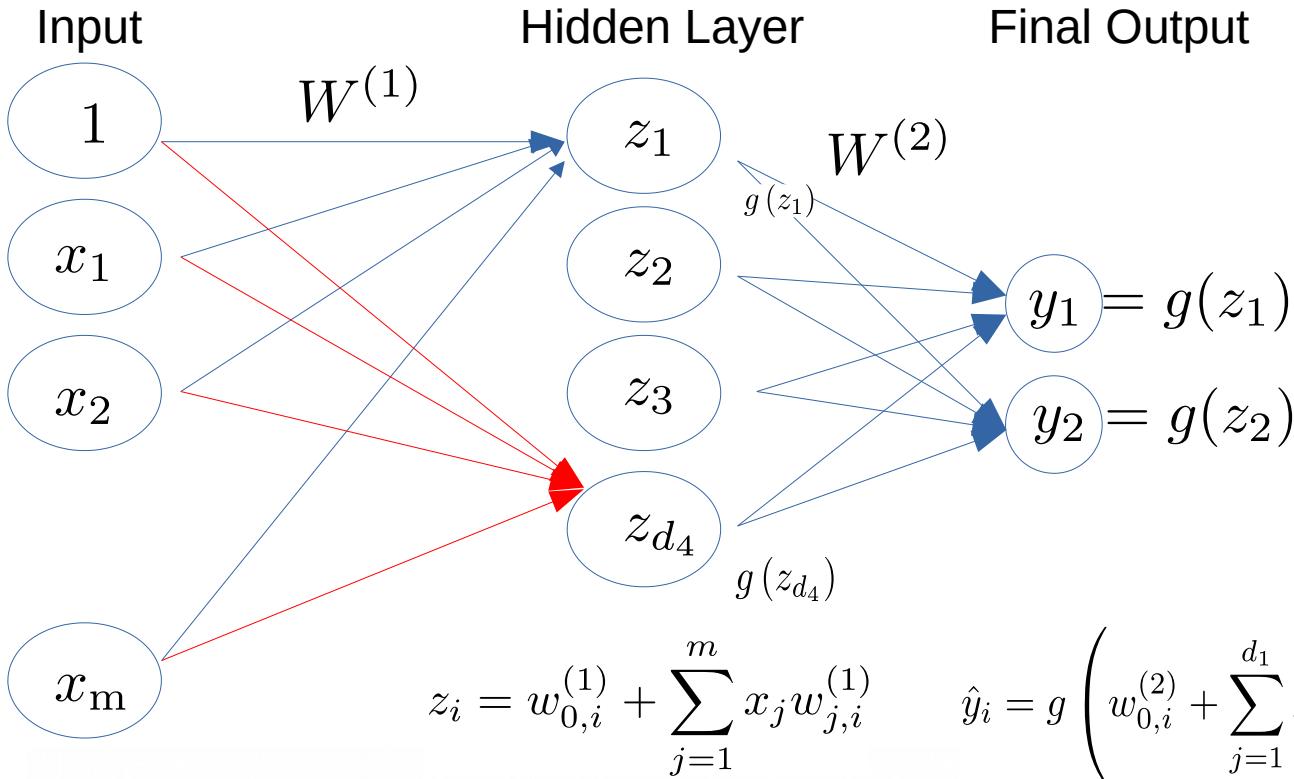
BUILDING A NN WITH PERCEPTRONS: A MULTI-OUTPUT PERCEPTRON



SINGLE LAYER NN

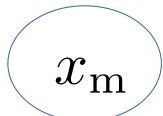
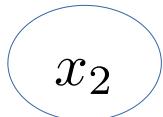
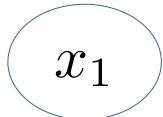
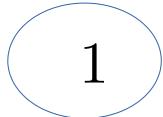


SINGLE LAYER NN

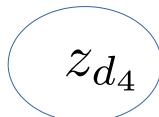
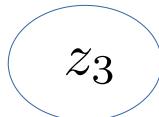
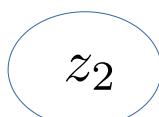


SINGLE LAYER NN

Input



Hidden Layer



X

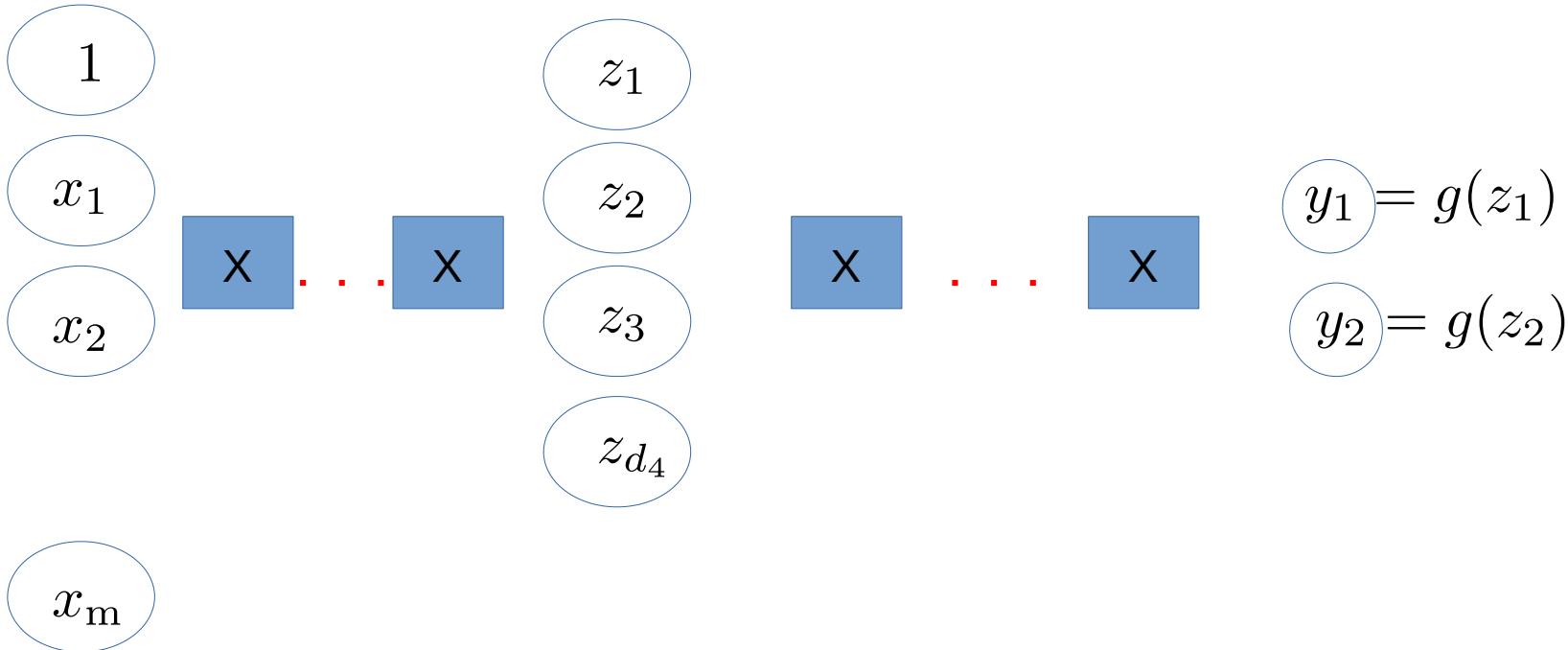
X

Final Output

$$y_1 = g(z_1)$$

$$y_2 = g(z_2)$$

FULLY CONNECTED DNN



EXPRESSIVENESS OF ANN

- **Boolean functions:**
 - Every Boolean function can be represented by a network with a single hidden layer.
 - Might require exponential (in number of inputs) hidden units.
- **Continuous functions:**
 - Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer [Cybenko 1989; Hornik et al. 1989].
- Deep NN are in practice superior to other ML methods in presence of large data sets.

UNIVERSAL FCT. APPROXIMATOR

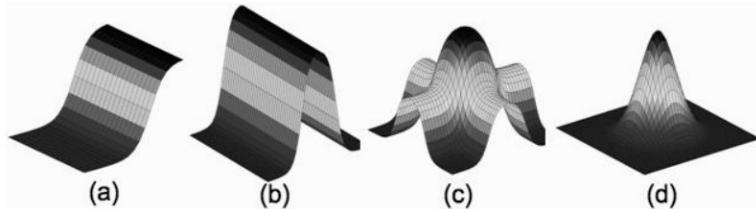


FIGURE 4.9 The learning of the MLP can be shown as the output of a single sigmoidal neuron (a), which can be added to others, including reversed ones, to get a hill shape (b). Adding another hill at 90° produces a bump (c), which can be sharpened to any extent we want (d), with the bumps added together in the output layer. Thus the MLP learns a local representation of individual inputs.

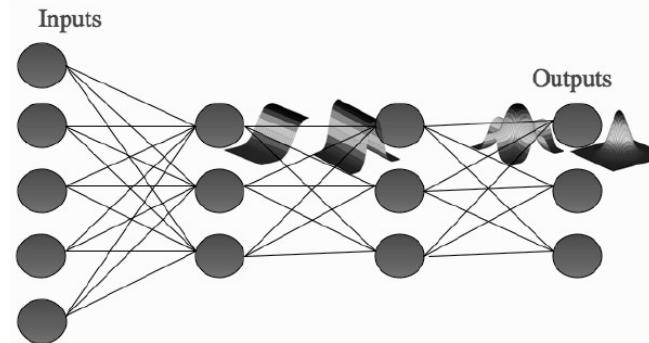
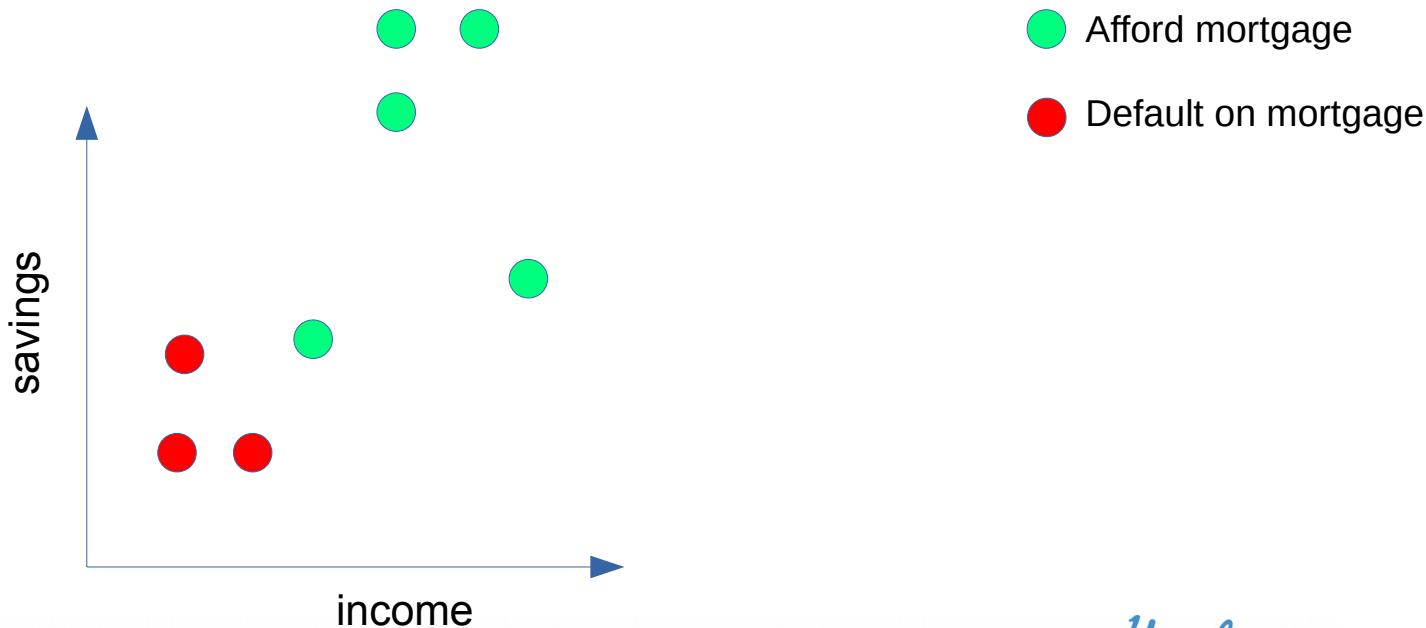


FIGURE 4.10 Schematic of the effective learning shape at each stage of the MLP.

CLASSIFICATION PROBLEM

- Can I afford a loan for a house?

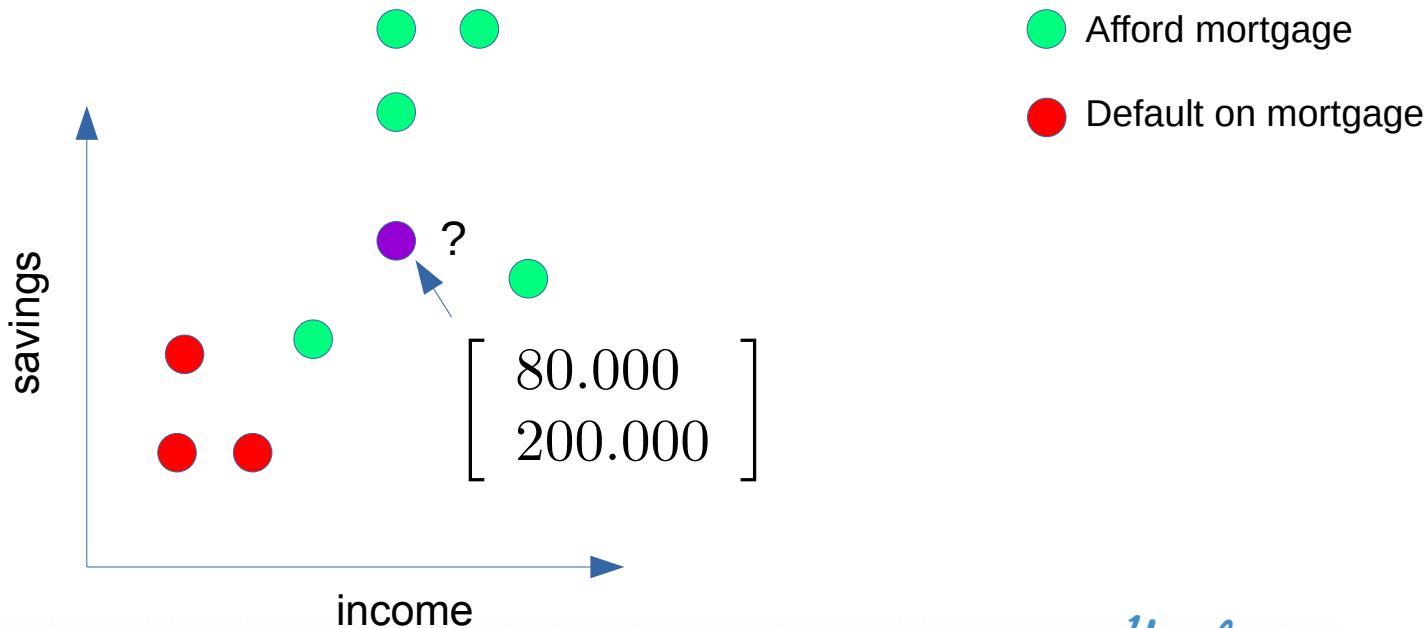
- X: income
- Y: savings



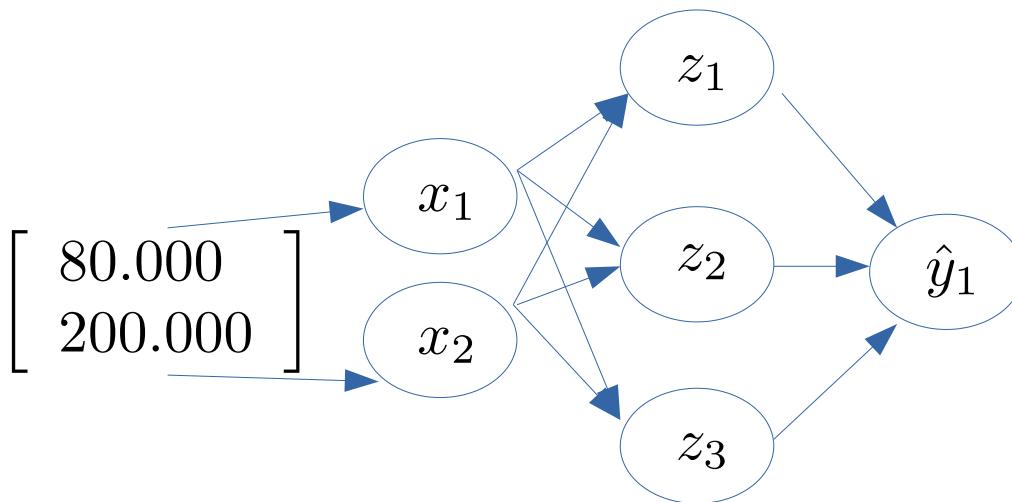
CLASSIFICATION PROBLEM

- Can I afford a loan for a house?

- X: income
- Y: savings

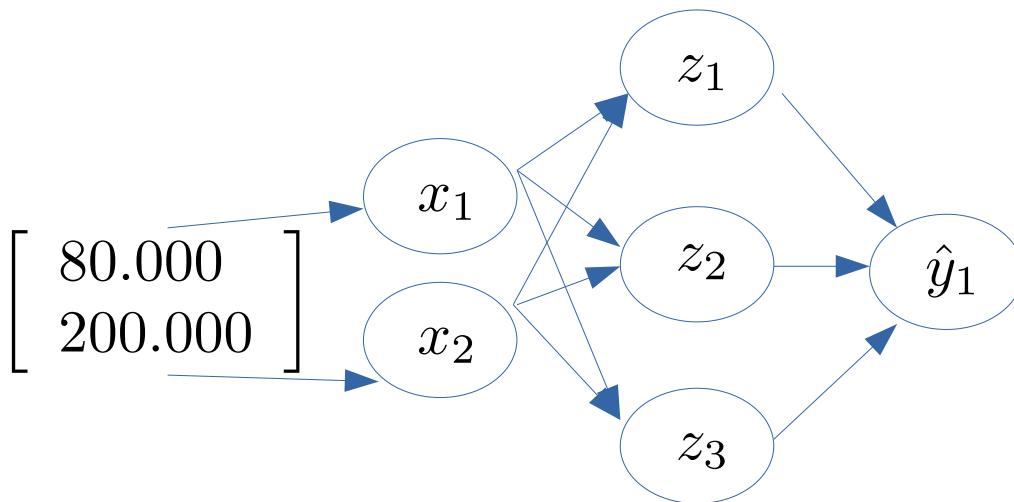


TRY TO DO A PREDICTION



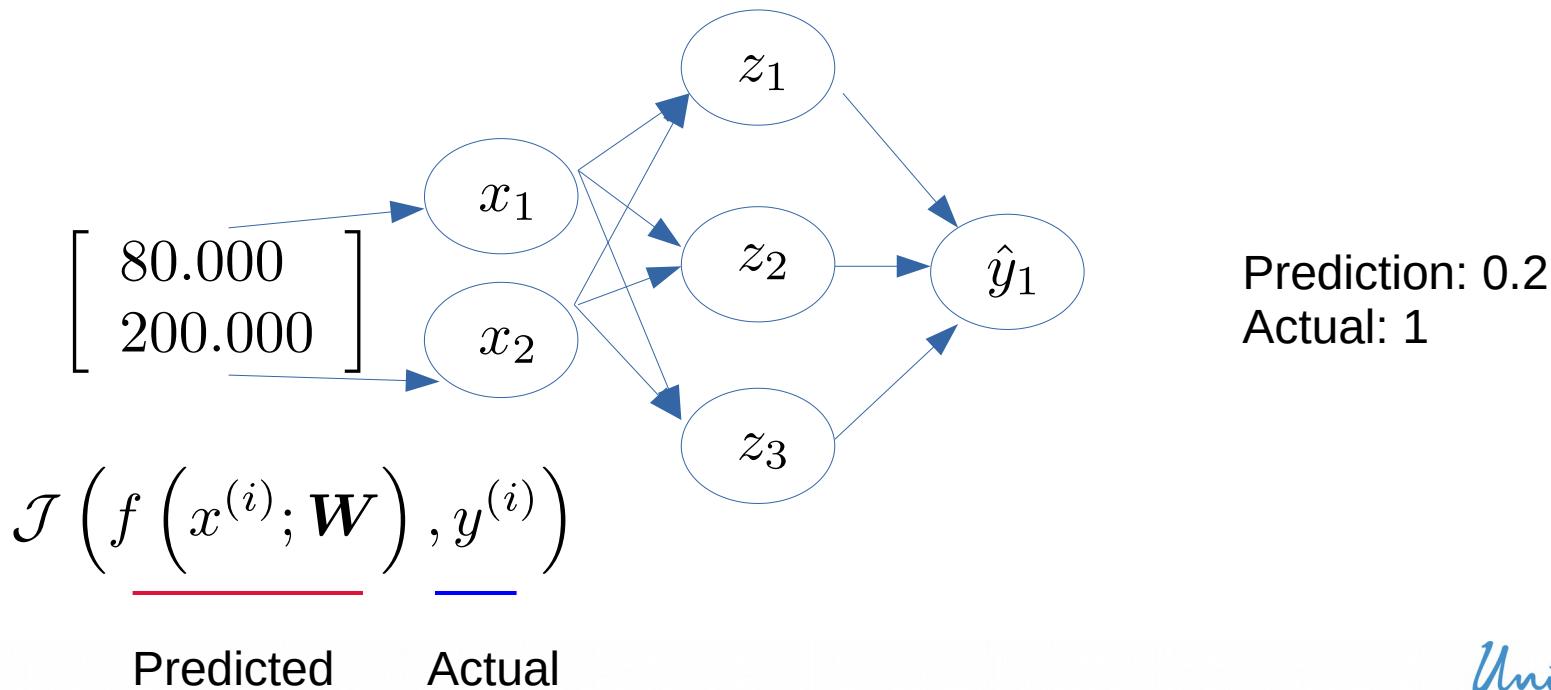
Prediction: 0.2

TRY TO DO A PREDICTION



Prediction: 0.2
Actual: 1

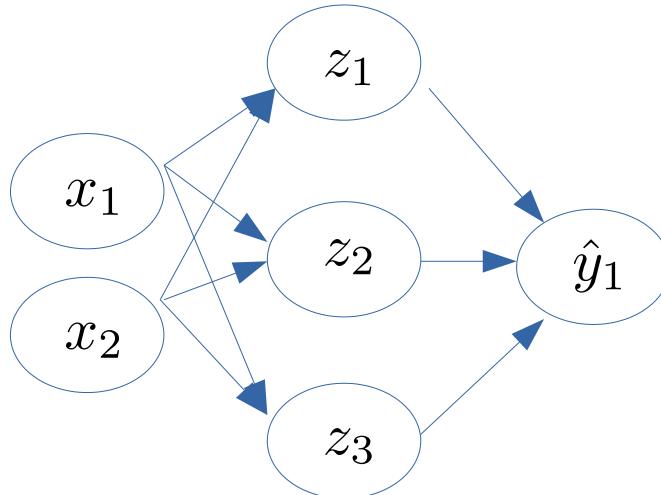
RECALL: QUANTIFYING THE “LOSS”



EMPIRICAL LOSS

The empirical loss measures the total loss over our entire data set.

$$\begin{bmatrix} 80.000, 200.000 \\ 120.000, 400.000 \\ 10.000, 12.000 \\ \dots, \dots \end{bmatrix}$$

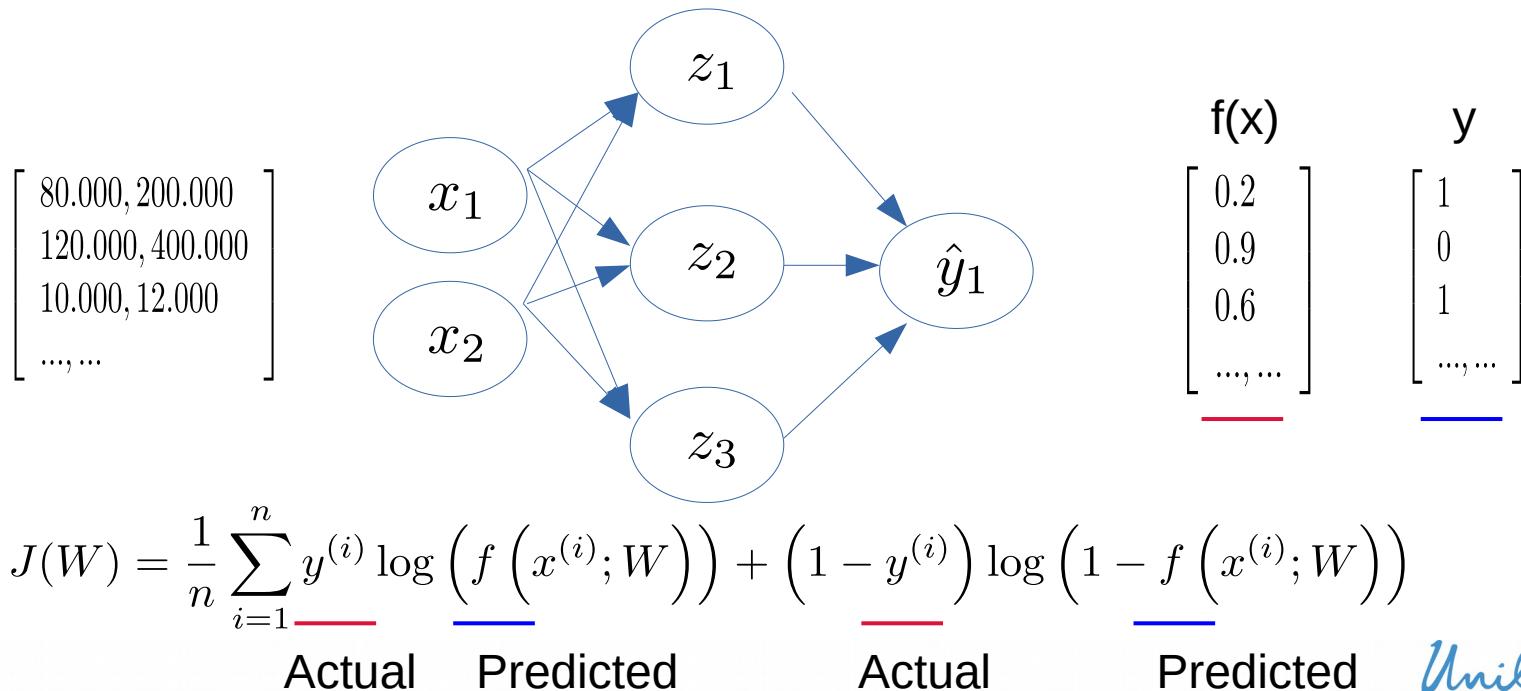


$$J(W) = \frac{1}{n} \sum_{i=1}^n \mathcal{J} \left(f \left(\underline{x^{(i)}}, \underline{W} \right), \underline{y^{(i)}} \right)$$

Predicted Actual

BINARY CROSS ENTROPY LOSS

Our example was a classification problem with output (0 or 1)



CROSS ENTROPY LOSS FUNCTION

- Classification: maximum likelihood principle.
 - Consider a set of m examples $X = \{x^{(1)}, \dots, x^{(m)}\}$ drawn independently from the true but unknown data-generating distribution $p_{\text{data}}(x; \theta)$.
- Let $p_{\text{model}}(x; \theta)$ be a parametric family of probability distributions over the same space indexed by θ . In other words, $p_{\text{model}}(x; \theta)$ maps any configuration x to a real number estimating the true probability $p_{\text{data}}(x; \theta)$.
- Maximum likelihood estimator for the parameters is given by

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(x^{(i)}; \boldsymbol{\theta})\end{aligned}$$

CROSS ENTROPY LOSS FUNCTION

- Numerically more stable: $\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$
- Because the arg max does not change when we rescale the cost function, we can divide by m to obtain a version of the criterion that is expressed as an expectation with respect to the empirical distribution \hat{p}_{data} defined by the training data:

$$\theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta)$$

CROSS ENTROPY LOSS FUNCTION

- One way to interpret maximum likelihood estimation is to view it as minimizing the dissimilarity between the empirical distribution \hat{p}_{data} , defined by the training set and the model distribution, with the degree of dissimilarity between the two measured by the KL divergence.

- The KL divergence is given by

$$D_{\text{KL}} (\hat{p}_{\text{data}} \parallel p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

- The term on the left is a function only of the data-generating process, not the model.
 - This means when we train the model to minimize the KL divergence, we need only minimize

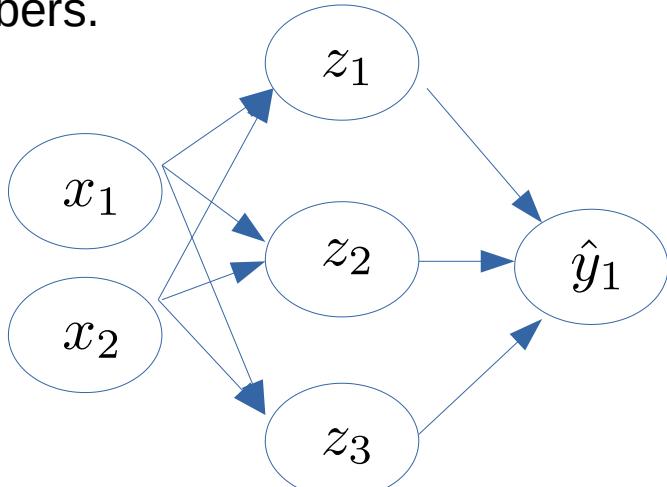
$$-\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$$

- this of course the same as the maximization in equation.
- Any loss consisting of a negative log-likelihood is a cross-entropy between the empirical distribution defined by the training set and the probability distribution defined by model.

MEAN SQUARED ERROR (MSE)

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

$$\begin{bmatrix} 80.000, 200.000 \\ 120.000, 400.000 \\ 10.000, 12.000 \\ \dots, \dots \end{bmatrix}$$



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

Actual Predicted

$$\begin{array}{c|c} f(x) & y \\ \hline 450.000 & 470.000 \\ 250.000 & 220.000 \\ 190.000 & 250.000 \\ \dots, \dots & \dots, \dots \end{array}$$

Loan requested Loan required

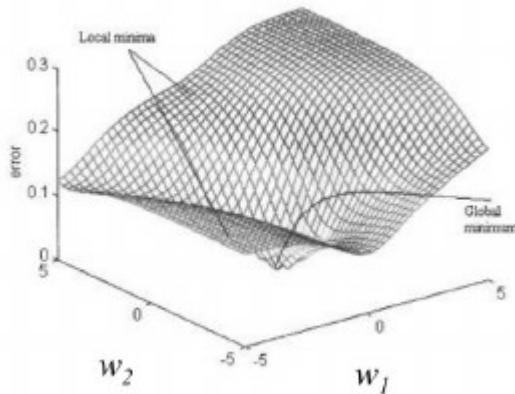
NETWORK TRAINING

We want to find the network weights that achieve the lowest loss!

$$\begin{aligned}\mathbf{W}^* &= \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L} \left(f \left(x^{(i)} ; \mathbf{W} \right) , y^{(i)} \right) \\ \mathbf{W}^* &= \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})\end{aligned}$$

GRADIENT DESCENT IN WEIGHT SPACE

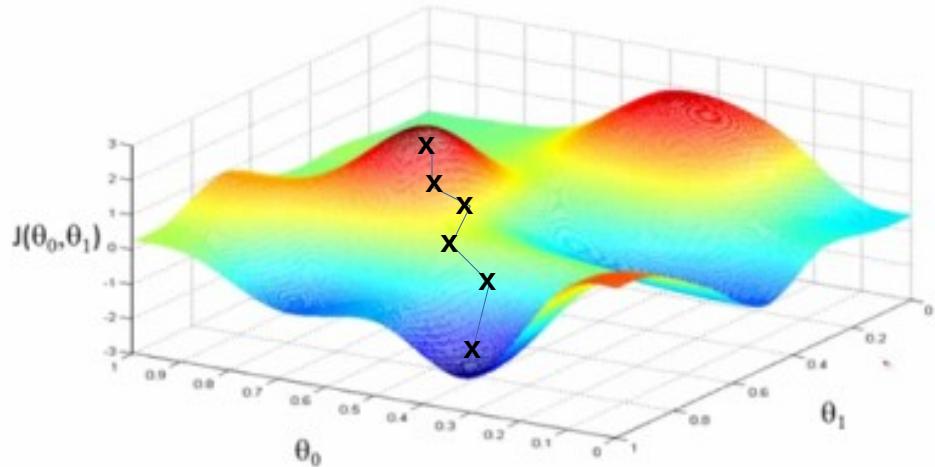
Goal: Given $(x_d, y_d)_{d \in D}$ find w to minimize $J(w) = \frac{1}{2} \sum_{d \in D} (f_w(x_d) - y_d)^2$



This error measure defines a Surface over the hypothesis (i.e. weight) space

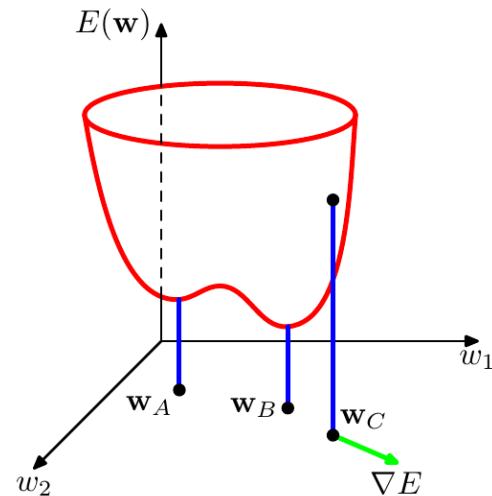
GRADIENT DESCENT IN WEIGHT SPACE

- $W^* = \underset{W}{\operatorname{argmin}} J(W)$
- Randomly pick an initial (w_0, w_1)
- Compute gradient
- Take small steps in the opposite direction of gradient.
- Repeat until convergence



RECALL PARAMETER OPTIMIZATION

- For either of these problems, the error function $J(w)$ is nasty ($E(w)$ in the figure)
- Nasty = non-convex
- Non-convex = has **local minima**



DESCENT METHODS IN GENERAL

- The typical strategy for optimization problems of this sort is a descent method:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \Delta\mathbf{w}^{(\tau)}$$

- These come in many flavors
 - Gradient descent $\nabla J(\mathbf{w}^{(\tau)})$
 - Stochastic gradient descent $\nabla J_n(\mathbf{w}^{(\tau)})$
 - Newton-Raphson (second order) ∇^2
- All of these can be used here, stochastic gradient descent is particularly effective
 - Redundancy in training data, escaping local minima.

GRADIENT DESCENT ALGORITHM

Algorithm

- I. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(W)}{\partial W}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

All that matters
to train a NN



Learning rate

GRADIENT DESCENT ALGORITHM

Algorithm

- I. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$ Can be computationally expensive
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

STOCHASTIC GRADIENT DESCENT

Algorithm

- I. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$  Can be noisy
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

STOCHASTIC GRADIENT DESCENT

Algorithm

- I. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(W)}{\partial W} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(W)}{\partial W}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

MINI-BATCHES WHILE TRAINING

- More accurate estimation of gradient
 - Smoother convergence
 - Allows for larger learning rates
- Mini-batches lead to fast training!
 - Can parallelize computation + achieve significant speed increases on GPU's
- Note: a complete pass over all the patterns in the training set is called an **epoch**.

COMPUTING GRADIENTS: ERROR BACKPROPAGATION

- How does a small change in one weight (e.g., w_2) affect the final loss $J(W)$?



COMPUTING GRADIENTS: ERROR BACKPROPAGATION

- How does a small change in one weight (e.g., w_2) affect the final loss $J(W)$?
- Chain rule



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

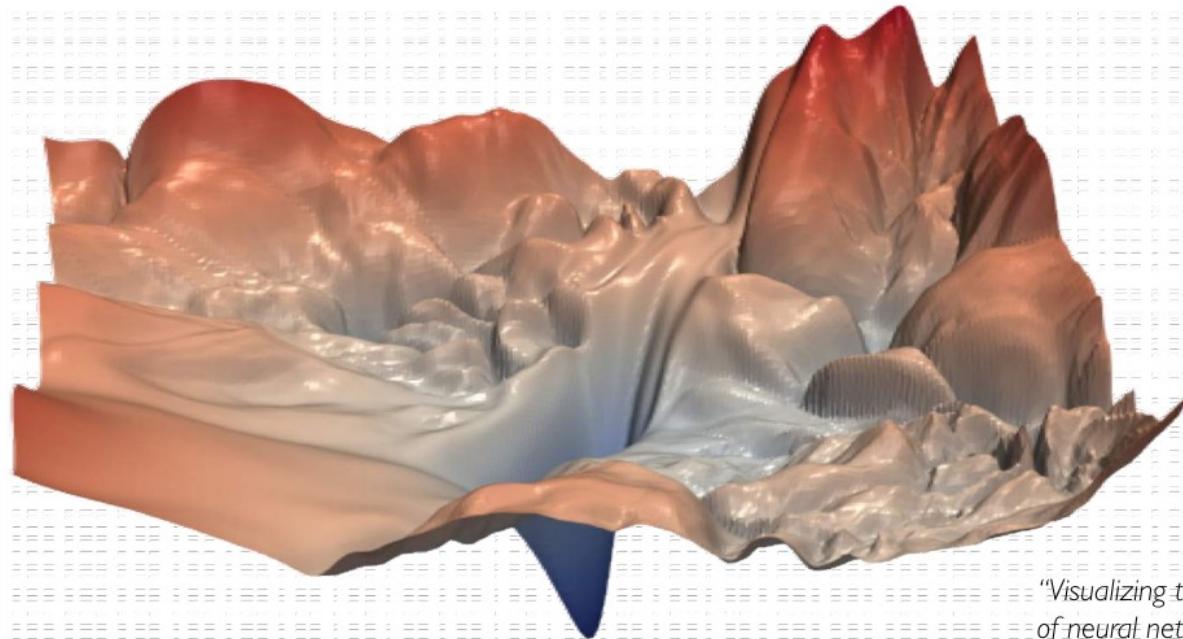
COMPUTING GRADIENTS: ERROR BACKPROPAGATION

- How does a small change in one weight (e.g., w_2) affect the final loss $J(W)$?
- Chain rule
- **Repeat this for every weight in the network using gradients from later layers**



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

TRAINING NEURAL NETWORKS



"Visualizing the loss landscape
of neural nets". Dec 2017.

See <https://papers.nips.cc/paper/7875-visualizing-the-loss-landscape-of-neural-nets.pdf>

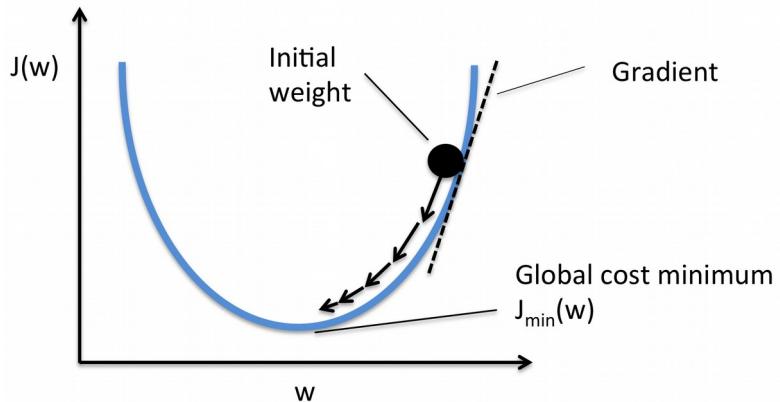
LOSS FUNCTION: CAN BE DIFFICULT TO OPTIMIZE

- Remember:
 - Optimization through gradient descent:
 - How can we set the learning rate?

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

SETTING THE LEARNING RATE

- Small learning rate converges slowly and gets stuck in false local minima
 - Design an adaptive learning rate that “adapts” to the landscape.



FEW VARIANTS OF SGD

Method	Formula
Learning Rate	$w^{(t+1)} = w^{(t)} - \eta \cdot \nabla \ell(w^{(t)}, z) = w^{(t)} - \eta \cdot \nabla w^{(t)}$
Adaptive Learning Rate	$w^{(t+1)} = w^{(t)} - \eta_t \cdot \nabla w^{(t)}$
Momentum [Qian 1999]	$w^{(t+1)} = w^{(t)} + \mu \cdot (w^{(t)} - w^{(t-1)}) - \eta \cdot \nabla w^{(t)}$
Nesterov Momentum [Nesterov 1983]	$w^{(t+1)} = w^{(t)} + v_t; \quad v_{t+1} = \mu \cdot v_t - \eta \cdot \nabla \ell(w^{(t)} - \mu \cdot v_t, z)$
AdaGrad [Duchi et al. 2011]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A_{i,t} + \epsilon}}; \quad A_{i,t} = \sum_{\tau=0}^t (\nabla w_i^{(\tau)})^2$
RMSProp [Hinton 2012]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A'_{i,t} + \epsilon}}; \quad A'_{i,t} = \beta \cdot A'_{t-1} + (1 - \beta) (\nabla w_i^{(t)})^2$
Adam [Kingma and Ba 2015]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot M_{i,t}^{(1)}}{\sqrt{M_{i,t}^{(2)} + \epsilon}}; \quad M_{i,t}^{(m)} = \frac{\beta_m \cdot M_{i,t-1}^{(m)} + (1 - \beta_m) (\nabla w_i^{(t)})^m}{1 - \beta_m^t}$

WEIGHT INITIALIZATION

- Before the training process starts: all weights vectors must be initialized with some numbers.
- There are many initializers of which random initialization is one of the most widely known ones (e.g., with a normal distribution).
 - Specifically, one can configure the mean and the standard deviation, and once again seed the distribution to a specific (pseudo-)random number generator.
 - which distribution to use, then?
 - random initialization itself can become problematic under some conditions: you may then face the vanishing gradients and exploding gradients problems.
- What to do against these problems?
 - e.g. Xavier & He initialization (available in Keras)
 - They are different in the way how they manipulate the drawn weights to arrive at approximately 1. By consequence, they are best used with different activation functions.
 - Specifically, He initialization is developed for ReLU based activating networks and by consequence is best used on those. For others, Xavier (or Glorot) initialization generally works best.

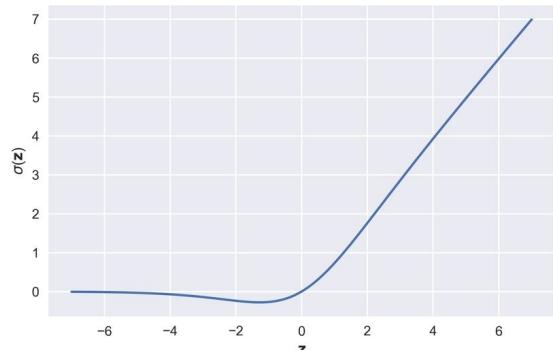
VANISHING GRADIENTS

- Deep learning community often deals with two types of problems during training: vanishing gradients (and exploding) gradients.
 - Vanishing gradients
 - the backpropagation algorithm, which chains the gradients together when computing the error backwards, will find really small gradients towards the left side of the network (i.e., farthest from where error computation started).
 - This problem primarily occurs e.g. with the Sigmoid and Tanh activation functions, whose derivatives produce outputs of $0 < x' < 1$, except for Tanh which produces $x' = 1$ at $x = 0$.
 - Consequently, when using Tanh and Sigmoid, you risk having a suboptimal model that might possibly not converge due to vanishing gradients.
 - ReLU does not have this problem – its derivative is 0 when $x < 0$ and is 1 otherwise.
 - It is computationally faster. Computing this function – often by simply maximizing between $(0, x)$ – takes substantially fewer resources than computing e.g. the sigmoid and tanh functions. By consequence, ReLU is the de facto standard activation function in the deep learning community today.

SWISH ACTIVATION FUNCTION

- Nevertheless, it does not mean that it cannot be improved.
 - Swish activation function.
 - Instead, it does look like the de facto standard activation function, with one difference: the domain around 0 differs from ReLU.
- Swish is a smooth function. That means that it does not abruptly change direction like ReLU does near $x = 0$.
 - Swish is non-monotonic. It thus does not remain stable or move in one direction, such as ReLU.
 - It is in fact this property which separates Swish from most other activation functions, which do share this monotonicity.
- In applications - Swish could be better than ReLu.

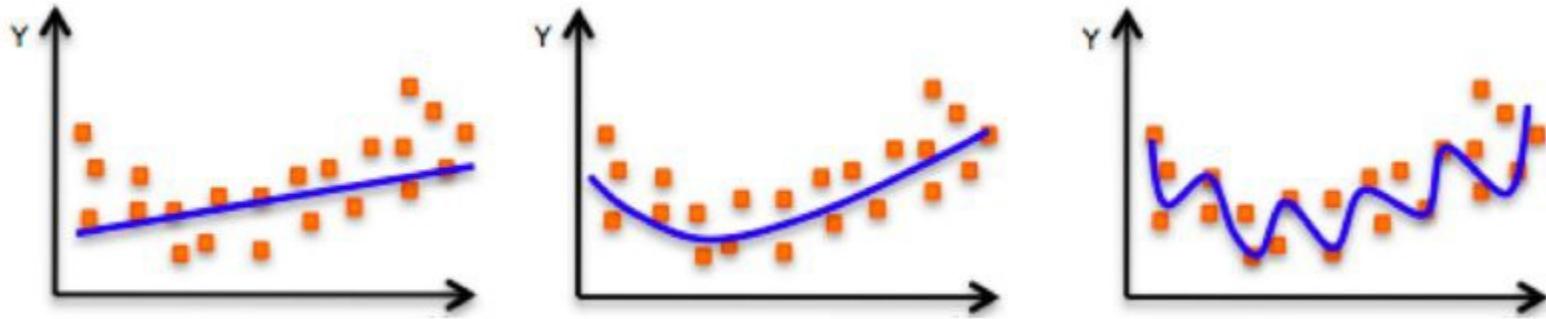
$$\begin{aligned} f(x) &= x * \text{sigmoid}(x) \\ &= x * (1 + e^{-x})^{-1} \end{aligned}$$



INTERMEZZO – ACTION REQUIRED

- Let's look at this notebook:
→ [GradientDescent_and_StochasticGradientDescent.ipynb](#)

NOTES ON OVERFITTING



Underfitting

Model does not have
capacity to fully learn the data

Ideal Fit

Overfitting

Too complex, extra parameters,
does not generalize well

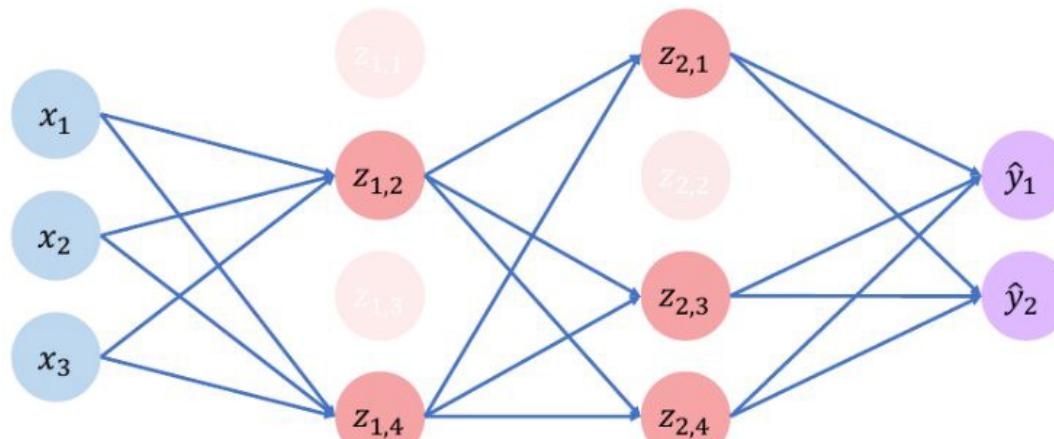
EARLY STOPPING

- Stop training before we have a chance to overfit



REGULARIZATION IN NN: DROPOUT

- During training, randomly set some activations to 0
 - Typically 'drop' 50 % of activations in layer
 - Forces network to not rely on any node



REGULARIZATION IN NN: DROPOUT

- It is an efficient way of performing model averaging with neural networks.
- Can be interpreted as some sort of bagging.
- Now, we assume that the model's role is to output a probability distribution. In the case of bagging, each model i produces a probability distribution $p^{(i)}(y | x)$.
- The prediction of the ensemble is given by the arithmetic mean of all these distributions:

$$\frac{1}{k} \sum_{i=1}^k p^{(i)}(y | x)$$

- In the case of dropout, each sub-model defined by mask vector μ defines a probability distribution $p(y | x, \mu)$.
- The arithmetic mean over all masks is given by $\sum p(\mu) p(y | x, \mu)$ where $p(\mu)$ is the probability distribution that was used to sample μ at training time.

REMARK: BATCH NORMALIZATION

- <https://arxiv.org/abs/1502.03167>
- Batch normalization is used to stabilize and perhaps accelerate the learning process.
- It does so by applying a transformation that maintains the mean activation close to 0 and the activation standard deviation close to 1.
 - Suppose we built a neural network with the goal of classifying gray-scale images. The intensity of every pixel in a gray-scale image varies from 0 to 255. Prior to entering the neural network, every image will be transformed into a 1 dimensional array. Then, every pixel enters one neuron from the input layer. If the output of each neuron is passed to a sigmoid function, then every value other than 0 (i.e. 1 to 255) will be reduced to a number close to 1. Therefore, it's common to normalize the pixel values of each image before training. Batch normalization, on the other hand, is used to apply normalization to the output of the hidden layers.

A GEOMETRIC INTERPRETATION

- In 3D, the following mental image may prove useful. Imagine two sheets of colored paper: **one red and one blue**.
- Put one on top of the other.
- Crumple them together into a small ball. That crumpled paper ball is your input data, and each sheet of paper is a class of data in a classification problem.
- What a neural network (or any other machine-learning model) is meant to do is figure out a **transformation of the paper ball** that would uncrumple it, so as to make the two classes cleanly separable again.
- With deep learning, this would be implemented as a series of simple transformations of the 3D space, such as those you could apply on the paper ball with your fingers, one movement at a time.



Figure 2.9 Uncrumpling a complicated manifold of data

BUILDING AN MLP FROM SCRATCH

Computing the **computational complexity** of this algorithm is very easy.

- The recall phase loops over the neurons, and within that loops over the inputs, so its complexity is **$O(mn)$** .
- The training part does this same thing, but does it for T iterations, so costs **$O(T mn)$** .

The Perceptron Algorithm

- Initialisation
 - set all of the weights w_{ij} to small (positive and negative) random numbers
- Training
 - for T iterations or until all the outputs are correct:
 - * for each input vector:
 - compute the activation of each neuron j using activation function g :
 - update each of the weights individually using:

$$y_j = g \left(\sum_{i=0}^m w_{ij} x_i \right) = \begin{cases} 1 & \text{if } \sum_{i=0}^m w_{ij} x_i > 0 \\ 0 & \text{if } \sum_{i=0}^m w_{ij} x_i \leq 0 \end{cases} \quad (3.4)$$

· update each of the weights individually using:

$$w_{ij} \leftarrow w_{ij} - \eta(y_j - t_j) \cdot x_i \quad (3.5)$$

Recall

- compute the activation of each neuron j using:

$$y_j = g \left(\sum_{i=0}^m w_{ij} x_i \right) = \begin{cases} 1 & \text{if } w_{ij} x_i > 0 \\ 0 & \text{if } w_{ij} x_i \leq 0 \end{cases} \quad (3.6)$$

PERCEPTRON IMPLEMENTATION

```
for data in range(nData): # loop over the input vectors
    for n in range(N): # loop over the neurons
        # Compute sum of weights times inputs for each neuron
        # Set the activation to 0 to start
        activation[data][n] = 0
        # Loop over the input nodes (+1 for the bias node)
        for m in range(M+1):
            activation[data][n] += weight[m][n] * inputs[data][m]

        # Now decide whether the neuron fires or not
        if activation[data][n] > 0:
            activation[data][n] = 1
        else
            activation[data][n] = 0
```

```
# Compute activations
activations = np.dot(inputs,self.weights)

# Threshold the activations
return np.where(activations>0,1,0)
```

PERCEPTRON IMPLEMENTATION (II)

- The weight update for the entire network can be done in one line (where eta is the learning rate, η):

```
self.weights -= eta*np.dot(np.transpose(inputs),self.activations-targets)
```

The Multi-layer Perceptron Algorithm

- Initialisation

- initialise all weights to small (positive and negative) random values

- Training

- repeat:

- * for each input vector:

Forwards phase:

- compute the activation of each neuron j in the hidden layer(s) using:

$$h_\zeta = \sum_{i=0}^L x_i v_{i\zeta} \quad (4.4)$$

$$a_\zeta = g(h_\zeta) = \frac{1}{1 + \exp(-\beta h_\zeta)} \quad (4.5)$$

- work through the network until you get to the output layer neurons, which have activations (although see also Section 4.2.3):

$$h_\kappa = \sum_j a_j w_{j\kappa} \quad (4.6)$$

$$y_\kappa = g(h_\kappa) = \frac{1}{1 + \exp(-\beta h_\kappa)} \quad (4.7)$$

Backwards phase:

- compute the error at the output using:

$$\delta_o(\kappa) = (y_\kappa - t_\kappa) y_\kappa (1 - y_\kappa) \quad (4.8)$$

- compute the error in the hidden layer(s) using:

$$\delta_h(\zeta) = a_\zeta (1 - a_\zeta) \sum_{k=1}^N w_{\zeta k} \delta_o(k) \quad (4.9)$$

- update the output layer weights using:

$$w_{\zeta\kappa} \leftarrow w_{\zeta\kappa} - \eta \delta_o(\kappa) a_\zeta^{\text{hidden}} \quad (4.10)$$

- update the hidden layer weights using:

$$v_\iota \leftarrow v_\iota - \eta \delta_h(\kappa) x_\iota \quad (4.11)$$

- * (if using sequential updating) randomise the order of the input vectors so that you don't train in exactly the same order each iteration

- until learning stops (see Section 4.3.3)

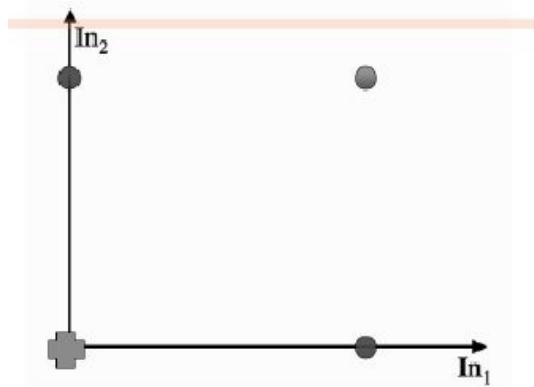
- Recall

- use the Forwards phase in the training section above



EXAMPLE: PERCEPTRON LEARNING: THE “OR” DATA-SET

In ₁	In ₂	t
0	0	0
0	1	1
1	0	1
1	1	1



- Data for the OR logic function and a plot of the four data points
- Action required – see *Multi-layer_Perceptron.ipynb*

RECIPE FOR USING MLP

- Select inputs and outputs for your problem
 - Before anything else, you need to think about the problem you are trying to solve, and make sure that you have data for the problem, both input vectors and target outputs.
 - At this stage you need to choose what features are suitable for the problem and decide on the output encoding that you will use — standard neurons, or linear nodes.
 - These things are often decided for you by the input features and targets that you have available to solve the problem.
 - Later on in the learning it can also be useful to re-evaluate the choice by training networks with some input feature missing to see if it improves the results at all.

RECIPE FOR USING MLP (II)

■ Normalize inputs

- Rescale the data by subtracting the mean value from each element of the input vector, and divide by the variance (or alternatively, either the maximum or minus the minimum, whichever is greater).

■ Split the data into training, testing, and validation sets

- You cannot test the learning ability of the network on the same data that you trained it on, since it will generally fit that data very well (often too well, overfitting and modeling the noise in the data as well as the generating function).
- Recall: we generally split the data into three sets, one for training, one for testing, and then a third set for validation, which is testing how well the network is learning during training.

RECIPE FOR USING MLP (III)

■ Select a network architecture

- You already know how many input nodes there will be, and how many output neurons.
- You need to consider whether you will need a hidden layer at all, and if so how many neurons it should have in it.
- You might want to consider more than one hidden layer.
- The more complex the network, the more data it will need to be trained on, and the longer it will take.
- It might also be more subject to over-fitting.
- The usual method of selecting a network architecture is to try several with different numbers of hidden nodes and see which works best.

RECIPE FOR USING MLP (IV)

- Train a network
 - The training of the NN consists of applying the MLP algorithm to the training data.
 - This is usually run in conjunction with early stopping, where after a few iterations of the algorithm through all of the training data, the generalization ability of the network is tested by using the validation set.
 - The NN is very likely to have far too many degrees of freedom for the problem, and so after some amount of learning it will stop modeling the generating function of the data, and start to fit the noise and inaccuracies inherent in the training data. At this stage the error on the validation set will start to increase, and learning should be stopped.
- Test the network
 - Once you have a trained network that you are happy with, it is time to use the test data for the first (and only) time. This will enable you to see how well the network performs on some data that it has not seen before, and will tell you whether this network is likely to be usable for other data, for which you do not have targets.

ACTION REQUIRED – MLP

- The implementation is a batch version of the algorithm, so that weight updates are made after all of the input vectors have been presented.
- The central weight update computations for the algorithm can be implemented as:

```
deltao = (targets-self.outputs)*self.outputs*(1.0-self.outputs)
deltah = self.hidden*(1.0-self.hidden)*(np.dot(deltao,np.transpose(self.))
weights2))

updatew1 = np.zeros((np.shape(self.weights1)))
updatew2 = np.zeros((np.shape(self.weights2)))

updatew1 = eta*(np.dot(np.transpose(inputs),deltah[:, :-1]))
updatew2 = eta*(np.dot(np.transpose(self.hidden),deltao))
self.weights1 += updatew1
self.weights2 += updatew2
```

EVALUATING A CLASSIFIER

- How can we assess the prediction quality of a classifier?
- Initially, we'll consider the case of binary classification (and extend it later to multi-class classification).
- Confusion matrix shows the performance of a classifier.

		Predicted	
		0 (No)	1 (Yes)
Actual	0 (No)	True Negatives (TN)	False Positives (FP)
	1 (Yes)	False Negatives (FN)	True Positives (TP)

ACTION REQUIRED – MLP (II)

- There are a few improvements that can be made to the algorithm, and there are some important things that need to be considered:
 - how many training data points are needed
 - how many hidden nodes should be used
 - how much training the network needs
- We will look at the improvements first, and then move on to practical considerations.
- The first thing that we can do is check that this MLP can indeed learn the logic functions, especially the XOR.
- See *Multi-layer_Perceptron.ipynb*

RECALL: DIFFERENT ACTIVATIONS

- In the algorithm described above, we used sigmoid neurons in the hidden layer and the output layer.
- This is fine for classification problems, since there we can make the classes be 0 and 1.
- However, we might also want to perform regression problems, where the output needs to be from a continuous range, not just 0 or 1.
- The sigmoid neurons at the output are not very useful in that case. We can replace the output neurons with linear nodes that just sum the inputs and give that as their activation.
- This does not mean that we change the hidden layer neurons; they stay exactly the same, and we only modify the output nodes.
- They are not models of neurons anymore, since they don't have the characteristic fire/don't fire pattern.
- Even so, they enable us to solve regression problems, where we want a real number out, not just a 0/1 decision.

DIFFERENT ACTIVATION FCT

```
# Different types of output neurons
if self.outtype == 'linear':
    return outputs
elif self.outtype == 'logistic':
    return 1.0/(1.0+np.exp(-self.beta*outputs))
elif self.outtype == 'softmax':
    normalisers = np.sum(np.exp(outputs),axis=1)*np.ones((1,np.shape(outputs))[0])
    return np.transpose(np.transpose(np.exp(outputs))/normalisers)
else:
    print "error"
```

Action required – see *Multi-layer_Perceptron.ipynb*

TEST – IRIS DATA SET



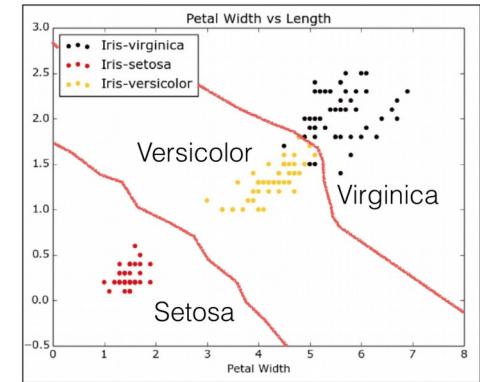
Iris Versicolor



Iris Setosa



Iris Virginica



Action required – see *Multi-layer_Perceptron.ipynb*

TEST – MNIST

- <http://yann.lecun.com/exdb/mnist/>
- The MNIST database of handwritten digits, available from this page, has a training set of 60,000 examples, and a test set of 10,000 examples. It is a subset of a larger set available from NIST. The digits have been size-normalized and centered in a fixed-size image.
- MNIST: Modified National Institute of Standards and Technology database.
- Action required – see *Multi-layer_Perceptron.ipynb*.

THURSDAY, JULY 23RD, 2020

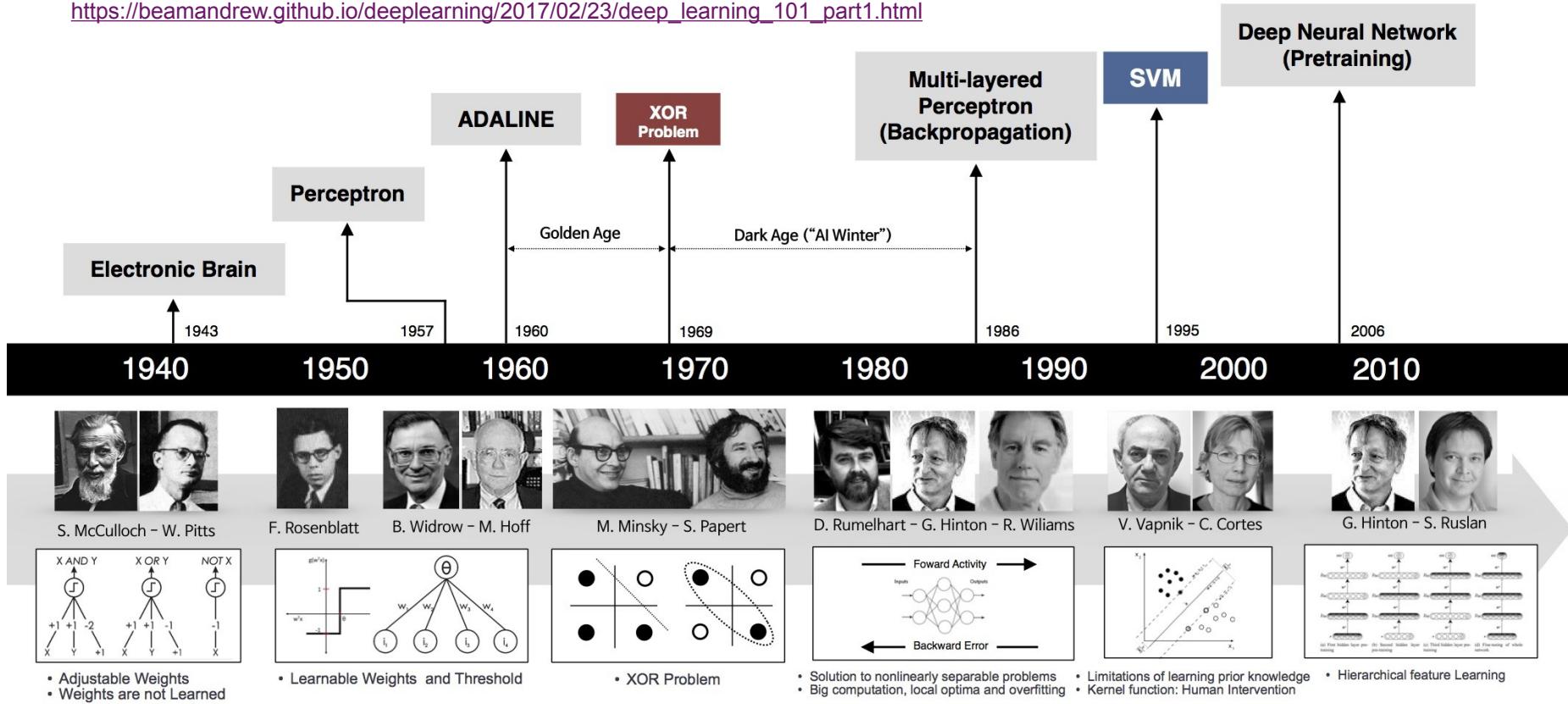


ROAD-MAP – THU (9.00-12.00, 1.30-3.30)

- Lecture (5 x 45-50 min):
 - Deep Learning cont'd
 - More advanced topics:
 - Recurrent neural networks and beyond.
 - Reinforcement learning.
 - Detour – think out of the box: "Deep Equilibrium Nets"
- Throughout lectures – hands-on:
 - Basics on Tensorflow & Keras
 - Examples related to the day's topics in Tensorflow

A TIMELINE OF DEEP LEARNING

https://beamandrew.github.io/deeplearning/2017/02/23/deep_learning_101_part1.html



RECAP – DAY 1

- Deep learning basics
 - Perceptron
 - Multi-layer Perceptron
 - Fully connected ANN
 - Activation functions
 - Stochastic Gradient Descent
- Examples (Regression & Classification)
 - Analytical Functions
 - Iris Data Set
 - MNIST Data Set

KERAS & TENSORFLOW BASICS

- [tensorflow.org](https://www.tensorflow.org)
- Keras API:
https://www.tensorflow.org/guide/keras/sequential_model
- Fun data sets to play with: <https://www.kaggle.com/datasets>
- Some “clean” data to play with: <https://archive.ics.uci.edu/ml/index.php>
- Help for debugging – Tensorboard: <https://www.tensorflow.org/tensorboard>

A GENTLE FIRST EXAMPLE

- Lets look at the notebook: **Gentle_DNN.ipynb**.
- This Notebook contains all the basic functionality that we learned yesterday from a theoretical point of view.
- 2 simple examples, one regression, and one classification.

A SEMI-COMPREHENSIVE TF TOUR

- **TF_tour.ipynb**
- 5 examples (incl. Kaggle data set from Lending Club)
- Tensorboard

SOME PERSONAL TAKE-AWAY

- Swish activation is the “best” if you need smooth and deep.
- Multiple of 2 for network (training speed).
- Smaller learning rate with deeper networks.
- Batch normalization for speed.
- Glorot initialization.
- Custom layers for custom models
(https://www.tensorflow.org/tutorials/customization/custom_layers).

A BREAK MOUNTAINS



BEYOND VANILLA DNN

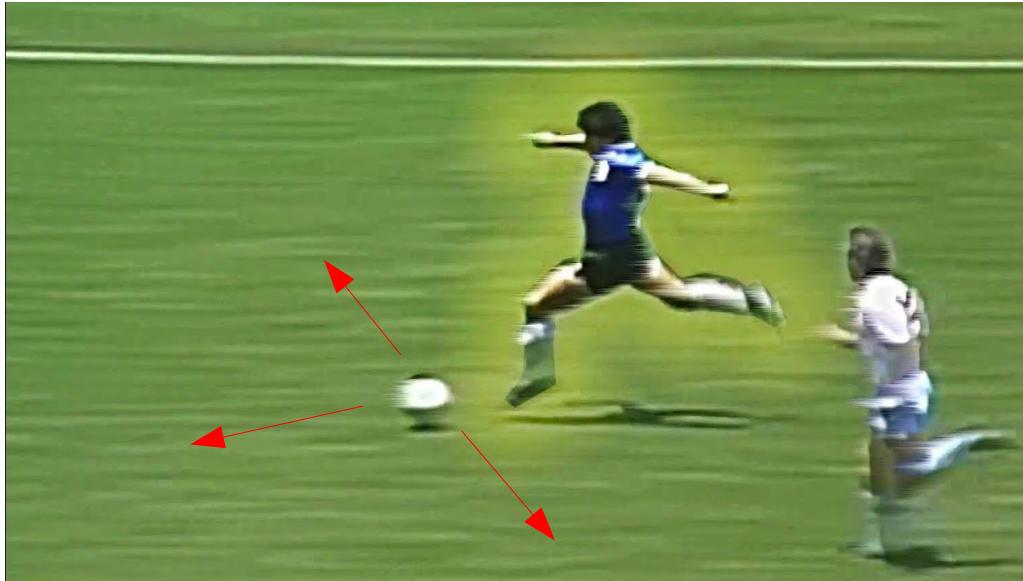
- Not all applications are plain vanilla deep neural nets.
- There exist situations where more intricate architectures are needed.
- Examples:
 - Time-series comparisons, such as estimating how closely related two documents or two stock tickers are.
 - Sequence-to-sequence learning, such as decoding an English sentence into French
 - Sentiment analysis, such as classifying the sentiment of tweets or movie reviews as positive or negative.
 - Time-series forecasting, such as predicting the future weather at a certain location, given recent weather data.

EXAMPLE



Given a picture of a ball, can we predict where it will go?

EXAMPLE



Given a picture of a ball, can we predict where it will go?

A SEQUENCE MODELING PROBLEM: PREDICT THE NEXT WORD

“Today, we hare having a class on deep learning”

A SEQUENCE MODELING PROBLEM: PREDICT THE NEXT WORD

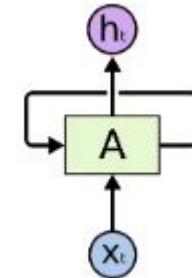
“Today, we hare having a class on deep learning”

given these words

predict the next word

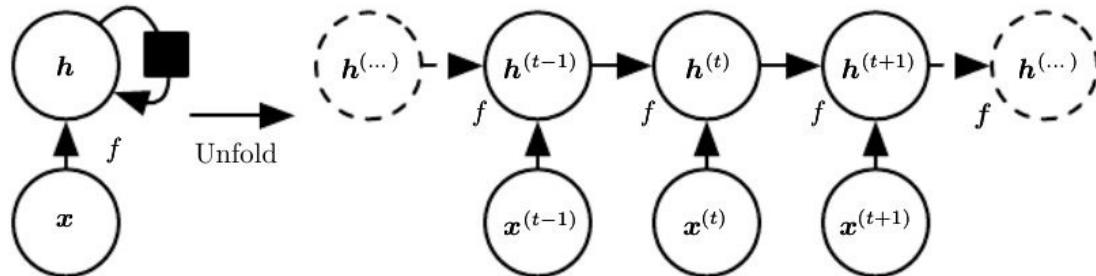
RECURRENT NEURAL NETS

- To model sequences, we need to
 - Handle variable-length sequences
 - Track long-term dependencies
 - Maintain information about the order
 - Share parameters across the sequence
- Recurrent Neural Networks (RNN) are an approach to sequence modeling problems (Rumelhart et al. (1986)).



RNN

- RNNs, are a family of neural networks for processing sequential data.
- A RNN is a neural network that is specialized for processing a sequence of values $x^{(1)}, \dots, x^{(T)}$.
- Unfold the computational graph of a dynamical system:



$$s^{(t)} = f(s^{(t-1)}; \theta)$$

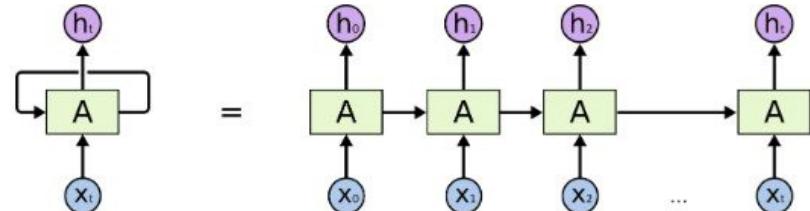
$$s^{(3)} = f(s^{(2)}; \theta)$$

$$= f(f(s^{(1)}; \theta); \theta)$$

RNN

- Recurrent layers use their own output as input
- Introduce history or time dependency in NNs.
- The only way to efficiently train them is to unroll them

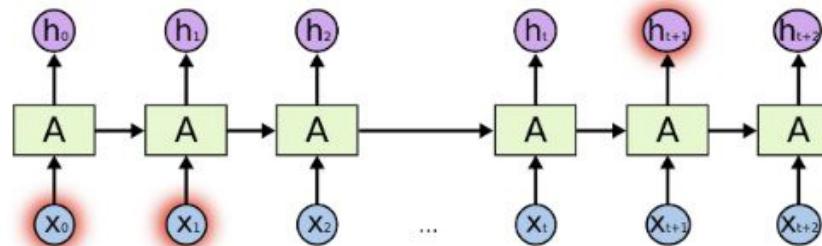
$$\mathbf{h}^{(t)} = f \left(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \theta \right)$$



An unrolled recurrent neural network.

RNN

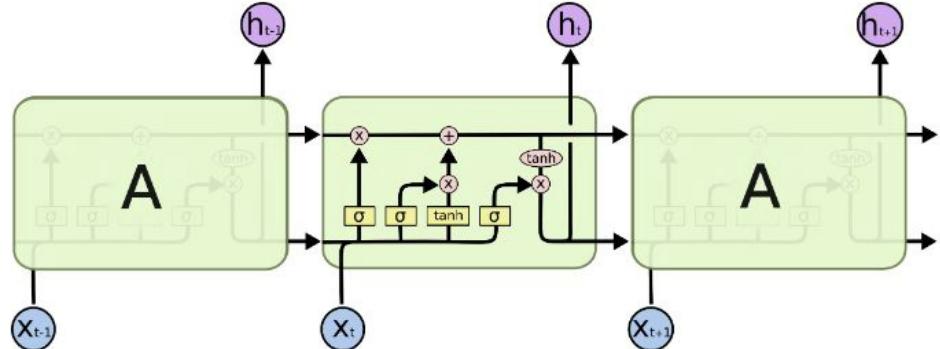
- Recurrent blocks suffer from two problems:
 - Long-term dependencies do not work well.
 - Difficult to connect two distant parts of the input.
 - Magnitude of the signal can get amplified at each recurrent connection.
 - At every time iteration, the gradient can either vanish or explode.
 - Very hard to train them.



LONG SHORT-TERM MEMORY (LSTM)

<http://www.bioinf.jku.at/publications/older/2604.pdf>

- Hochreiter & Schmidhuber (1997)
- LSTM layers are improved versions of the recurrent layers
 - They can learn long-term dependencies.
 - They can forget.
- They have an internal state and a structure which is composed of four actual layers
 - Layers labelled with σ are gates which can block or let information flow.

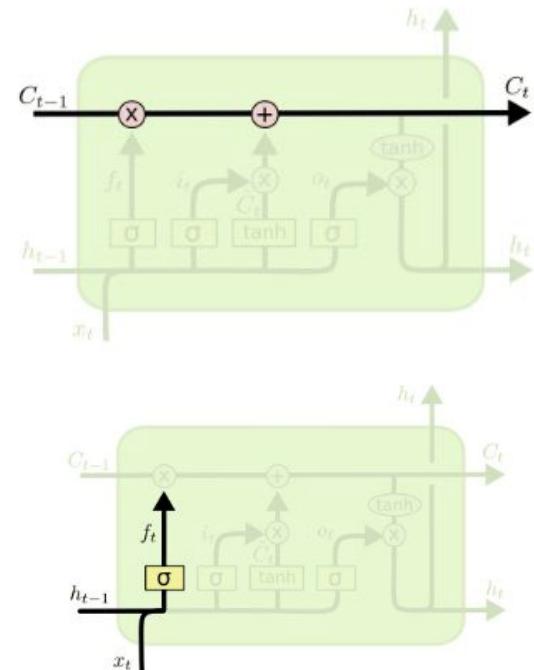


The repeating module in an LSTM contains four interacting layers.

LSTM FORGET GATE

<http://www.bioinf.jku.at/publications/older/2604.pdf>

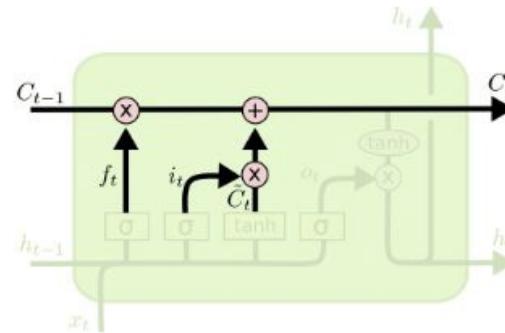
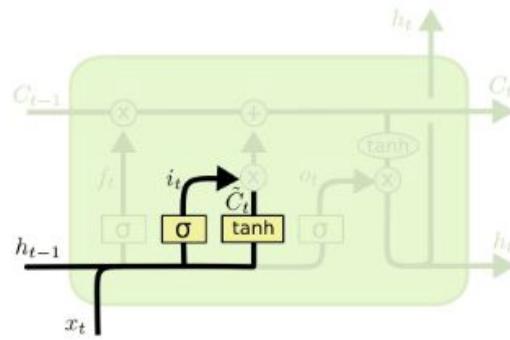
- LSTMs follow two paths
 - They update their internal state
 - They give an output based on the internal state and the input.
- A gate layer σ decides if we should forget an old part of the internal state
 - Something which has to be replaced by new information.



LSTM NEW STATE

<http://www.bioinf.jku.at/publications/older/2604.pdf>

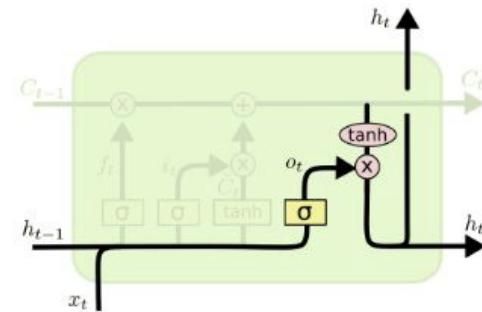
- Once that the layer decided what to forget, it computes
 - What has to replace it, i_t , based on the input and the old state.
 - What has to be used to replace it, the candidate value C_t
- The new state C_t can be computed based on the new information.



LSTM OUTPUT

<http://www.bioinf.jku.at/publications/older/2604.pdf>

- Based on the new state and the input, the layer can produce a result
 - this is the output.
 - the same value is also passed to the next iteration
- Why is this so important?
 - Many translation algorithms and voice interpreters are based on small variations of this layer.
- Action required: **RNN_intro.ipynb**

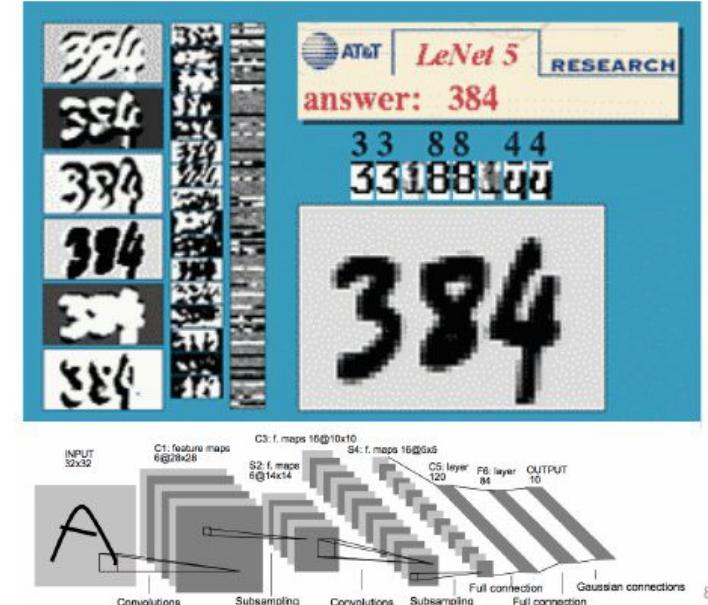


(see also <https://www.tensorflow.org/guide/keras/rnn>)

CONVOLUTIONAL NEURAL NETS

- Possibly the most successful types of networks.
- Uses sequences of convolutional layers.
- Can be interleaved with pooling operations or fully-connected layers.
- Train faster than MLPs.
- Can be used for 2D, 3D or higher dimensions (though 2D are the most common).
- Used for image recognition, object detection, sound analysis, etc.
- There exist more intricate architectures.
- Yann LeCun (1998)

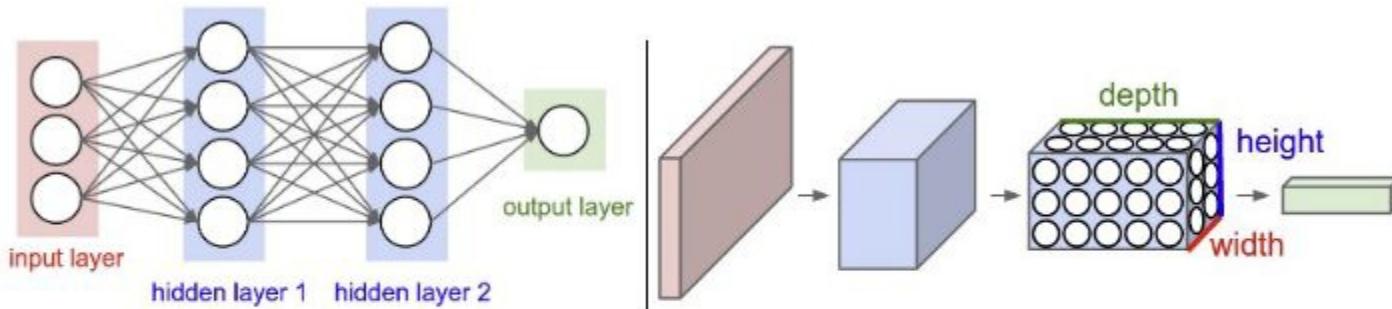
<http://yann.lecun.com/exdb/publis/pdf/lecun-98.pdf>



Unil

UNIL | Université de Lausanne

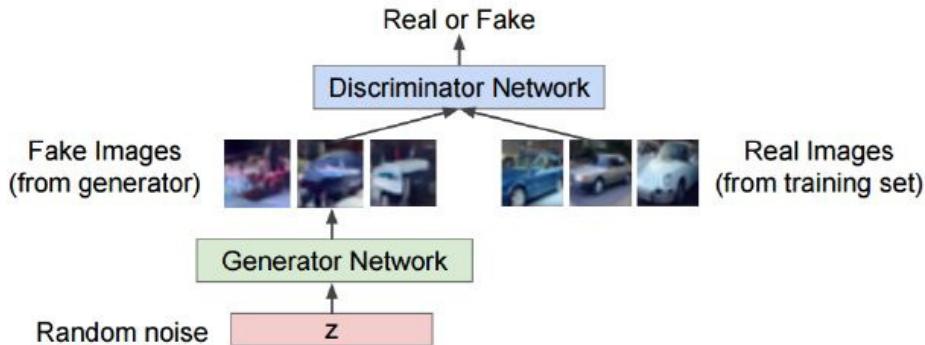
CONVOLUTIONAL NEURAL NETS



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

GENERATIVE ADVERSARIAL NETS

- GANs were introduced by Goodfellow et al. (2014)
- <https://arxiv.org/abs/1406.2661>.
- The idea is to train a network to generate samples which are indistinguishable from real ones (from the training set).
 - The input is a random noise sample (latent space).
- Another network is trained at the same time to distinguish between real and fake samples.



Training Data
(CelebA)



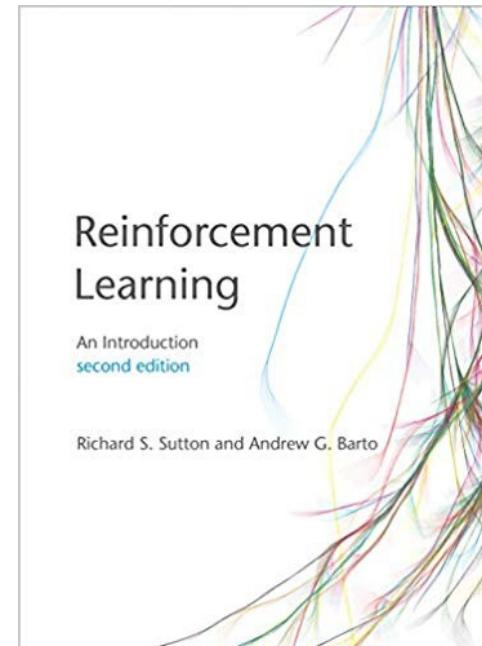
Sample Generator
(Karras et al, 2017)

A BREAK WITH MOUNTAINS



REINFORCEMENT LEARNING

- <https://www.alphagomovie.com>
- https://youtu.be/8tq1C8spV_g



RECALL – CLASSES OF LEARNING

- Supervised Learning
- Data: (x, y)
 x is data, y is label
- Goal:
Learn function to map $x \rightarrow y$
- Example (Classification):



This thing is an apple

RECALL – CLASSES OF LEARNING

- Unsupervised Learning
- Data: x
 x is data, no labels
- Goal:
Learn underlying structure
- Example:



This thing is like the other thing

REINFORCEMENT LEARNING

- Reinforcement Learning (RL)
- Data: state-action pairs
- Goal: Act optimally in a given environment.
Maximize future rewards over many time steps



- Example:
 - Eat this thing because it will keep you alive

RL KEY CONCEPTS



Agent

Agent – takes actions

RL KEY CONCEPTS



Agent



Environment

Environment: the world in which the agent exists and operates.

RL KEY CONCEPTS



Agent

Action: a_t

Actions

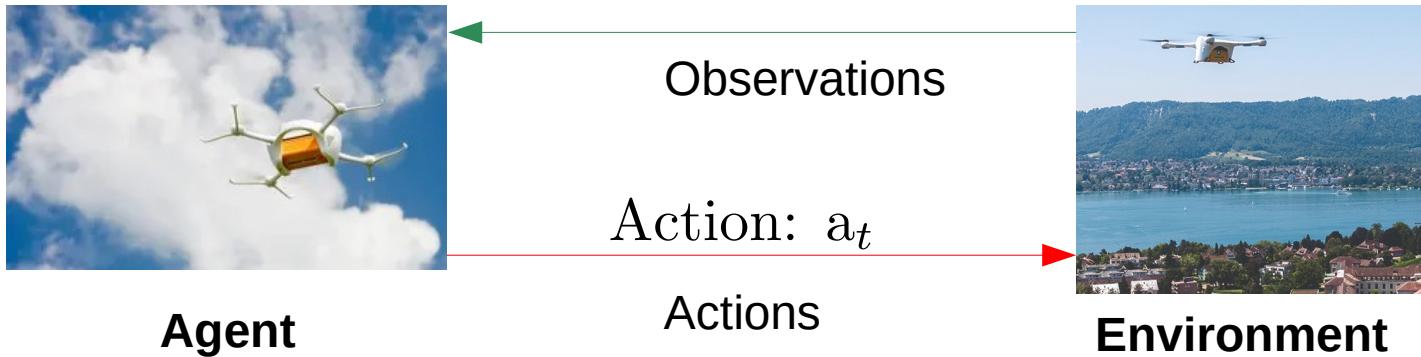


Environment

Action: a (possible) move the agent can make in the environment.

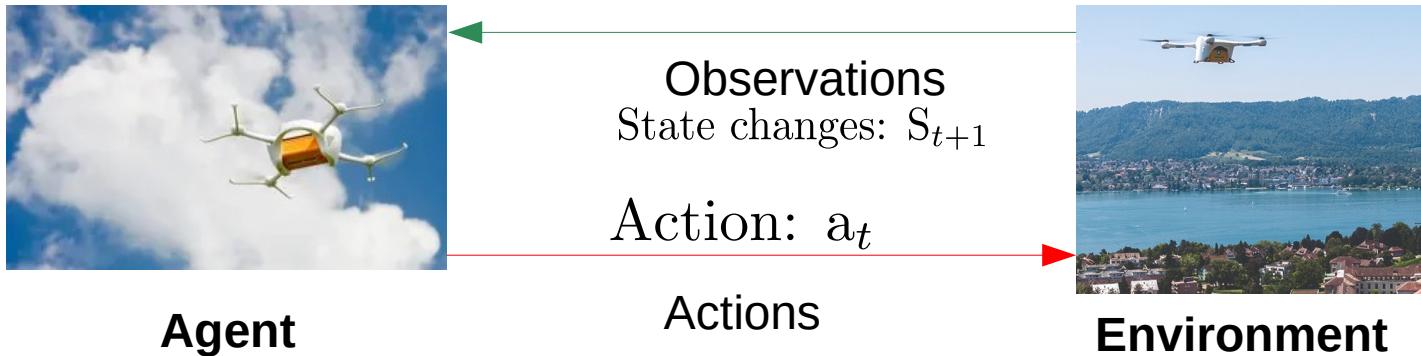
Actions for now: discrete (limited amount of possible actions for now).

RL KEY CONCEPTS



Observations: of the environment after taking actions.

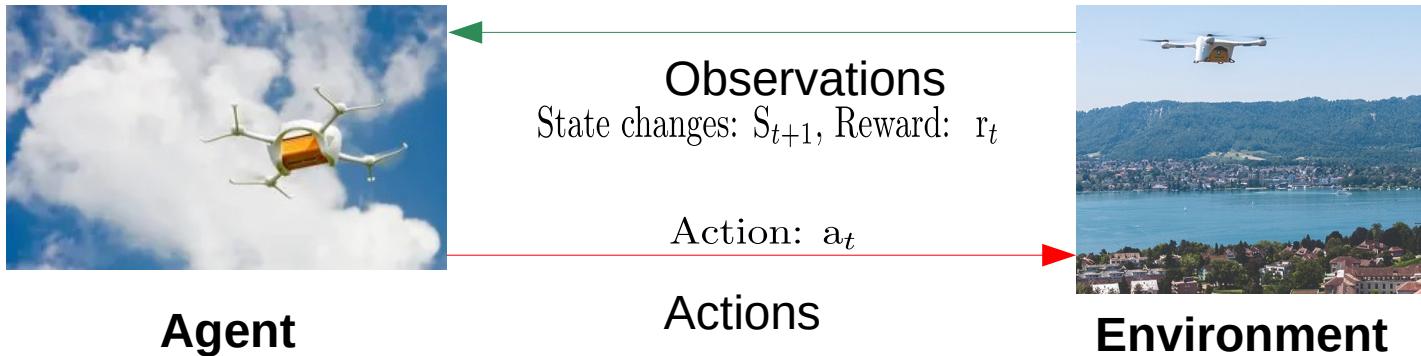
RL KEY CONCEPTS



State: a situation which the agent perceives.

“A concrete situation where the agent is in”, measure e.g., by a sensor.

RL KEY CONCEPTS



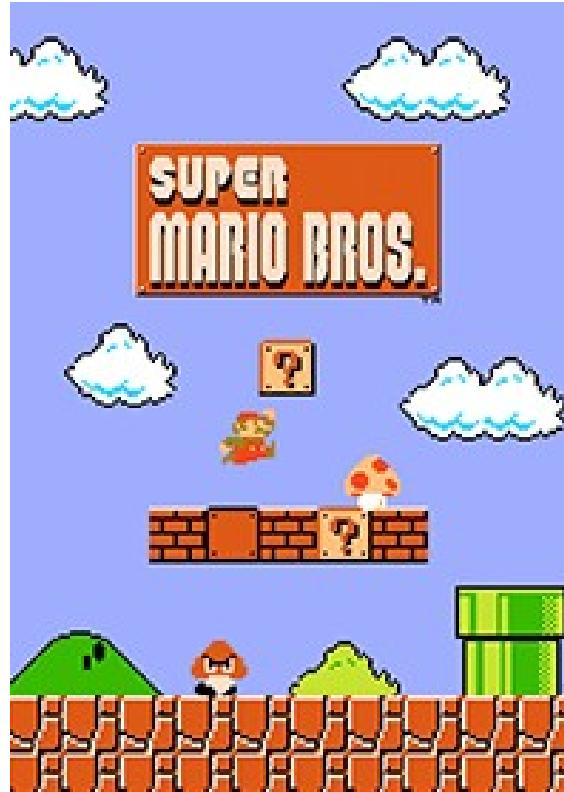
Reward: feedback that measures the success or failure of the agent's action in the environment.

“Good action” : high reward.

“Poor action” : low reward.

Rewards: can be immediate, or they can be delayed!

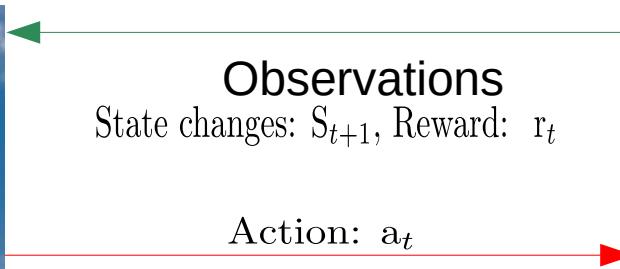
EXAMPLE: GETTING A REWARD



RL KEY CONCEPTS



Agent



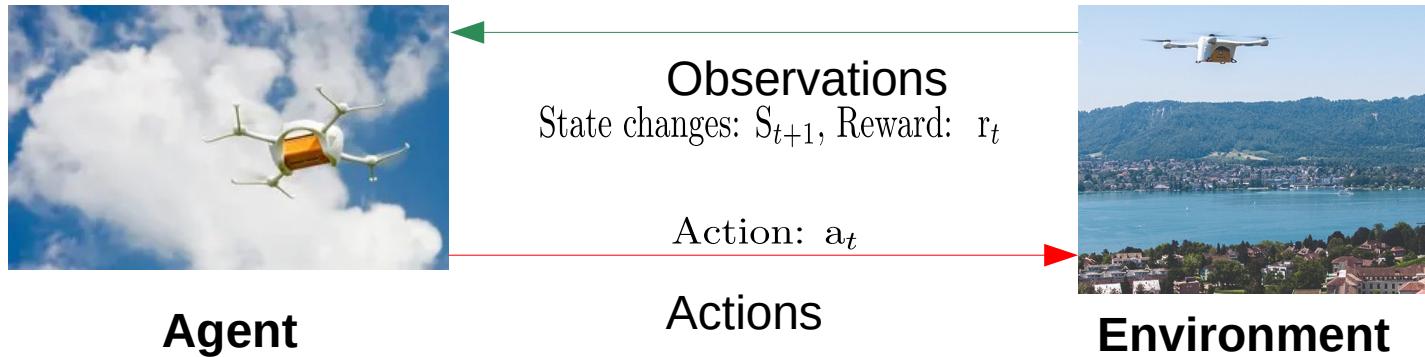
Environment

$$R_t = \sum_{i=t}^{\infty} r_i$$



Total reward

RL KEY CONCEPTS



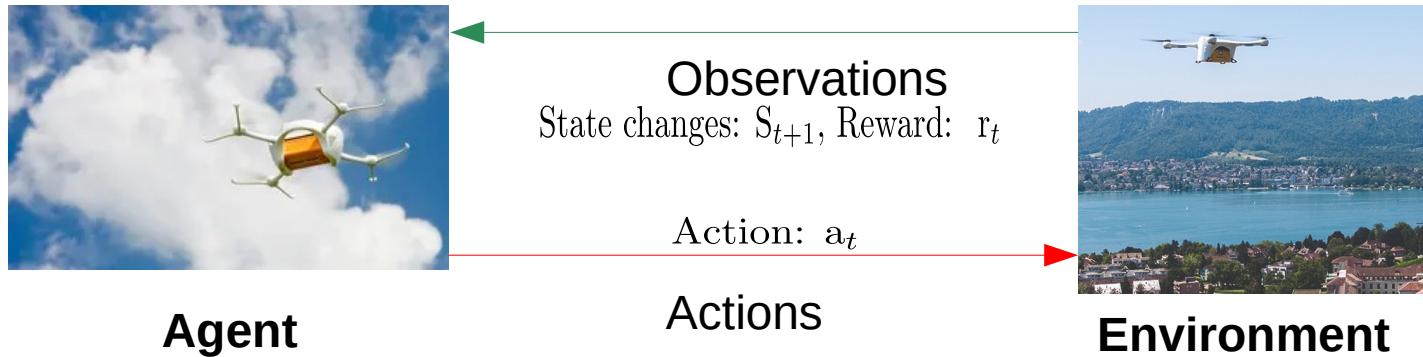
$$R_t = \sum_{i=t}^{\infty} r_i = r_t + r_{t+1} + \dots + r_{t+n} + \dots$$



Total reward

Problem: Want to avoid infinitely large rewards!

RL KEY CONCEPTS



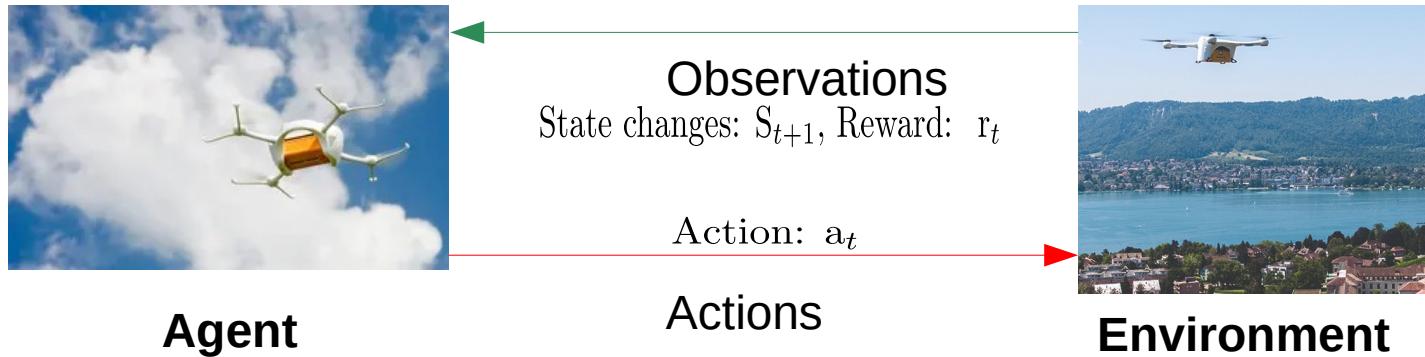
$$R_t = \sum_{i=t}^{\infty} \gamma^i r_i$$



Discounted total reward

Problem: Want to avoid infinitely large rewards → use discount factor!

RL KEY CONCEPTS



$$R_t = \sum_{i=t}^{\infty} \gamma^i r_i = \gamma^t r_t + \gamma^{t+1} r_{t+1} + \dots + \gamma^{t+n} r_{t+n} + \dots$$

γ : discount factor

Discounted total reward

Problem: Want to avoid infinitely large rewards → use discount factor!

EXAMPLE: LEARN TO DRIVE

<https://wayve.ai/blog/learning-to-drive-in-a-day-with-reinforcement-learning>



Self-driving car

Action: Steering angle

Next State: Next camera input that the car sees.

THE Q-FUNCTION

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

Total reward, R_t is the discounted sum of all rewards obtained from time t .

$$Q(s, a) = \mathbb{E}[R_t]$$

The Q-function captures the expected total future reward an agent in state, s, can receive by executing a certain action, a.

HOW TO TAKE ACTIONS GIVEN A Q-FUNCTION?

$$Q(s, a) = \mathbb{E} [R_t]$$

(state, action)

Ultimately, the agent needs a policy $\pi(s)$, to infer the best action to take at its state s .

Strategy: the policy should choose an action that maximizes the future reward:

$$\pi^*(s) = \operatorname{argmax}_a Q(s, a)$$

RL ALGORITHMS

Value Learning

Find $Q(s, a)$

$a = \underset{a}{\operatorname{argmax}} Q(s, a)$

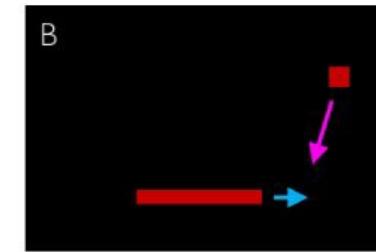
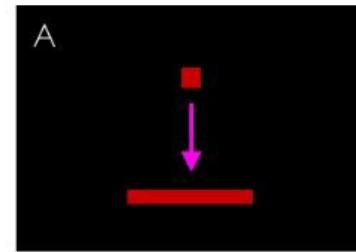
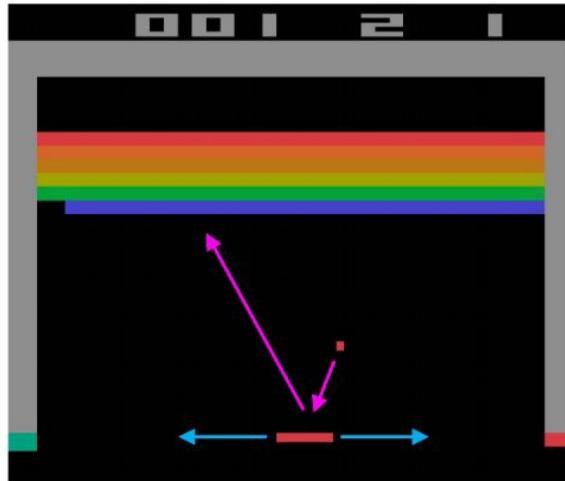
Policy Learning

Find $\pi(s)$

Sample $a \sim \pi(s)$

LOOKING CLOSER INTO THE Q-FUNCTION

Example: Atari Breakout

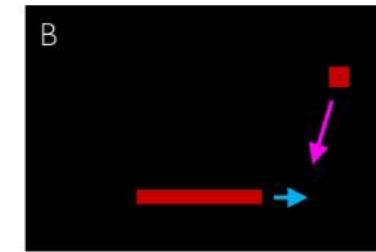
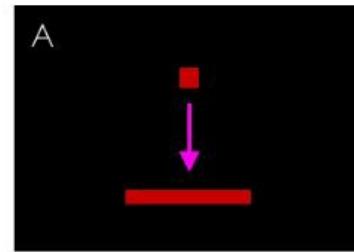
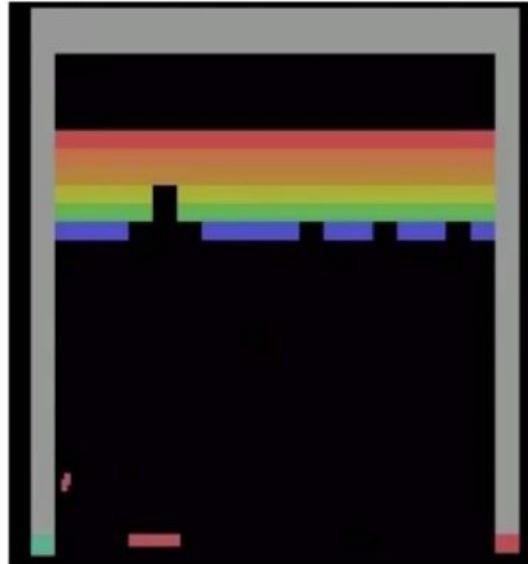


It can be very difficult for humans to Accurately estimate Q-values.

Which (s,a) pair has a higher Q-value?

LOOKING CLOSER INTO THE Q-FUNCTION

Example: Atari Breakout - Middle

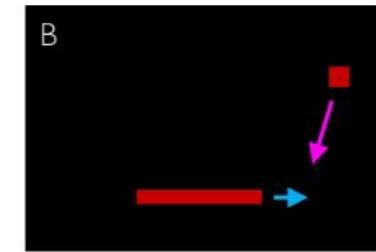
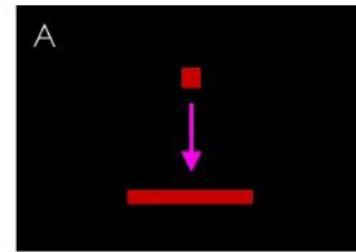


It can be very difficult for humans to Accurately estimate Q-values.

Which (s,a) pair has a higher Q-value?

LOOKING CLOSER INTO THE Q-FUNCTION

Example: Atari Breakout - Side



It can be very difficult for humans to Accurately estimate Q-values.

Which (s,a) pair has a higher Q-value?

ATARI BREAKOUT

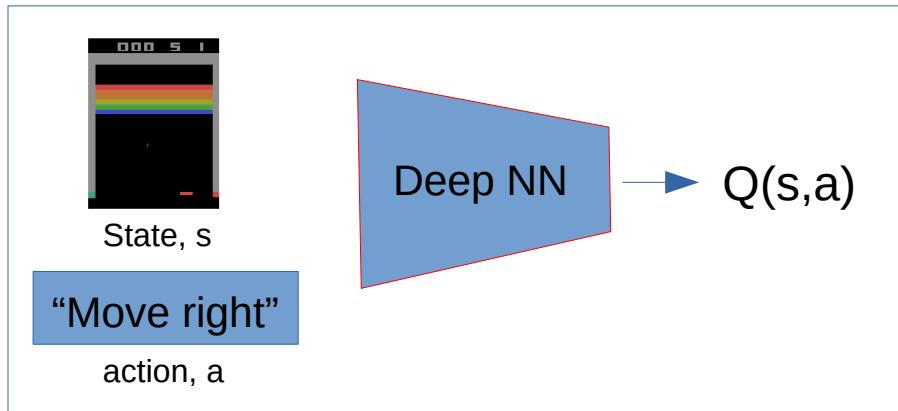
- <https://www.youtube.com/watch?v=V1eYniJ0Rnk>

MISSPECIFIED REWARD FUNCTION

- <https://www.youtube.com/watch?v=tOIHko8ySg>

EXAMPLE: RF & DNN

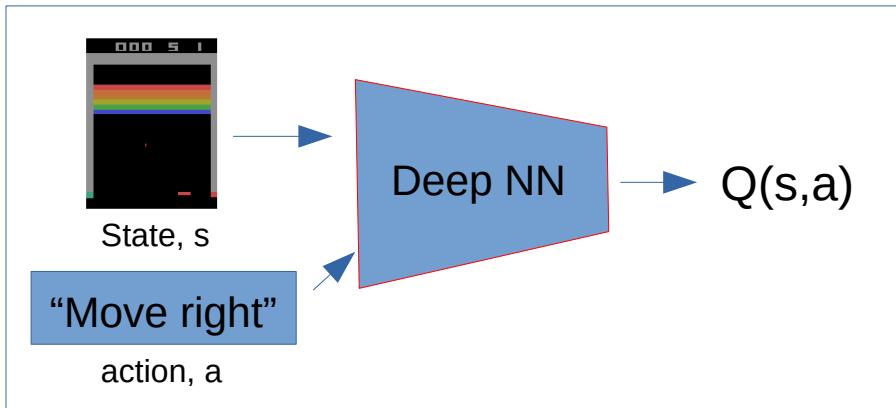
How can we use deep neural networks to model Q-functions?



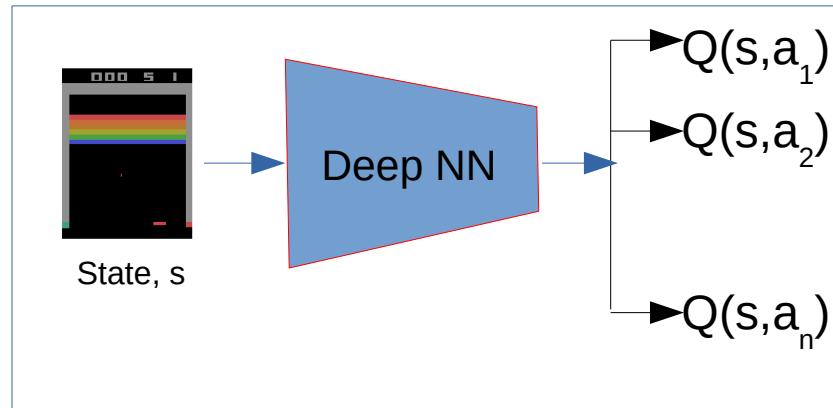
Problem: for a given state, we need to compute Q for all possible a
→ computationally “expensive”

EXAMPLE: RF & DNN

How can we use deep neural networks to model Q-functions?



Problem: for a given state, we need to compute Q for all possible a
→ computationally “expensive”.



Alternative approach.
→ Multiple outputs.
→ Take max.

EXAMPLE: RF & DNN

How can we use deep neural networks to model Q-functions?



$$\mathcal{L} = \mathbb{E} \left[\left\| \left(r + \gamma \max_{a'} Q(s', a') \right) - Q(s, a) \right\|^2 \right]$$

EXAMPLE: RF & DNN

How can we use deep neural networks to model Q-functions?



Value you observed when taking a particular action

target

$$\mathcal{L} = \mathbb{E} \left[\left\| \left(r + \gamma \max_{a'} Q(s', a') \right) - Q(s, a) \right\|^2 \right]$$

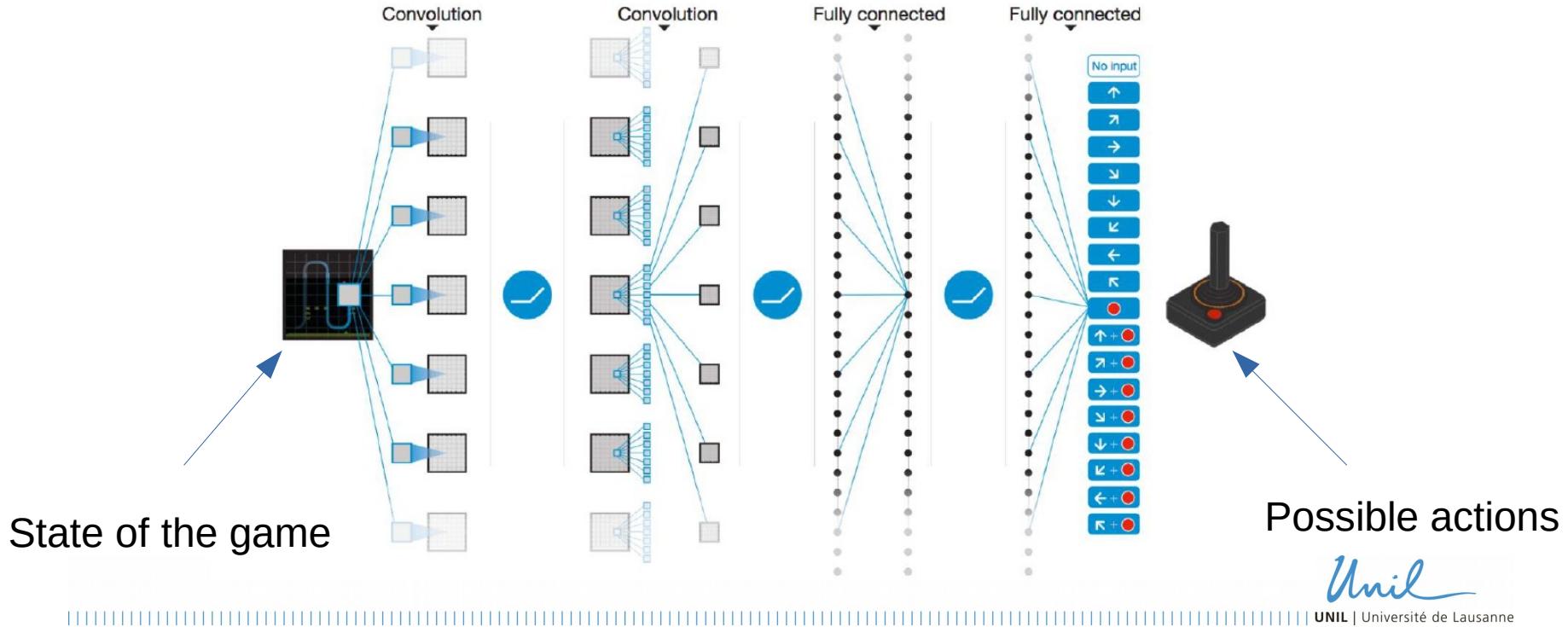
predicted

Output of NN

→ Minimize Cost function

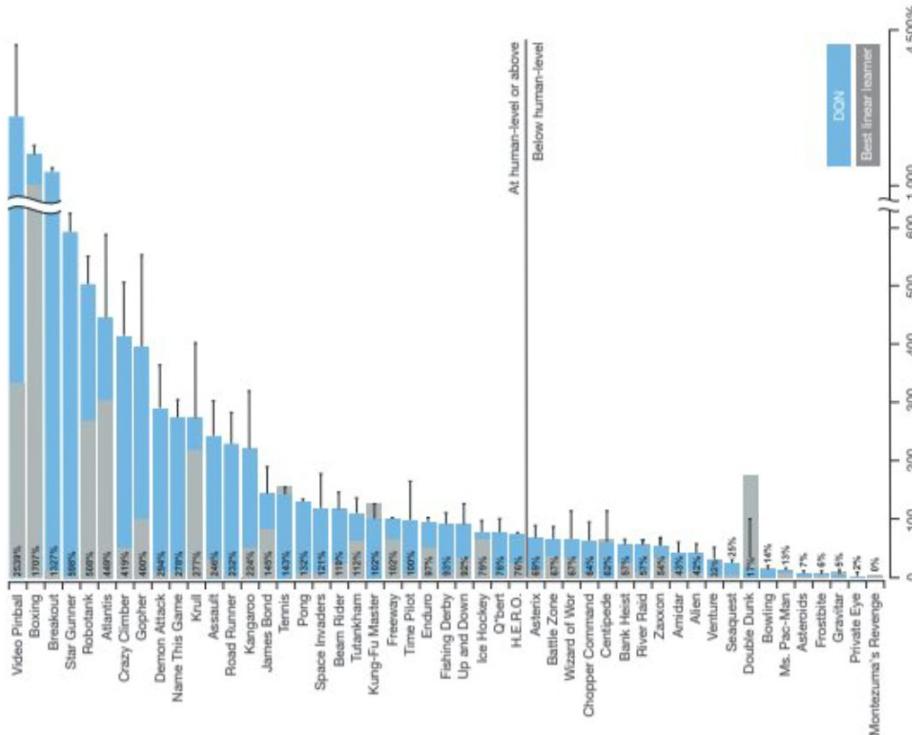
TRAIN A DQN FOR THE ATARI TEST CASE

<http://deepmind.com/>



ATARI: SUPER HUMAN PERFORMANCE

Mnih et al. (2015) – <https://web.stanford.edu/class/psych209/Readings/MnihEtAlHassabis15NatureControlDeepRL.pdf>



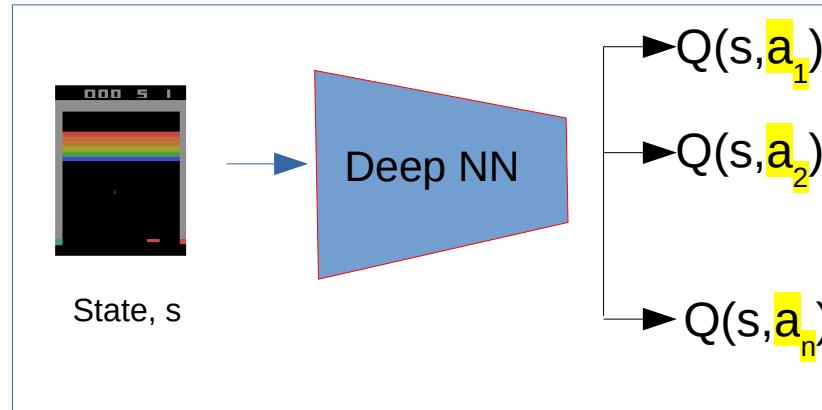
SOME PROBLEMS WITH DQN

- Trouble with Complexity:
 - Can model scenarios where the action space is discrete and small
 - Cannot handle continuous action spaces
- Trouble with Flexibility:
 - Cannot learn stochastic policies since policy is deterministically computed from the Q function (take argmax).
 - To overcome, consider a new class of RL training algorithms:
Policy gradient methods.

IMPORTANT:
Imagine you want to predict
steering wheel angle of a car!

THE BASIC IDEA OF POLICY GRADIENT

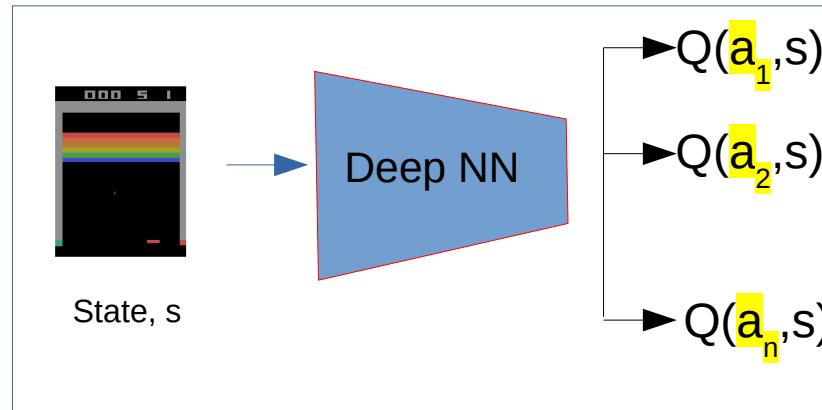
DQN (before): Approximating Q and inferring the optimal policy



THE BASIC IDEA OF POLICY GRADIENT

DQN (before): Approximating Q and inferring the optimal policy

Policy Gradient: Directly optimize the policy!

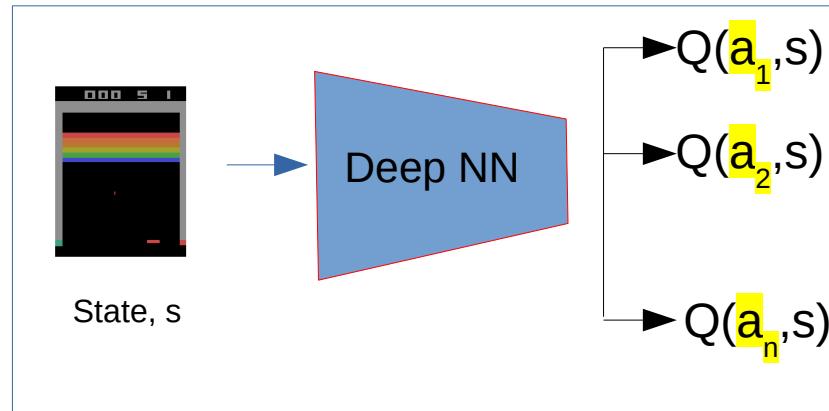


THE BASIC IDEA OF POLICY GRADIENT

DQN (before): Approximating Q and inferring the optimal policy

Policy Gradient: Directly optimize the policy!

Achieve this e.g.
via softmax.

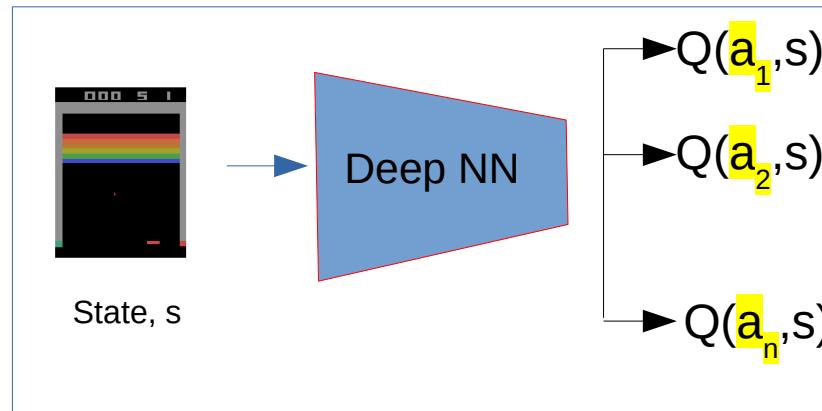


$$\sum_{a_i \in A} \pi(a_i | s) = 1$$

THE BASIC IDEA OF POLICY GRADIENT

DQN (before): Approximating Q and inferring the optimal policy

Policy Gradient: Directly optimize the policy!



Achieve this e.g.
via softmax.

$$\sum_{a_i \in A} \pi(a_i | s) = 1$$

$$\pi(a | s) = P(\text{action} | \text{state})$$

CONNECTION TO DYNAMIC PROGRAMMING

- The solution is approached in the limit as $j \rightarrow \infty$ by iterations on at every coordinate of the “grid” of observations.

$$V_{j+1}(x) = \max \{r(x, u) + \beta V_j(\tilde{x})\}$$

s.t.

$$\tilde{x} = g(x, u)$$

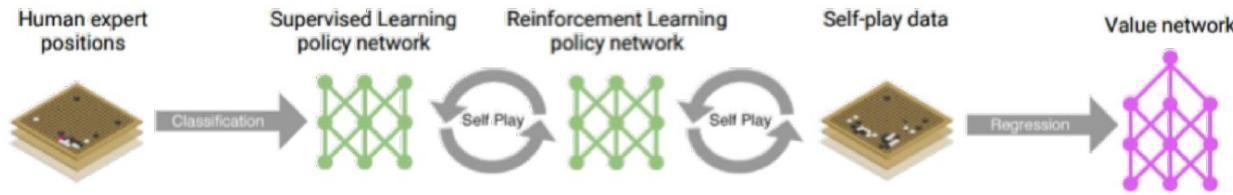
x : grid point/observation/state variable.
Describes your problem.

↑
`old solution' on which we
Interpolate (i.e., “predict”).

RL WITH AVAILABLE TOOLS

- Keras & TF resources
 - <https://keras.io/examples/rl>
- Deep Policy Gradient
 - https://keras.io/examples/rl/ddpg_pendulum/
 - If time permits, lets look at the code: ***ddpg_pendulum.ipynb***
- Atari Breakout
 - https://colab.research.google.com/github/keras-team/keras-io/blob/master/examples/rl/ipynb/deep_q_network_breakout.ipynb

ALPHA GO BEATS HUMAN PLAYER



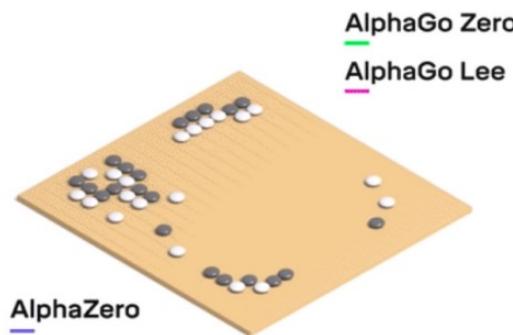
ALPHA GO BEATS HUMAN PLAYER



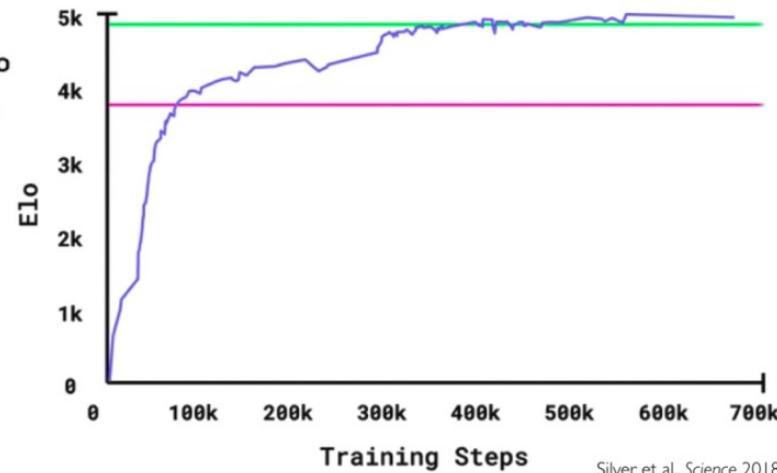
ALPHAZERO: RL FROM SELF-PLAY (2018)

<https://deepmind.com/blog/alphago-zero-learning-scratch/>

Go

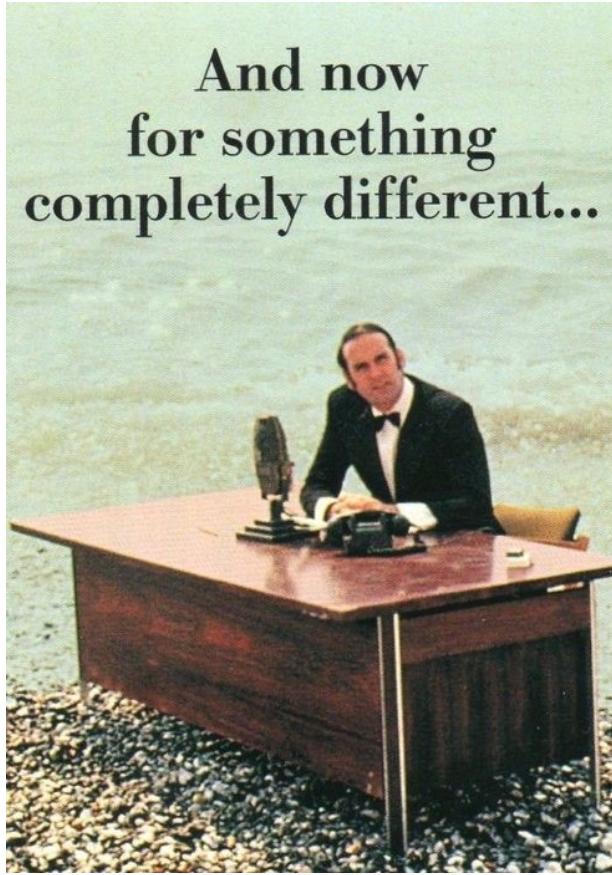


AlphaGo Zero
AlphaGo Lee



Silver et al., Science 2018.

**And now
for something
completely different...**

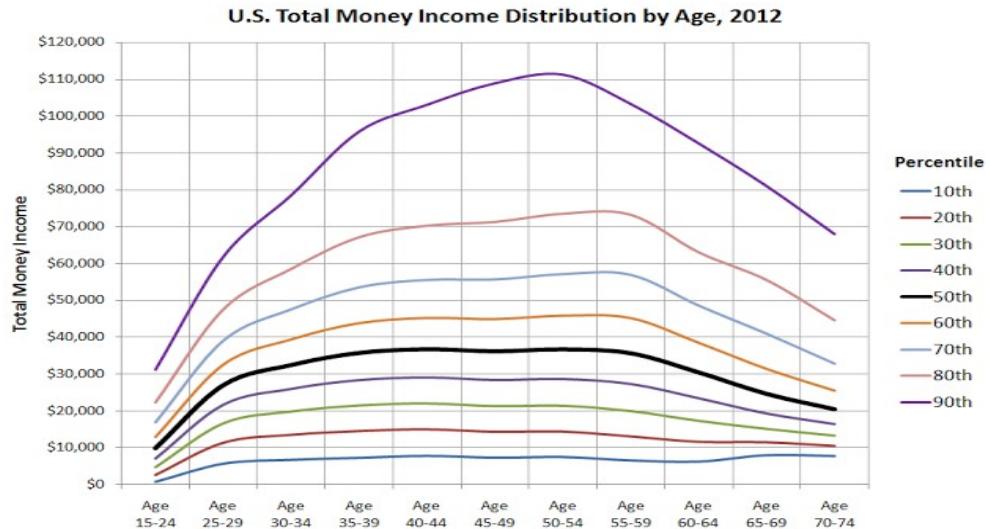


DYNAMIC STOCHASTIC ECONOMIC MODELS



- **Heterogeneity** a crucial ingredient in contemporary models:
 - to study e.g. cross-sectional consumption response to aggregate shocks.
 - to model, e.g., social security.
- Example OLG models:
 - How many age groups?
 - borrowing constraints?
 - aggregate shocks?
 - idiosyncratic shocks?
 - liquid / illiquid assets**?

→ **Models: heterogeneous & high-dimensional**



**see, e.g., Kaplan et al. (2018), Wong (2018),...

DEEP EQUILIBRIUM NETS

<https://github.com/sischei/DeepEquilibriumNets>

- “Economic” loss function:
$$l_\rho := \frac{1}{N_{\text{path length}}} \sum_{x_i \text{ on sim. path}} (\mathbf{G}(x_i, \mathcal{N}_\rho(x_i)))^2$$
- where we use \mathcal{N}_ρ to simulate a path.
- G is chosen such that the true equilibrium policy $f(x)$ is defined by $G(x, f(x)) = 0 \forall x$.
- $G(., .)$: implied error in the optimality conditions (unit-free Euler errors).
- Therefore, there is no need for labels to evaluate our loss function.
 - Unsupervised Machine Learning.

QUESTIONS?



THANK YOU FOR YOUR ATTENTION



KeepCalmAndPosters.com