# Machine Learning and Artificial Neural Networks

Anh Khoa Doan

# What you have seen previously

- So far: Uncertainty Quantification
  - Combines "knowledge by reasoning" (from numerical analysis) and "knowledge by data" (statistics)…
  - To get a better understanding (and prediction) of truth
- What we will see in the next two sessions
  - Emphasis on "knowledge by data"…
    - "Machine/Deep Learning"
  - … and one of the form of combining Bayesian philosophy with machine learning
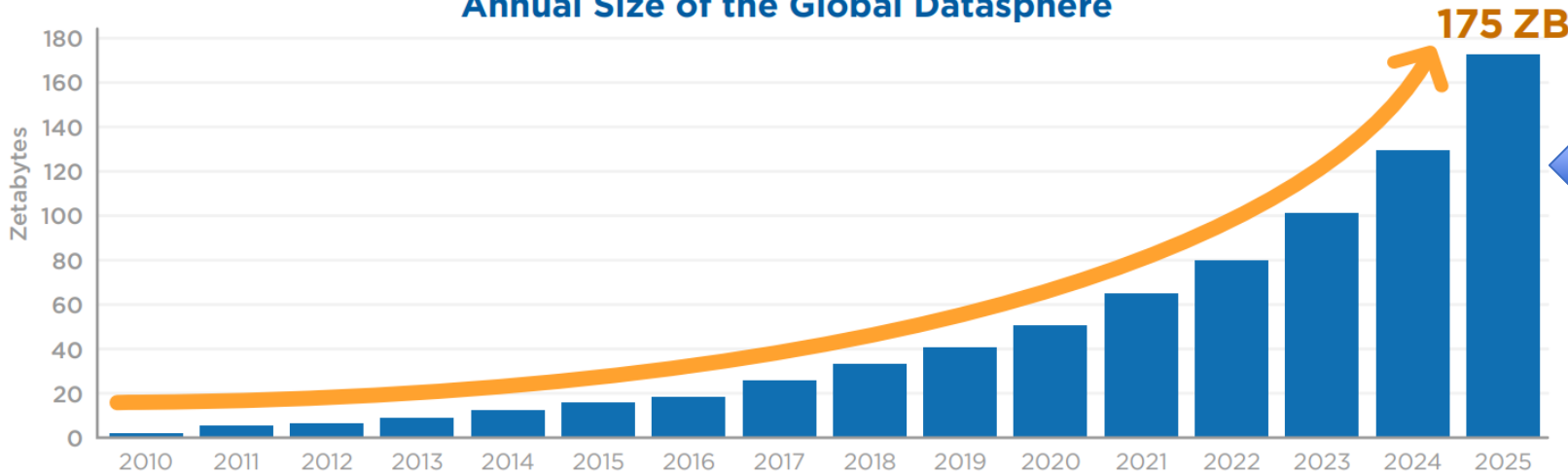
# Structure of the Lecture

- Introduction to Machine Learning (ML)
  1. Drivers behind current ML
  2. Position of ML in science
  3. Classification of ML methods

- Introduction to Neural Network
  1. Linear regression and computational graph
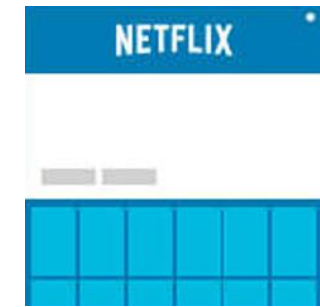  2. Gradient descent
  3. The neuron

**TU**Delft

# Data is becoming increasingly prevalent

- In recent time: exponential explosion in data

**Annual Size of the Global Datasphere**

175 ZB

(Bar chart showing Zetabytes on y-axis from 0 to 180, years 2010 to 2025 on x-axis, with exponential growth reaching 175 ZB by 2025)

Entire Netflix catalogue
**~500 million times**

NETFLIX

37 trillion

DVD
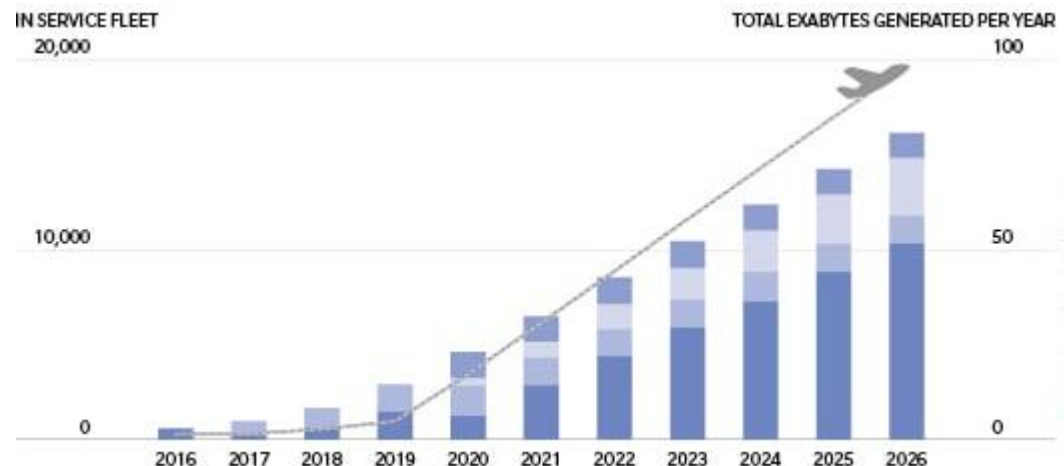
Stacked:
100 ladders to the moon

→ Driven by "digitalization"

# Data in aeronautical industry is also exploding

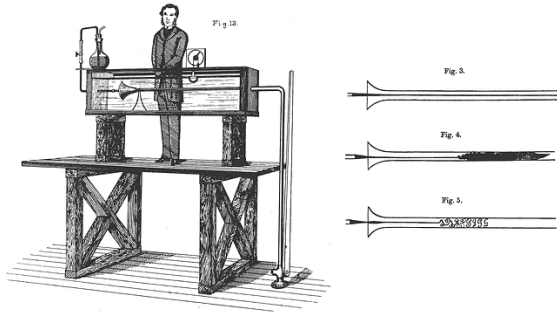- How much data is generated per flight?

Just from the engines:
GE (2018): "around 1TB per engine per flight"
Approximately 120,000 flights per day
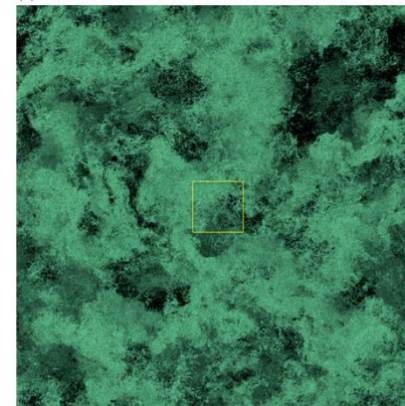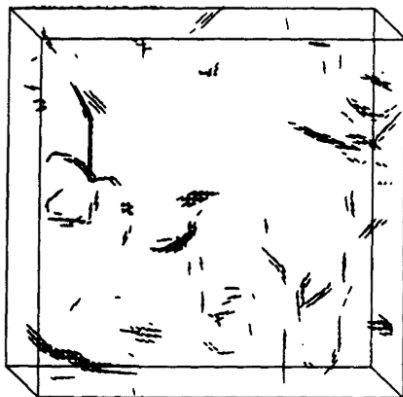→ 120 PB/day
→ 43.8EB/year

# On a smaller (lab-)scale

- But actually, how much data produced by
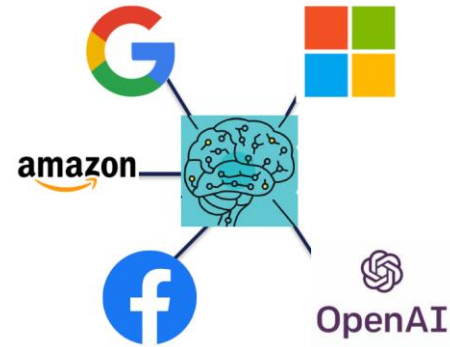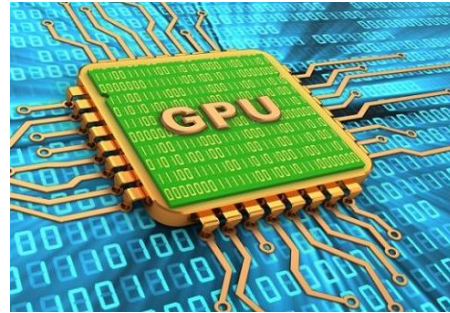  - One lab-experiment?





  - One CFD simulation?

# Recent success in ML

- "Big data" and "Machine Learning" techniques
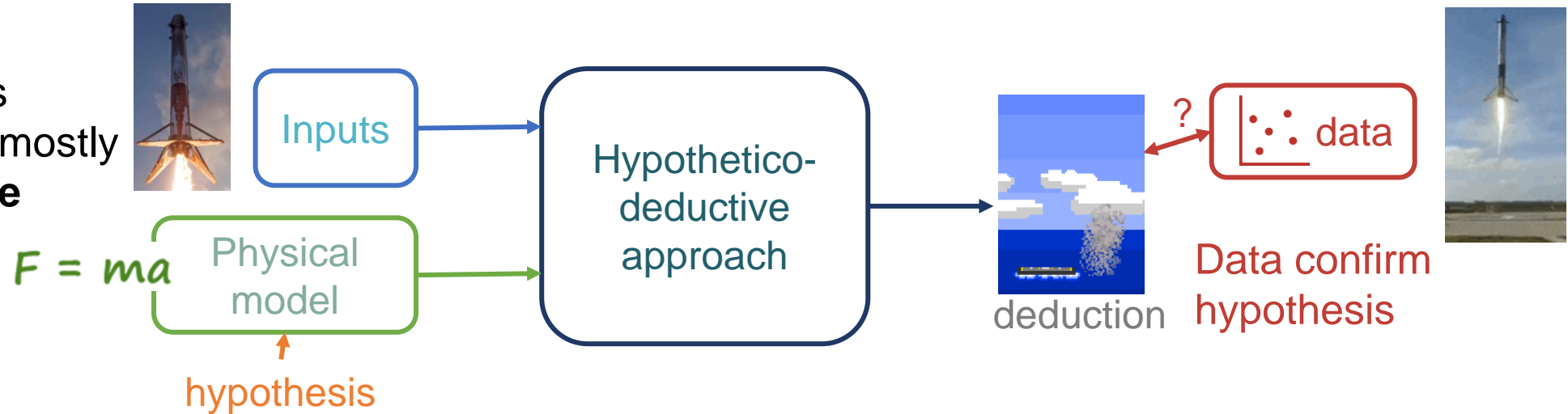  - Many success in exploiting/analysing "large dataset"



- Can we leverage "new" data-driven techniques (such as machine learning) for fluid mechanics research?
  - What makes this process "different"?
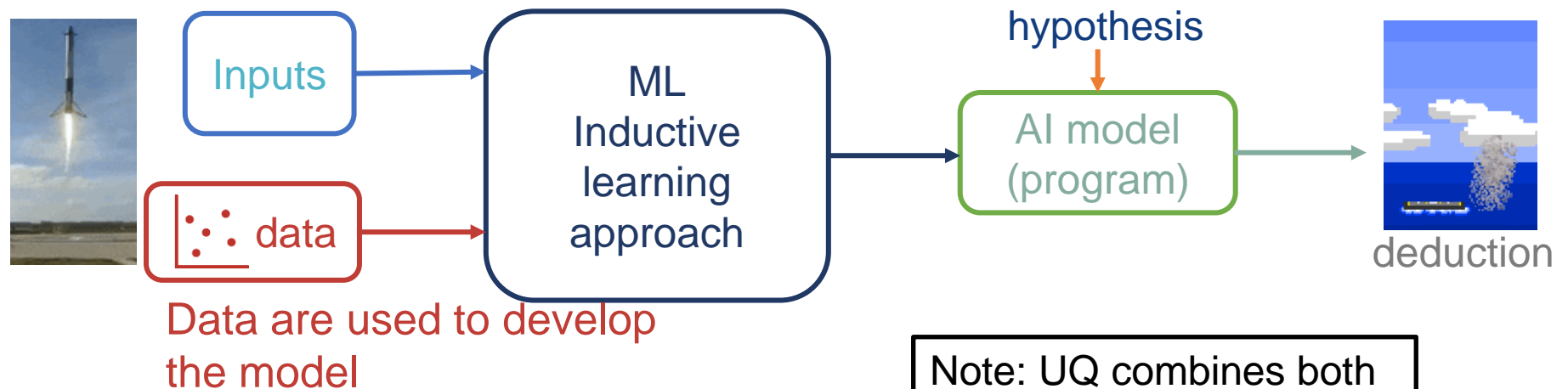
# Structure of the Lecture

- Introduction to Machine Learning (ML)
    1. Drivers behind current ML
    2. **Position of ML in science**
    3. Classification of ML methods

- Introduction to Neural Network
    1. Linear regression and computational graph
    2. Gradient descent
    3. The neuron

# The objective of ML is the obtention of the model

Scientific method is currently mostly **deductive**

$F = ma$

Inputs

Physical model

↑ hypothesis

Hypothetico-deductive approach

deduction

? → data

Data confirm hypothesis

Machine Learning is an **inductive** process

Inputs

data

ML Inductive learning approach

hypothesis ↓

AI model (program)

deduction

Data are used to develop the model

Note: UQ combines both approaches

# What is a good model?

- ## Fundamentally: Bayesian inference

$$\underline{P(H|E)} \propto \underline{P(E|H)} \cdot \underline{P(H)}$$
Posterior Likelihood Prior

$H$: Hypothesis
$E$: Evidence

"curve fitting" example

Try to find the prior (hypothesis) that provides the best $P(H|E)$?

$H: y = ax + b$

$E$

$P(E|H)$ large

$H: y = ax + b$

$E$

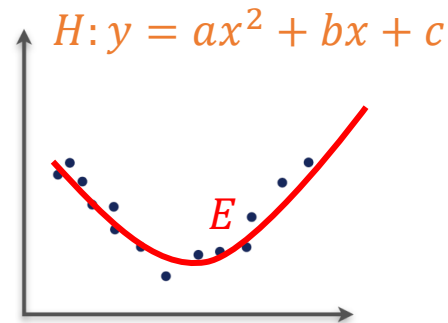$P(E|H)$ small

$H: y = ax^2 + bx + c$

$E$

Good prior is key

My prior hypothesis is supported by my evidence.
→Confidence in prior.

My prior hypothesis is not supported by my evidence.
→Need to review my prior

Prior ⟹ Likelihood
⟸ Posterior

[Dehaene, S. (2020). How We Learn: Why Brains Learn Better Than Any Machine... for Now. Penguin.]

# Current Data Science Landscape



2020: ML

Data Science

Data mining

Big
data

Statistical Analysis: Estimators, correlations,
unsupervised clustering

Machine Learning

Decision tree, artificial neural network, support
vector machines, supervised clustering,
reinforcement learning, …

Deep Learning

2010: Modern ML
Disciplines

# Machine learning in a nutshell: terminology

$u$: Inputs → ML model $\mathcal{F}(u; \phi)$

Prediction: $\hat{y} = \mathcal{F}(u; \phi)$

$y$: Target data

$\mathcal{L}$

$\mathcal{L}$: Loss function

Training/Learning
$$\frac{\partial \mathcal{L}}{\partial \phi}$$

How to pick the right ML approach?
Categorization of ML model based on the link of the *output* of the ML model with the target data/ loss function $\mathcal{L}$

# Structure of the Lecture

- Introduction to Machine Learning (ML)
  1. Drivers behind current ML
  2. Position of ML in science
  3. **Classification of ML methods**

- Introduction to Neural Network
  1. Linear regression and computational graph
  2. Gradient descent
  3. The neuron

# ML can be categorized on the role to perform

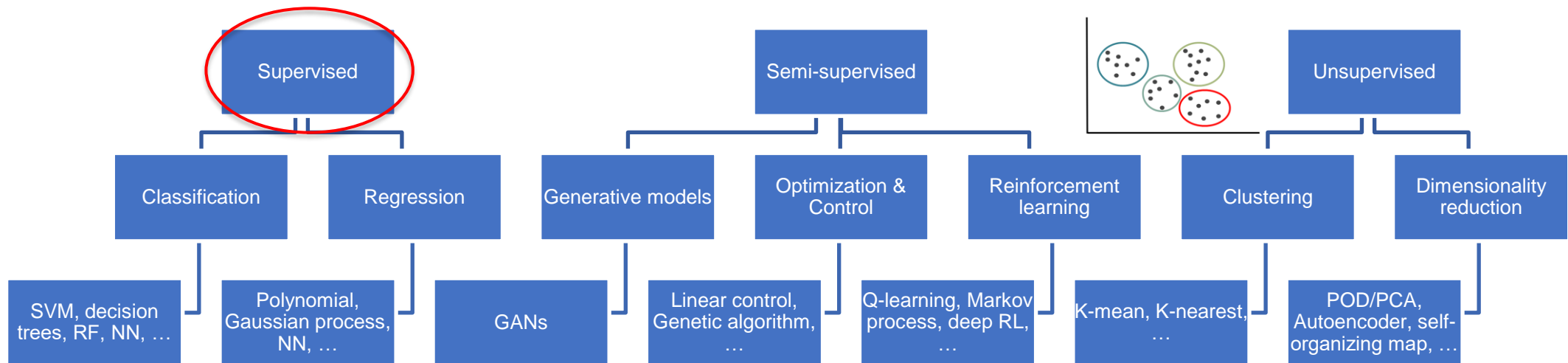Available labelled (input,target) data

Loss function directly related to output of ML model:
*Getting as close to the target data is the objective*

Partially labelled data

Loss function indirectly related to output of ML model

Labelled (input,target) data not available

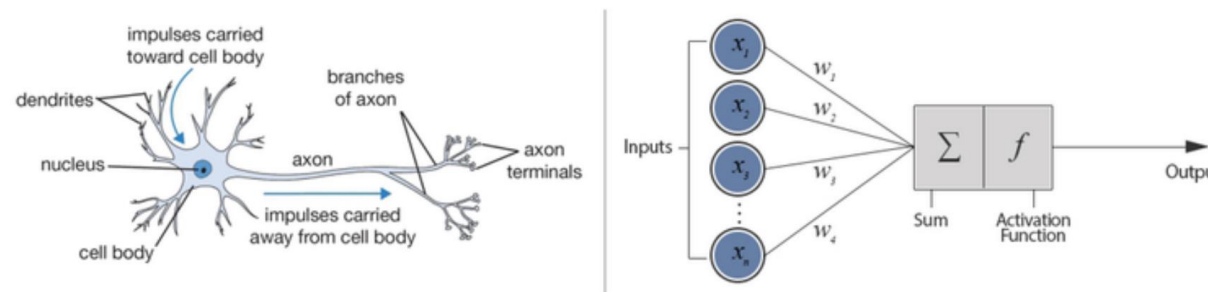Loss function not related to the role of ML model as no data available

# Structure of the Lecture

- Introduction to Machine Learning (ML)
    1. Drivers behind current ML
    2. Position of ML in science
    3. Classification of ML methods

- **Introduction to Neural Network**
    1. Linear regression and computational graph
    2. Gradient descent
    3. The neuron

# Neural network can be used as a basis for "any model"
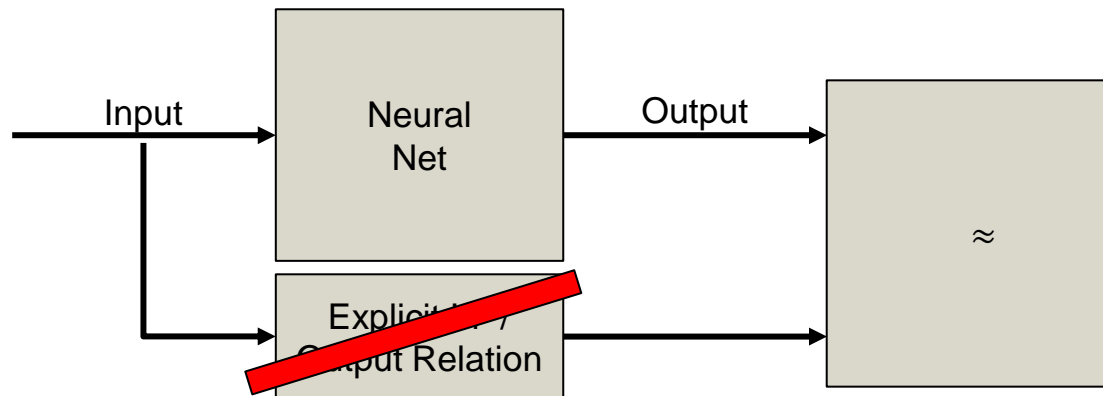
- ## Neural networks
  - Inspired to reproduce the biological process of learning (late 1940s/early 1950s)



  - Universal approximator
    - *"Can approximate any kind of nonlinear function"*
  - Very strong expressivity
    - *"Can represent a large variety of functions"*

[Hornik, Neural Networks, 1991.]
[Anthony & Barnett, Neural Network Learning: Theoretical Foundations, CUP, 2009]
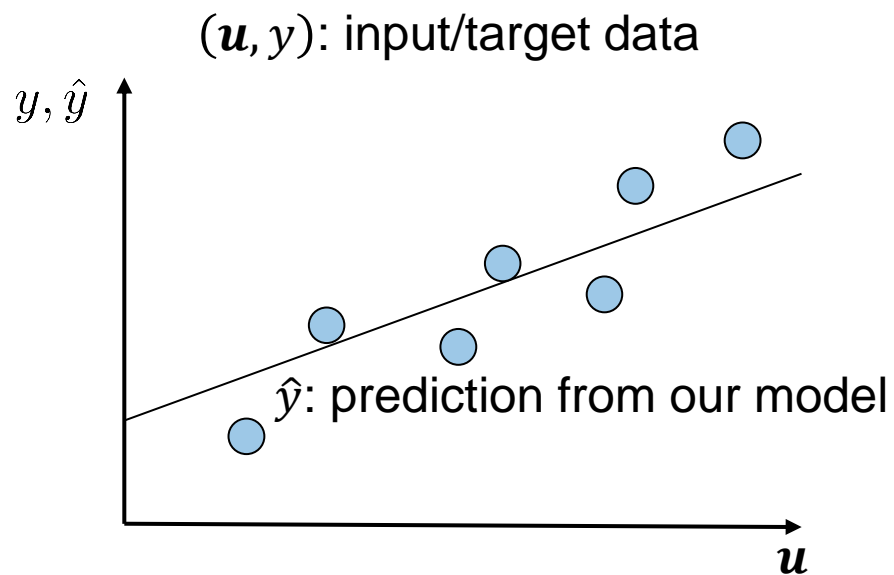
# Objective of AI/ML is to obtain a model

- Strength of neural networks: Universal approximators
  - Neural networks can approximate "any function" given a large enough number of neurons
  - BUT: no means of knowing beforehand what kind of network to use for that nor the appropriate weights
  - Neural networks are (generally) trained on input/target data
    - Approximate the underlying function existing in the data

# Structure of the Lecture

- Introduction to Machine Learning (ML)
  1. Small history of AI
  2. What is AI (and ML)?
  3. Position of ML in science
  4. Classification of ML methods

- **Introduction to Neural Network**
  1. **Linear regression and computational graph**
  2. **Gradient descent**
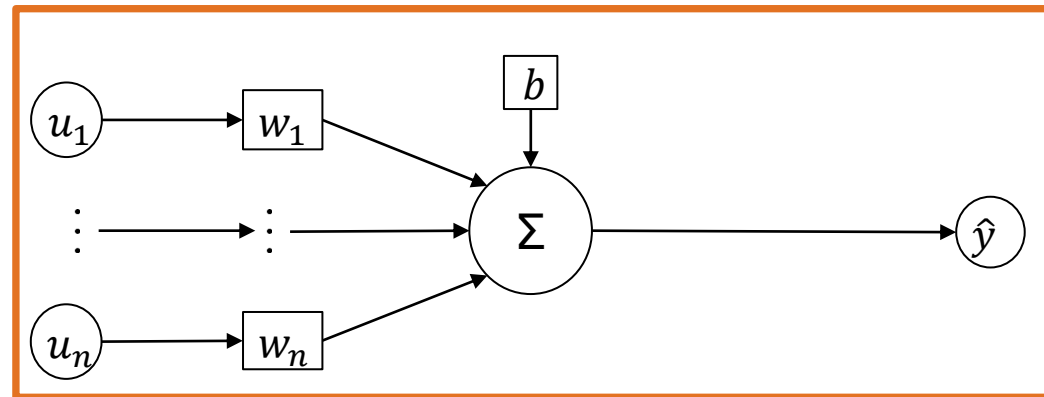  3. The neuron

# Linear regression

$(\boldsymbol{u}, y)$: input/target data

$y, \hat{y}$

$\hat{y}$: prediction from our model

$\boldsymbol{u}$

Simplest regression:
$$y \approx \hat{y} = f(\boldsymbol{u}; \boldsymbol{w}, b) = \boldsymbol{u} \cdot \boldsymbol{w} + b$$

$$= \sum_i u_i w_i + b$$

$\boldsymbol{w}$: weights of the model

$b$: bias of the model

# Computational Graph: some terminology
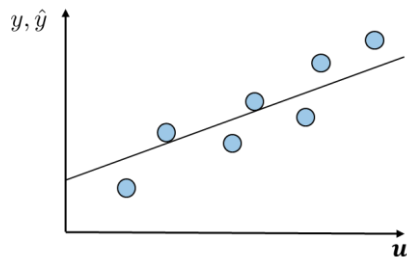
**Forward pass**

0.1
0.55
0.5
$b$

$u_1$ → $w_1$
0.055

1
0.7
$u_2$ → $w_2$
0.7

$\Sigma$ → $\hat{y}$ 1.855

0.3
2
$u_n$ → $w_n$
0.6

→ Information flow

**Loss function**

$y, \hat{y}$

0.1
0.55
0.5
$b$

$u_1$ → $w_1$

1
0.7
$u_2$ → $w_2$

$\Sigma$ → $\hat{y}$ 1.855 ↔ $y$ 2   Training data

0.3
2
$u_n$ → $w_n$

$$L(\hat{y}, y) = \frac{1}{2}\left\|y - \hat{y}\right\|^2$$

0.01051125

# Training of the neuron/graph

- "Training a model" is solving an optimization problem:



$$\underset{\boldsymbol{w},b}{\mathrm{argmin}}\, L(\hat{y}, y)$$

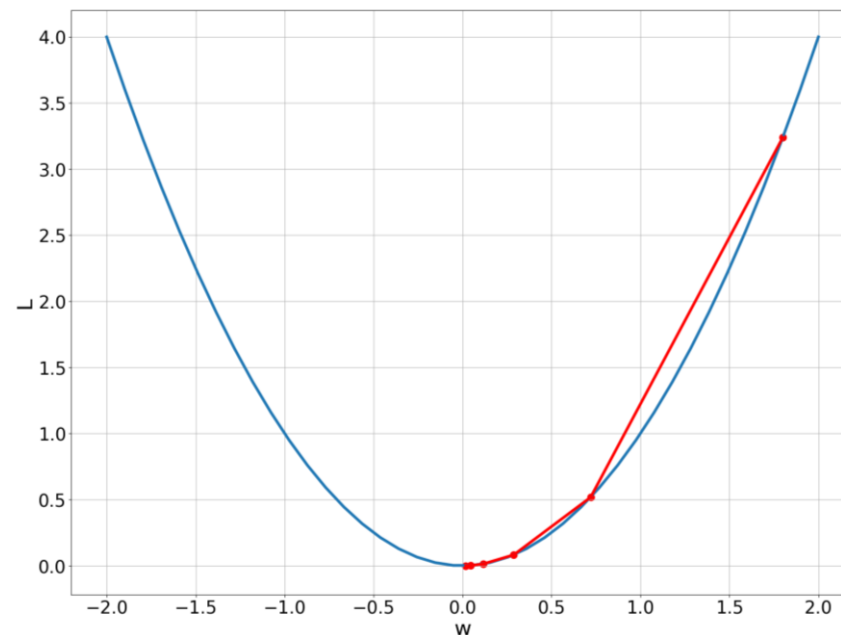$$\text{Subject to } \hat{y} = \boldsymbol{w} \cdot \boldsymbol{u} + b$$

→ Many ways to solve this.

→ Focus on numerical gradient-based iterative optimizer

# Gradient descent

- Minimization of $L$ via gradient descent:
- Computation of $-\nabla L$ (steepest descent)
- Update weights $\boldsymbol{w}$ in that direction by factor $\alpha$ ("learning rate")
- Stop when optimization criteria are met ($N$ iterations or threshold $L < \epsilon$)

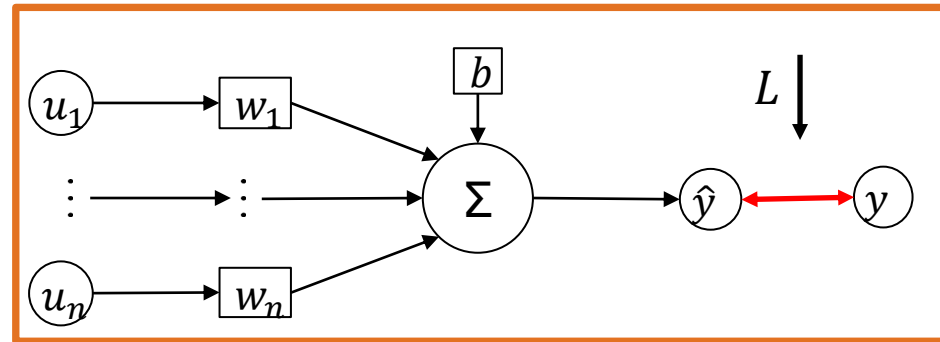$$\nabla L = \begin{pmatrix} \dfrac{\partial L}{\partial w_1} \\ \cdots \\ \dfrac{\partial L}{\partial w_n} \end{pmatrix}_{\boldsymbol{w}} \qquad \boldsymbol{w}_{new} = \boldsymbol{w}_{old} - \alpha \nabla L$$

# Gradient Descent and Computational Graph

- To solve our optimization problem, we need the gradient of
$L = \frac{1}{2}(\hat{y} - y)^2$



With the chain rule:

$$\frac{\partial L}{\partial w_i} \approx \frac{\delta L}{\delta w_i} = \frac{dL}{d\hat{y}} \cdot \frac{\delta \hat{y}}{\delta w_i} \qquad \frac{\partial L}{\partial b} \approx \frac{\delta L}{\delta b} = \frac{dL}{d\hat{y}} \cdot \frac{\delta \hat{y}}{\delta b}$$

$$\nabla L = \begin{pmatrix} \dfrac{\partial L}{\partial w_1} \\ \vdots \\ \dfrac{\partial L}{\partial w_n} \\ \dfrac{\partial L}{\partial b} \end{pmatrix}_{\boldsymbol{w,b,y}}$$
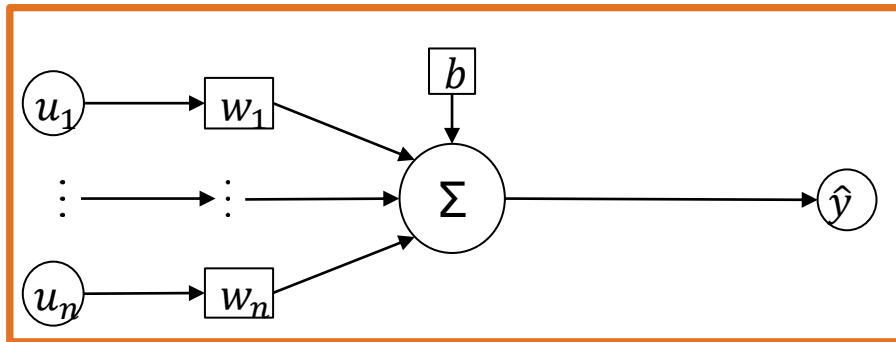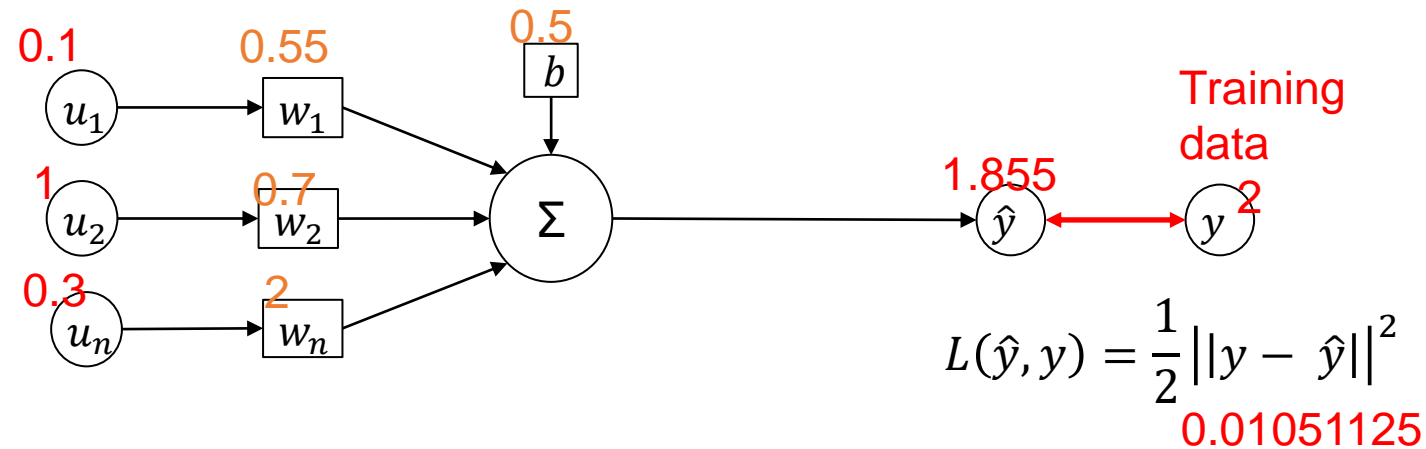
# Gradient Descent and Computational Graph

$$\nabla L = \begin{pmatrix} \dfrac{\partial L}{\partial w_1} \\ \vdots \\ \dfrac{\partial L}{\partial w_n} \\ \dfrac{\partial L}{\partial b} \end{pmatrix}_{\boldsymbol{w},b,\boldsymbol{y}}$$

$$\frac{\partial L}{\partial w_i} \approx \frac{\delta L}{\delta w_i} = \textcolor{red}{\frac{dL}{d\hat{y}}} \cdot \frac{\delta \hat{y}}{\delta w_i} \qquad\qquad \frac{\partial L}{\partial b} \approx \frac{\delta L}{\delta b} = \textcolor{red}{\frac{dL}{d\hat{y}}} \cdot \frac{\delta \hat{y}}{\delta b}$$

$$\textcolor{red}{L = \frac{1}{2}(\hat{y} - y)^2 \rightarrow \frac{dL}{d\hat{y}} = \hat{y} - y = \Delta y}$$

# Gradient Descent and Computational Graph

$$\nabla L = \begin{pmatrix} \dfrac{\partial L}{\partial w_1} \\ \vdots \\ \dfrac{\partial L}{\partial w_n} \\ \dfrac{\partial L}{\partial b} \end{pmatrix}_{\boldsymbol{w},b,\boldsymbol{y}}$$

$$\frac{\partial L}{\partial w_i} \approx \frac{\delta L}{\delta w_i} = \frac{dL}{d\hat{y}} \cdot \frac{\delta \hat{y}}{\delta w_i}$$

$$\frac{\partial L}{\partial b} \approx \frac{\delta L}{\delta b} = \frac{dL}{d\hat{y}} \cdot \frac{\delta \hat{y}}{\delta b}$$

$$\hat{\boldsymbol{y}} = f(u; \boldsymbol{w}, b) = \boldsymbol{u} \cdot \boldsymbol{w} + b = \sum_i u_i w_i + b$$

$$\frac{\delta \hat{y}}{\delta w_i} = \frac{\partial}{\partial w_i}\left(\sum_i u_i w_i + b\right) = u_i \qquad \frac{\delta \hat{y}}{\delta b} = \frac{\partial}{\partial b}\left(\sum_i u_i w_i + b\right) = 1$$



$$\nabla L = \begin{pmatrix} \dfrac{\partial L}{\partial w_1} \\ \vdots \\ \dfrac{\partial L}{\partial w_n} \\ \dfrac{\partial L}{\partial b} \end{pmatrix}_{\boldsymbol{w},b,\boldsymbol{y}} = \begin{pmatrix} \Delta y\, u_i \\ \vdots \\ \Delta y\, u_n \\ \Delta y \end{pmatrix}_{\boldsymbol{w},b,\boldsymbol{y}}$$

# Gradient Descent and Computational Graph

0.1   0.55   0.5

$u_1$ → $w_1$

1   0.7

$u_2$ → $w_2$

0.3   2

$u_n$ → $w_n$

$b$

$\Sigma$

1.855   Training data   2

$\hat{y}$ ↔ $y$

$$L(\hat{y}, y) = \frac{1}{2}\left|\left|y - \hat{y}\right|\right|^2$$
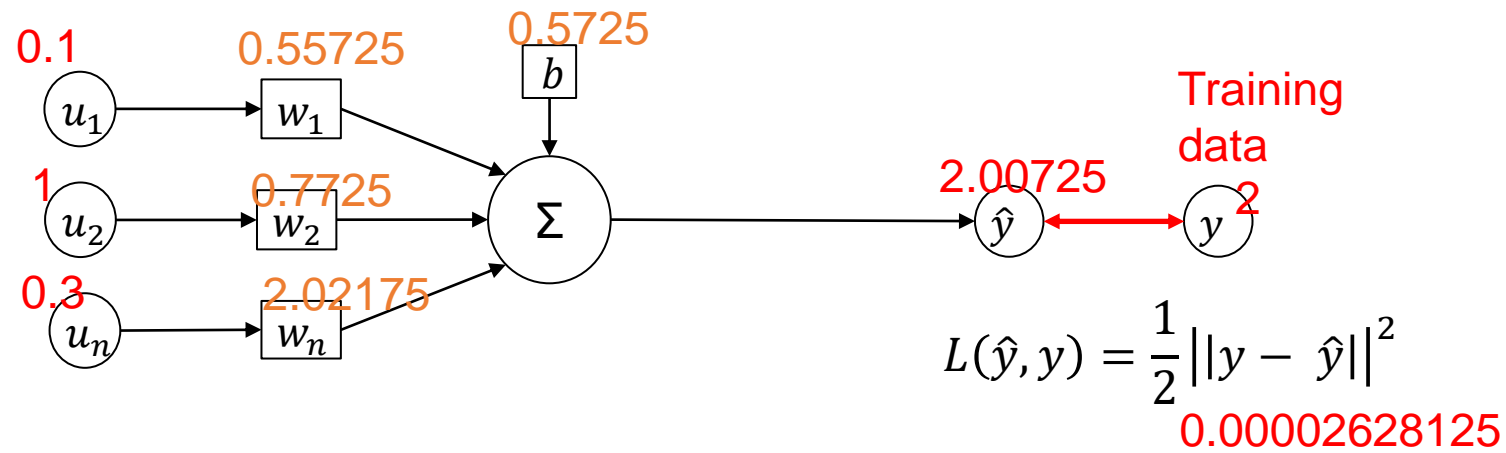
0.01051125

$$\nabla L = \begin{pmatrix} \Delta y\, u_i \\ \vdots \\ \Delta y\, u_n \\ 1 \end{pmatrix}_{\boldsymbol{w,b,y}} = -\begin{pmatrix} 0.0145 \\ 0.145 \\ 0.0435 \\ 0.145 \end{pmatrix}$$

$$\boldsymbol{w}_{new} = \boldsymbol{w}_{old} - \alpha \nabla L$$

$$\begin{pmatrix} 0.55 \\ 0.7 \\ 2 \\ 0.5 \end{pmatrix} + (0.5)\begin{pmatrix} 0.0145 \\ 0.145 \\ 0.0435 \\ 0.145 \end{pmatrix} = \begin{pmatrix} 0.55725 \\ 0.7725 \\ 2.02175 \\ 0.5725 \end{pmatrix}$$

# Gradient Descent and Computational Graph

0.1
0.55725
0.5725
$u_1$ → $w_1$ → $b$

1
0.7725
$u_2$ → $w_2$ → $\Sigma$

0.3
2.02175
$u_n$ → $w_n$

2.00725
$\hat{y}$ ↔ $y$ — Training data — 2

$$L(\hat{y}, y) = \frac{1}{2} \left|\left| y - \hat{y} \right|\right|^2$$

0.00002628125

# Limitation of gradient descent

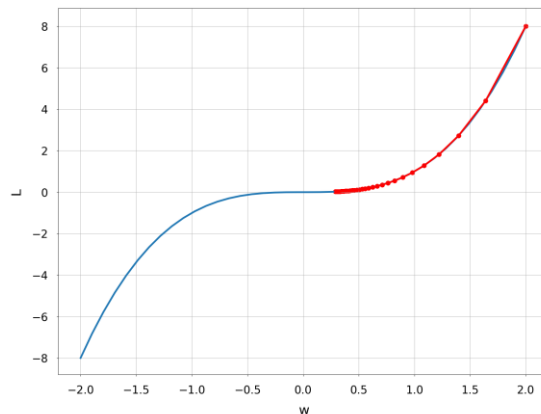- If the loss function is not well-behaved, the gradient descent may not converge appropriately:
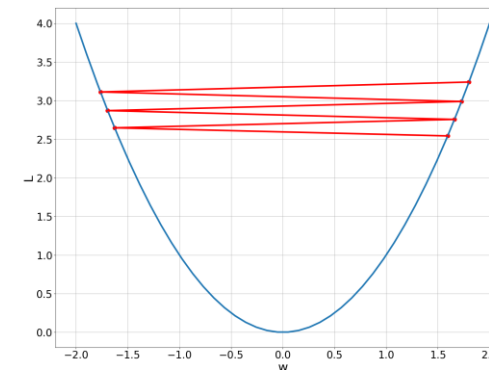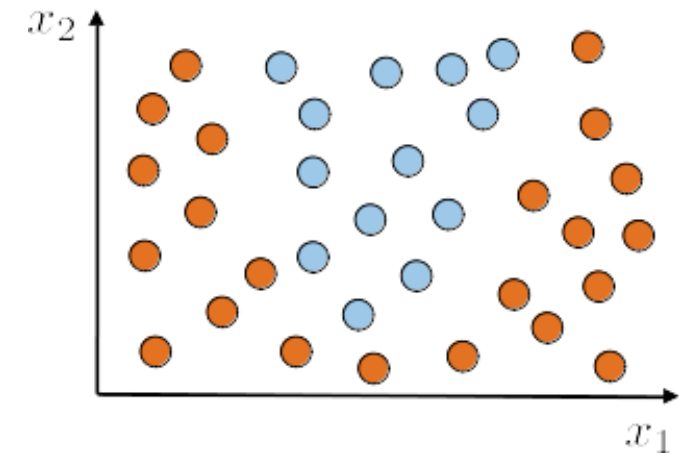
Local minimum
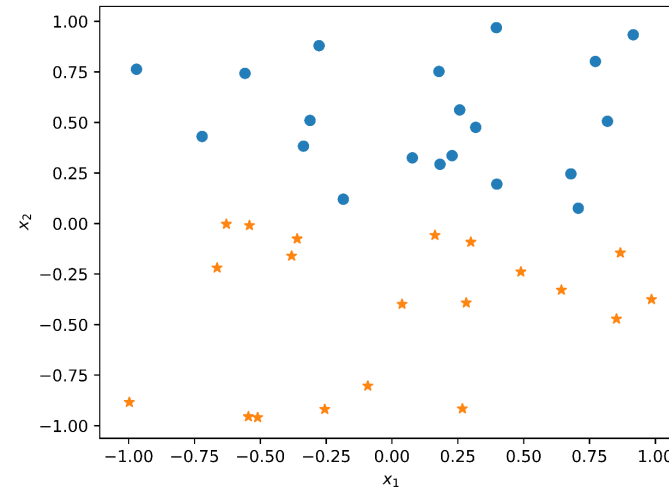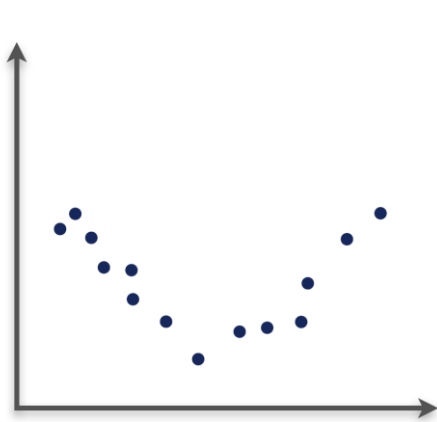


Jumping out of minimum



Vanishing gradient



Oscillating

# Structure of the Lecture

- **Introduction to Machine Learning (ML)**
  1. Small history of AI
  2. What is AI (and ML)?
  3. Position of ML in science
  4. Classification of ML methods

- **Introduction to Neural Network**
  1. Linear regression and computational graph
  2. Gradient descent
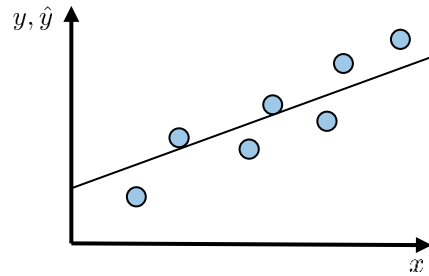  3. **Logistic regression and the neuron**

# More complex tasks cannot be achieved with just linear regressors

- Up to now: linear distribution – linear regression

- What about non-linear distribution? Or classification problem?



→ Need for nonlinear behaviour

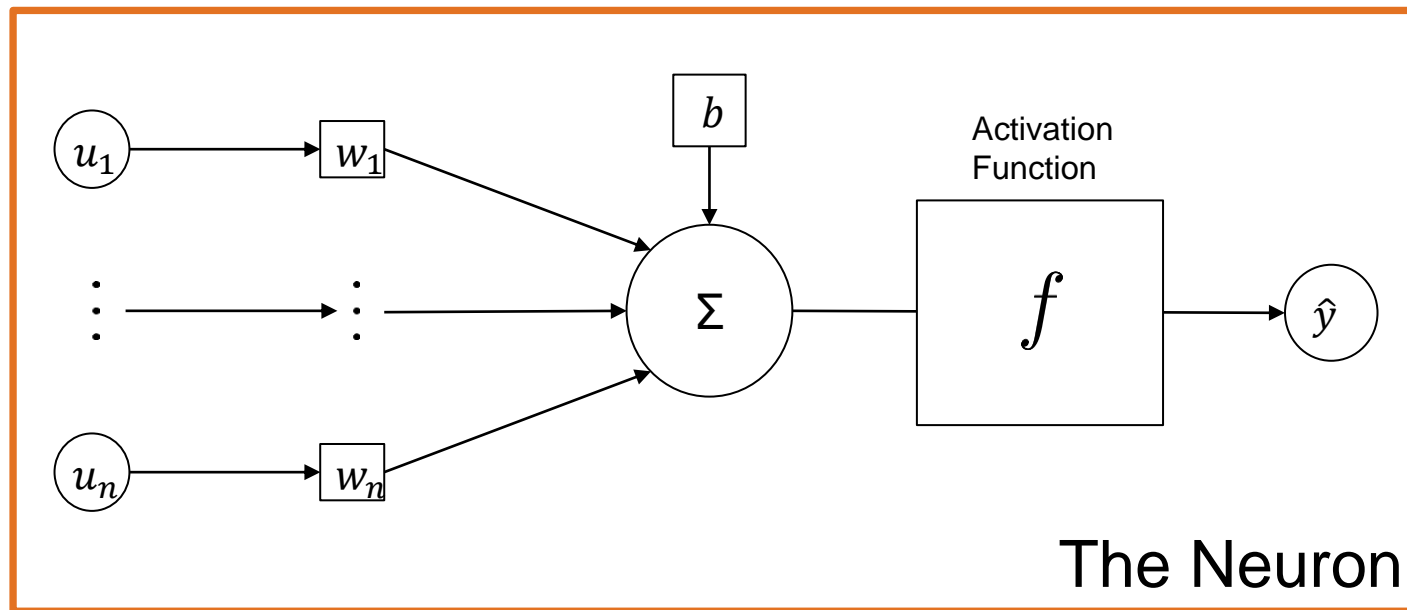# The neuron introduces nonlinearity through its activation function
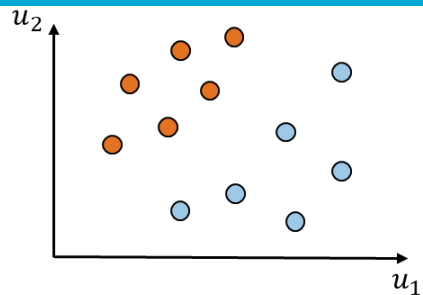


Standar linear regression:

$$\hat{y} = \sum_{i}^{n} w_i u_i + b_i = \boldsymbol{w} \cdot \boldsymbol{u} + \boldsymbol{b}$$

Nonlinearity introduced with activation function :
$$\hat{y} = f(\boldsymbol{w} \cdot \boldsymbol{u} + \boldsymbol{b})$$

The Neuron

# Classification requires step-like function

$u_2$

$u_1$

## Binary classification

Binary Step

$$f(x) = \begin{cases} 1, x > 0 \\ 0, x < 0 \end{cases}$$

## Linear regression

Identity

$$f(x) = x$$

$u_1$ $w_1$

$b$

$\Sigma$

Activation Function

$\hat{y}$

$u_2$ $w_2$

$u_2$

$u_1$

$u_2$

$y = 1$

$y = 0$

$u_1$

# Classification requires step-like function
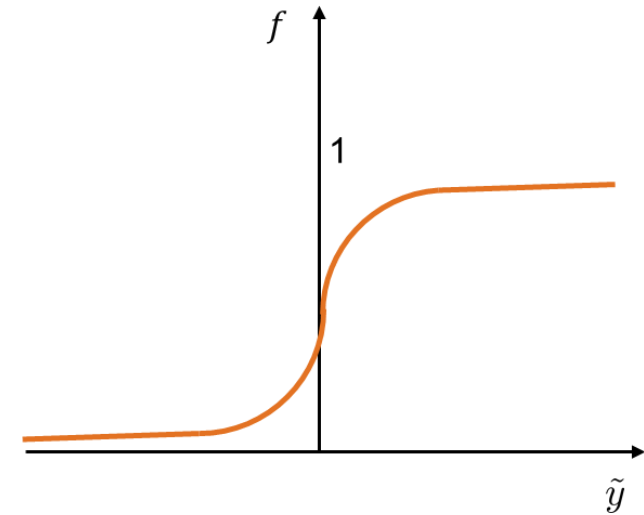


**Binary classification**

Binary Step

**Undefined derivative!!**

$$f(x) = \begin{cases} 1, x > 0 \\ 0, x < 0 \end{cases}$$

**Sigmoid Function**

$$f(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

$$f'(x) = f(1 - f)$$

# Logistic regression

- Used for binary classification (0/1)
- Consider $\boldsymbol{u} \in \mathbb{R}^n$ as input vector and $y$ as the target class
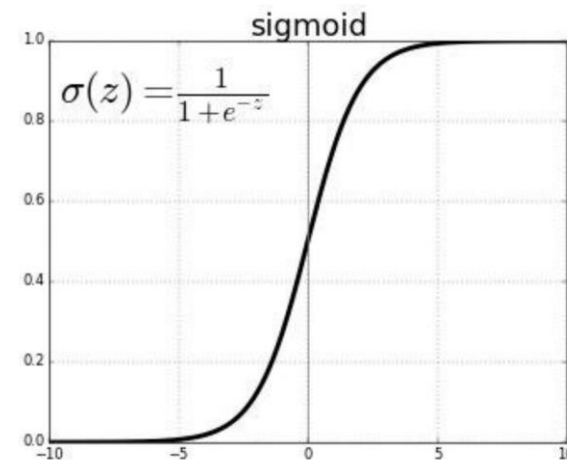
With parameters $\boldsymbol{w} \in \mathbb{R}^n$ and $b \in \mathbb{R}$

Linear output: $\hat{y} = \boldsymbol{u} \cdot \boldsymbol{w} + b$ → Cannot be used for binary classification

Logistic output: $\hat{a} = \sigma(\boldsymbol{u} \cdot \boldsymbol{w} + b)$

- Loss function: $\mathcal{L}(y, \hat{a}) = -(y \log(\hat{a}) + (1 - y) \log(1 - \hat{a}))$
- Cost function:

$$J(\boldsymbol{w}, b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(y^{(i)}, \hat{a}^{(i)})$$

sigmoid

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

# Note on the log-loss function

- Let's start again from a linear regression, but with some error

$$y \approx \boldsymbol{w}^T \boldsymbol{x} + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

- Assuming that $\boldsymbol{x}$ is also normally distributed, we can estimate likelihood

$$p(y|\boldsymbol{x}, \boldsymbol{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y - \boldsymbol{w}^T \boldsymbol{x})^2}{2\sigma^2}\right]$$

Log likelihood of the data:

$$\log \prod_i p(y^{(i)} | \boldsymbol{x}^{(i)}, \boldsymbol{w}) = \sum_i \left[-\frac{1}{2}\log(2\pi\sigma^2) - \frac{(y^{(i)} - \boldsymbol{w}^T x^{(i)})^2}{2\sigma^2}\right]$$

Maximising log-likelihood ↔ Minimising MSE

# Note on the log-loss function

- **Consider our logistic model. Given input $\boldsymbol{x}^{(i)}$, outputs $a_i$:**
  - probability $a\left(x^{(i)}\right)$ for class 1
  - probability $1 - a\left(x^{(i)}\right)$ for class 0
- **We can estimate the likelihood for the entire dataset:**

$$p(y|\boldsymbol{X}, \boldsymbol{w}) = \prod_i p\left(y^{(i)}|\boldsymbol{x}^{(i)}, \boldsymbol{w}\right) = \prod_i a_i^{y^{(i)}} (1 - a_i)^{1 - y^{(i)}}$$

- **And its negative log:**

$$\sum_i \left(-y^{(i)} \log a_i - \left(1 - y^{(i)}\right) \log(1 - a_i)\right)$$

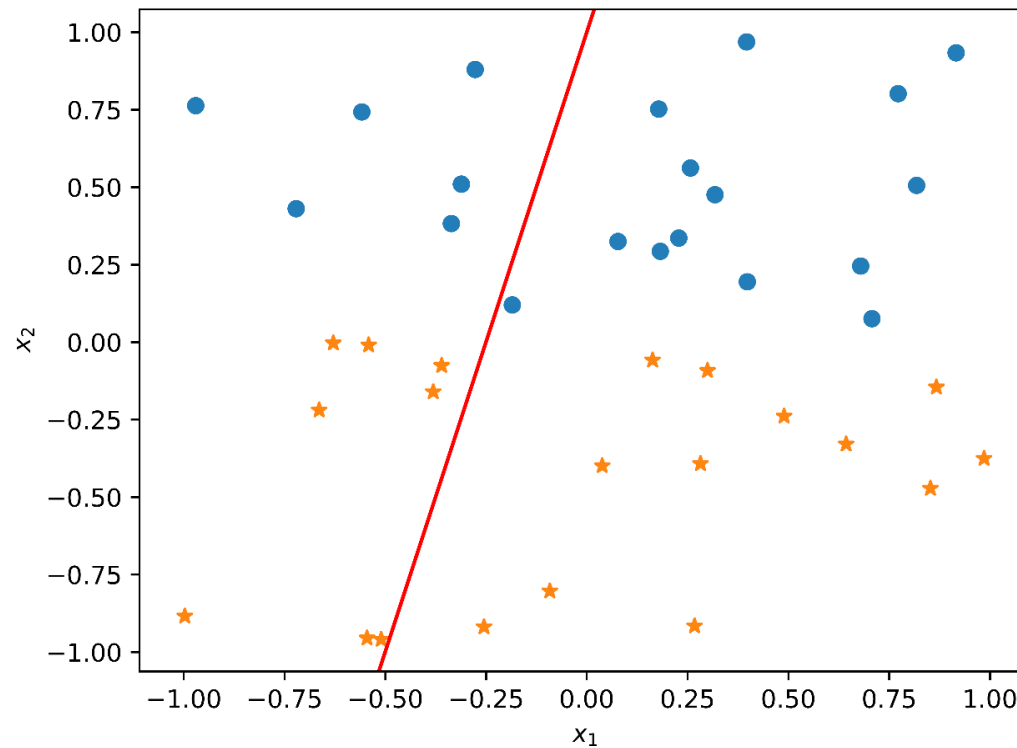$\rightarrow$ To be minimized

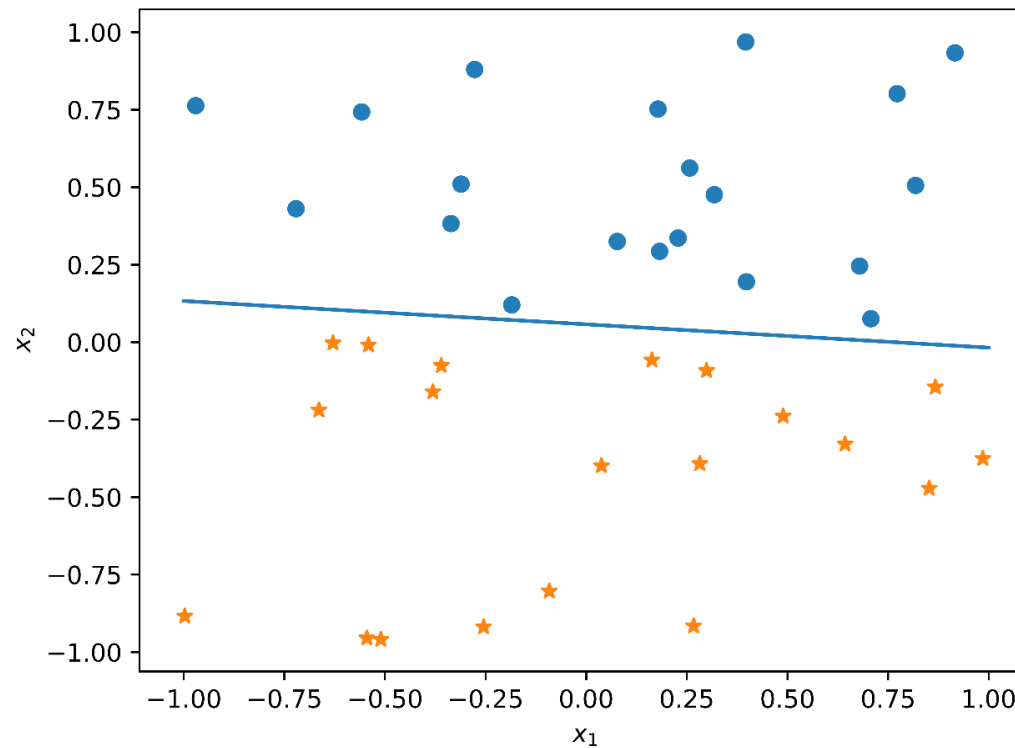# Classification: example

Initial data

# Classification: example

Initialization

# Classification: example

After training

# Summary – What we have seen so far

- Aim of ML: data-driven model obtention

- Place of ML and ML problem

- Classification of ML tool


- Linear regression and the neuron

- Logistic regression