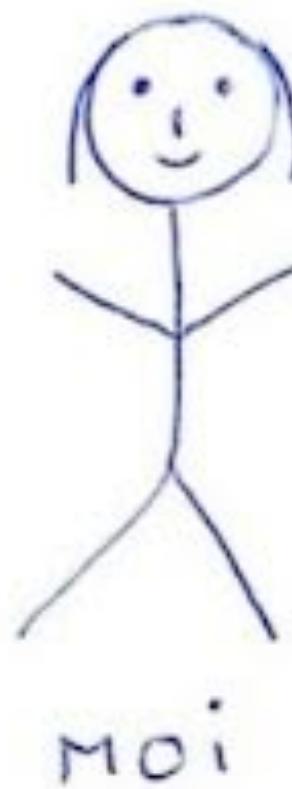


From Basic Machine Learning Models to Advanced Kernel Learning



Scott Pesme
Post-doc Inria Grenoble



Michaël Arbel
Researcher Inria Grenoble



Julien Mairal
Researcher Inria Grenoble

We work on the “theoretical aspects” of machine learning

Course structure



Scott Pesme
Post-doc Inria Grenoble

Basic Machine Learning Models



Michaël Arbel
Chercheur Inria Grenoble

Advanced Kernel Learning



Julien Mairal
Chercheur Inria Grenoble

2 x (~1h15 lecture) + 15 minutes break

Final grade: 50% exam + 50% homework

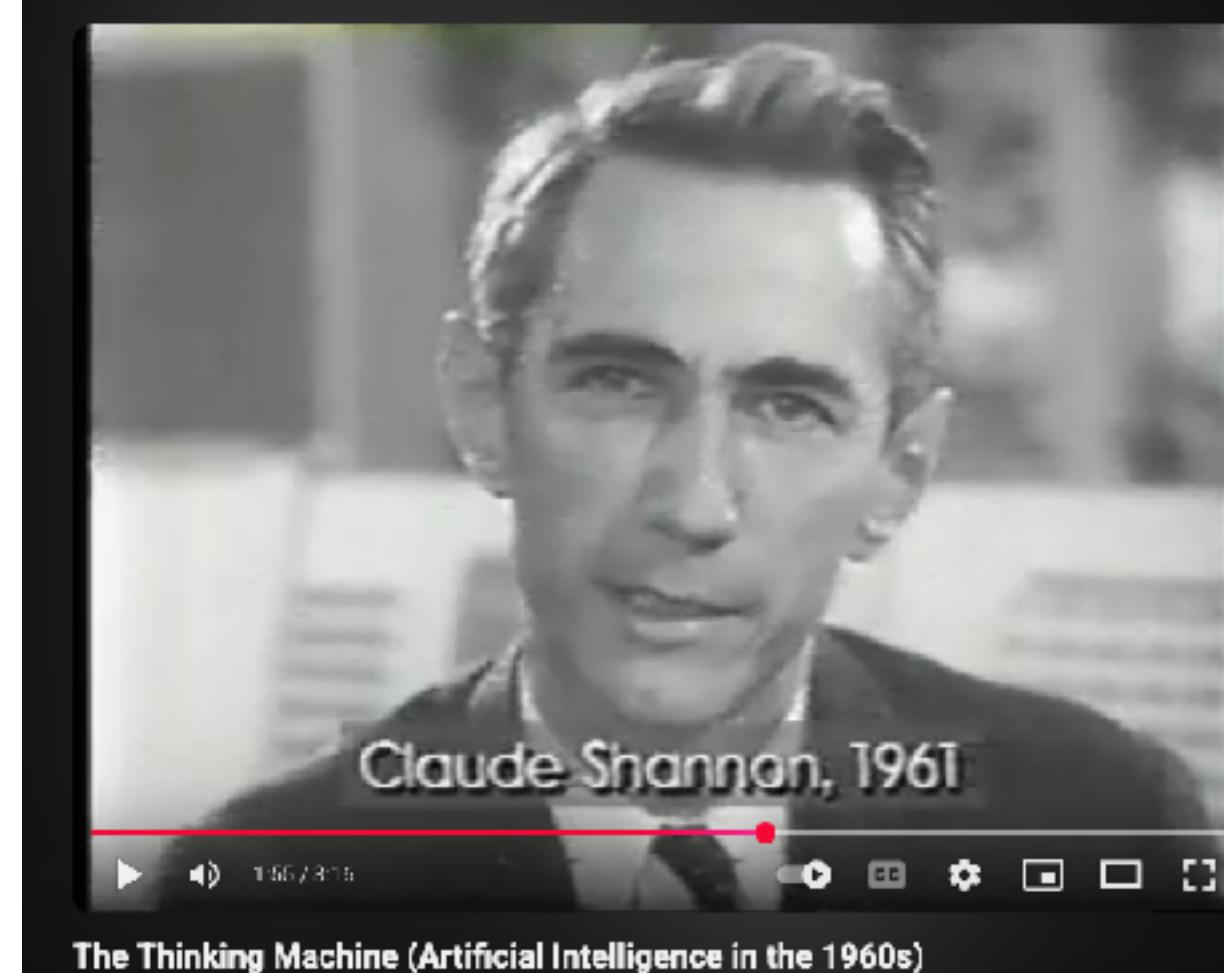
What to expect?

- Understanding of basic ML methods:
 - *Linear regression, linear classification, unsupervised learning, neural networks*
 - *Optimisation methods (gradient descent)*
- A deeper dive into kernel methods

You will mostly need pre-requisites from linear algebra!

Introduction

A bit of AI history: the idea of thinking machines is not new!

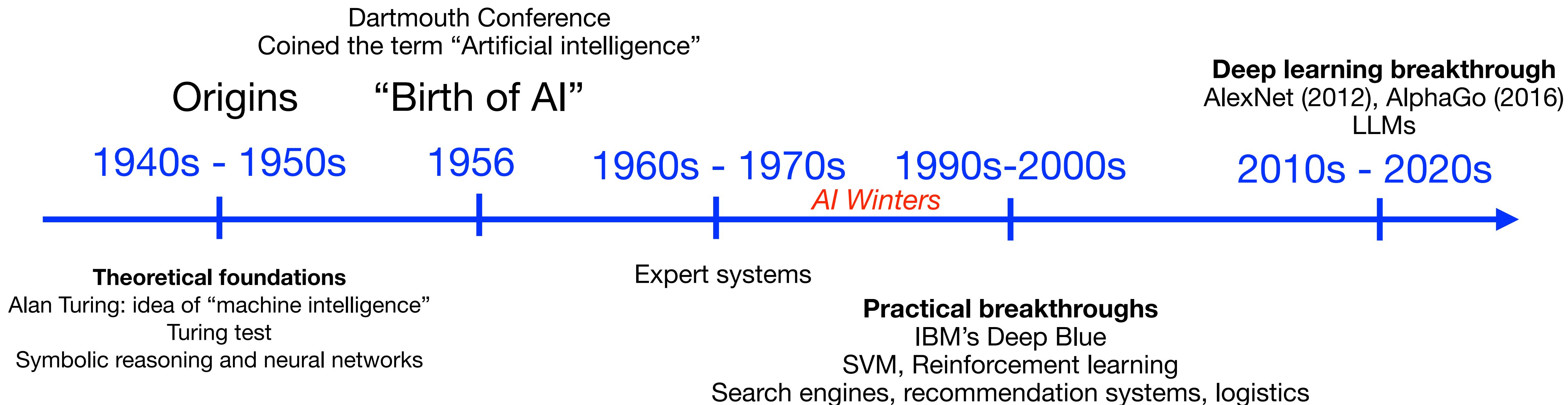
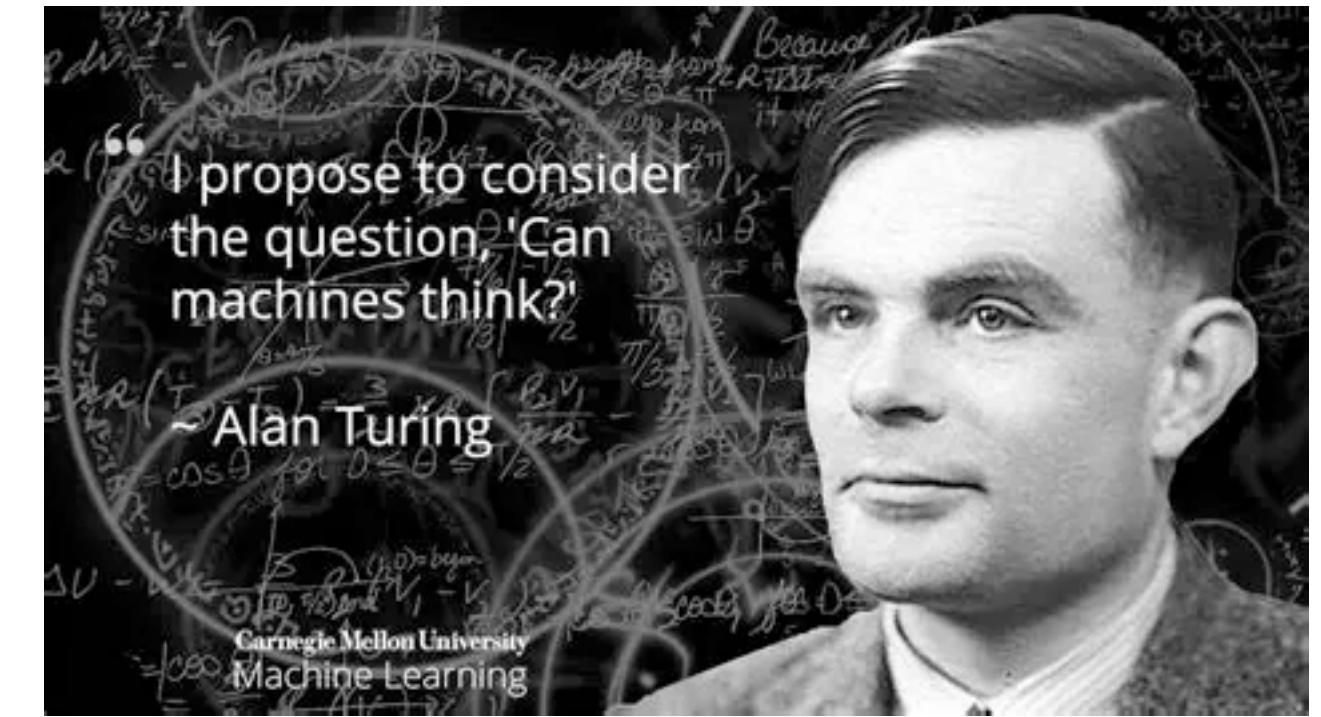


<https://www.youtube.com/watch?v=aygSMgK3BEM>

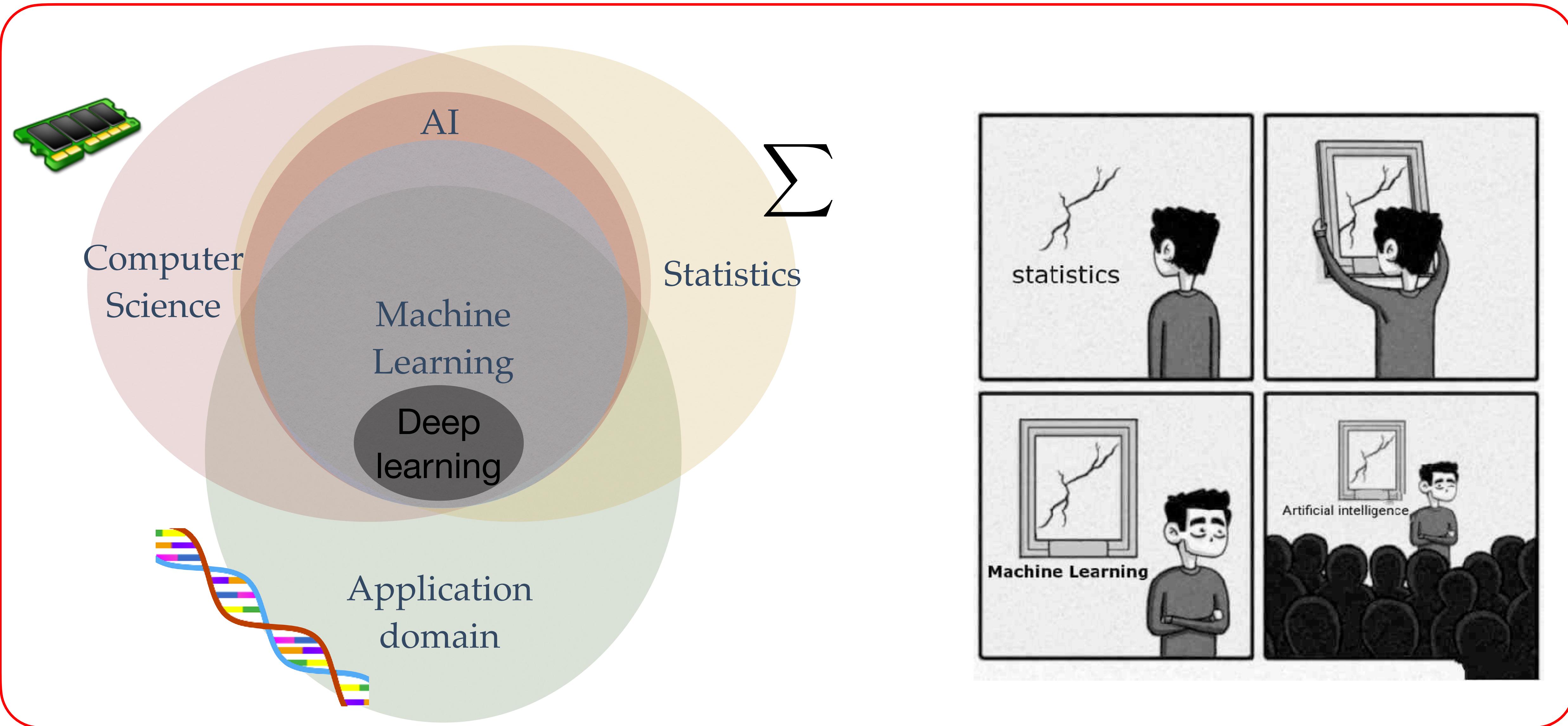
What's going to happen to us if machines can think? Can they?

*I don't really know, but you come back in 4 or 5 years I'll probably say
“Sure they can think”.*

A bit of AI history



Some vocabulary



Back to the introduction

What is Machine Learning?

Field which develops
“machines” / “algorithms” / “functions”
that learn from data

How does it work?

Training phase:

Neural network at school

Data



‘Dogkey’ ✗✓

After training:

Proficient neural network



‘Lagopède
Alpin’

It learns ‘on its own’ to extract the important features

Machine learning vs expert systems

After training:

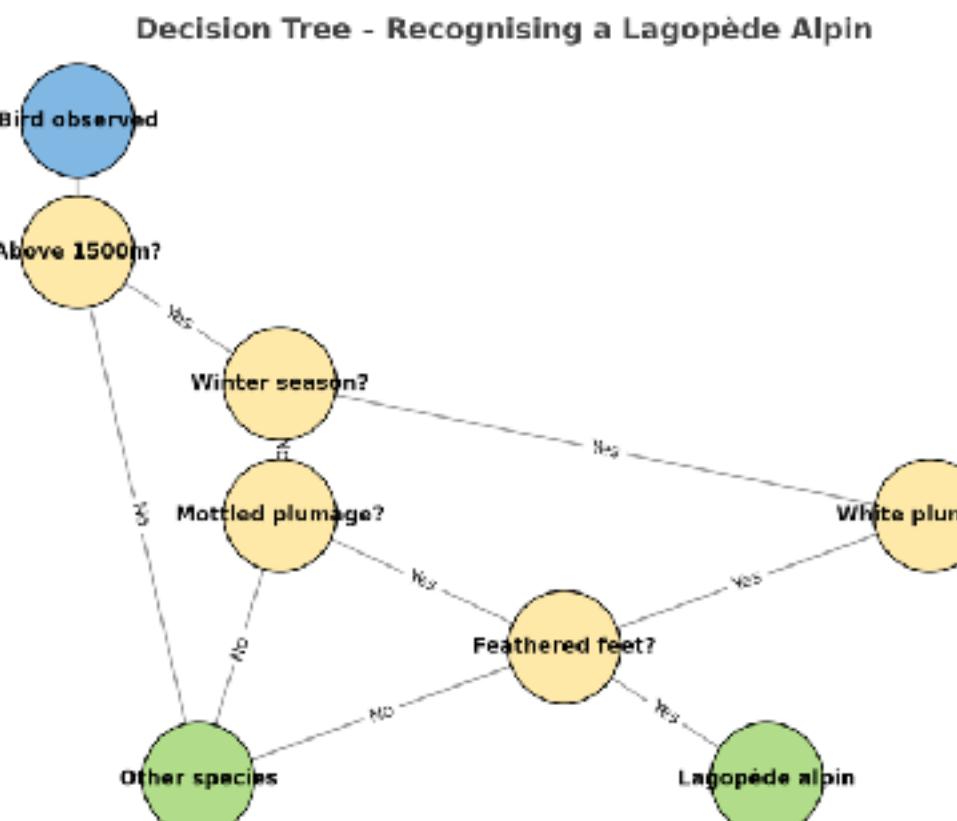
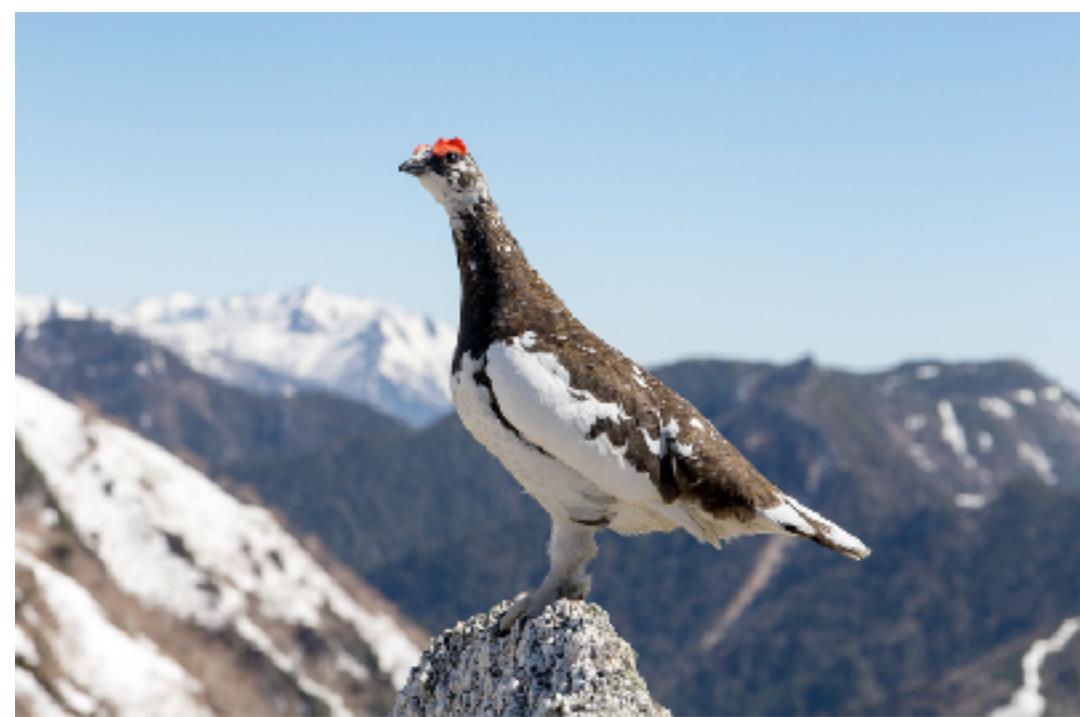


Proficient neural network



‘Lagopède
Alpin’

It learns ‘on its own’ to extract the important features



‘Lagopède
Alpin’

An “expert” gives a list of “if/then” rules

Pros and cons of two different “AI strategies”

Machine learning

Pros:

- No need to find the rules
- Adapts to data
- Works very well!!

Cons:

- Opacity (what the heck is going on??)
- Needs a “lot” of data

Expert systems

Pros:

- Transparent and understandable

Cons:

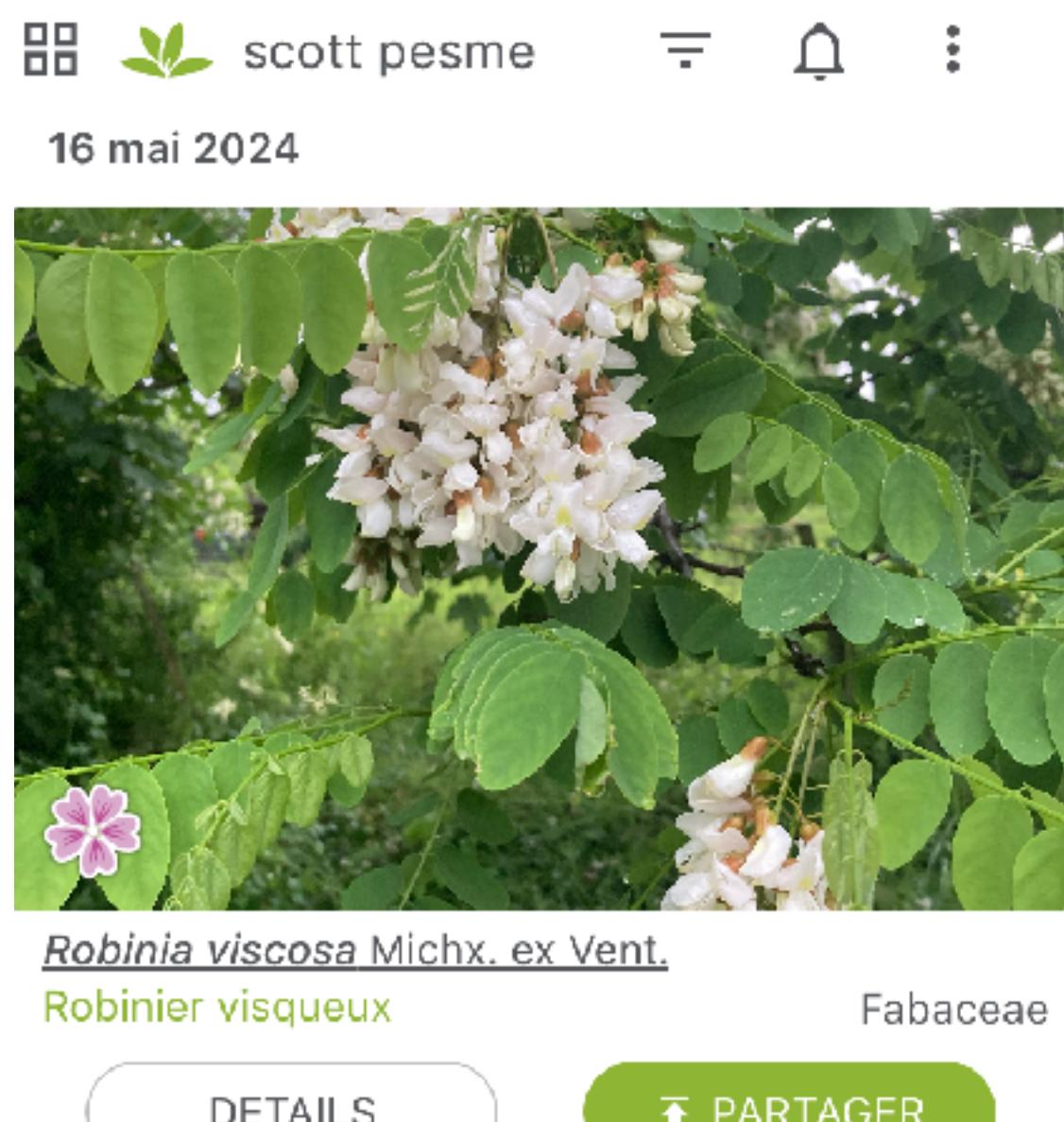
- Finding the rules is hard (impossible?)
- Doesn't work very well

Where is machine learning?

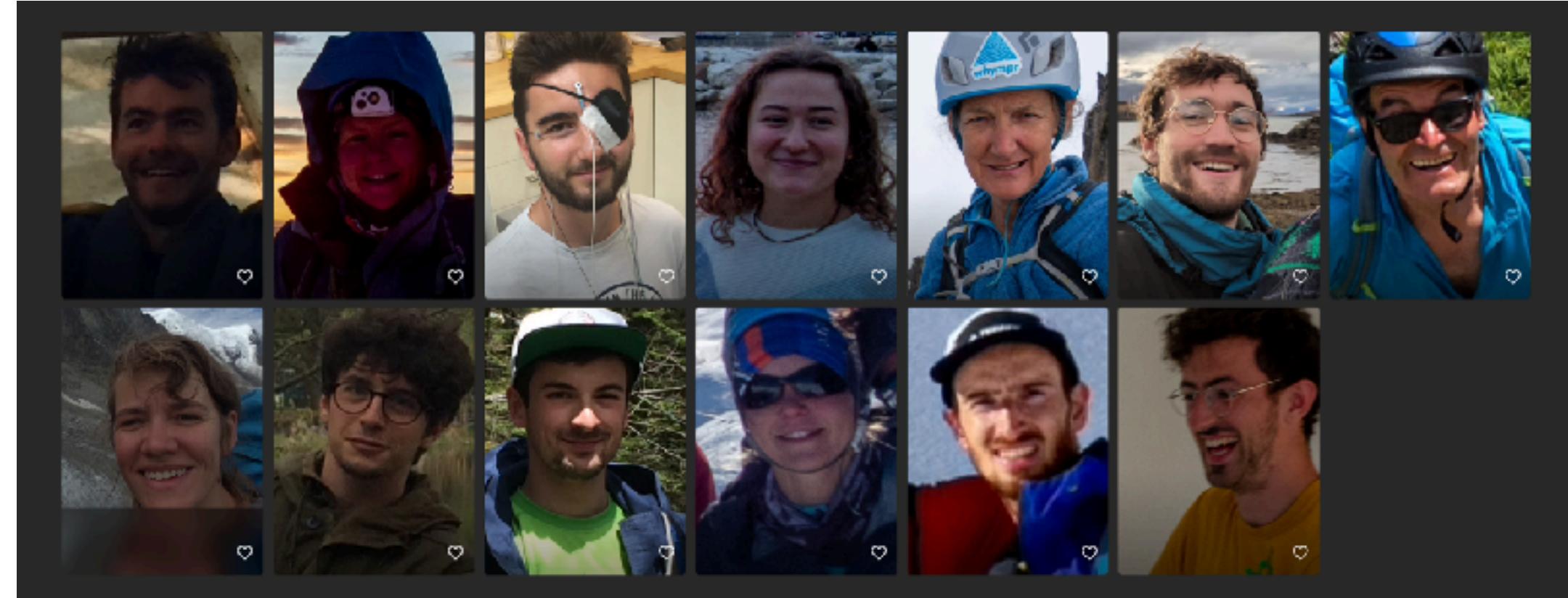
Text generation:



PlantNet:



Apple's photo app:



Drug discovery:

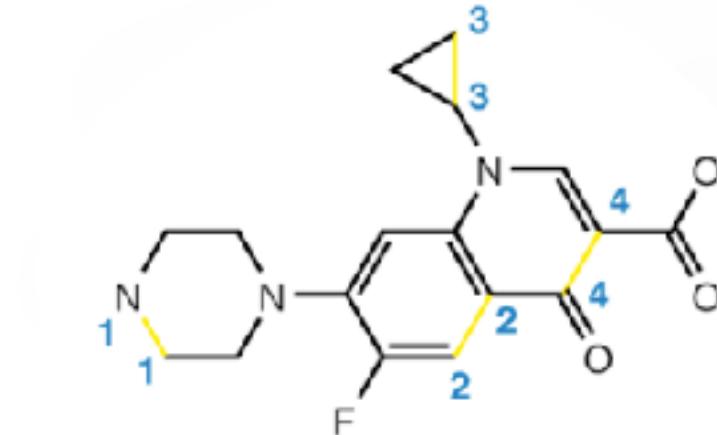


Image generation:



But also:

Deep fakes



Scams

Mass surveillance

Learning a mapping / function from input to output

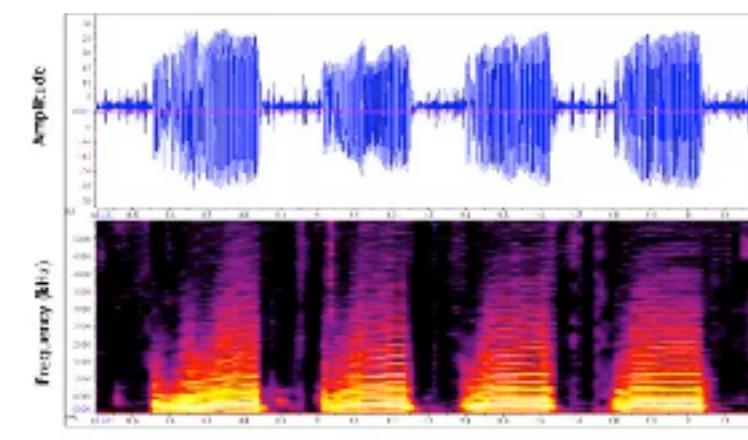
Input $x \in \mathbb{R}^d$

Text *What is the capital
of France?*

Image



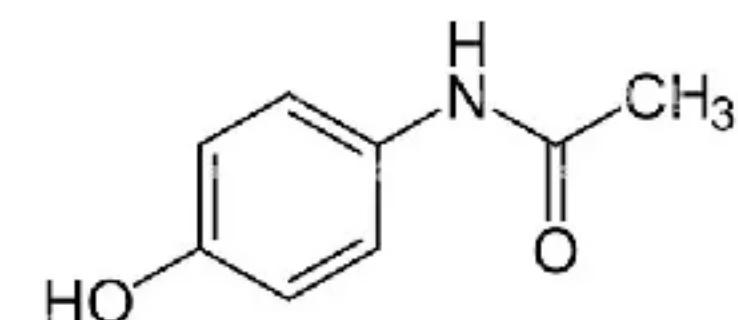
Audio



Video



Graph



“Feature” vector (Height, age, income)

Trainable
function

$$f_w(x) = \hat{y}$$



Output $y \in \mathbb{R}^p$

Text

Image

Audio

Video

Graph

“Feature” vector

How is the input represented “inside the machine”?

Input

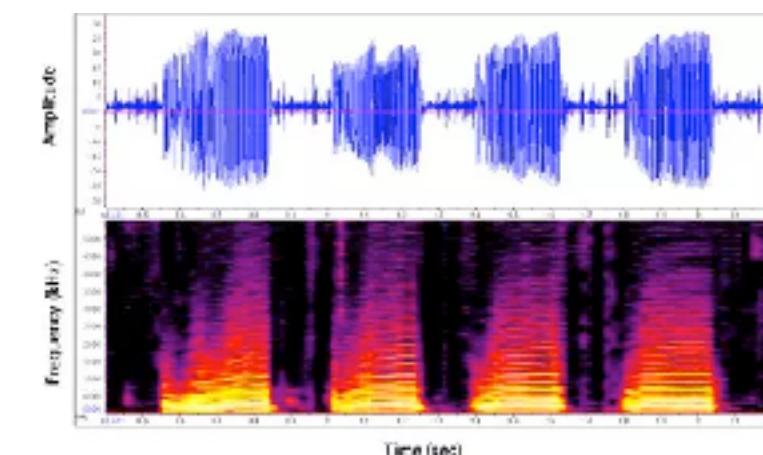
Text

*What is the capital
of France?*

Image



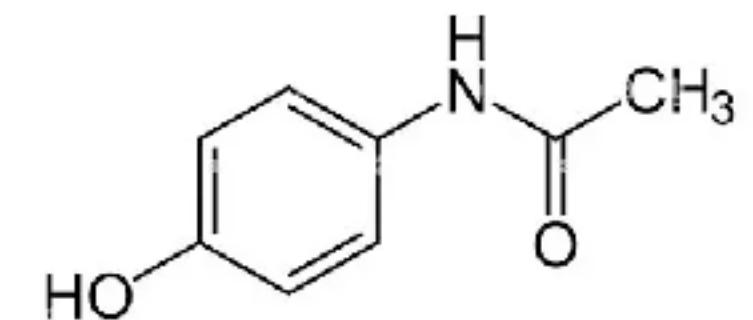
Audio



Video



Graph



“Feature” vector (Height, age, income)

“Digital” representation

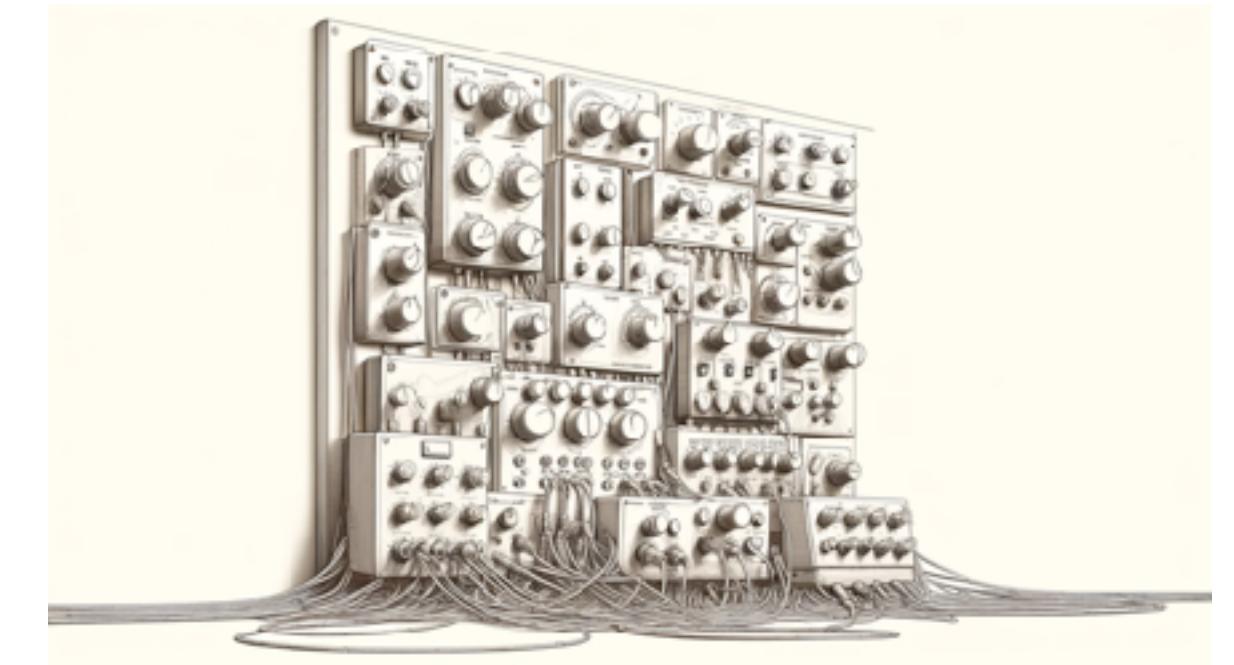
$x = ???$

How do we get a machine to learn?

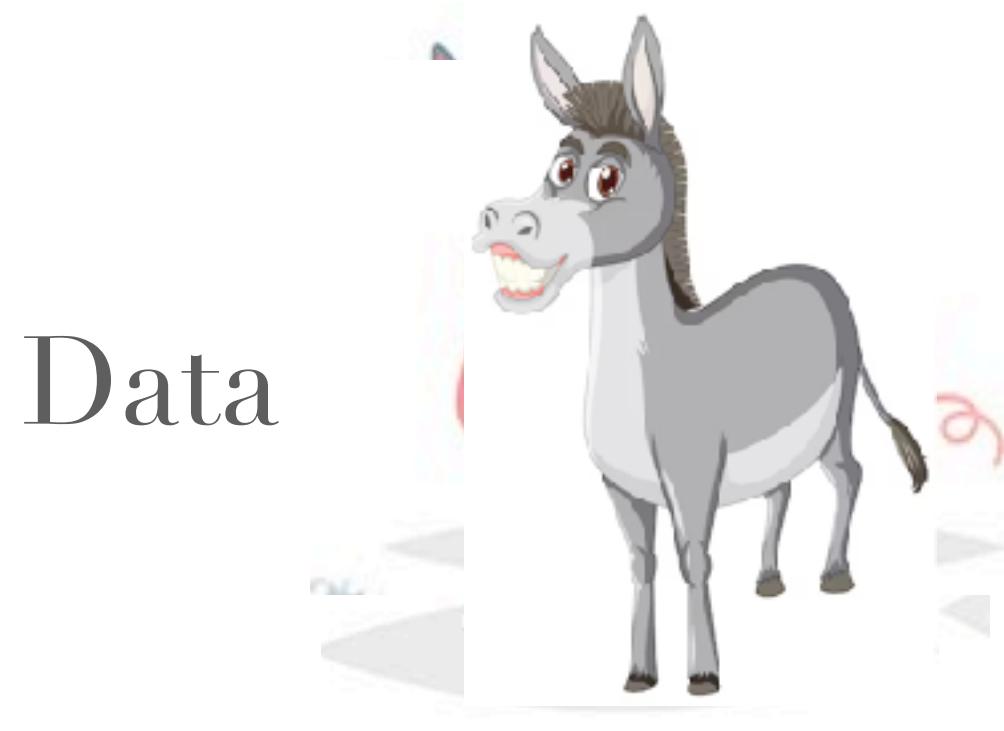
The neural network



looks rather like:



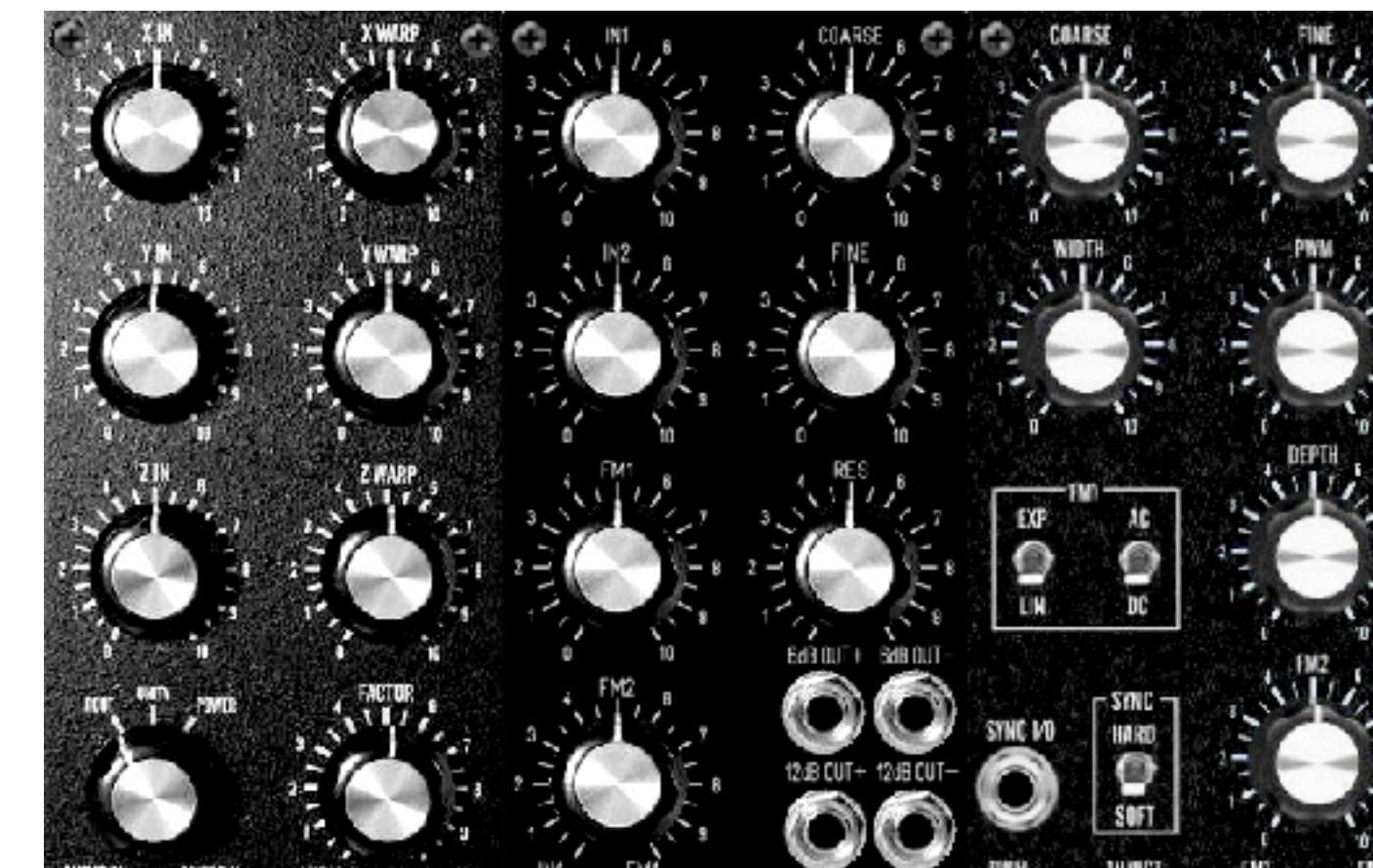
Training:



Data



Neural network at school

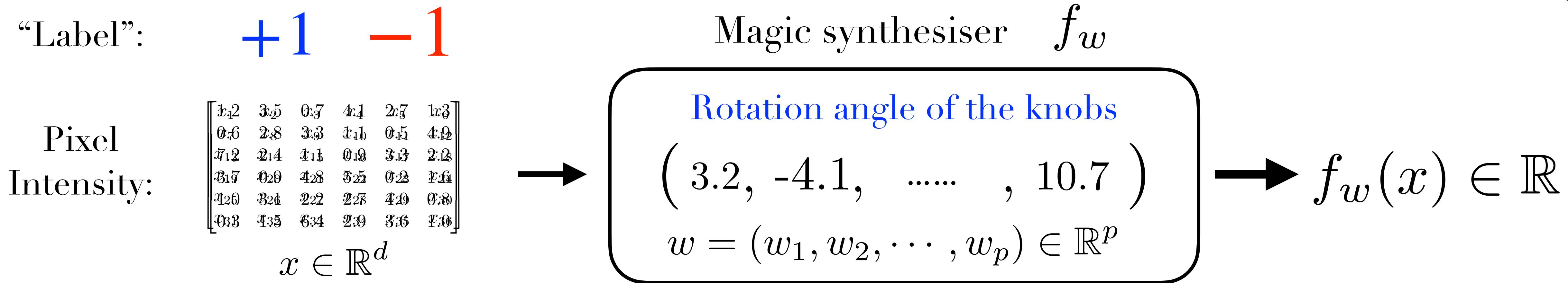
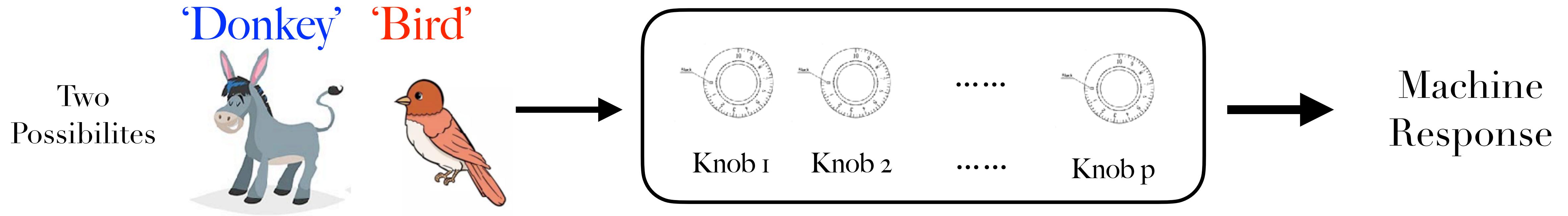


“*Vâilee*”



Training a neural network corresponds to turning knobs

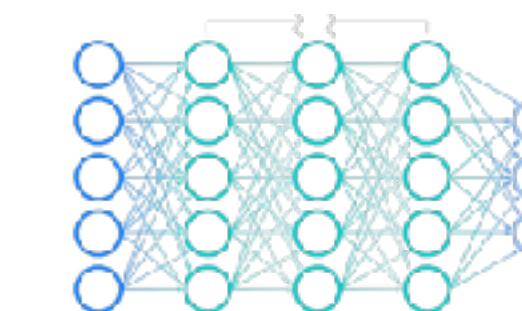
Same story, but with math



The magical synthesiser: neural networks

$$x \in \mathbb{R}^d \rightarrow f_w(x) = w_7 \max(0, w_1 x_1 + w_2 x_2 + w_3 x_3) + w_8 \max(0, w_4 x_4 + w_5 x_5 + w_6 x_6)$$

$$\rightarrow f_w(x) = W_L \phi(W_{L-1} \phi(\cdots \phi(W_1 x + b_1) \cdots) + b_{L-1}) + b_L.$$



The learning procedure, but with math

'Donkey'

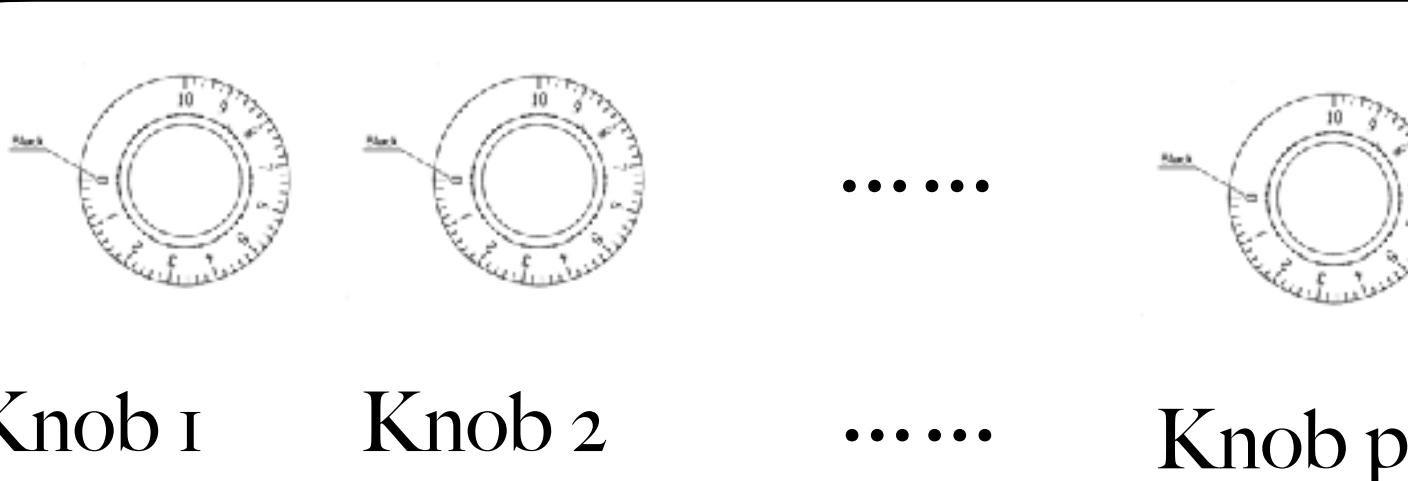


Label +1

$$\begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 \\ x_7 & x_8 & x_9 & x_{10} & x_{11} & x_{12} \\ x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} \\ x_{19} & x_{20} & x_{21} & x_{22} & x_{23} & x_{24} \\ x_{25} & x_{26} & x_{27} & x_{28} & x_{29} & x_{30} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} \end{bmatrix}$$

$$x \in \mathbb{R}^d$$

Magic synthetiser f_w



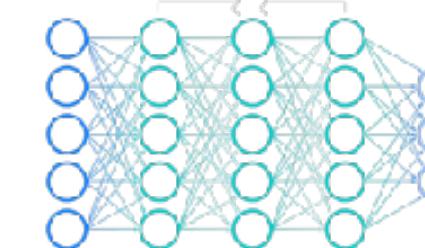
$$w = (w_1, w_2, \dots, w_p) \in \mathbb{R}^p$$

Machine Response

$$f_w(x) \in \mathbb{R}$$

Magic synthesisers: neural networks

$$x \rightarrow f_w(x) = w_7 \max(0, w_1 x_1 + w_2 x_2 + w_3 x_3) + w_8 \max(0, w_4 x_4 + w_5 x_5 + w_6 x_6)$$



How can we get it to learn?

Before training, it predicts $f_w(x)$ (-234.5), and we want it to predict "Donkey" (+1)

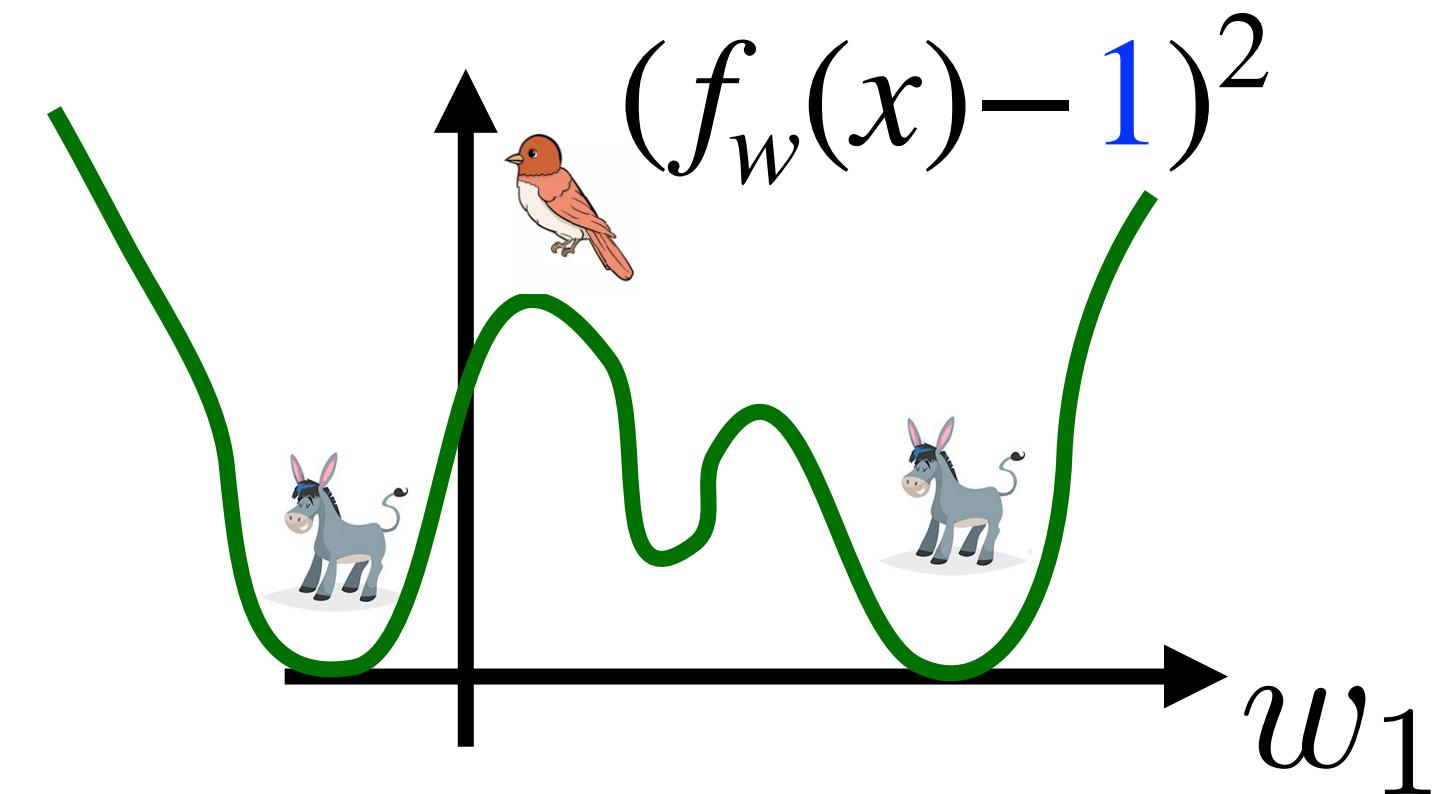
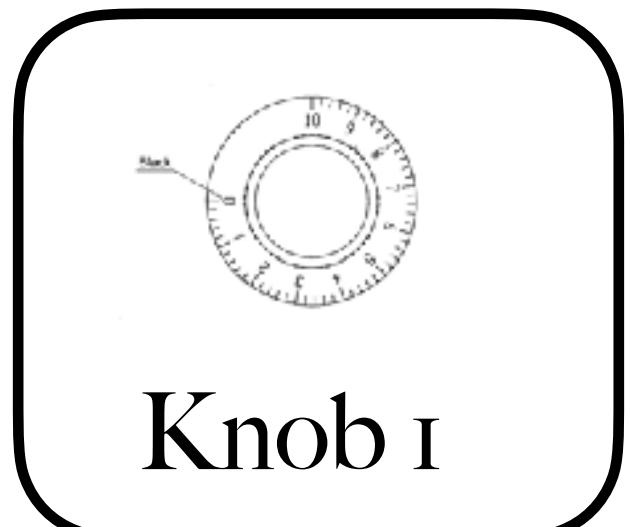
We can model its mistake by $L(w) = (f_w(x) - 1)^2$

We want to minimise the mistake! $\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$

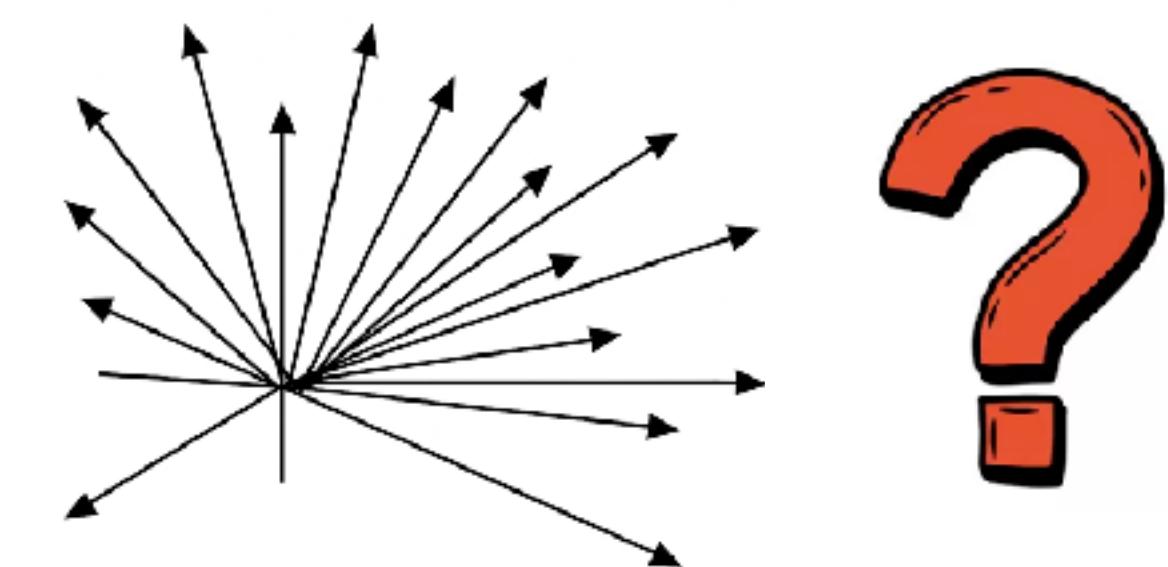
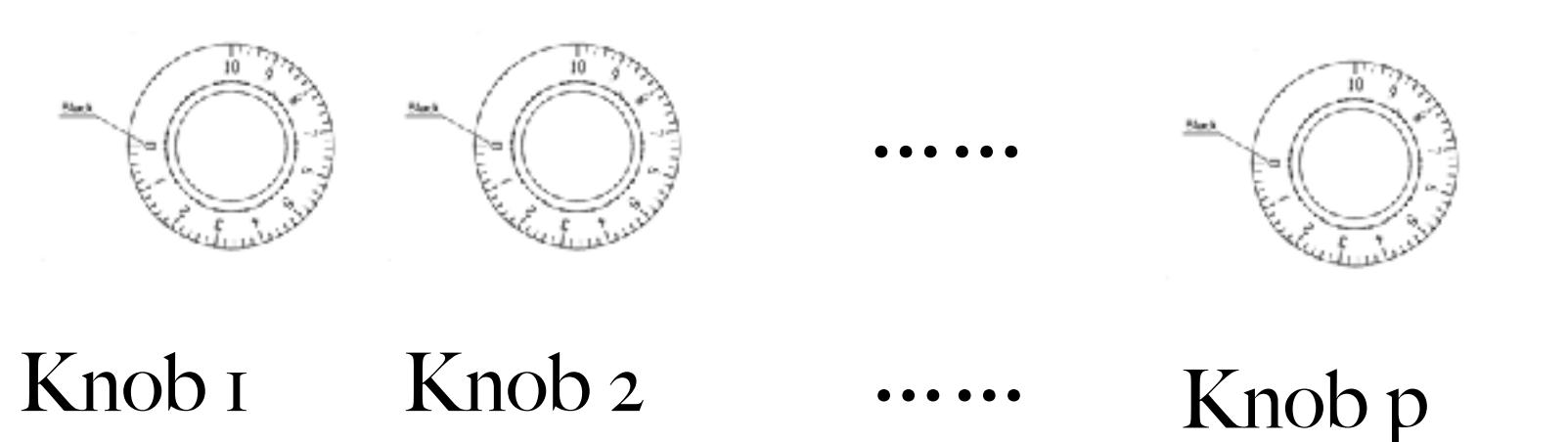
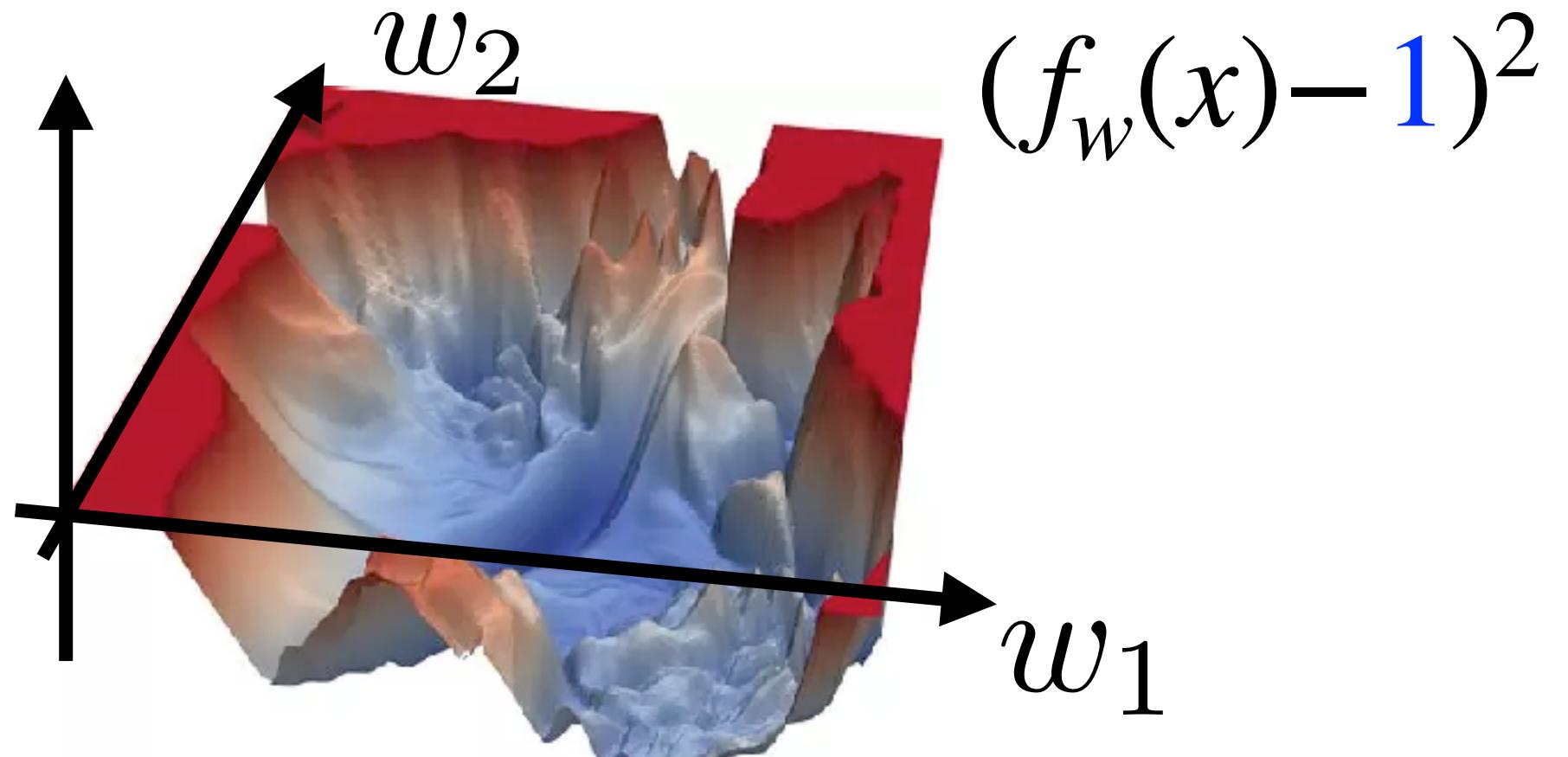
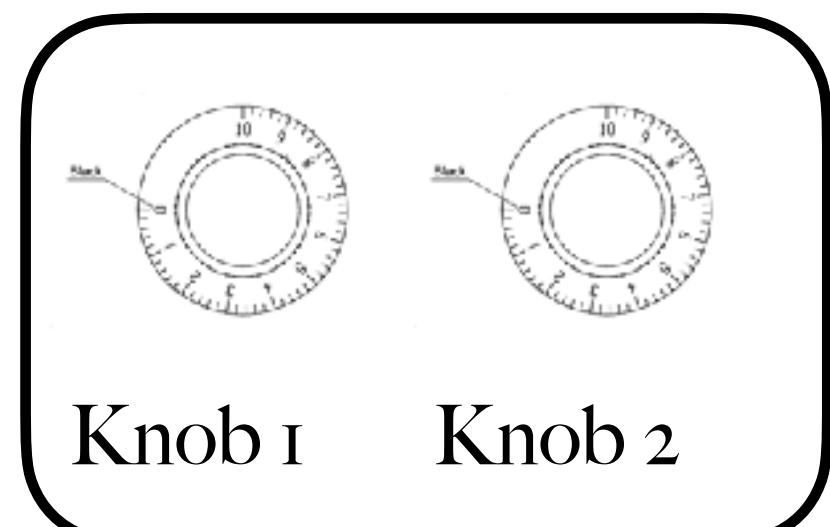
It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

$$w = (w_1)$$



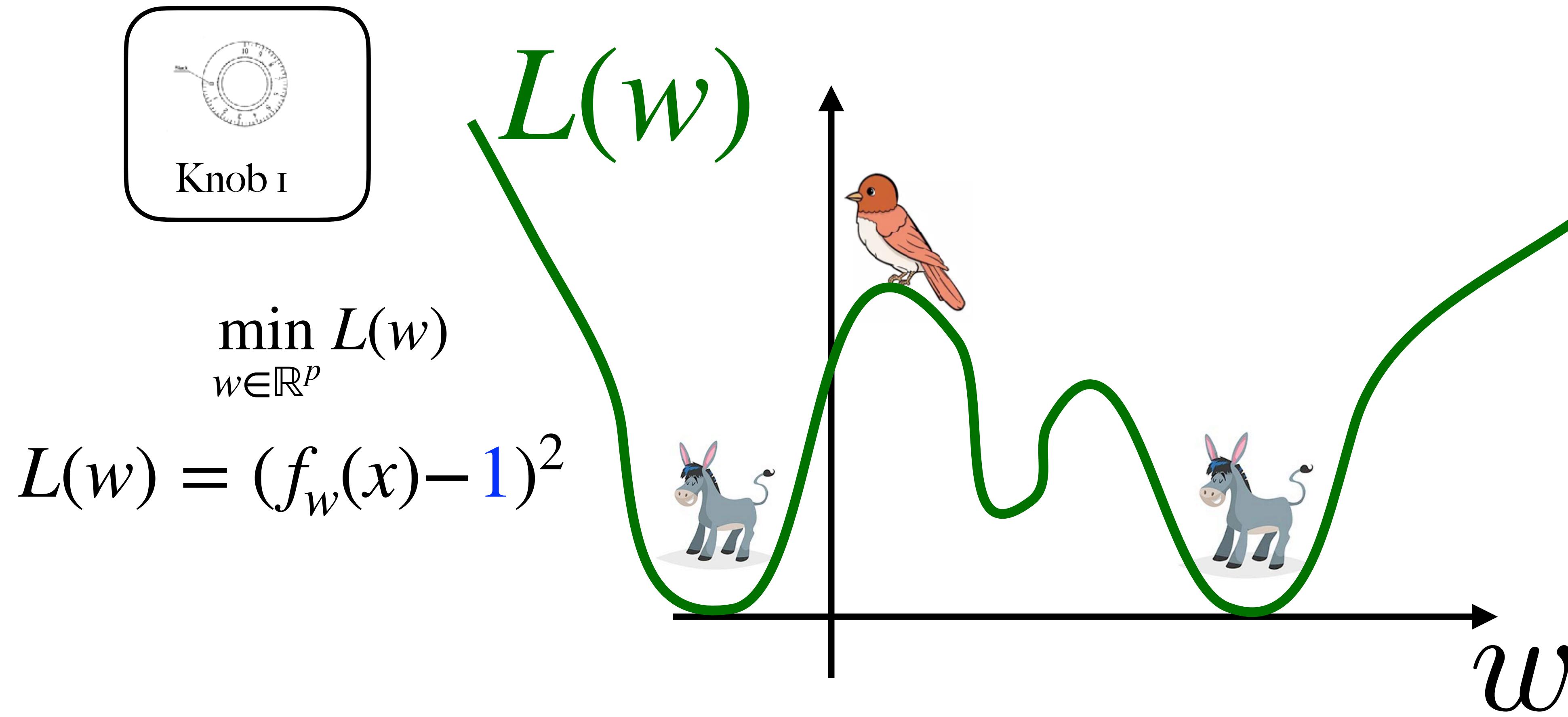
$$w = (w_1, w_2)$$



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

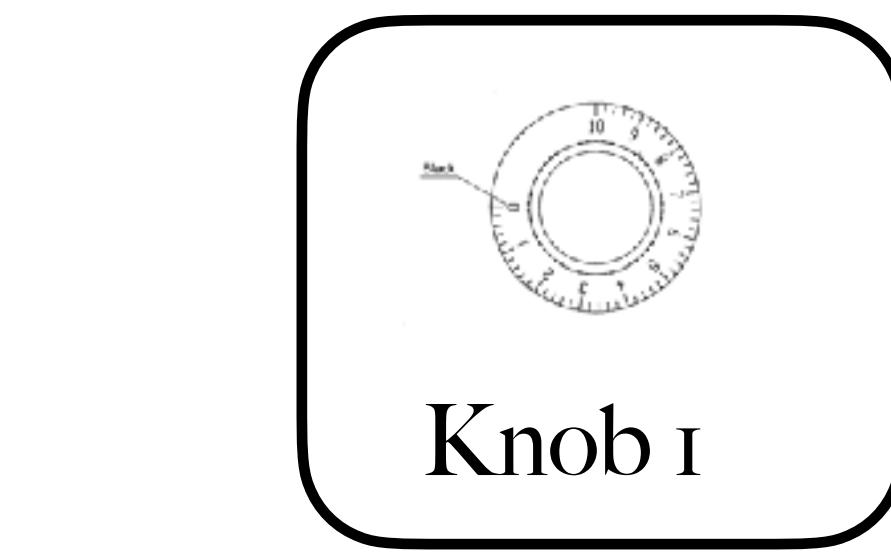
Back to 1D:



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

Back to 1D:

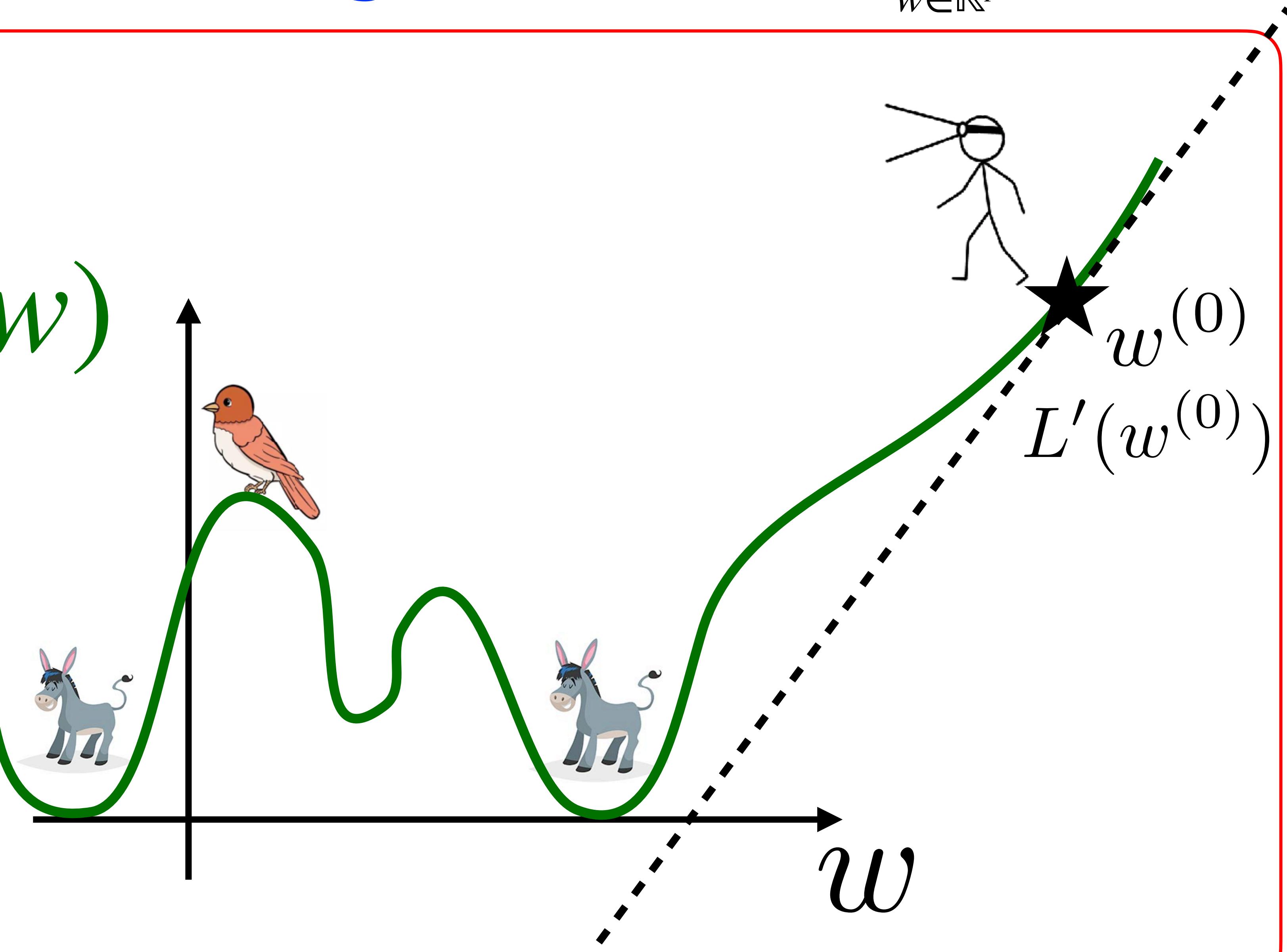


$$\min_{w \in \mathbb{R}^p} L(w)$$

$$L(w) = (f_w(x) - 1)^2$$

$$L(w)$$

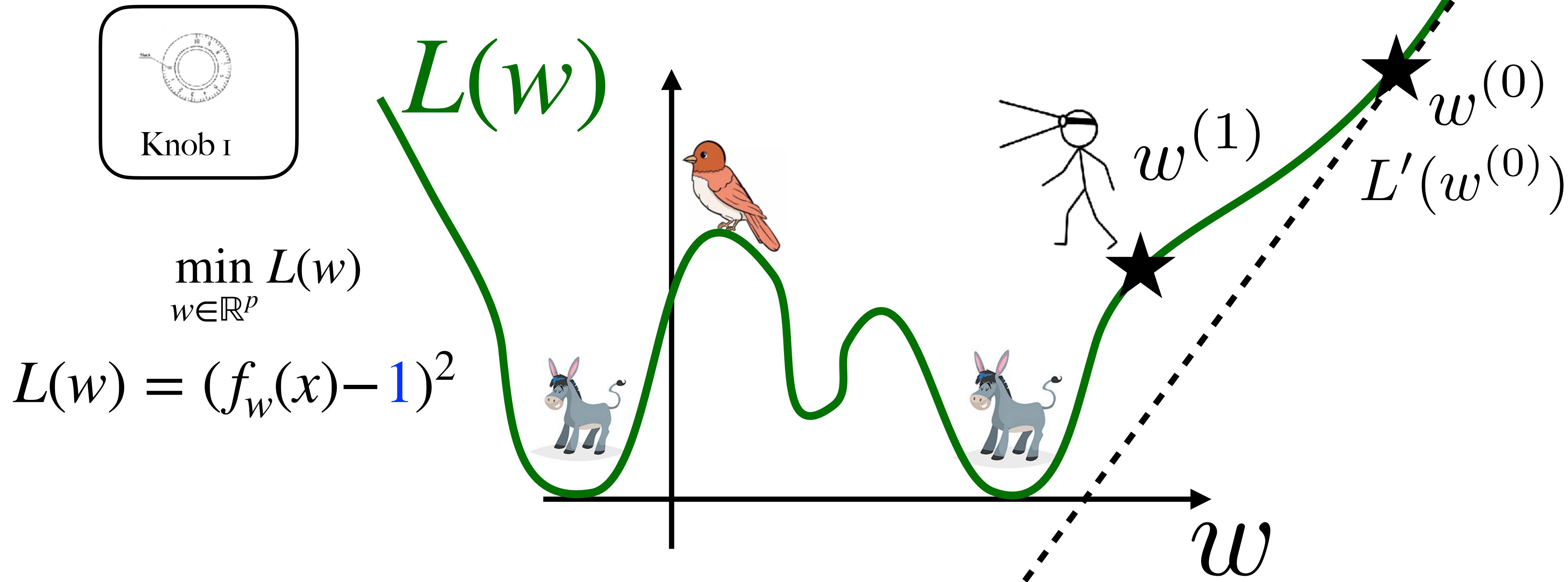
$$w^{(0)}$$
$$L'(w^{(0)})$$



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

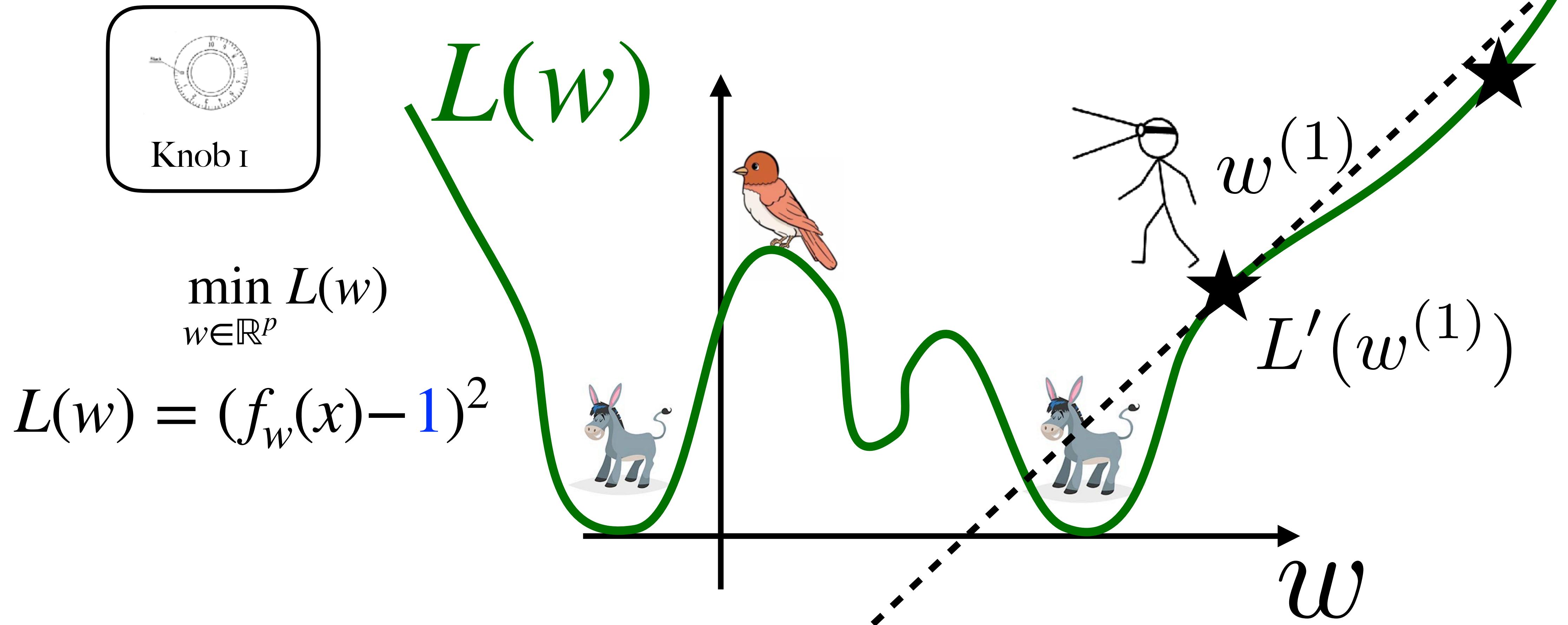
Back to 1D:



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

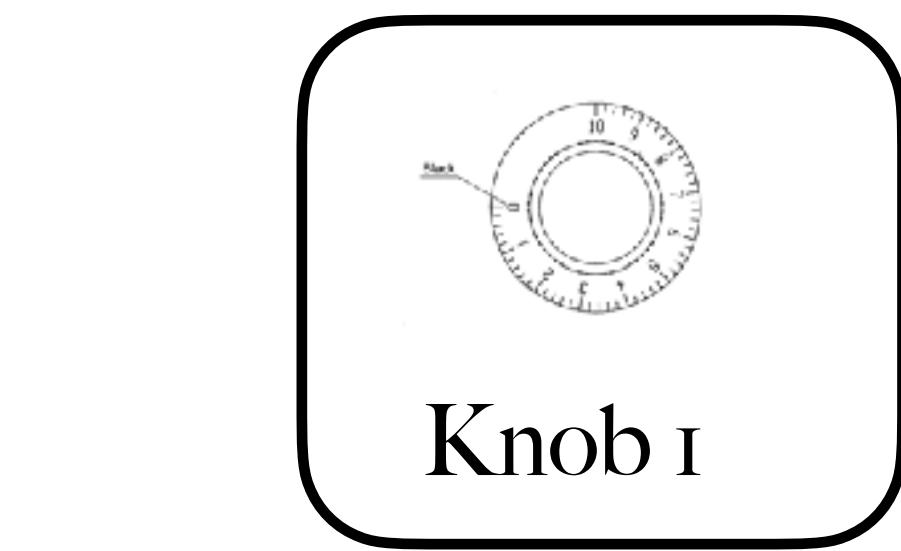
Back to 1D:



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

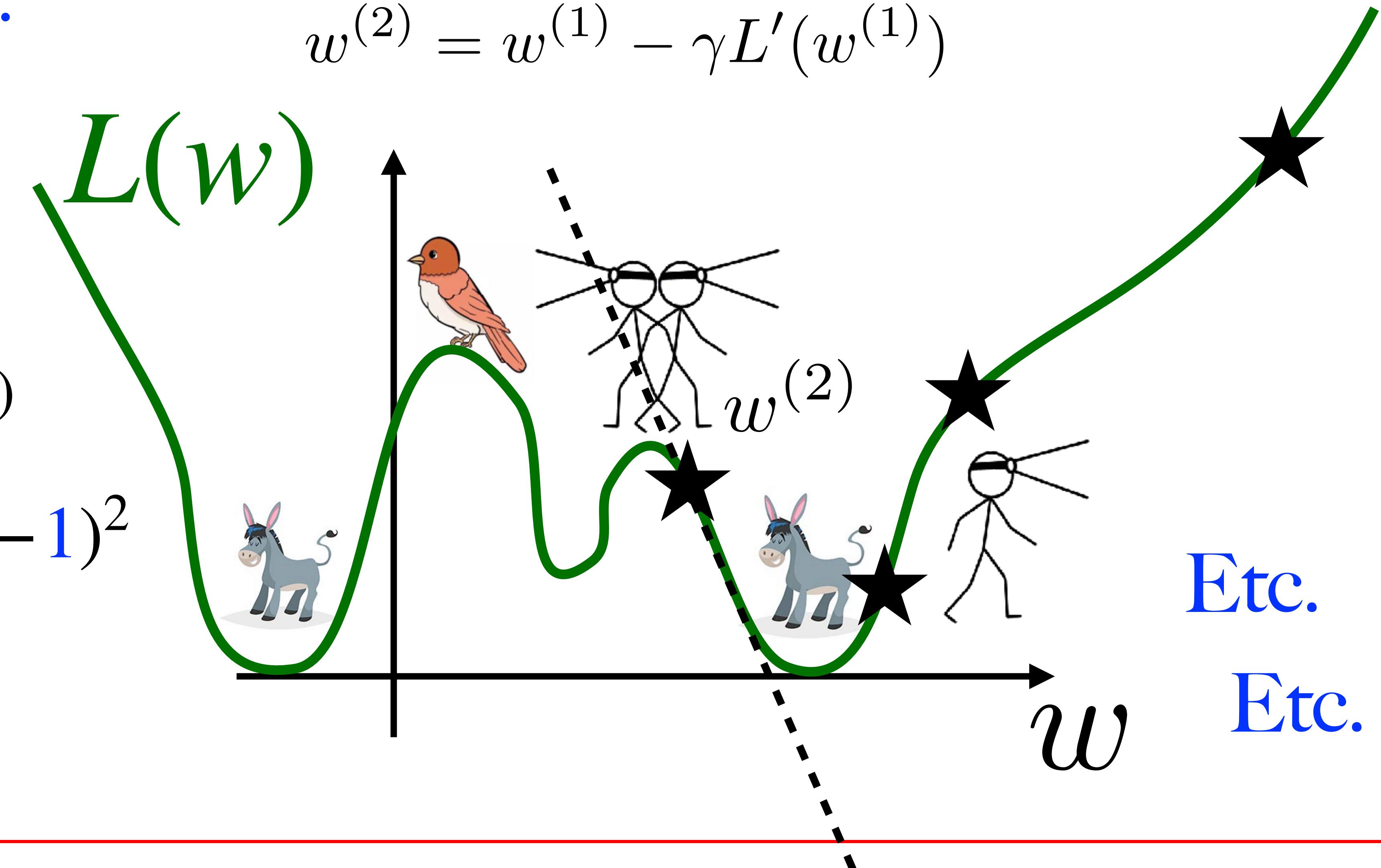
Back to 1D:



$$\min_{w \in \mathbb{R}^p} L(w)$$

$$L(w) = (f_w(x) - 1)^2$$

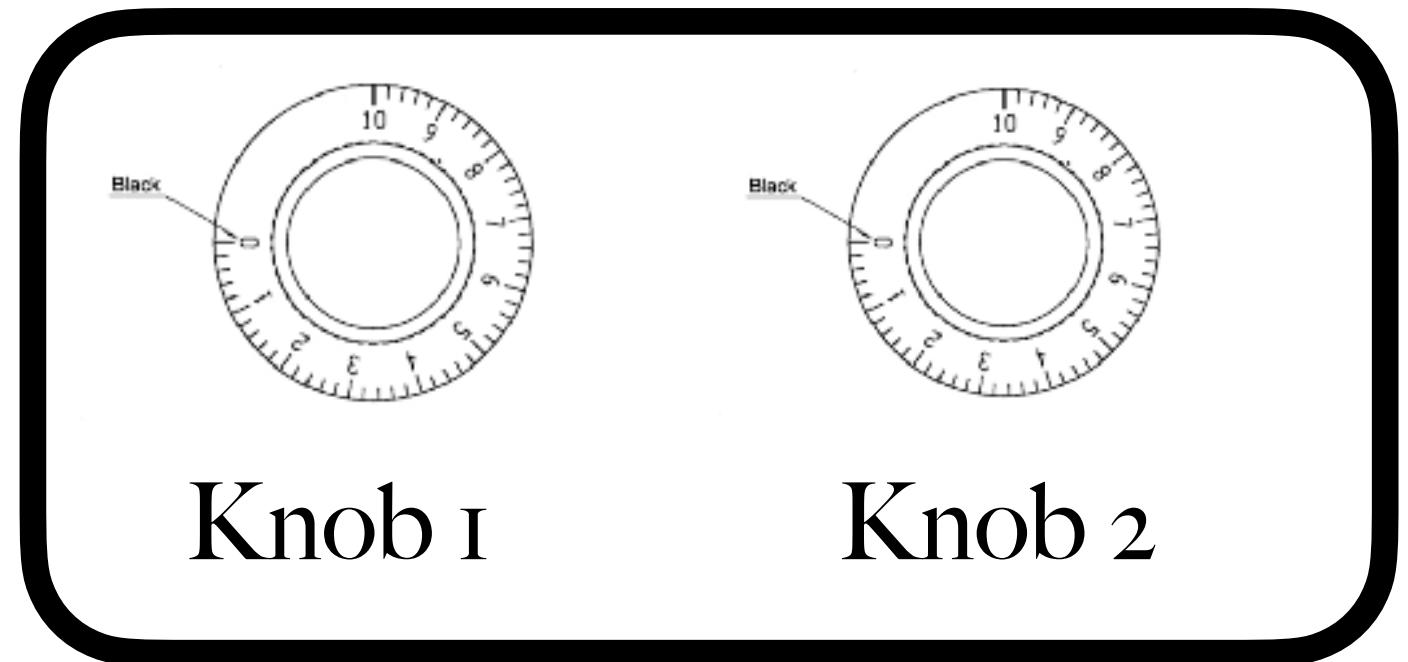
$$w^{(1)} = w^{(0)} - \gamma L'(w^{(0)})$$
$$w^{(2)} = w^{(1)} - \gamma L'(w^{(1)})$$



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$

En 2D:



Knob 1

Knob 2

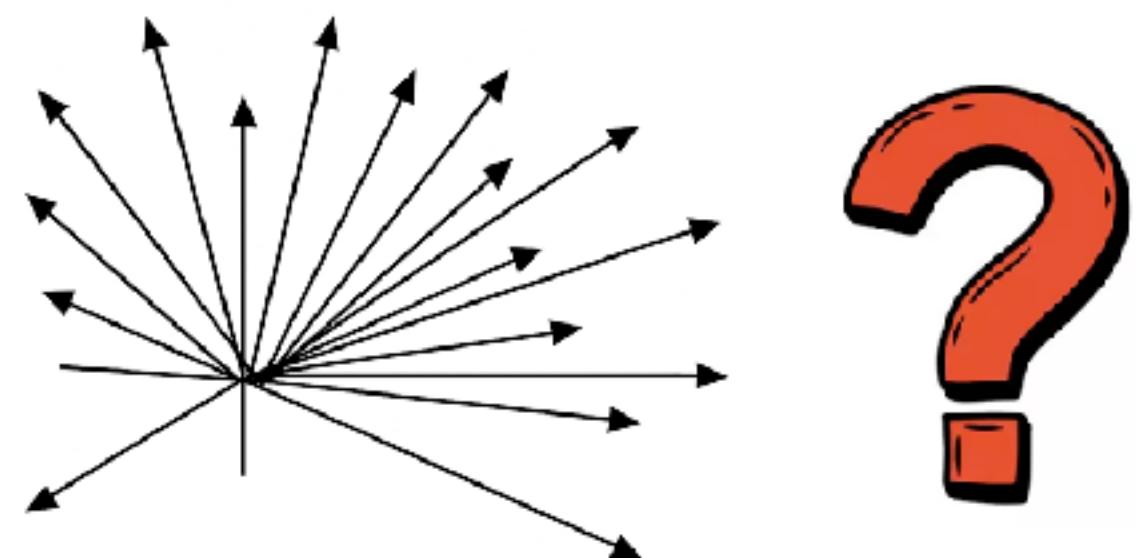
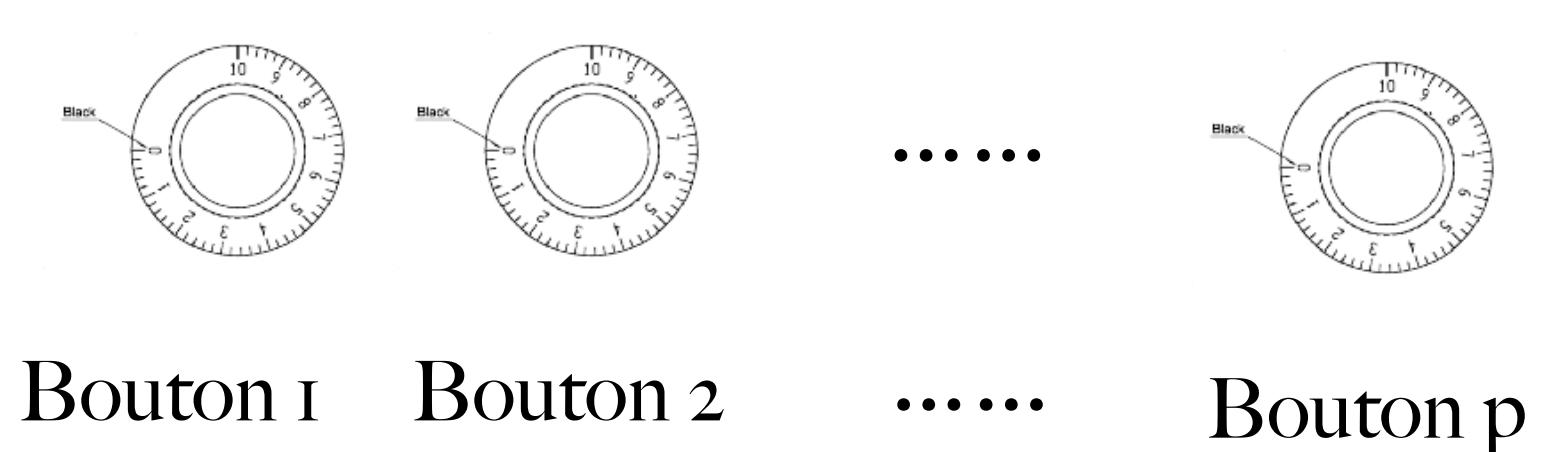
$$\min_{w \in \mathbb{R}^p} L(w)$$

$$L(w) = (f_w(x) - 1)^2$$



It boils down to minimising a function!

$$\min_{w \in \mathbb{R}^p} (f_w(x) - 1)^2$$



$$\min_{w \in \mathbb{R}^p} L(w)$$

$$L(w) = (f_w(x) - 1)^2$$

Gradient descent method:

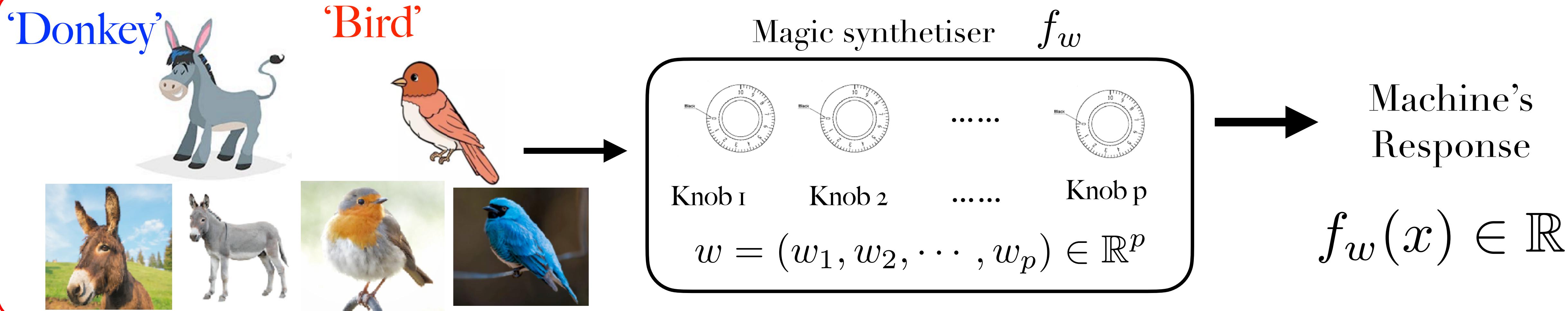
$$w^{(k+1)} = w^{(k)} - \gamma \nabla L(w^{(k)})$$

$$\text{New knob configuration} = \text{Old configuration} + \text{How much to turn them} \times \text{In which direction}$$

If “everything goes well”: $L(w^{(k)}) \xrightarrow{k \rightarrow \infty} 0$ $f_{w^{(k)}}(x) \xrightarrow{k \rightarrow \infty} 1$



But there is more than one image!



Data:

$$(x_1, y_1), \dots, (x_n, y_n)$$

1.2	3.5	0.7	4.1	2.7	1.3
0.6	2.8	3.3	1.1	0.5	4.9
7.2	2.4	1.1	0.9	3.3	2.2
3.7	0.9	4.8	5.5	0.2	1.6
1.0	3.1	2.2	2.7	4.0	0.8
0.3	1.5	6.4	2.9	3.6	1.0

+1

3.1	0.8	2.4	5.2	1.9	3.6
4.7	1.2	0.9	3.4	2.8	5.0
0.5	3.3	4.1	1.8	2.0	0.7
2.6	5.5	1.1	0.4	3.9	2.3
3.0	2.7	5.8	1.5	0.6	4.2
1.4	3.6	0.2	4.5	2.1	5.3

-1

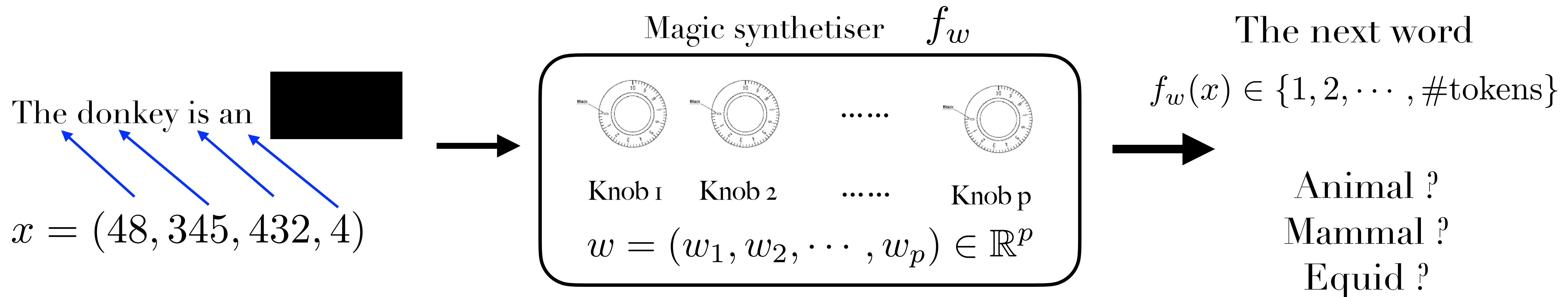
Loss to minimise:

$$\min_w \sum_{i=1}^n (y_i - f_w(x_i))^2$$

“Supervised learning”

What about ChatGPT?

The donkey is an equid. As a herbivore, it frequently consumes fibrous plants.



Data (sentences):

$$x_1, x_2, \dots, x_n$$

Fonction à minimiser:

$$\min_w \sum_{i=1}^n (\text{mot d'après} - f_w(x_i))^2$$

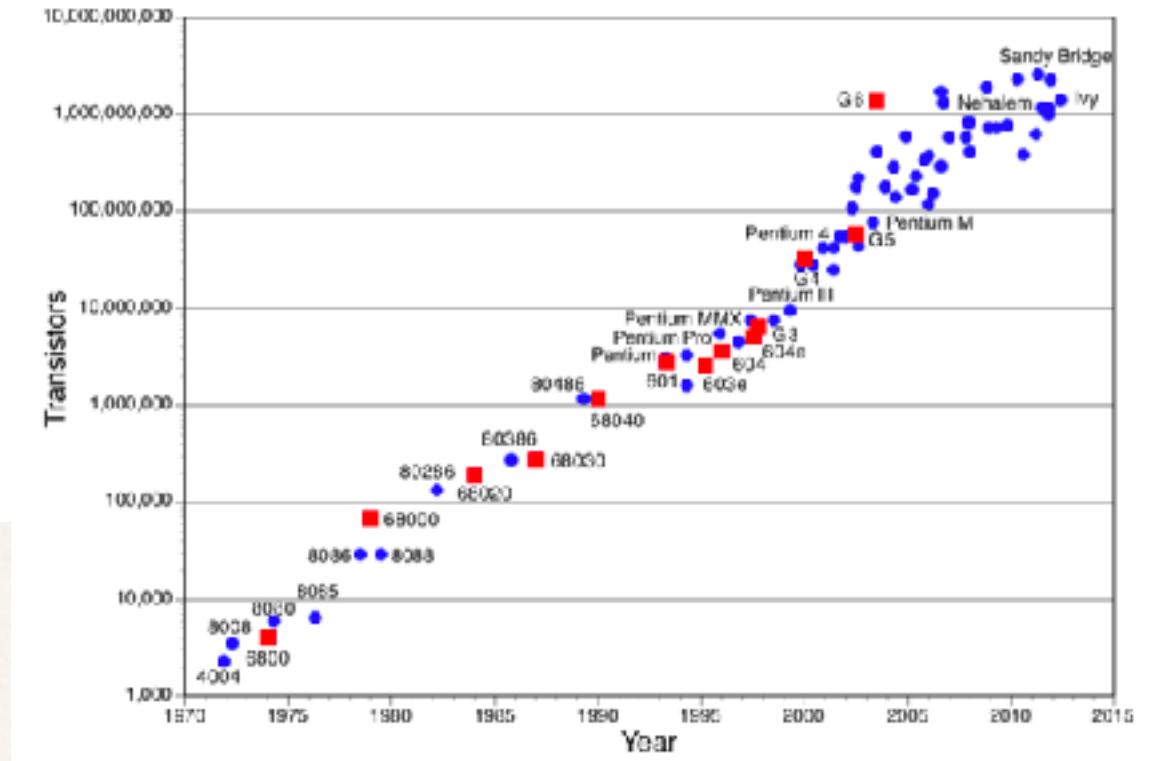
Apprentissage “semi-supervisé”

Taking a step back

Three ingredients behind the success of modern machine learning:

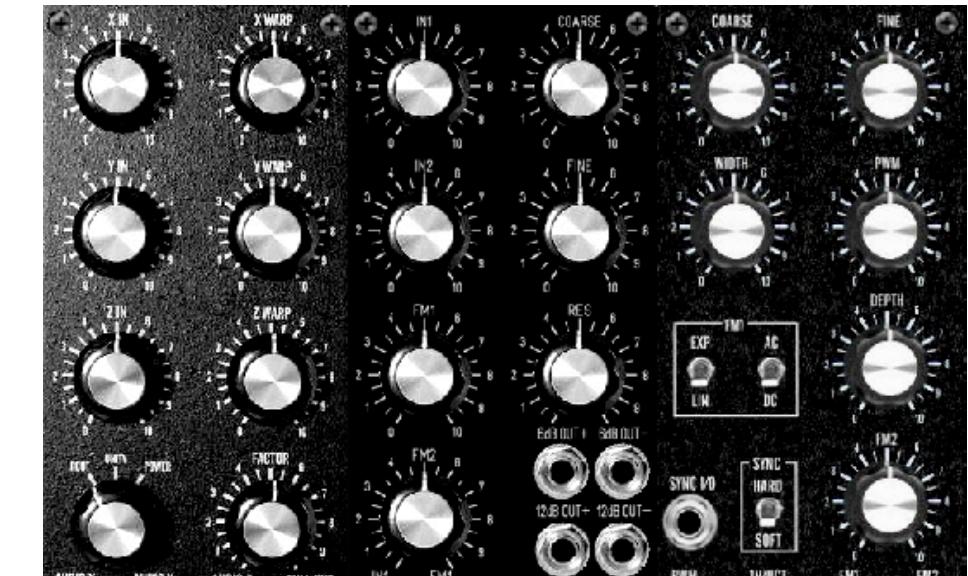
Data: internet (text et images)

1950s: 10^3 FLOPs
2024: 10^{25} FLOPs



Training Chat-GPT, some numbers

What's the capital
of Switzerland?



→ Bern

Size of the synthesiser:

A trillion (10^{12}) trainable knobs \approx Paris covered with knobs

Energy to turn the knobs: 10GWh energy consumption (€2M euros electricity bill)

Bringing up a child to its twenties:

To compare:
 $2000\text{Kal/day} \rightarrow 100\text{Watts}$
 $20\text{ MWh from 0 to 20 years old}$

therefore: Training ChatGPT $\approx 500 \times$
 $0 \rightarrow 20$

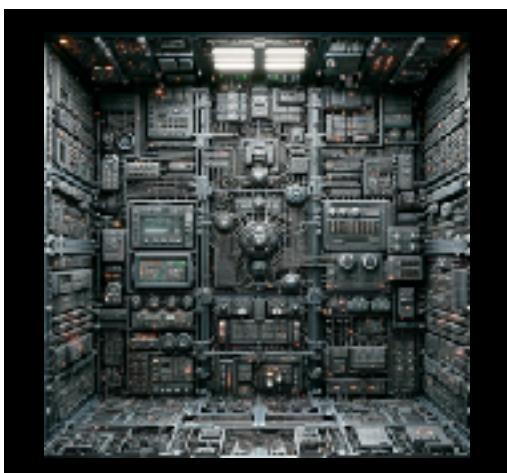


But not much compared to the $\approx 1\text{GWh}$ daily consumption due to queries!

Text seen by ChatGPT:

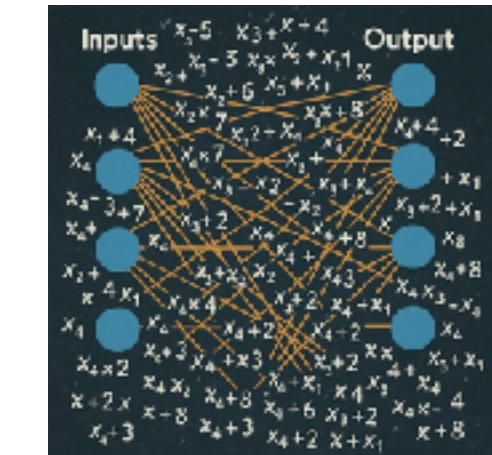
Trained on a 10 trillion (10^{13}) words \approx 10 million copies of War and Peace \approx 100 000 years to read (no sleeping)

Why don't we understand?

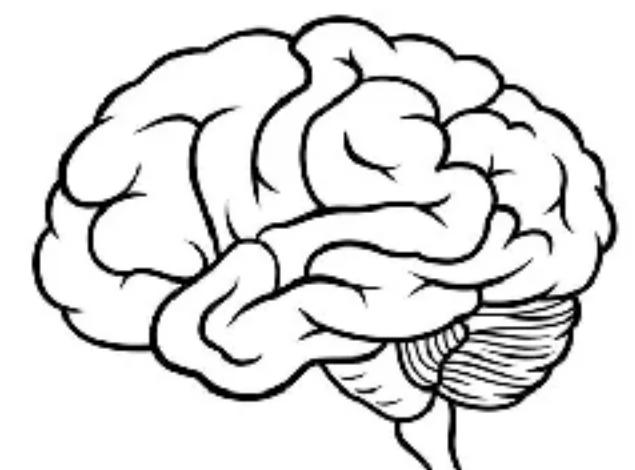


‘Lagopède
Alpin’

$$f_w(x) = W_L \phi(W_{L-1} \phi(\cdots \phi(W_1 x + b_1) \cdots) + b_{L-1}) + b_L.$$

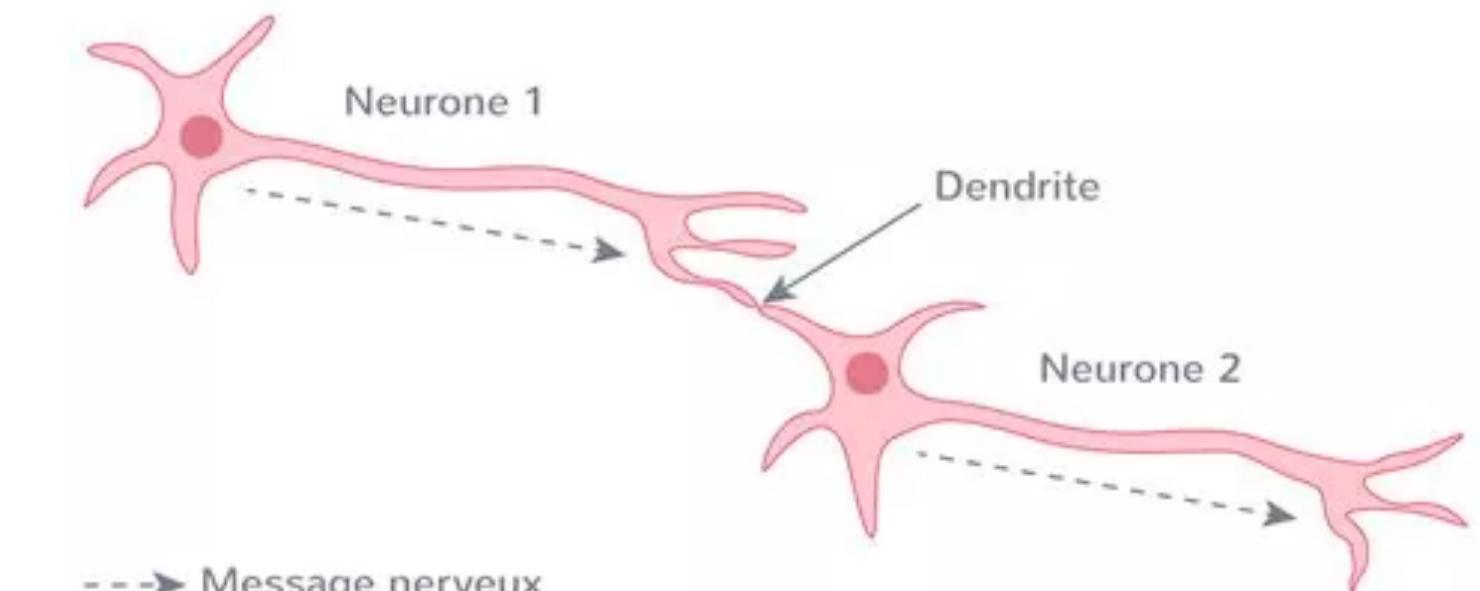


Just because we know exactly what happens at the ‘local scale’ doesn’t mean we understand the behaviour at the macroscopic level.

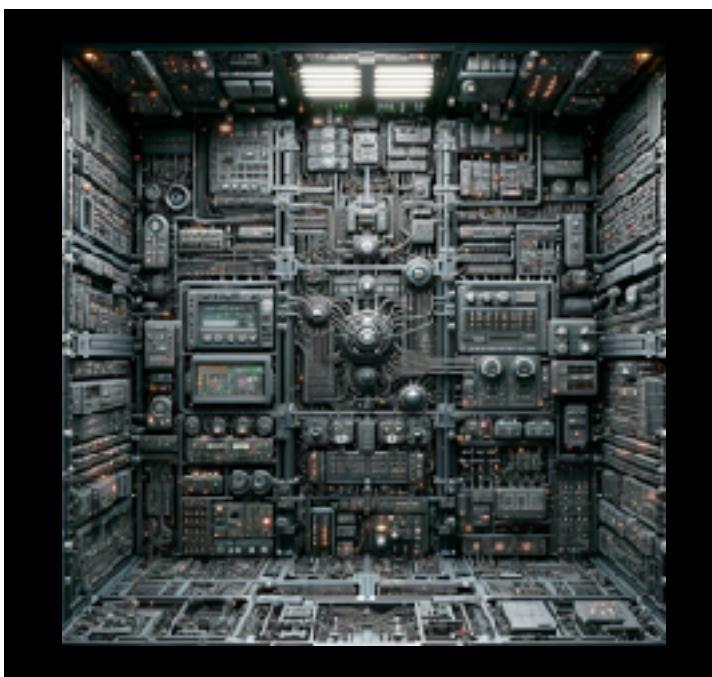


‘Lagopède
Alpin’

La propagation d'un message nerveux entre neurones



What are the issues?



‘Lagopède
Alpin’

It works brilliantly, but we don't understand why!

How does it make a decision? With which criteria?
When does it make a mistake? Can we fool it?
Is it biased? Can we extract sensible data?

Issues: security, ethics, reliability, sustainability...

Recent successes challenge many of our conceptions: Machines now navigate territories reserved to conscious beings

Does ChatGPT *understand* things?

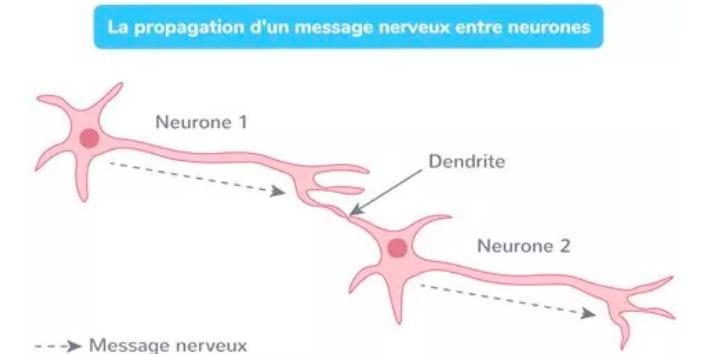
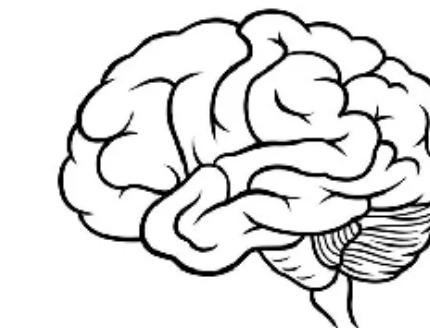
Is it *intelligent*?

Is it *creative*?

Could it be *conscious*?

How do we even define these words?

What about us, humans:



Is consciousness / intelligence only computational?

Is creativeness about interpolation or extrapolation?

A better understanding of these seemingly ‘intelligent’ machines will shed light on many aspects of our own.

So what can we do?

As we don't properly understand them,
Current neural network can be depicted as::

Chat-GPT:



Hard to follow what is going on

What theoreticians look at:



It already has similar properties (it flies!), and much simpler to analyse and understand.

“Nothing is more practical than a good theory”

Practical
deep learning



Transformers

$$\text{Atten}(h)(t) = h(t) + \sum_{i=1}^n W_i^1 \sum_{s=1}^T \sigma[(W_Q^i h(t))^T W_K^i h(s)] W_V^i h(s)$$
$$\text{FF}(h)(t) = f(h)(t).$$

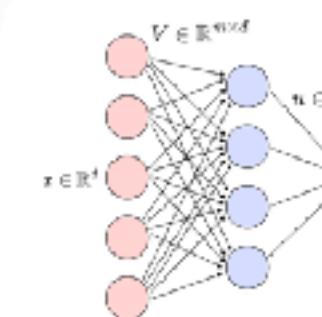
Intuition
through
experiments



Theoretical
deep learning



Simplified
network



Understanding
and
Improvements

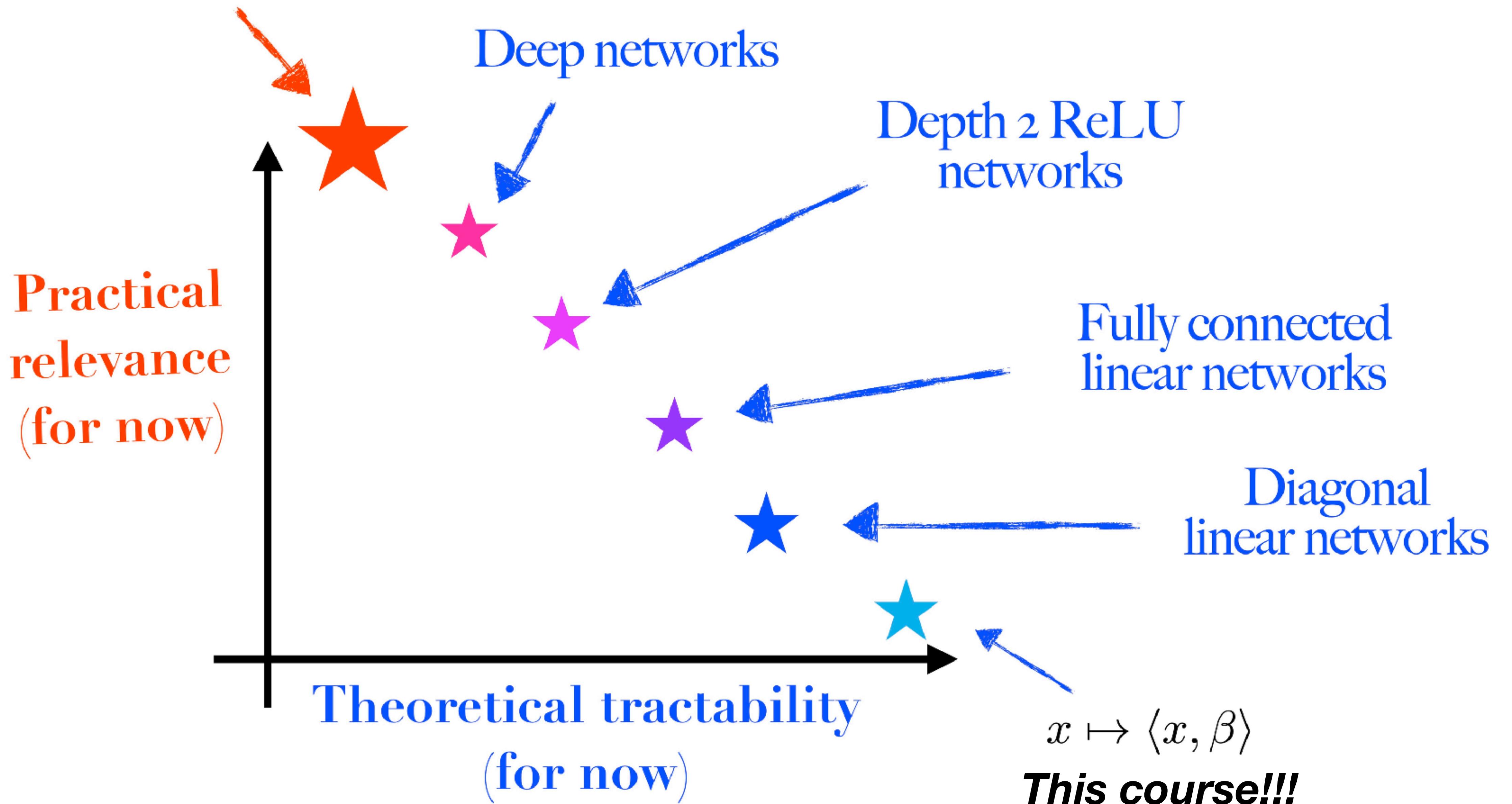


Future of
deep learning (??)



Architecture of
the future

Transformers



Up next: let's do some math!

What you know / don't know yet:



1. What is the gradient of $w \in \mathbb{R}^d \mapsto \langle w, x \rangle$
2. Let $A \in R^{d \times d}$, what is $w^\top A w$ equal to? a) $\sum_{i,j=1}^d A_{ij} w_i w_j$ b) $\sum_{i=1}^d A_{ij} w_i w_j$ c) $\|Aw\|^2$
3. Let $A \in R^{d \times d}$, what is the gradient of $w \in \mathbb{R}^d \mapsto w^\top A w$?
4. Which of the following are convex functions?
 - a) $w \mapsto \|w\|^2$
 - b) $w \mapsto w^\top A w$
 - c) $w \mapsto \exp(w)$
 - d) $w \mapsto \ln(w)$
5. Suppose $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and differentiable. Which minimal condition guarantees that w^\star is a global minimiser?
 - a) $\nabla f(w^\star) = 0$
 - b) $\nabla^2 f(w^\star) \succeq 0$
 - c) $\nabla f(w^\star) = 0$ and $\nabla^2 f(w^\star) \succeq 0$