

scai

SORBONNE CLUSTER FOR
ARTIFICIAL INTELLIGENCE



AI Introduction: From Linear Regression to Large Language Generative Models (part 1/2)

Raphael Cousin
Data Scientist
raphael.cousin@sorbonne-universite.fr
raphaelcousin.com



Plan

AI - Introduction:

- Definition
- Application
- History
- Context

Duration : ~10 min

Machine Learning:

- Data characteristics
- Parametric Models
- Loss function
- Training
- Not Parametric Model

Duration : ~30 min

Deep Learning:

- From Neuron to MLP
- CPU vs GPU
- The NN Layers

Duration : ~5 min

Data Science:

- Data Preprocessing
- Hyperparameters
- Model evaluation methodology

Duration : ~10 min

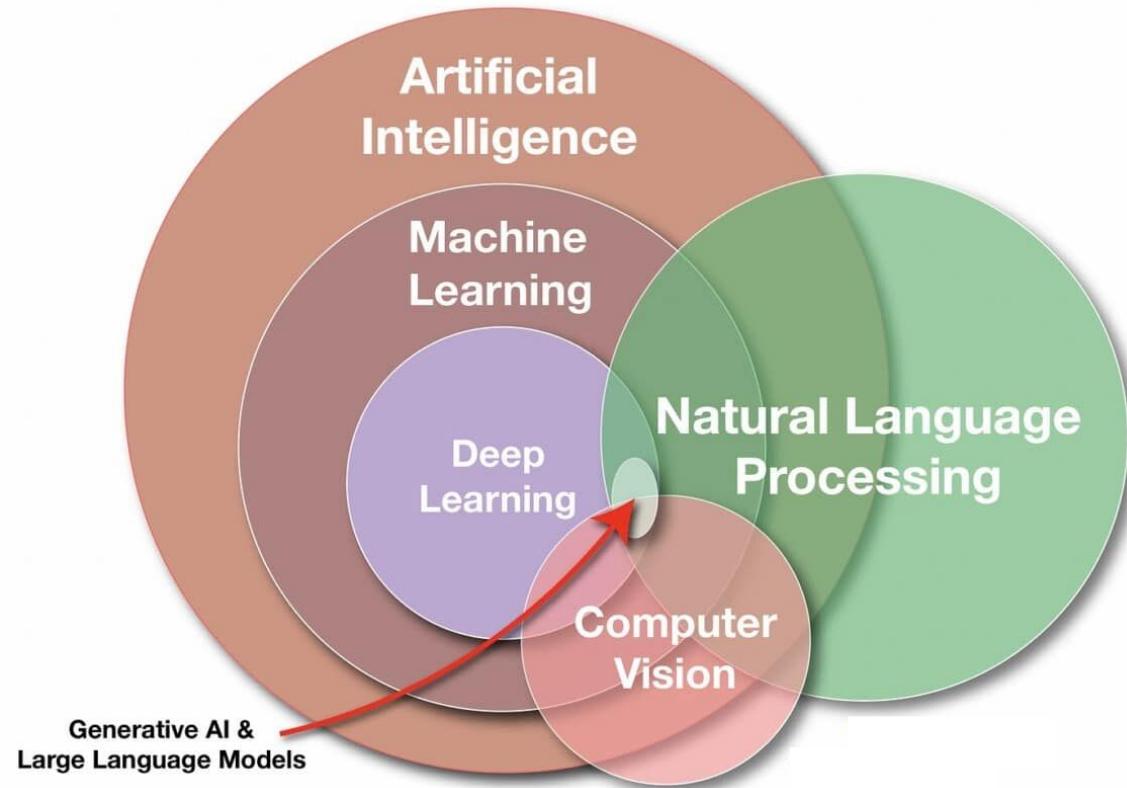
AI - Landing Overview

AI : Definitions

AI: Automation involving algorithms to reproduce cognitive capabilities such as reasoning, perception, and decision-making.

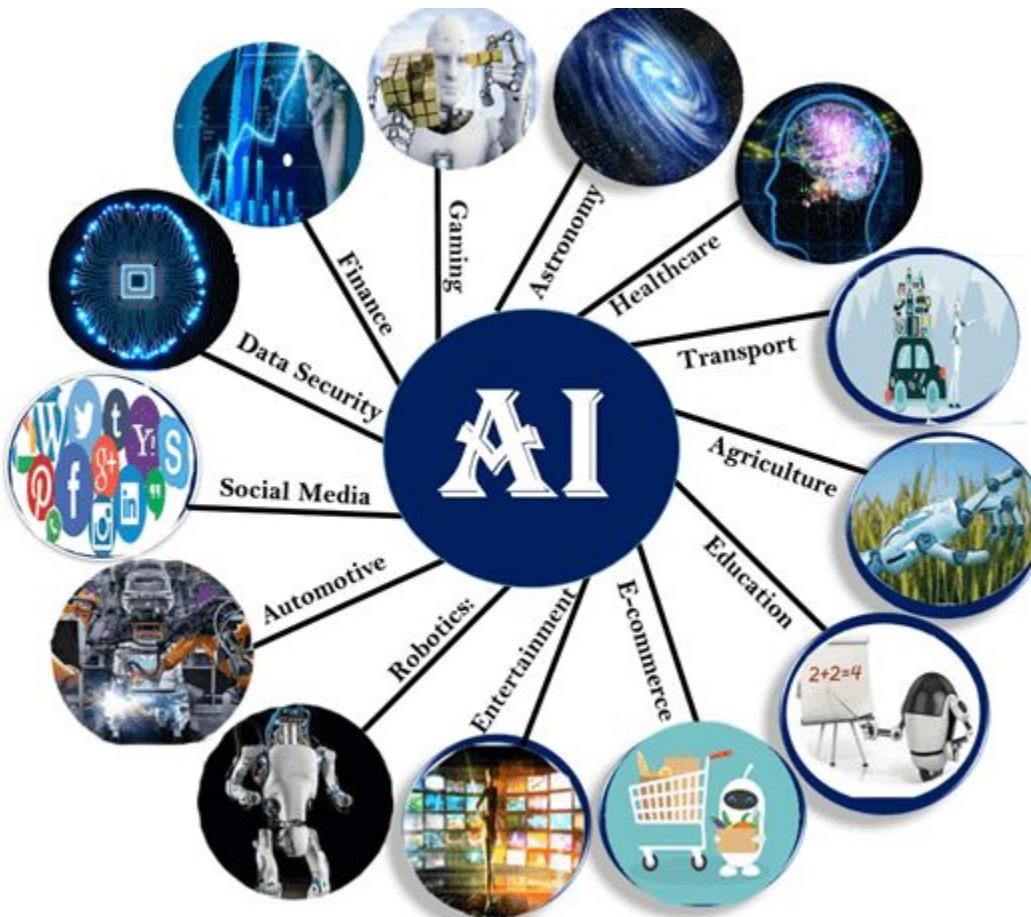
ML: Statistical algorithms that automatically learn patterns from data without being explicitly programmed for each task.

DL: Neural networks with multiple layers (millions of parameters) that can learn complex representations from raw data.

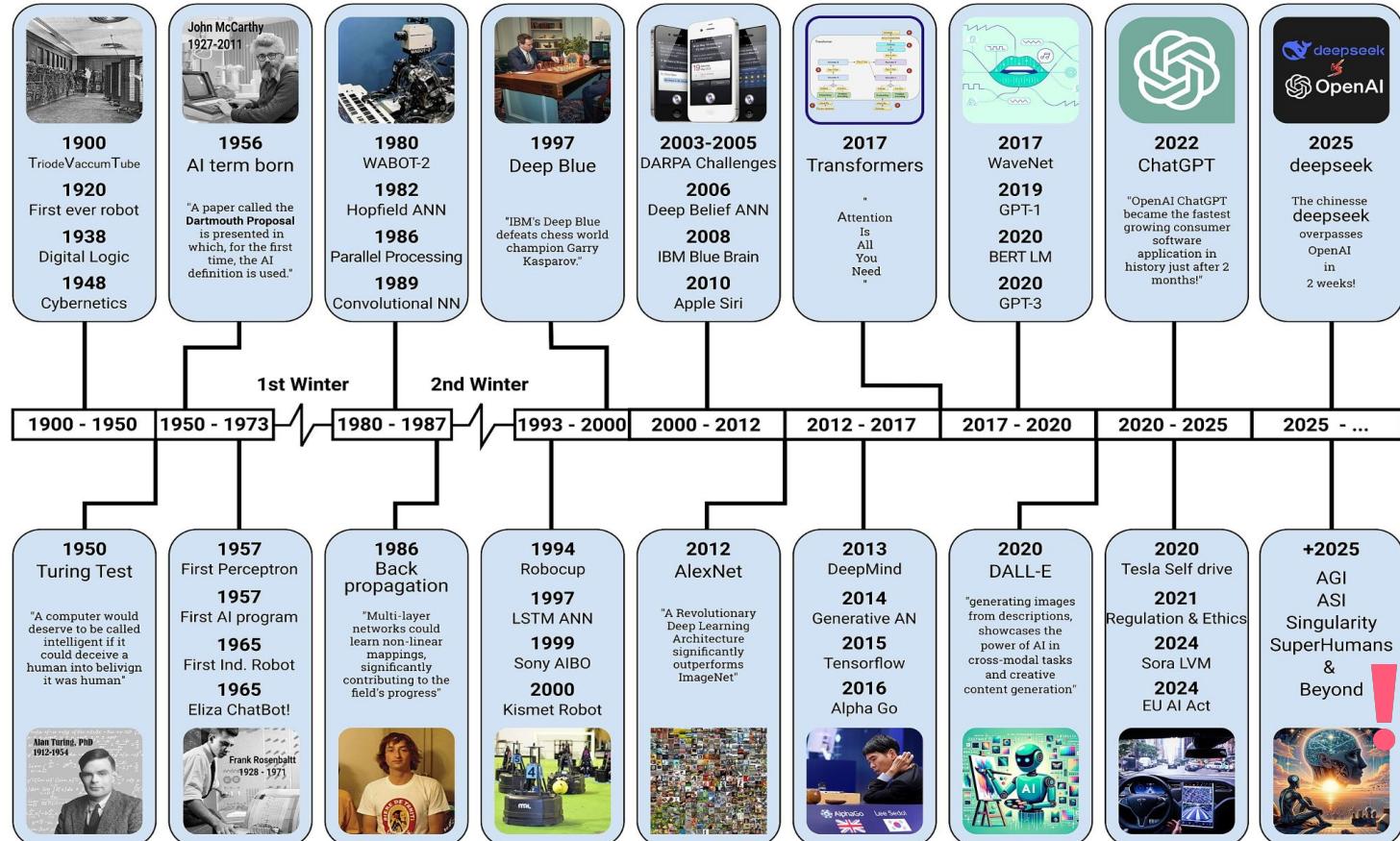


source: <https://medium.com/geekculture/ai-revolution-your-fast-paced-introduction-to-machine-learning-914ce9b6ddf>

AI : applications



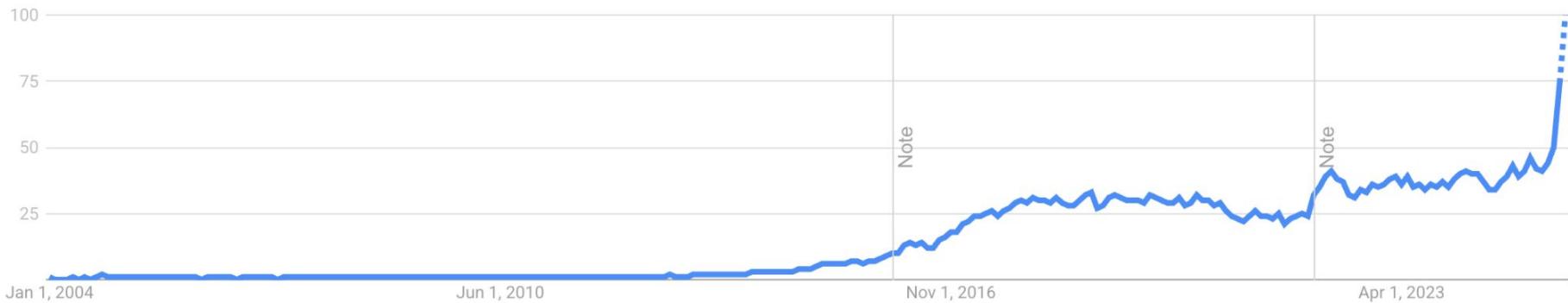
AI History



source: <https://commons.wikimedia.org/wiki/File:AI-History-Timeline-300dpi.jpg>

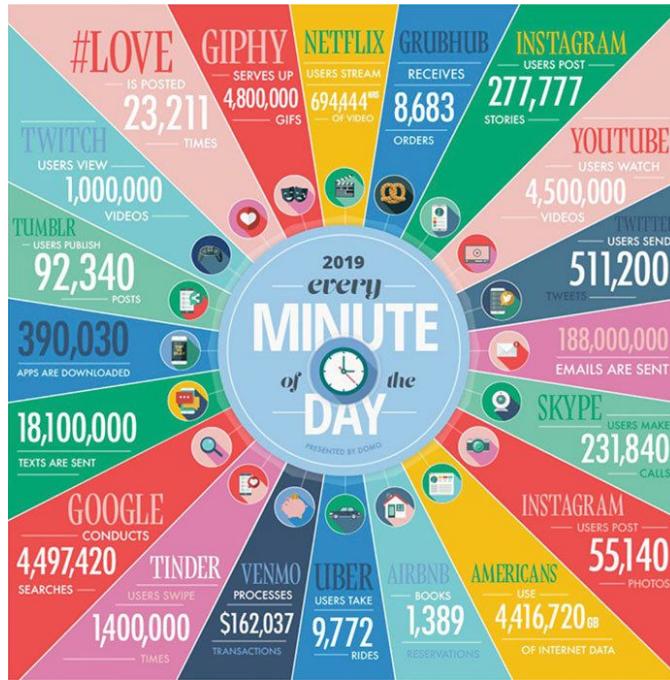
The Deep Learning Revolution

Interest over time ?



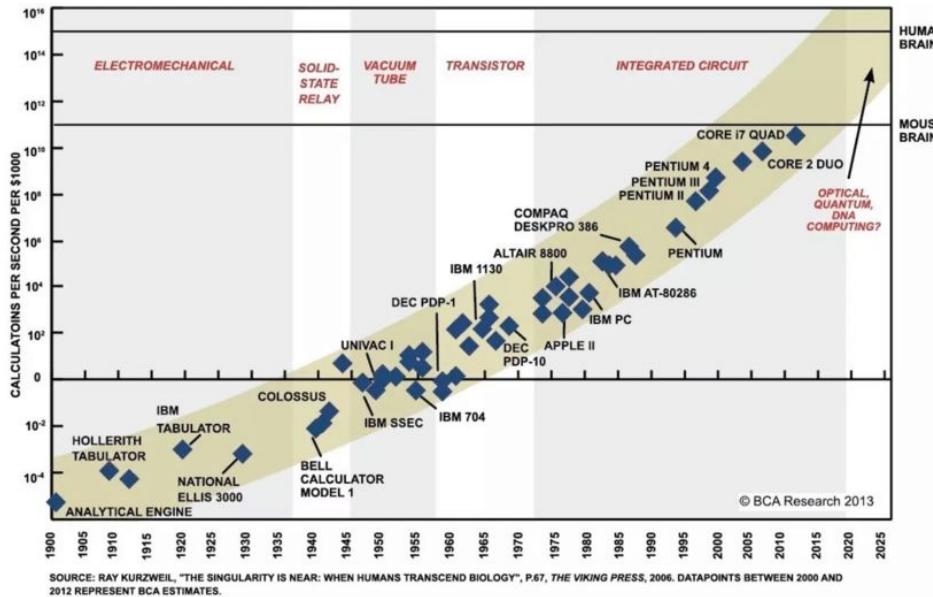
The convergence of three critical factors enabled the deep learning revolution:

1. Big Data



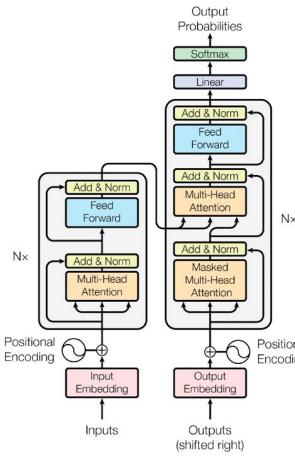
The internet era generated massive datasets: billions of images, text documents, videos, and user interactions. This data explosion provided the fuel needed to train complex models that require millions of examples to learn robust patterns.

2. Computational Power



GPUs originally designed for gaming proved perfect for neural network computations. A single modern GPU can perform trillions of operations per second, enabling training of models with billions of parameters in days rather than years.

3. Algorithmic Innovation



 PyTorch

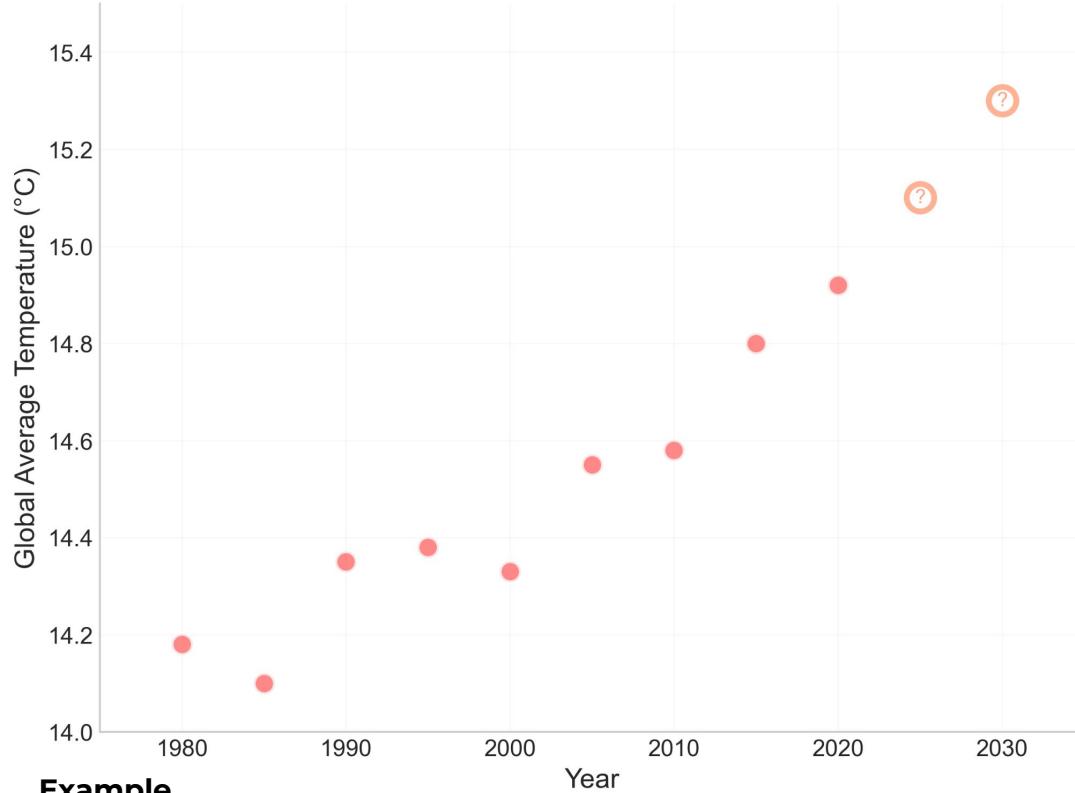
Breakthrough techniques like ReLU activation, batch normalization, dropout, and attention mechanisms solved critical training challenges. These innovations made it practical to train networks with dozens or even hundreds of layers. Open source code with big communities.

Machine Learning

The Machine Learning Objective

Machine Learning definition:

Systems that automatically learn patterns from data.

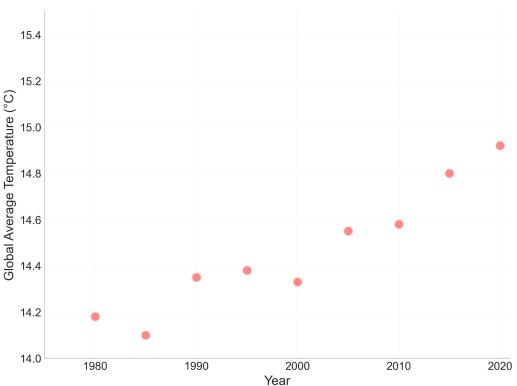


Example

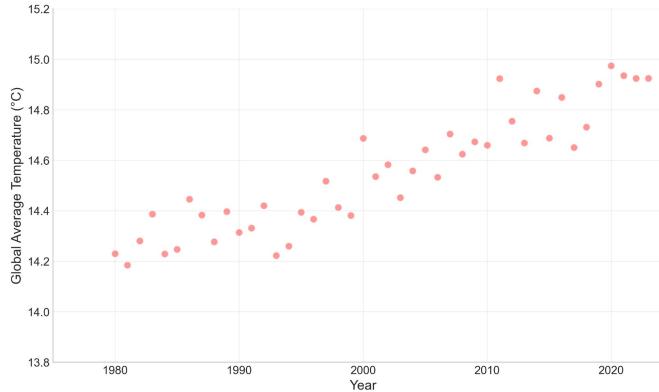
Using historical data to predict the future

Data Characteristics

Quantity: Sample size: number of observations

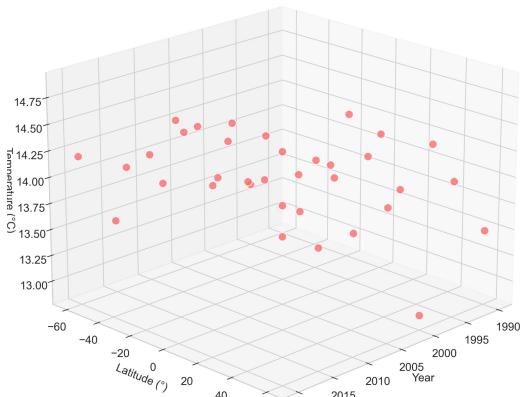


9 observations / 2 variables



44 observations / 2 variables

Dimension: number of variables per observation (features)



9 observations / 3 variables

Année	Latitude	Altitude	CO ₂	Humidité	°C
2000	45.0	200	370	65	13.8
2010	30.0	800	390	45	12.5
2020	60.0	50	415	80	14.2

3 observations / 6 variables

ML : Features and Target/Label

X : Explanatory variables (features/characteristics)

Y: Variable to explain (target)

Concretely, we want to find a relationship $Y = f(X)$

Example: $f(\text{year}, \text{latitude}, \text{altitude}) = \text{predicted temperature}$

Features								Label
date	lat	long	temp	humidity	cloud_coverage	wind_direction	atmp_pressure	rainfall
2021-09-09	49.71N	82.16W	74	20	3	N	18.6	.01
2021-09-09	32.71N	117.16W	82	42	6	SW	29.94	.23

Example

Regression, Classification

Regression: Predict continuous values

$$y \in \mathbb{R}^k$$

Examples: House price



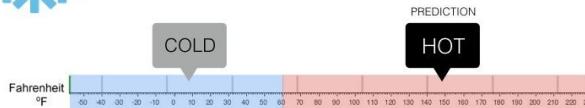
Regression

What is the temperature going to be tomorrow?



Classification

Will it be Cold or Hot tomorrow?



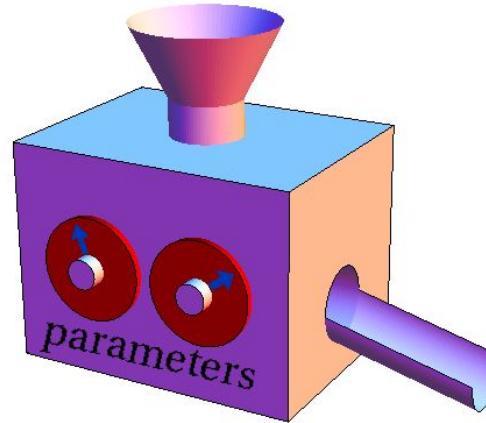
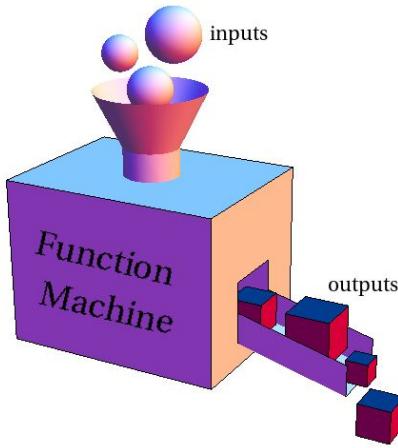
Classification: Predict discrete categories

$$y \in \{1, 2, \dots, K\}$$

Examples: Spam detection

In practice we will predict continuous value "probability" in $[0,1]$

Parametric Models characteristics



Parametric Models are functions defined by

Parameters: Number of values the model must learn

Parameter vector: $\theta \in \mathbb{R}^p$

p is the number of parameters

Input/Output Dimension:

$$f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^k$$

Input dimension: Dimension of features

Output dimension:

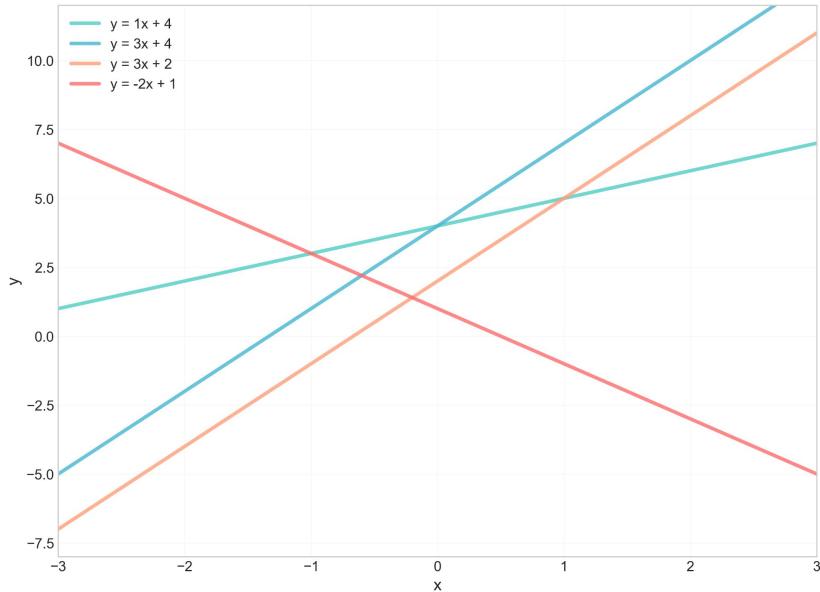
Dimension of features targets

Operations: Between parameters and input data

Operations: addition, multiplication, composition

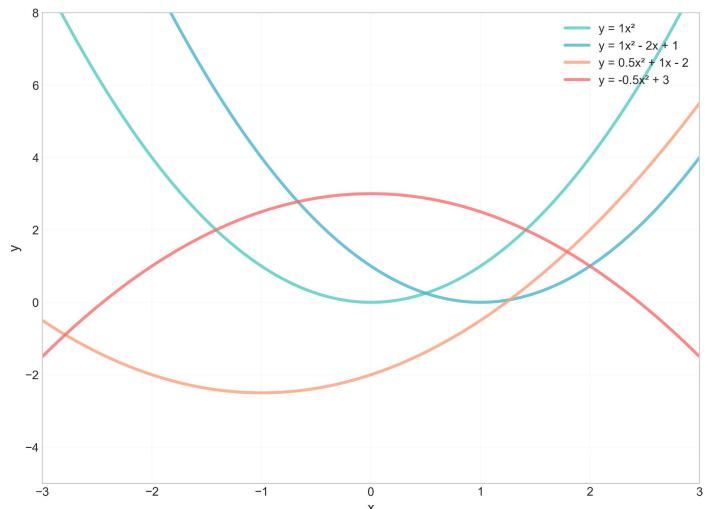
Parametric Models

Linear Functions: $y = ax + b$
(2 parameters: slope a , intercept b)



$\mathbb{R} \rightarrow \mathbb{R}$, 2 params

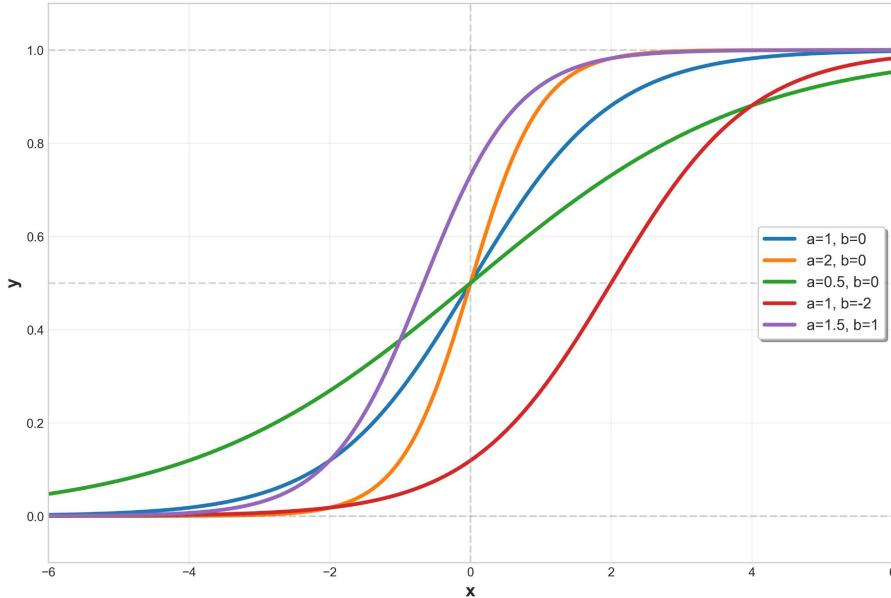
Polynomial Functions: $y = ax^2 + bx + c$
(3 parameters: a , b , c)



$\mathbb{R} \rightarrow \mathbb{R}$, 3 params

Parametric Models

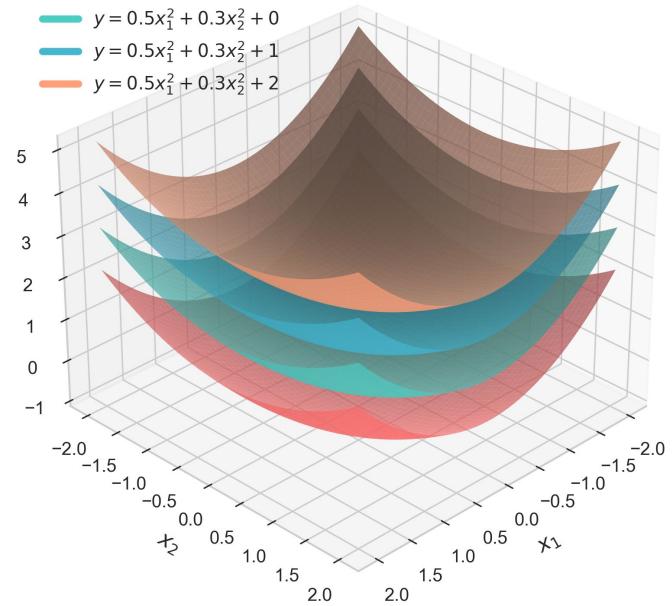
Logistic Functions: $y = 1/(1 + e^{-(ax + b)})$



Logistic $\mathbb{R} \rightarrow [0, 1]$, 2 params

Quadratic Functions: $y = ax_1^2 + bx_2^2 + cx_1 + dx_2 + e$
(5 parameters)

- $y = 0.5x_1^2 + 0.3x_2^2 + -1$
- $y = 0.5x_1^2 + 0.3x_2^2 + 0$
- $y = 0.5x_1^2 + 0.3x_2^2 + 1$
- $y = 0.5x_1^2 + 0.3x_2^2 + 2$



Quadratic function $\mathbb{R}^2 \rightarrow \mathbb{R}$, 5 params

Loss Functions characteristics

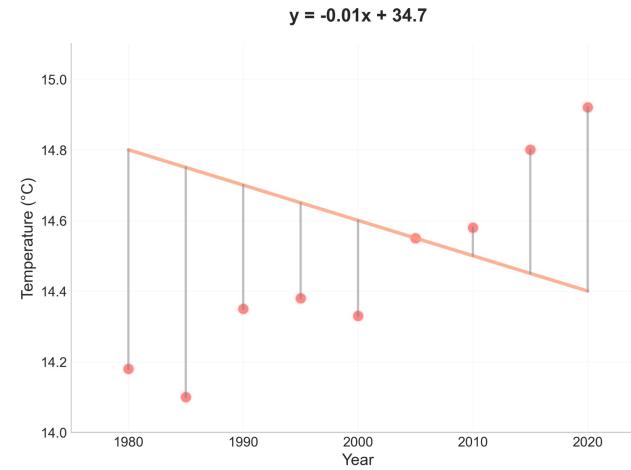
The Role of Loss Functions Loss functions quantify the error/distance between real targets $\hat{y} = f_\theta(x)$

$$l : \mathbb{R}^{k \times n} \times \mathbb{R}^{k \times n} \rightarrow \mathbb{R}$$
$$(Y, \hat{Y}) \mapsto l(Y, \hat{Y})$$

Example Loss Functions

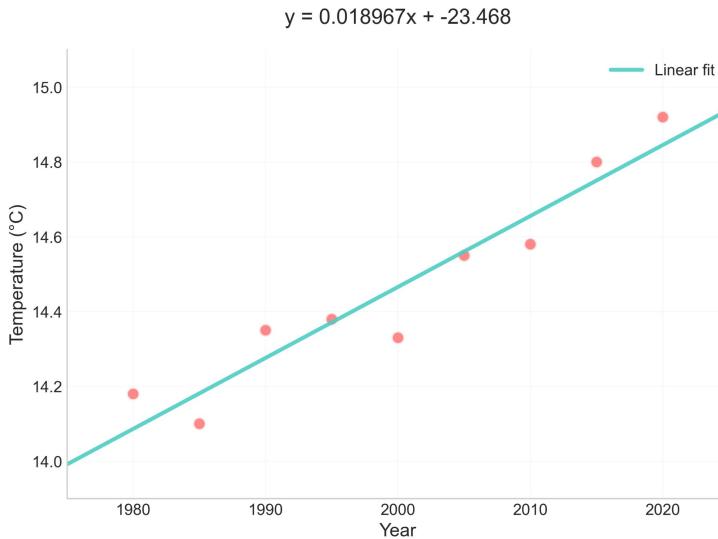
Mean Squared Error

$$\ell_{MSE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



The Fundamental Learning Problem

For a given parametric mode machine learning seeks to find the parameters that minimize the loss between real and predicted value :



$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \ell(\theta) = \arg \min_{\theta \in \mathbb{R}^p} l(Y, f_\theta(X))$$

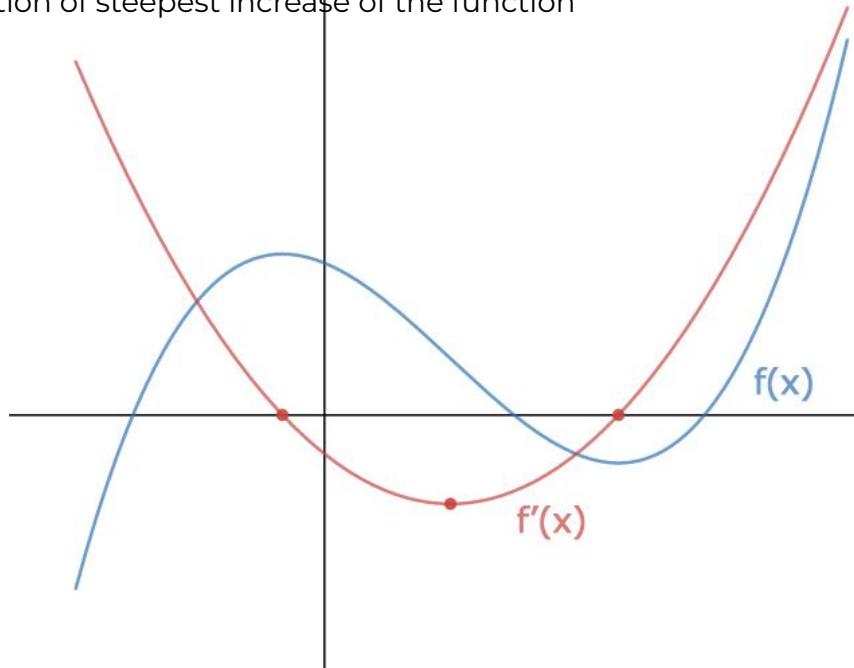
Derivatives and Gradients

Derivative (1D):

- Measures how much a function changes when its input changes

Gradient (Multi-dimensional):

- Extension of derivative to functions with multiple inputs $\nabla f = [\partial f / \partial x_1, \partial f / \partial x_2, \dots, \partial f / \partial x_n]$
- Points in the direction of steepest increase of the function



Gradient Descent and Optimization

0. Initialization

Initialize θ_0

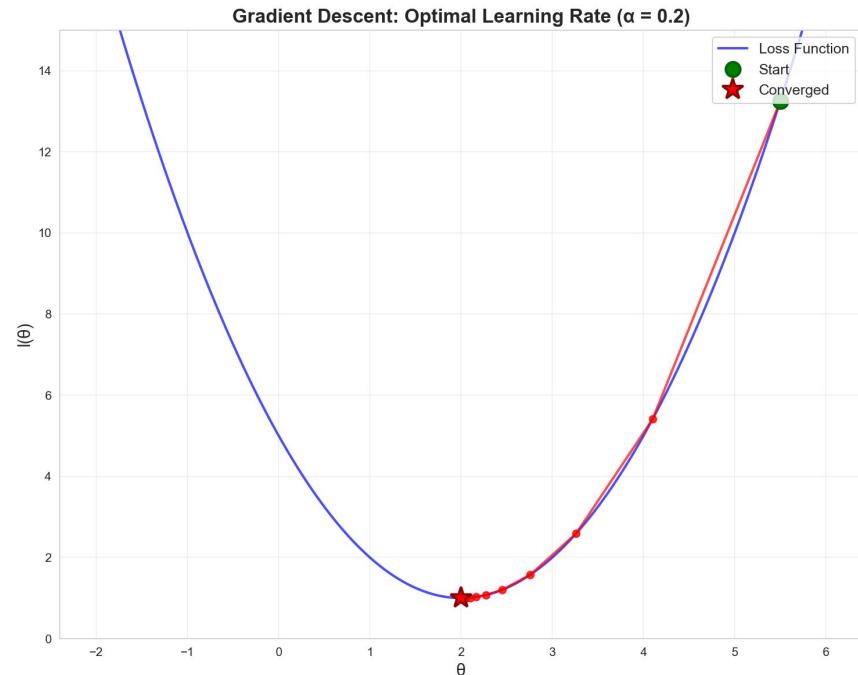
Initialize parameters randomly

1. Update Rule

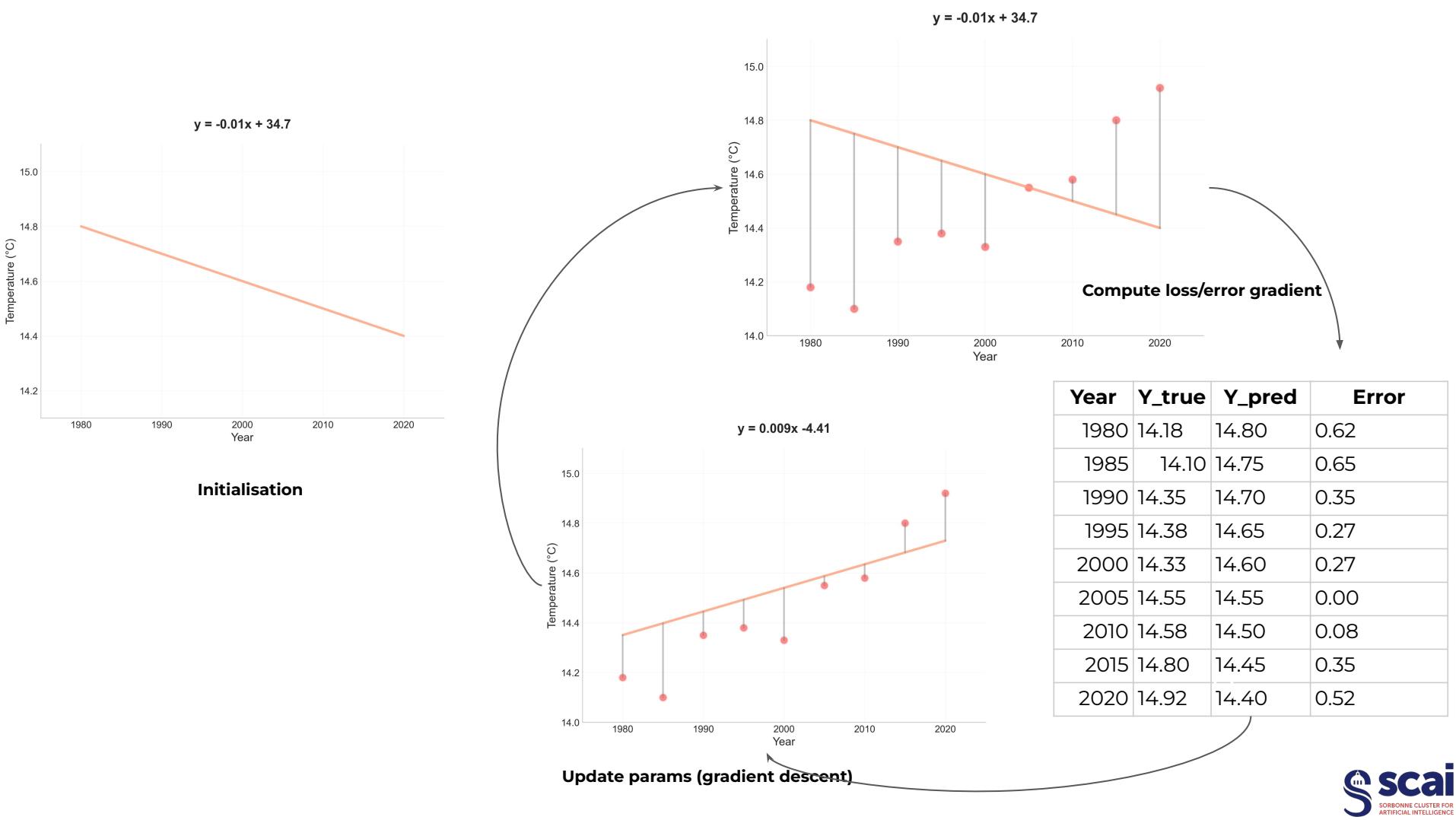
$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} l(\theta_t)$$

1. Iterate

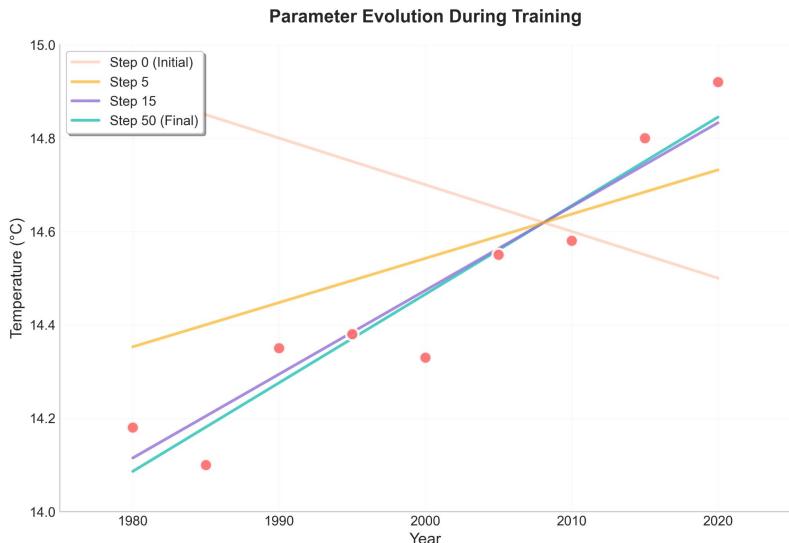
Repeat 1. while $l(\theta_{t+1}) < l(\theta_t)$



$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \ell(\theta) = \arg \min_{\theta \in \mathbb{R}^p} l(Y, f_\theta(X))$$

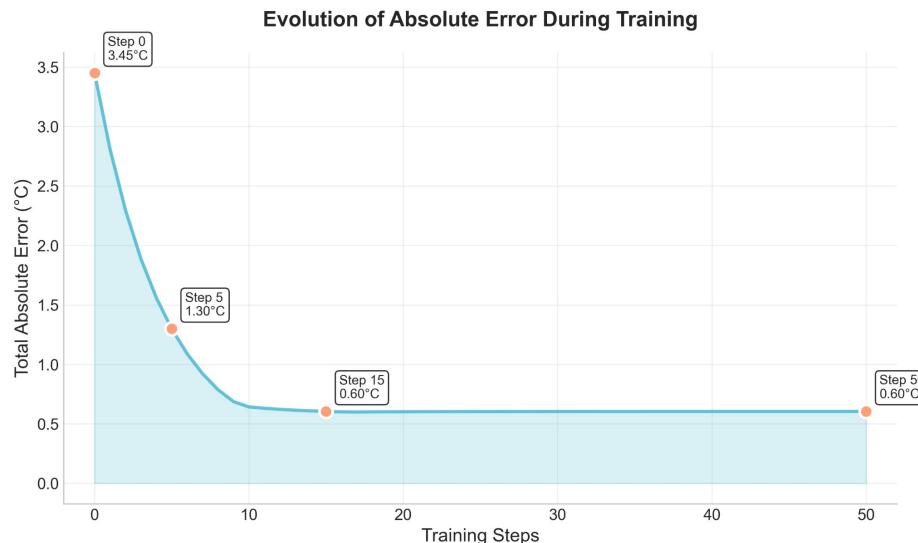


ML : Training example



Initial Model
 $a=-0.01, b=34.70$

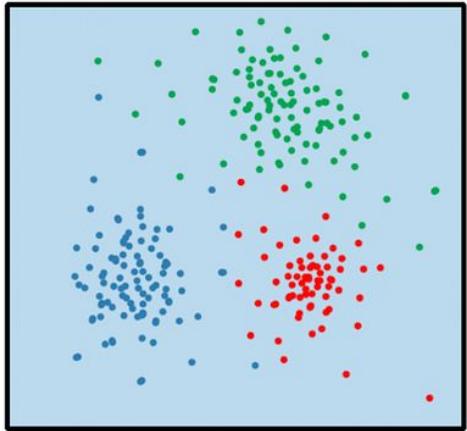
Trained Model
 $a=0.019, b=-23.47$



Initial model prediction for 2030
 $-0.01 * 2030 + 34.70 = 14.4$

Trained model prediction for 2030
 $0.019 * 2030 - 23.47 = 15.1$

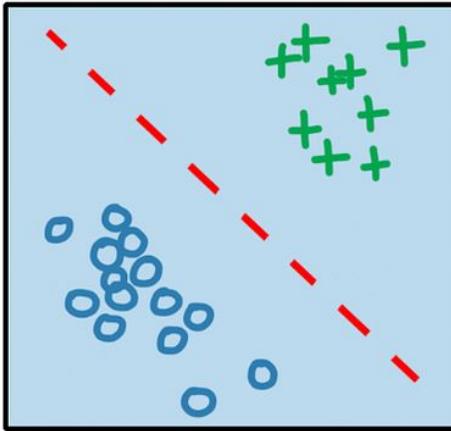
Types of Learning Problems



Unsupervised Learning:

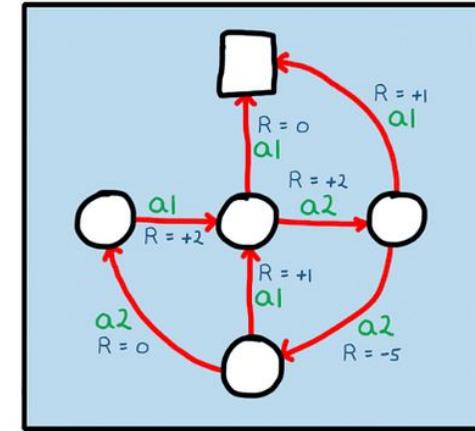
Learning structure from unlabeled data x

Clustering: Group similar data points



Supervised Learning:

Learning from labeled pairs
features target (x,y)



Reinforcement Learning:

Learning through interaction
and rewards in an environment

For most of the cases, we will reformulate the problem as an optimization problem

$$\arg \min_{f \in \mathcal{F}} \ell(f(X))$$

Non Parametric Models

Parametric Models:

Fixed number of parameters
regardless of dataset size

Learning: gradient descent

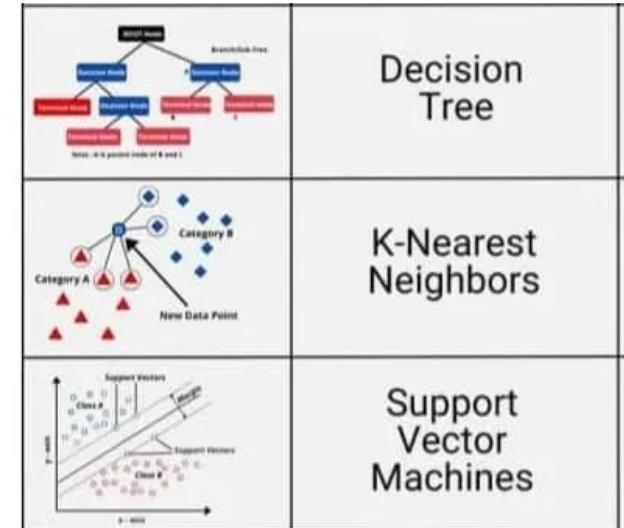
$$\arg \min_{\theta \in \mathbb{R}^p} \ell(Y, f_\theta(X))$$

Non-Parametric Models:

Number of parameters grows
with data

Learning: Various techniques
(not gradient descent)

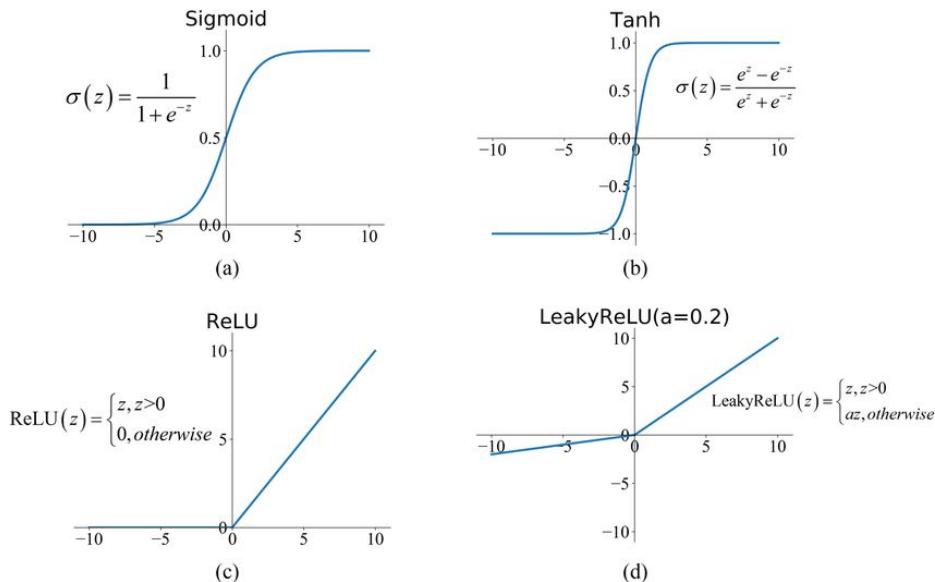
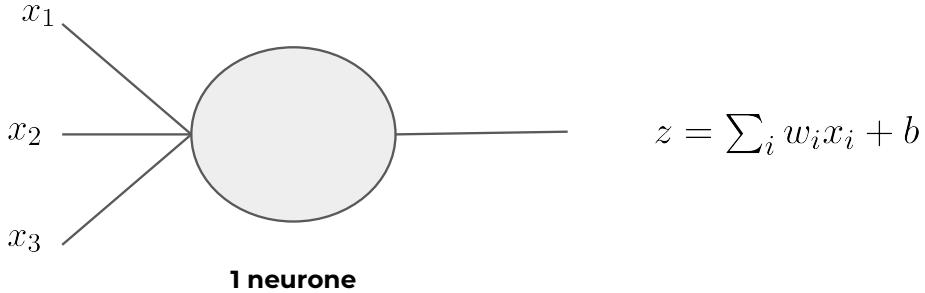
$$\arg \min_{f \in \mathcal{F}} \ell(Y, f(X))$$



The objective remains the same, the way we minimize will differ

Deep Learning

Single Neuron unit



Multi-Layer Perceptron (MLP)

A model like any other

Input: $X = (x_1, \dots, x_{r_0})$

Parameters:

$W^1, \dots, W^l \quad W^k \in \mathbb{R}^{r_{k-1} \times r_k}$

l Layers

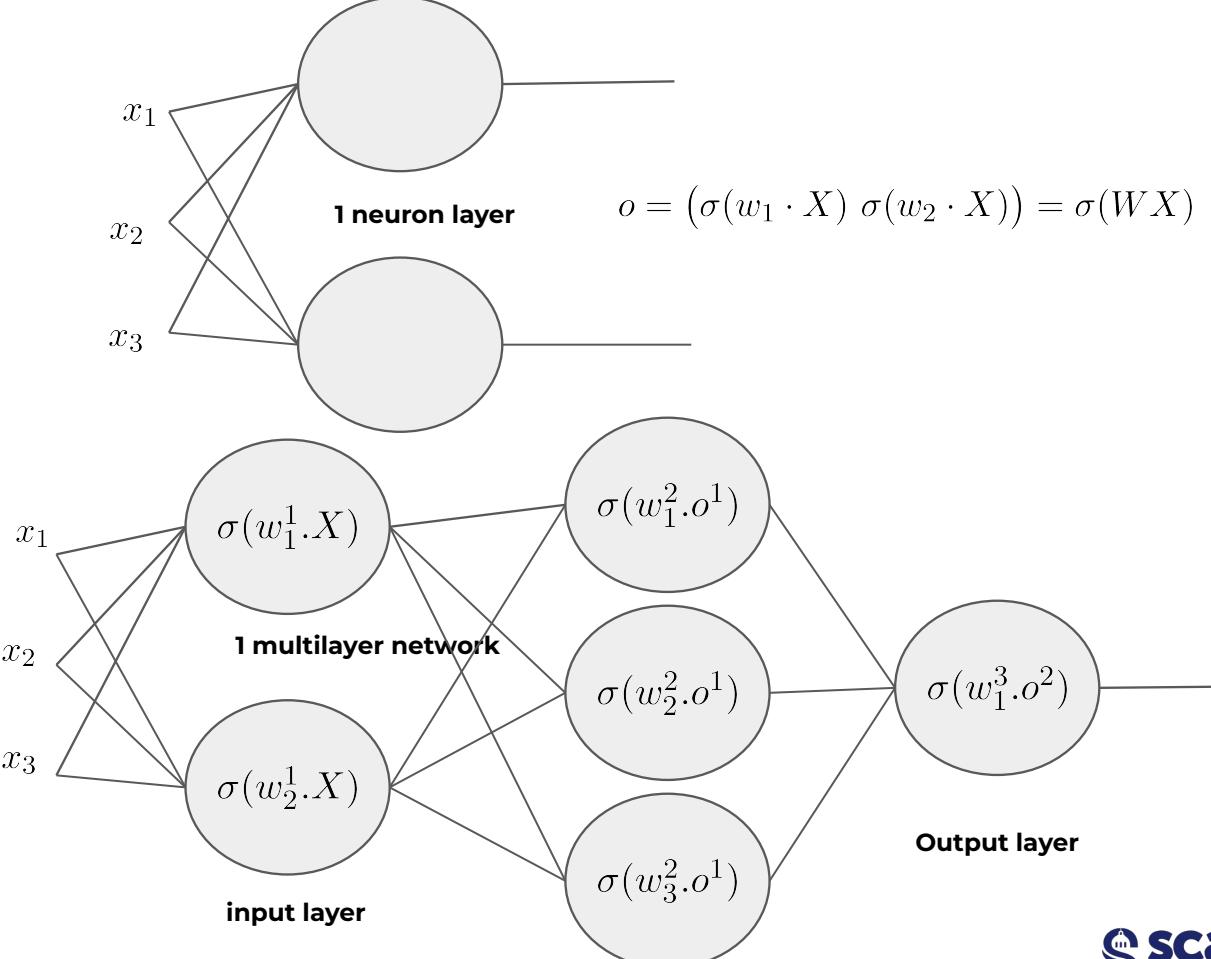
r_k neurons per layer

total = $\sum_{k=1}^l r_k \cdot (r_{k-1} + 1)$

Operation:

$f_\theta(X) = W^l \sigma(W^{l-1} \sigma(\dots \sigma(W^1 X)))$

Output: number of neuron in last layer



Training Neural Networks

Same objective and training as other parametric model :

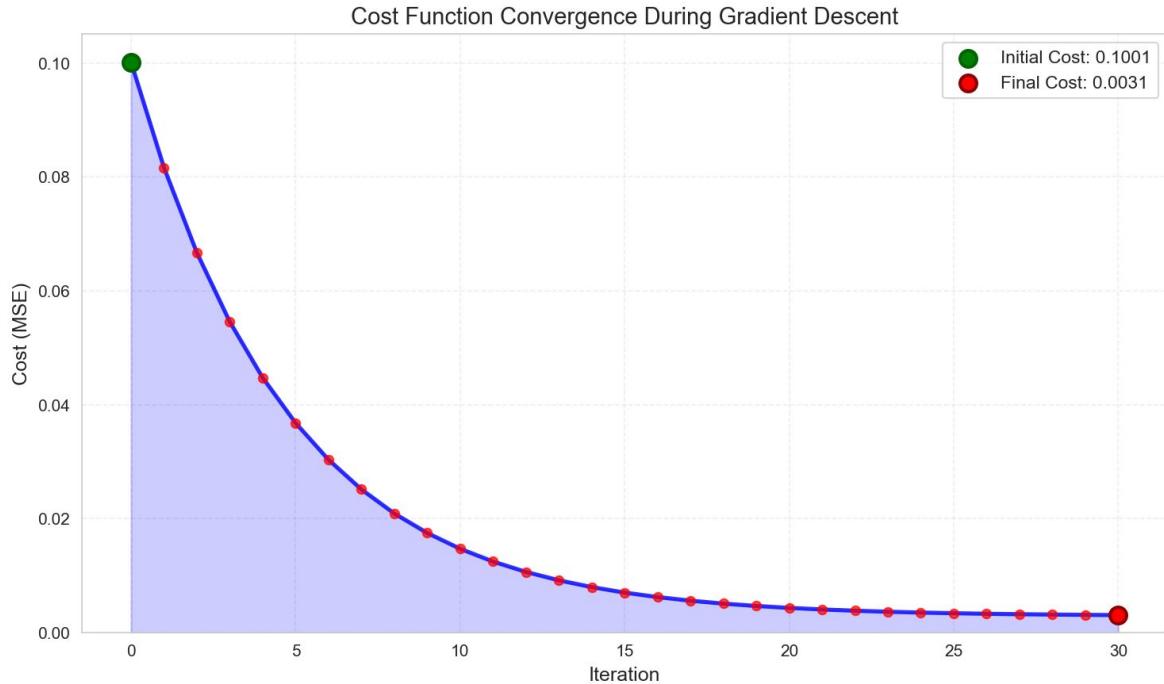
0. initialise the weights randomly

Train:

1. Compute the gradient of the loss
2. Update the parameters using gradient descent
3. Iterate

Predict:

Use the model to predict



$$\arg \min_{\theta \in \mathbb{R}^p} \ell(Y, f_\theta(X))$$

Backpropagation

Compute directly the gradient of the loss with NN can be taught

$$f_{\theta}(X) = W^l \sigma(W^{l-1} \sigma(\dots \sigma(W^1 X)))$$

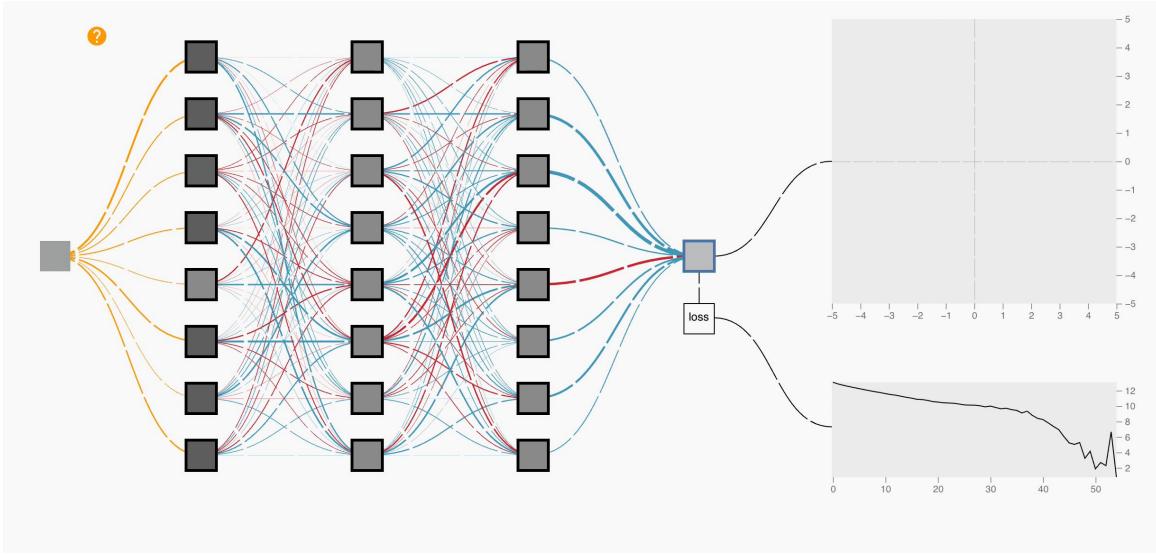
$$\nabla \ell(Y, f_{\theta}(X)) = \frac{\partial \ell(Y, f(X))}{\partial w_i^k} = ?$$

$$\forall k \in [0, l], \forall i \in [0, r_k]$$

Using chain rules

$$\frac{\partial z}{\partial x} = \sum_i \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

We are able to compute them efficiently



Step 1 : Forward Pass

For each Layer k, compute pre-activation and activation value

Step 2 : Backward Pass

Compute gradient using value from forward and next layer gradient value

Step 3: Apply Gradient Descent

Neural Networks became a standard



Key advantages:

- **Modularity:** very varied and adaptable architectures
- **Parallelization:** fast and efficient training
- **Performance:** modeling of complex information

Computer Vision

2012 — AlexNet wins ImageNet by a massive margin (15.3% vs 26.2% error)

Game Playing

2016 — AlphaGo defeats Lee Sedol

Vlad Tenev
@vladtenev

Another one bites the dust!

Last night, @AcerFur used Aristotle from @HarmonicMath to autonomously prove in @leanprover Erdos Problem #481, which had been open for over 45 years.

A systematic approach targeting all unsolved Erdos problems is in the works.

Speech Recognition

2012 — Deep neural networks replace GMM-HMMs as the standard for acoustic modeling

Machine Translation

2017 — Transformer architecture ("Attention Is All You Need") becomes the new paradigm

Data Science

Data: The Fuel of Deep Learning

Volume: Deep networks typically require thousands to millions of training examples.

Quality: Clean, accurate data are crucial - error will reduce strongly model performances.

Bias: Systematic distortions in the data will be learned and reproduced by the model.

Diversity: Training data must capture the full spectrum of real-world scenarios.



Data Preprocessing

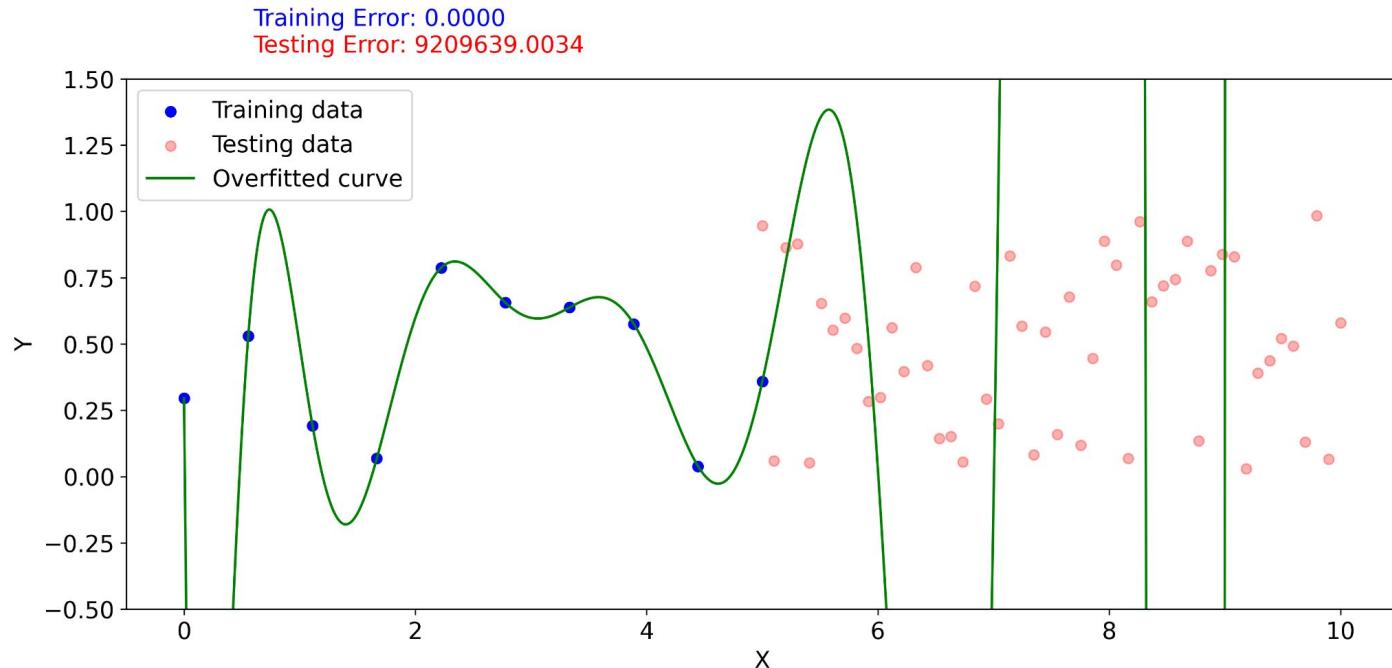
Steps in Data Preprocessing :

- Handle Inconsistencies
- Handle Missing Values
- Handle Categorical Values
- Handle Duplicates
- Handle Outliers
- Feature Engineering
- Scaling and Normalization
- Feature Selection and Dimensionality Reduction

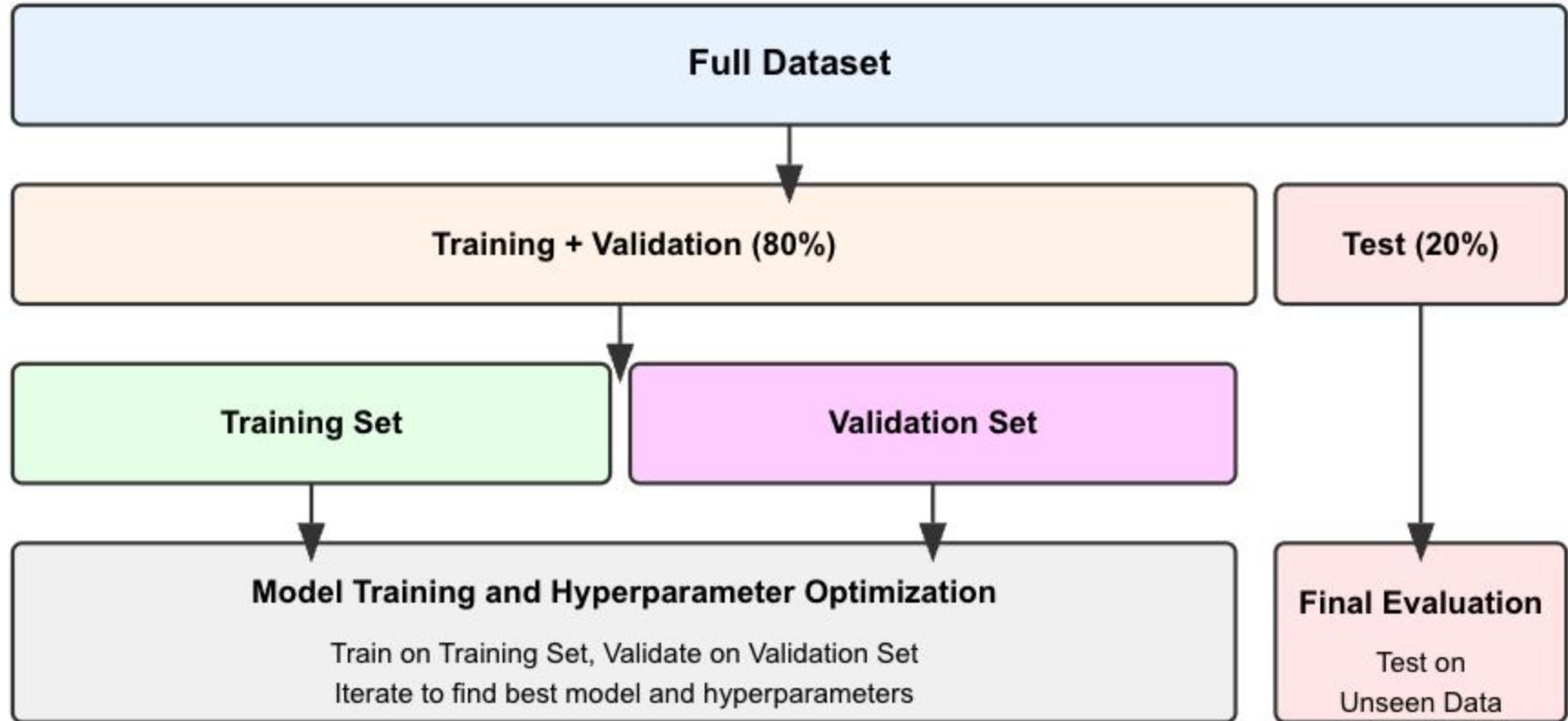


Is the model learning something?

$$\arg \min_{\theta \in \mathbb{R}^p} \ell(Y, f_\theta(X))$$

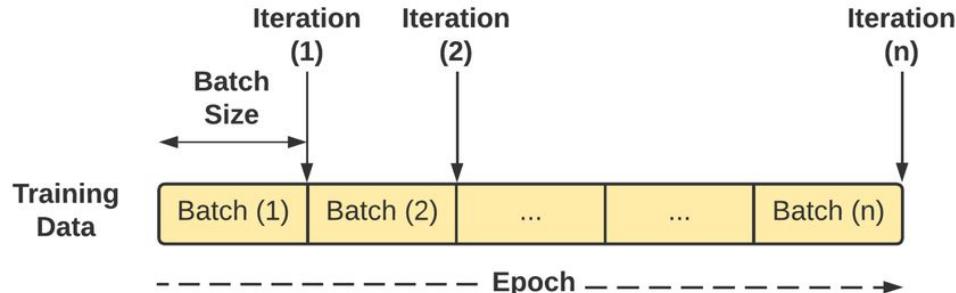


Model Evaluation

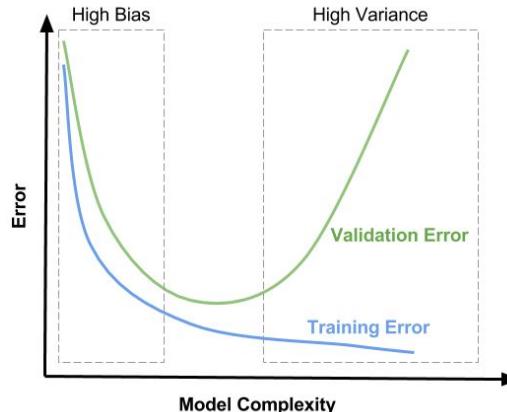


Science of hyperparameters

Hyperparameters: are settings chosen that control how a model learns



Batch size: Number of samples processed before updating model parameters.
Epoch: One complete pass through the entire training dataset.



Overfitting: model learns the training data too well, including its noise and peculiarities, resulting in poor generalization to unseen data.

Loss choice

The Role of Loss Functions

Loss functions quantify the error/distance between real y targets and our model outputs \hat{y}

Mean Squared Error (MSE): Binary Cross-Entropy:

$$\ell_{MSE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\ell_{BCE}(y, \hat{y}) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

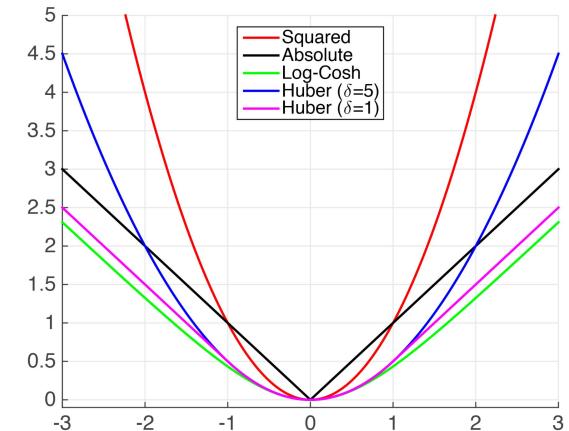
Mean Absolute Error (MAE): Cross-Entropy Loss:

$$\ell_{MAE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

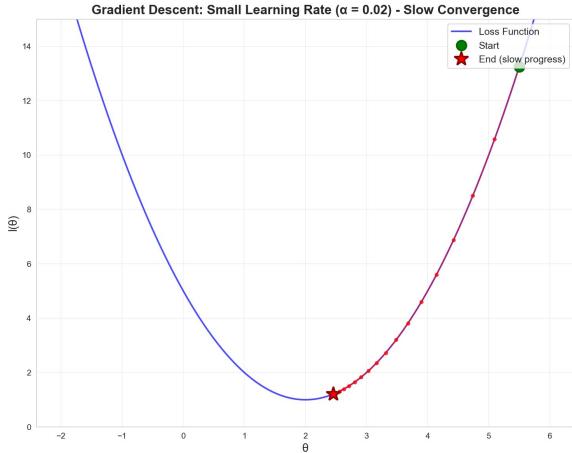
$$\ell_{CE}(y, \hat{y}) = -\sum_{i=1}^n \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic})$$

Custom Loss (eg: Focal Loss for unbalanced classification)

$$\ell_{Focal}(y, \hat{y}) = -\alpha(1 - \hat{y})^\gamma y \log(\hat{y}) - (1 - \alpha)\hat{y}^\gamma(1 - y) \log(1 - \hat{y})$$

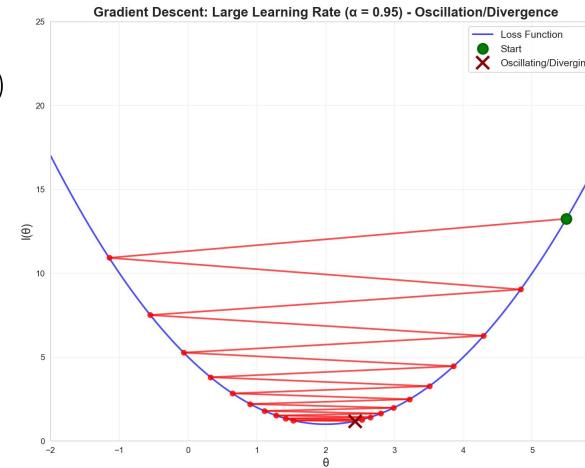


Learning rate and batch size

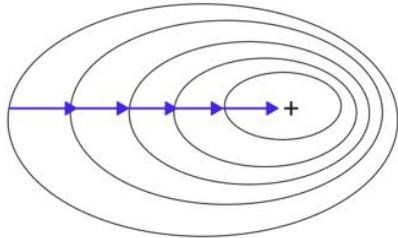


$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \ell(\theta) = \arg \min_{\theta \in \mathbb{R}^p} l(Y, f_\theta(X))$$

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} l(\theta_t)$$

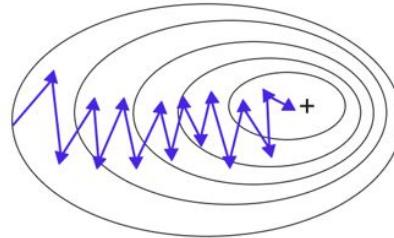


Batch Gradient Descent



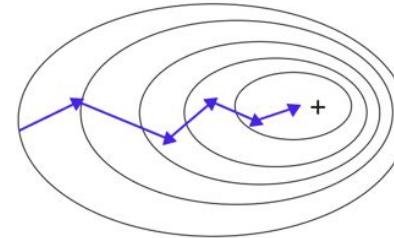
$$\nabla_{\theta} \mathcal{L} = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(f_{\theta}(x_i), y_i)$$

Stochastic Gradient Descent



$$\nabla_{\theta} \mathcal{L} \approx \nabla_{\theta} \ell(f_{\theta}(x_i), y_i)$$

Mini-Batch Gradient Descent



$$\nabla_{\theta} \mathcal{L} \approx \frac{1}{m} \sum_{i \in \mathcal{B}} \nabla_{\theta} \ell(f_{\theta}(x_i), y_i)$$

Optimizer choice

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla L(w_t)$$

Gradient Descent + Momentum:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla L(w_t)$$

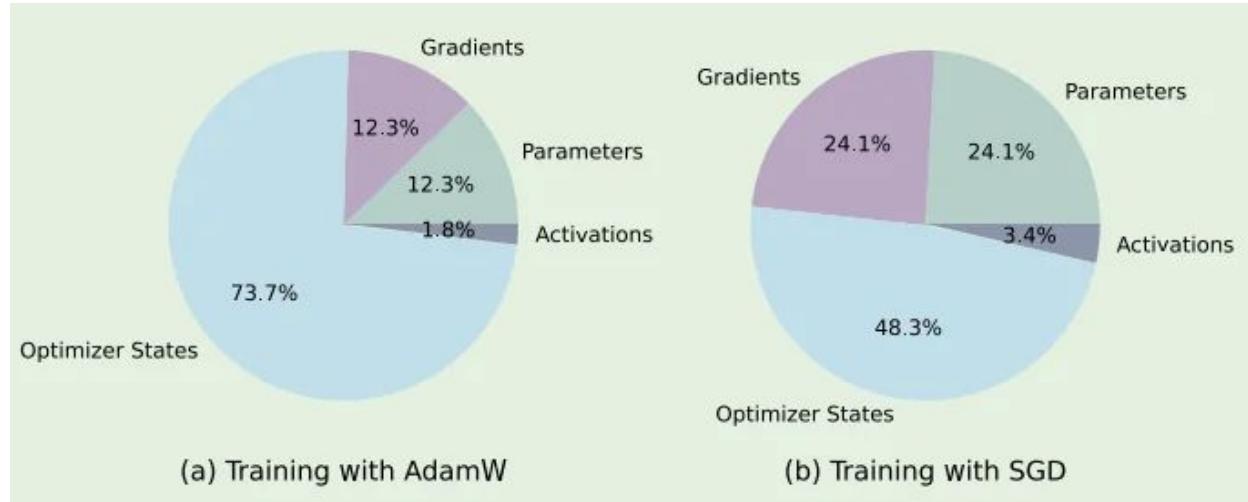
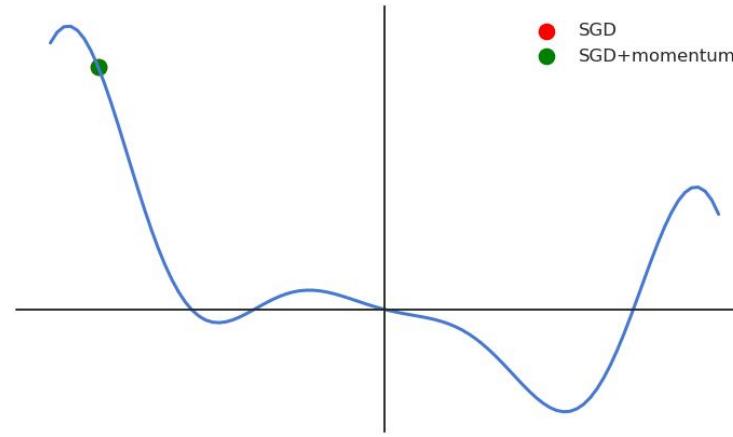
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla L(w_t))^2$$

Adam/Adamw:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$w_{t+1} = w_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$



Linear: $y = xW^T + b$

The fundamental building : Applies a weighted sum of inputs plus bias.

Embedding: $E(i) = W_i = [w_{i,1}, w_{i,2}, \dots, w_{i,d}]$

Maps discrete value to dense continuous vectors

Convolution: $y_{i,j,c} = \sum_{m,n,k} x_{i+s+m, j+s+n, k} \cdot W_{m,n,k,c} + b_c$

Detect spatial patterns like edges, textures, and shapes

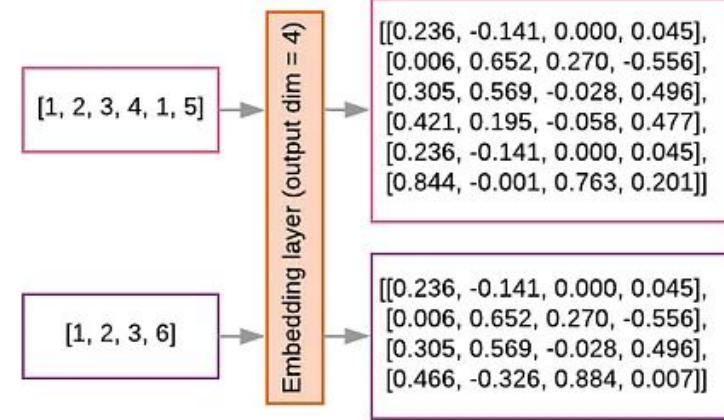
Recurrent: $y^{(t)} = \tanh(W_{xh}x^{(t)} + W_{hh}h^{(t-1)} + b)$

Processes sequences by combining current input with previous hidden state

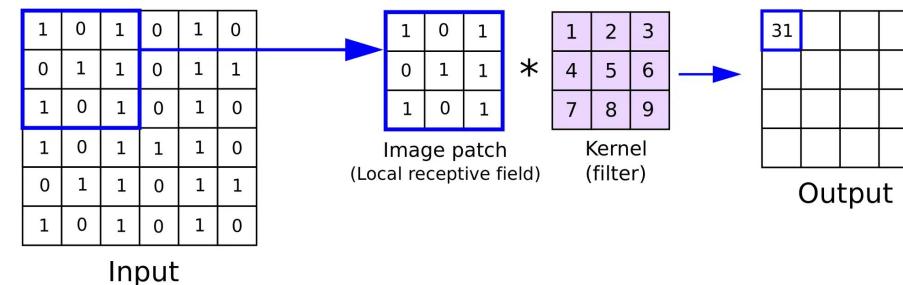
Attention: $y = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$

Computes weighted combinations of values based on query-key similarity

NN unit layers



Embedding : categorical data



CNN : image data

NN architectures

Number of Layers:

Depth of the network; more layers allow learning hierarchical representations but increase training difficulty.

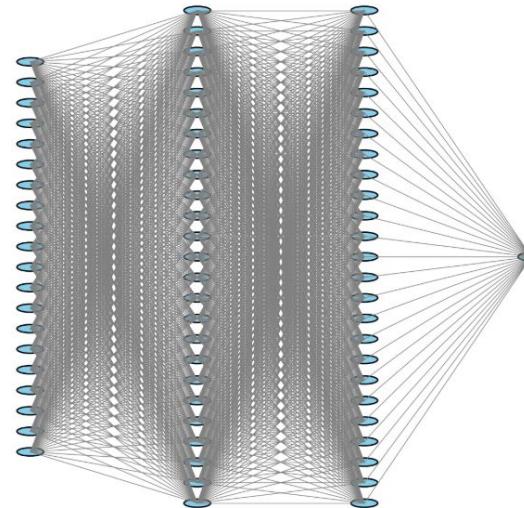
Number of neurons per Layers:

Width of each layer; wider layers capture more features but increase computational cost and risk of overfitting.

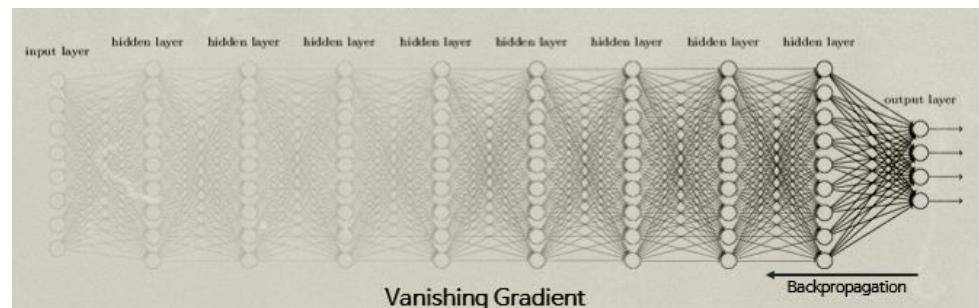
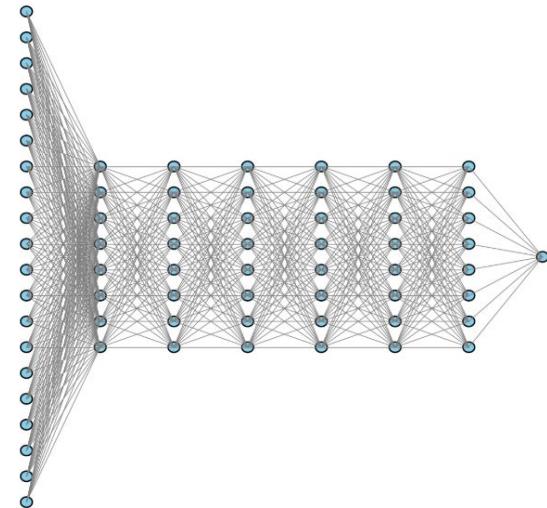
Vanishing gradient:

During backpropagation, gradients shrink exponentially through layers, making early layers learn extremely slowly or not at all.

WideMLP (Shallow & Wide)



DeepMLP (Deep & Narrow)



Architectural priors Components

Dropout:

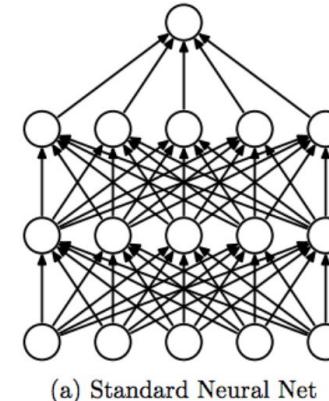
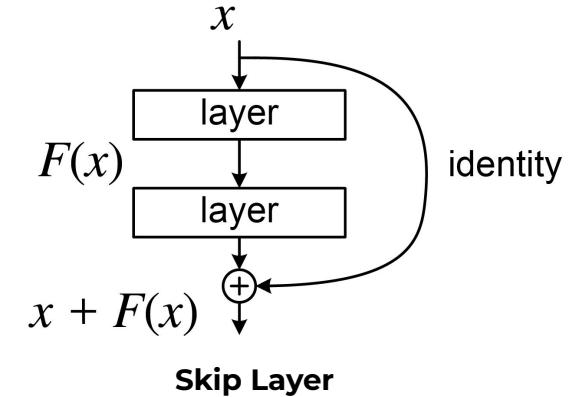
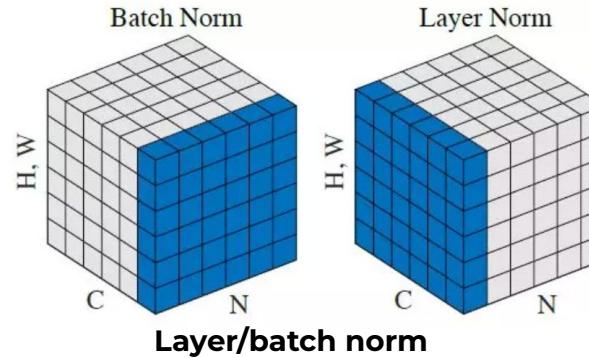
Randomly deactivates neurons during training, forcing the network to learn redundant representations and reducing overfitting.

Layer/Batch Normalization:

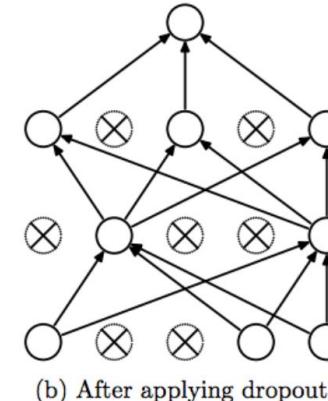
Normalizes activations across batches or features to stabilize training, allow higher learning rates, and reduce internal covariate shift.

Skip Layer:

Adds the input directly to the layer's output ($x + F(x)$), enabling gradient flow through deep networks and making it easier to learn identity mappings.



(a) Standard Neural Net



(b) After applying dropout.

Dropout

Data Science Librairies



Thanks

Metrics

