

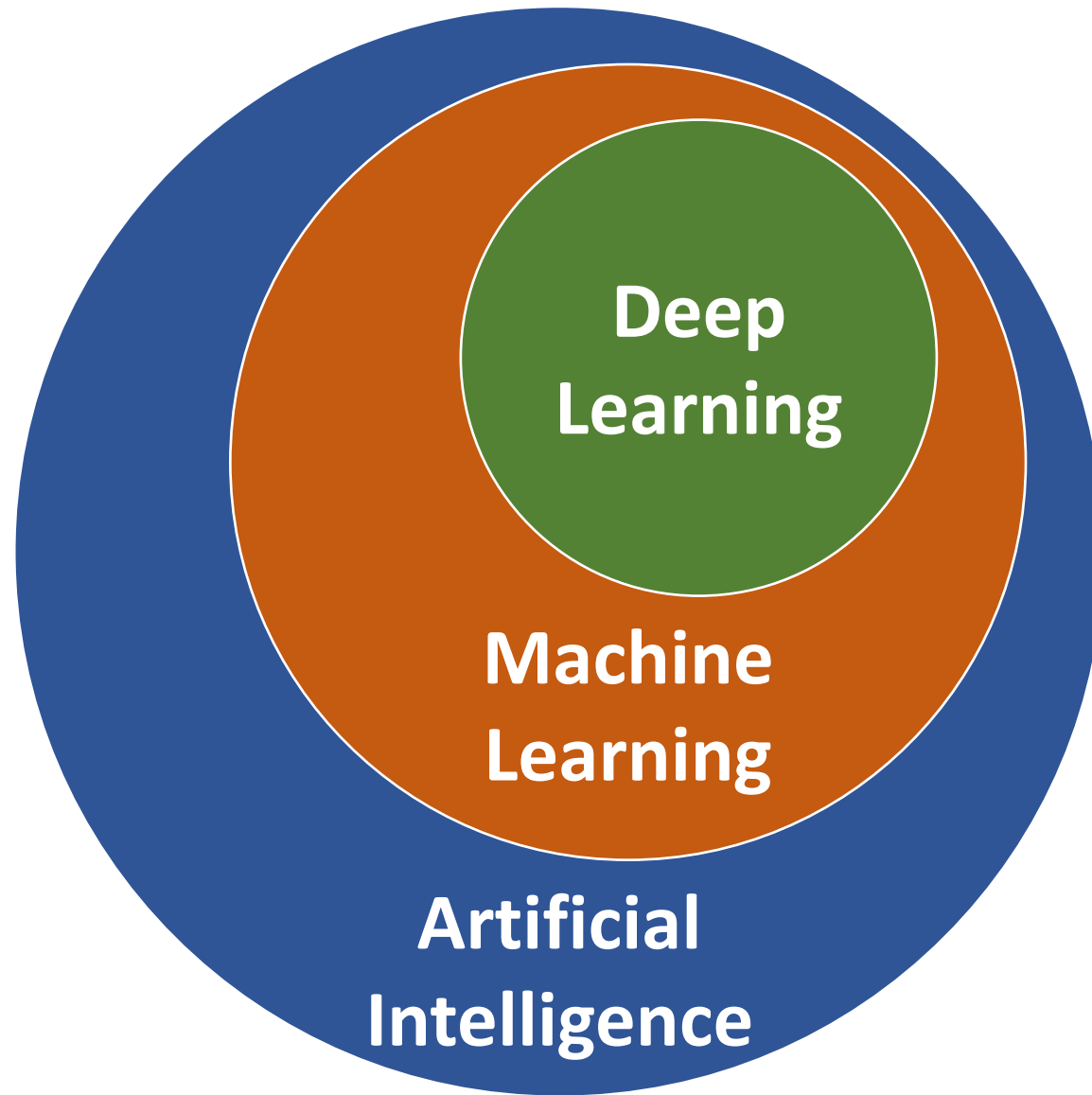


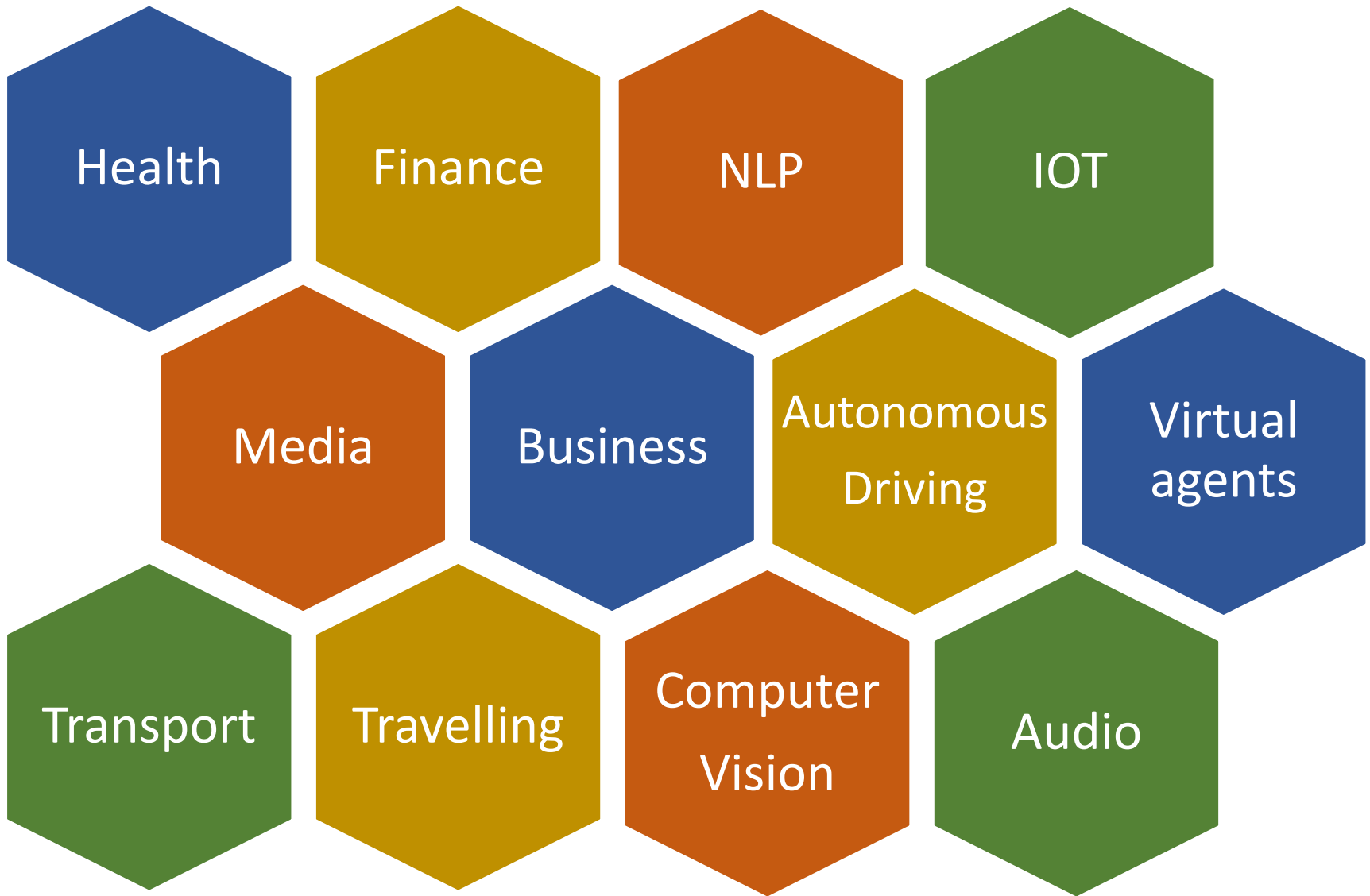
Hands-on Introduction to Deep Learning

Artificial Neural Networks



INSTITUT DU
DÉVELOPPEMENT ET DES
RESSOURCES EN
INFORMATIQUE
SCIENTIFIQUE





Application fields

IOT



- Personal assistant
- Smart connected object

Healthcare



- Medical Imaging
- Drug research

Finance

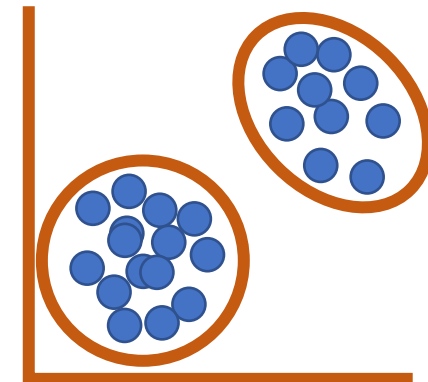
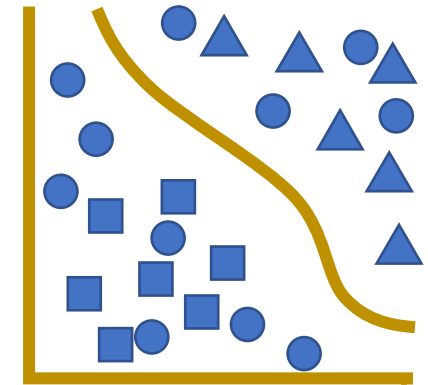
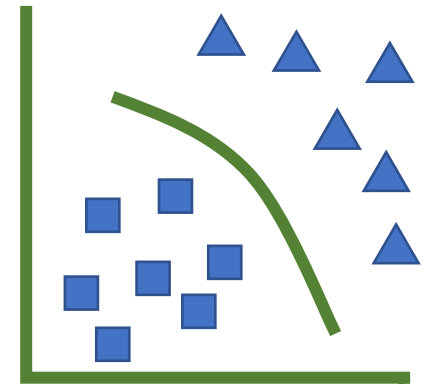
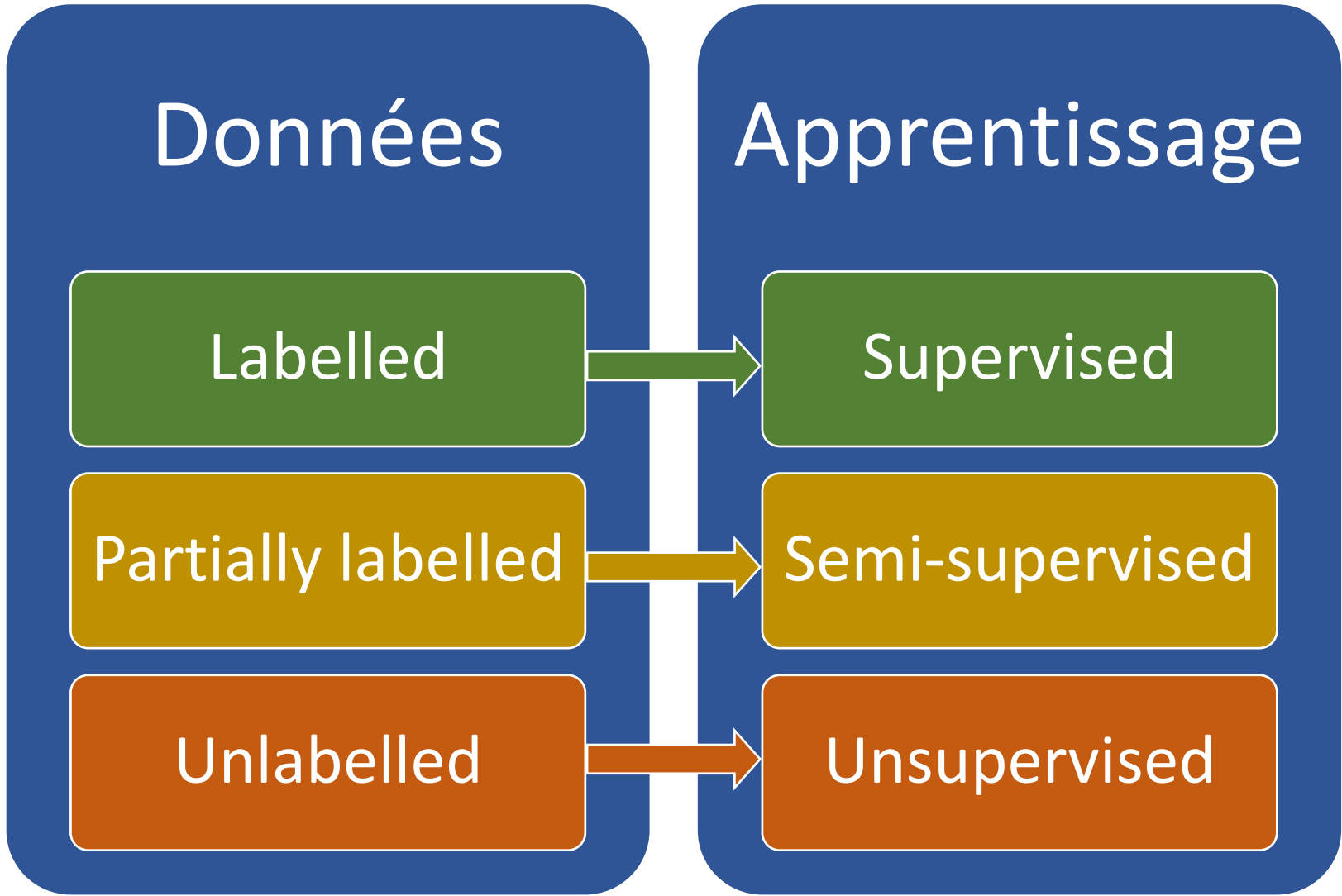


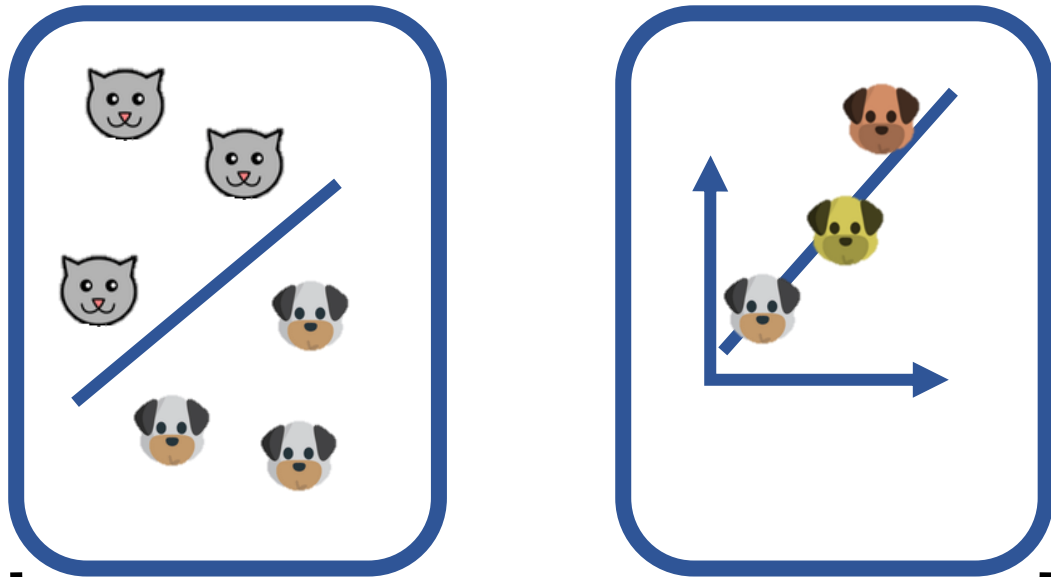
- Fraud detection
- Prediction of financial flows

Autonomous driving

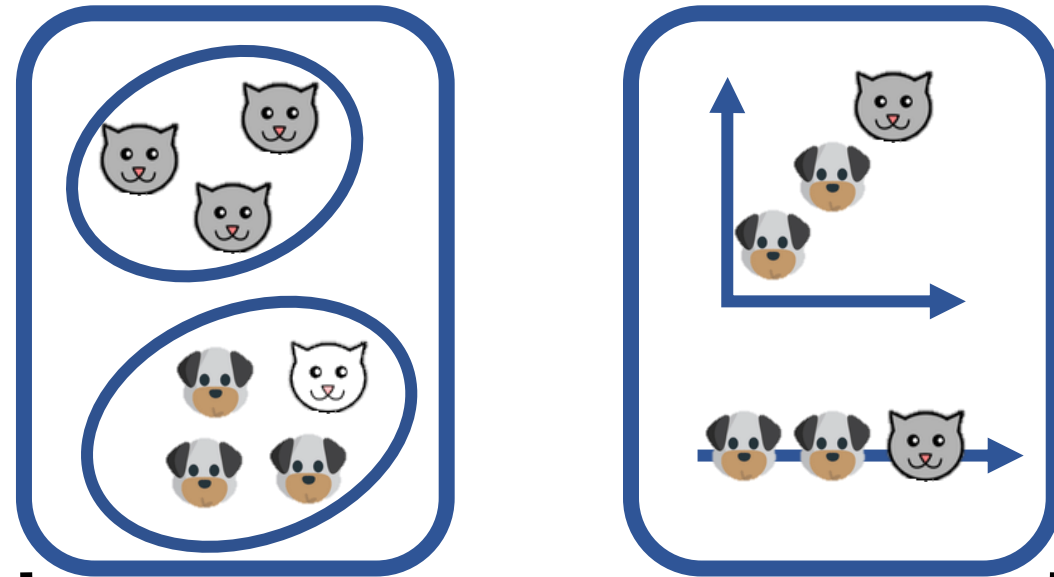


- Route optimization
- Autonomous driving

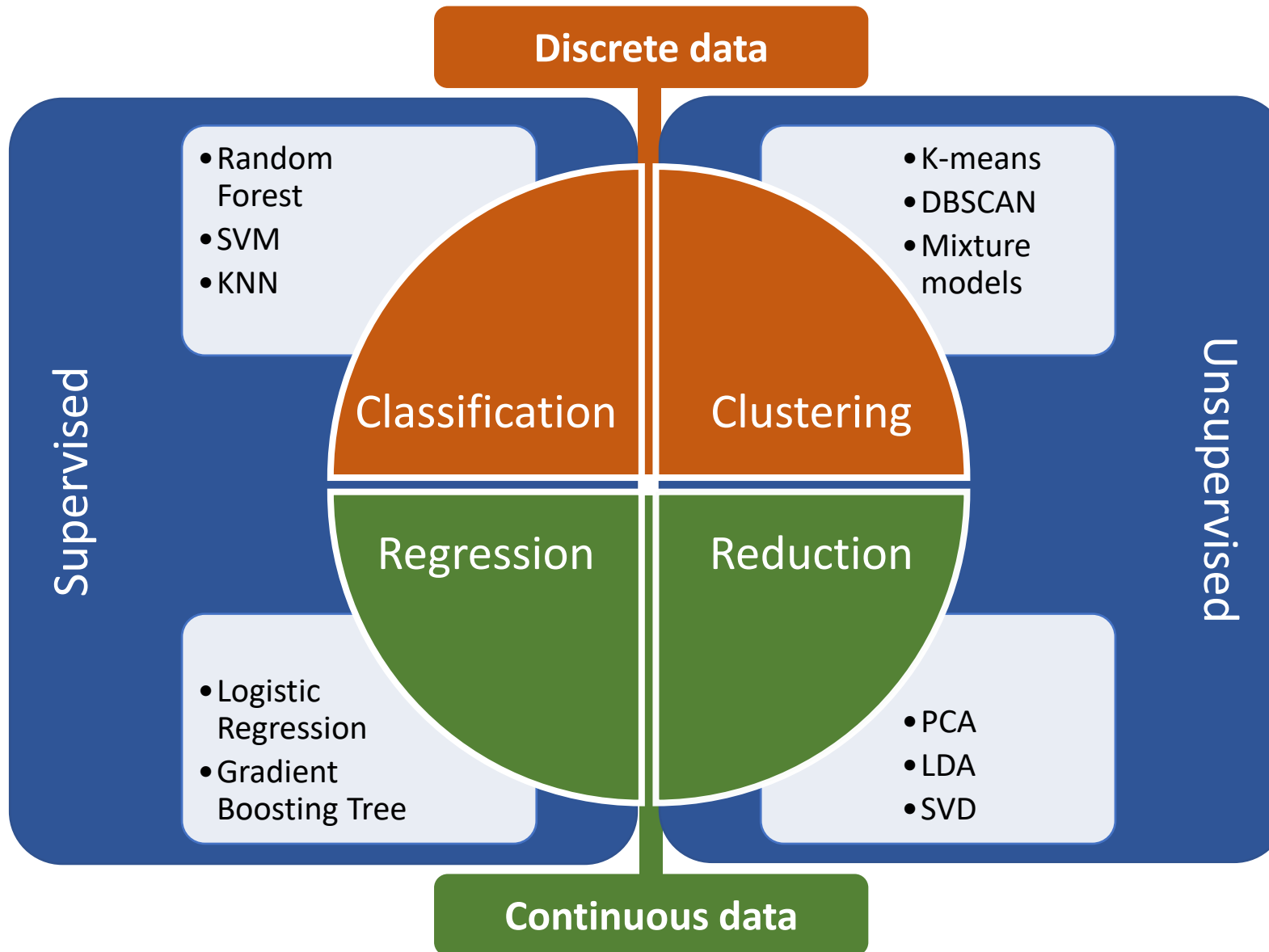


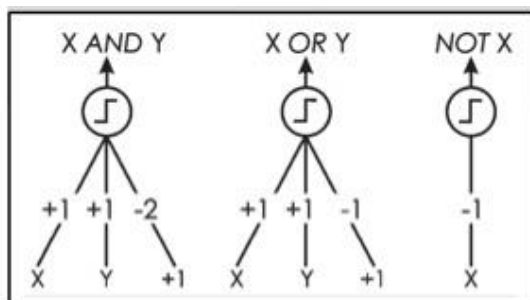


Supervised Learning

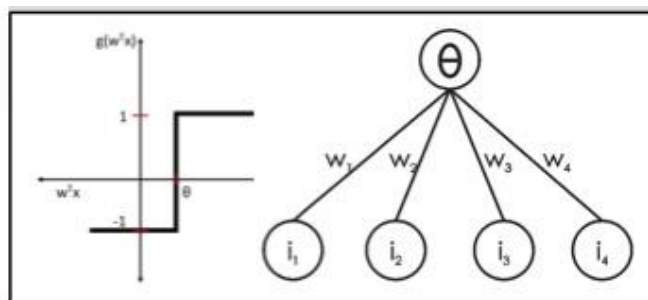


Unsupervised Learning

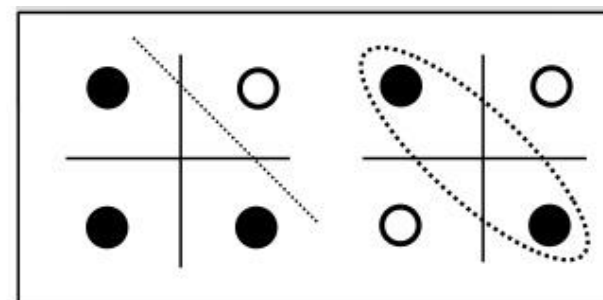




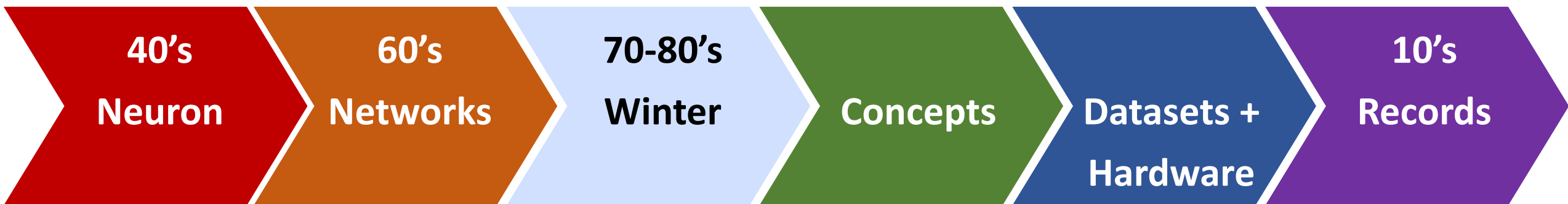
- Adjustable Weights
- Weights are not Learned

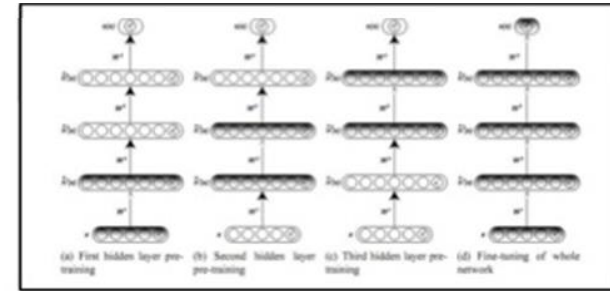
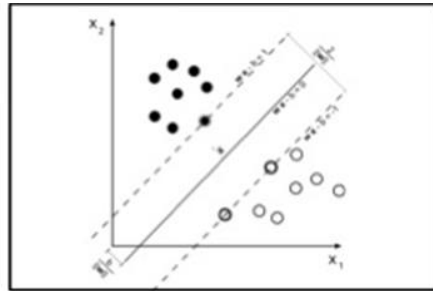
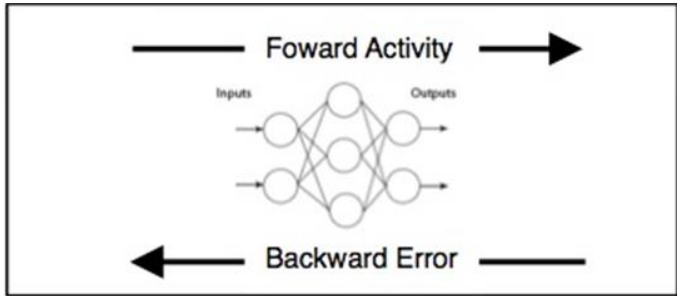


- Learnable Weights and Threshold

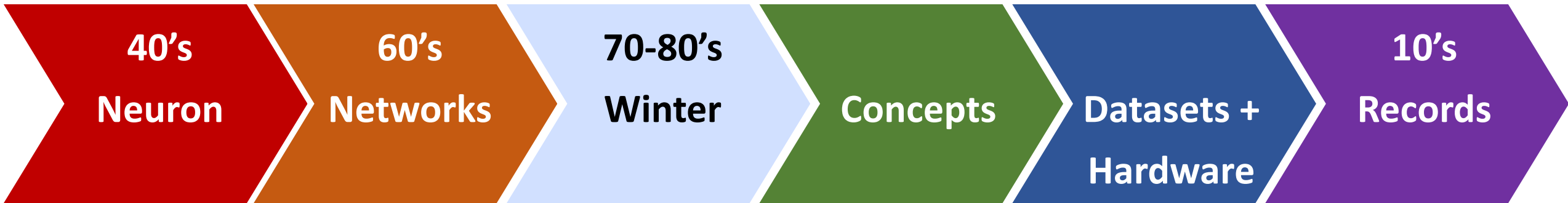


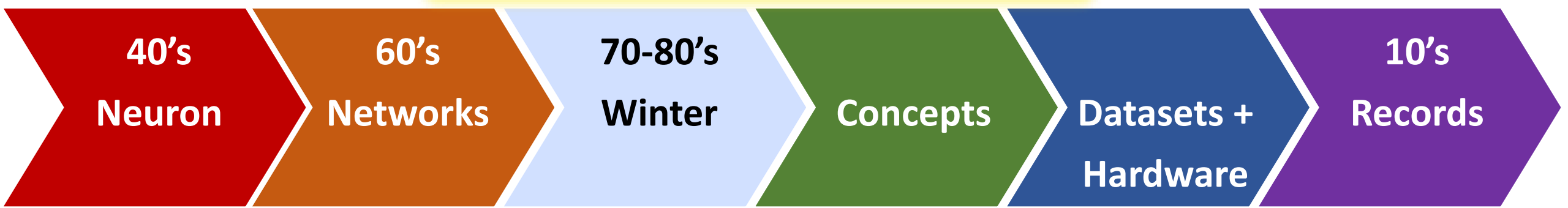
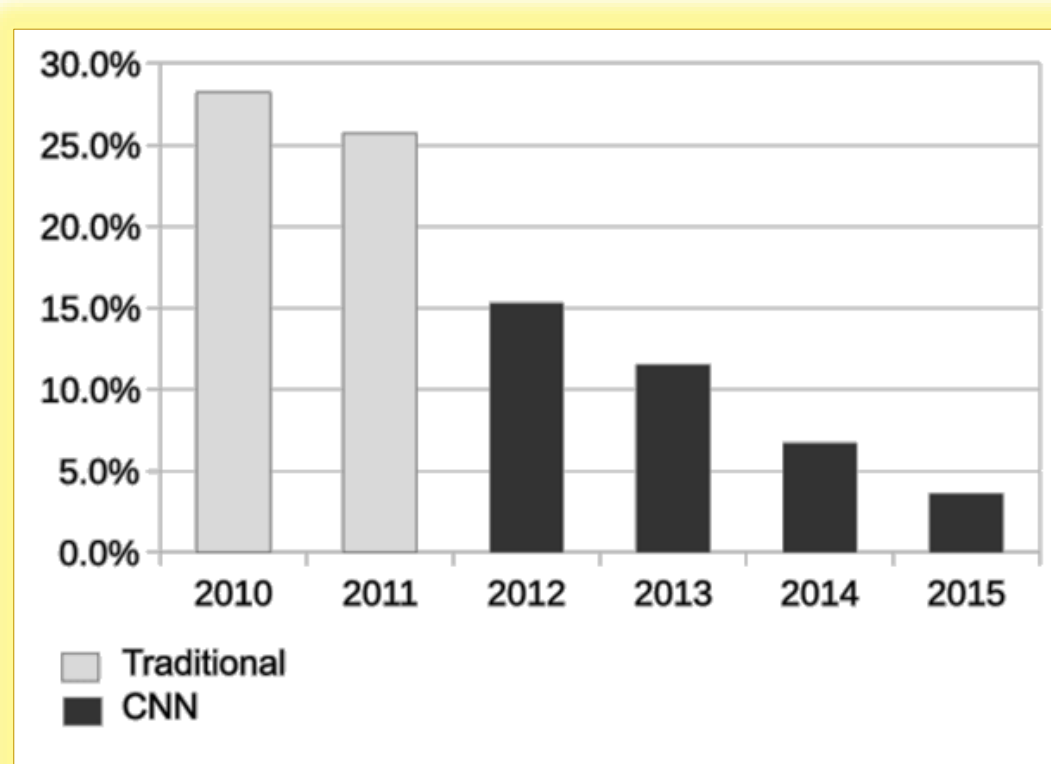
- XOR Problem



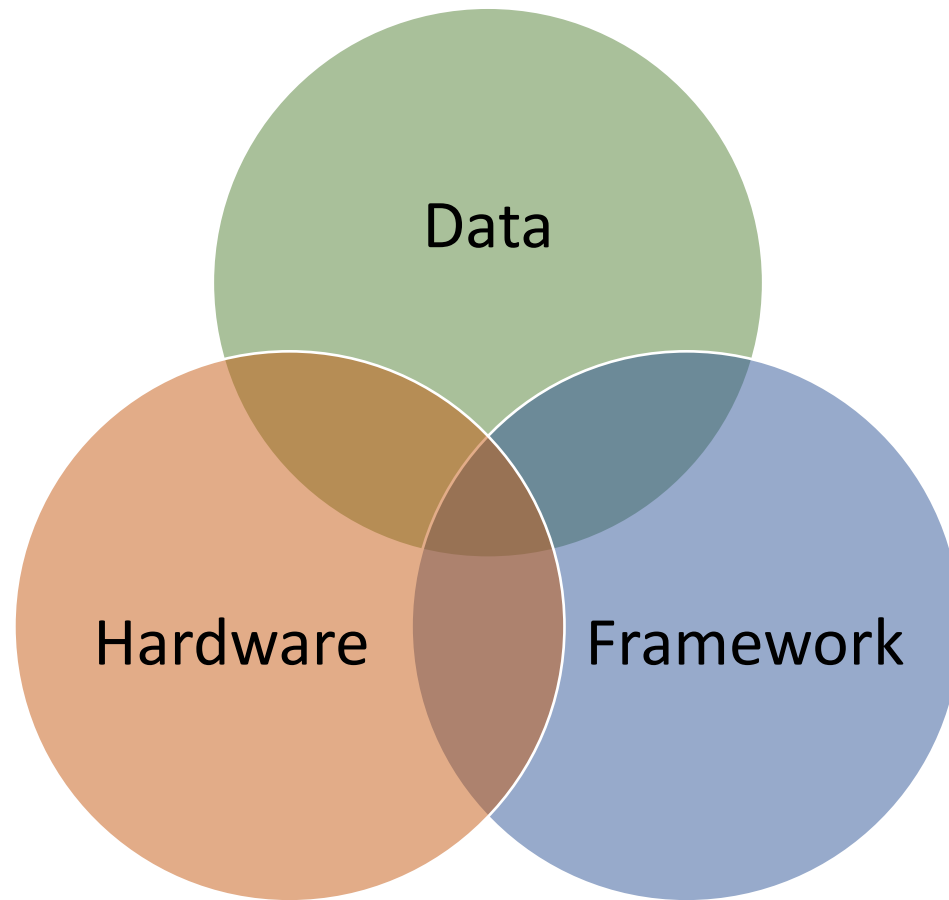


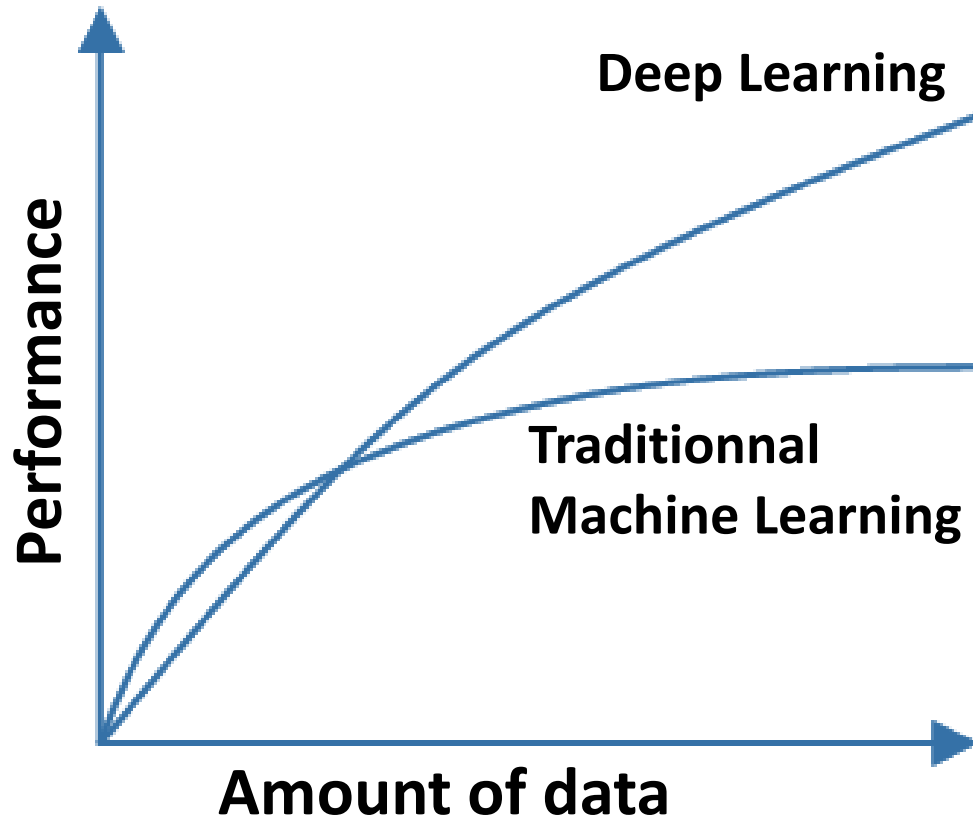
- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting
- Limitations of learning prior knowledge
- Kernel function: Human Intervention
- Hierarchical feature Learning





History of Deep Learning



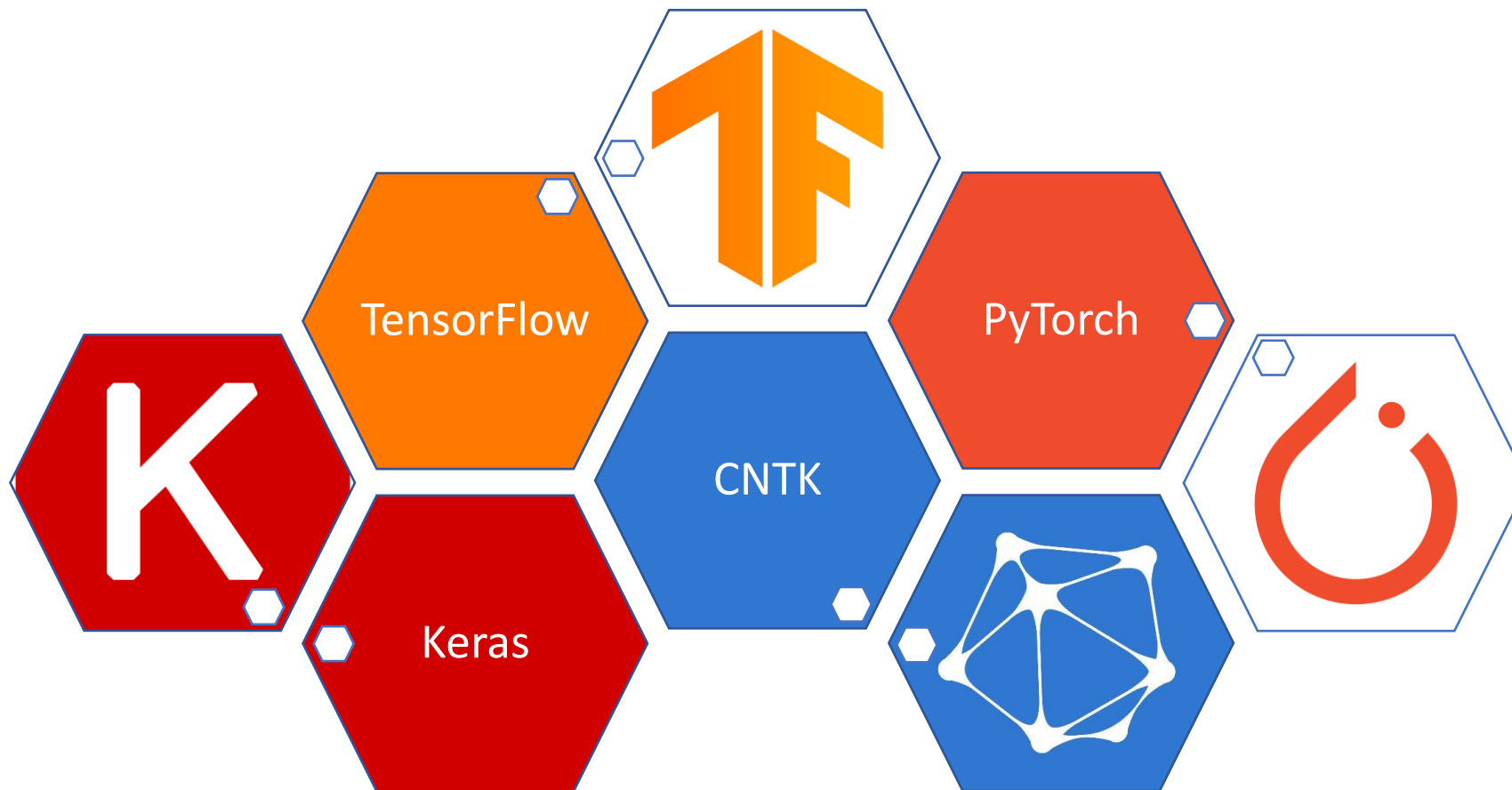


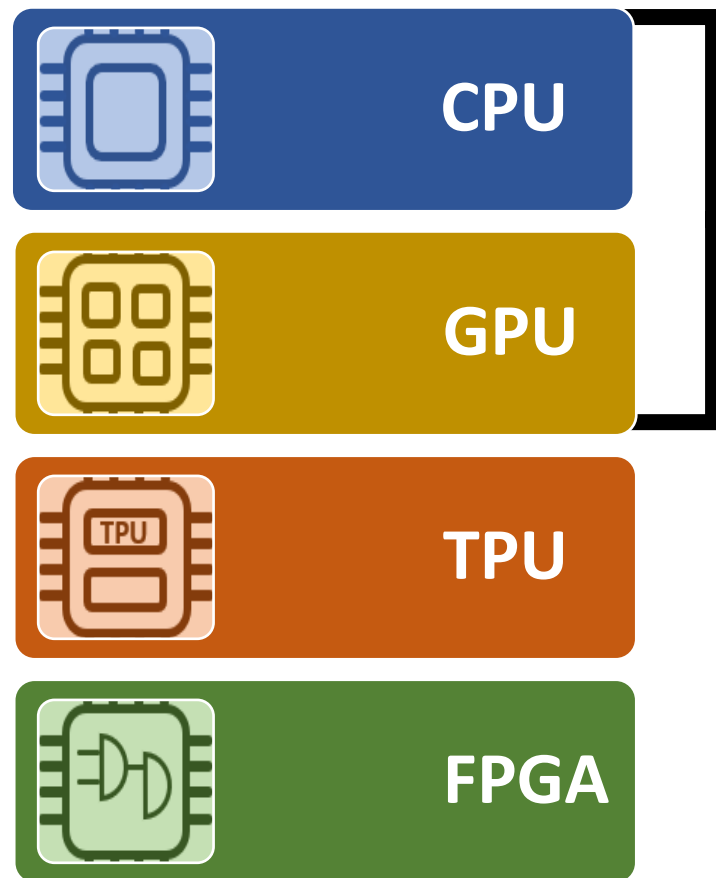
Amount

Open source

Collaboration

Competitions





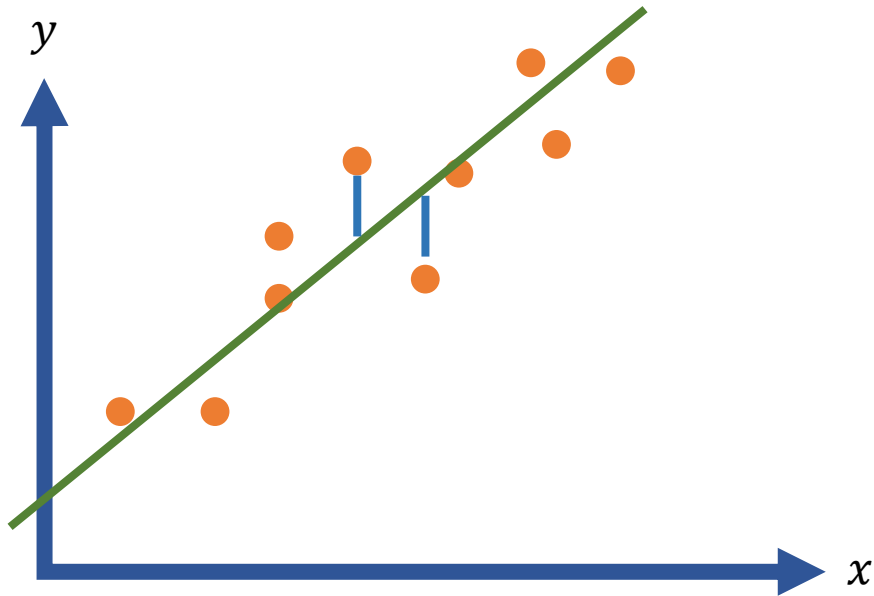
$$Y = X \cdot \Theta + N$$

$$\hat{Y} = X \cdot \hat{\Theta}$$

With :

$$\Theta = (a, b)$$

N , noise

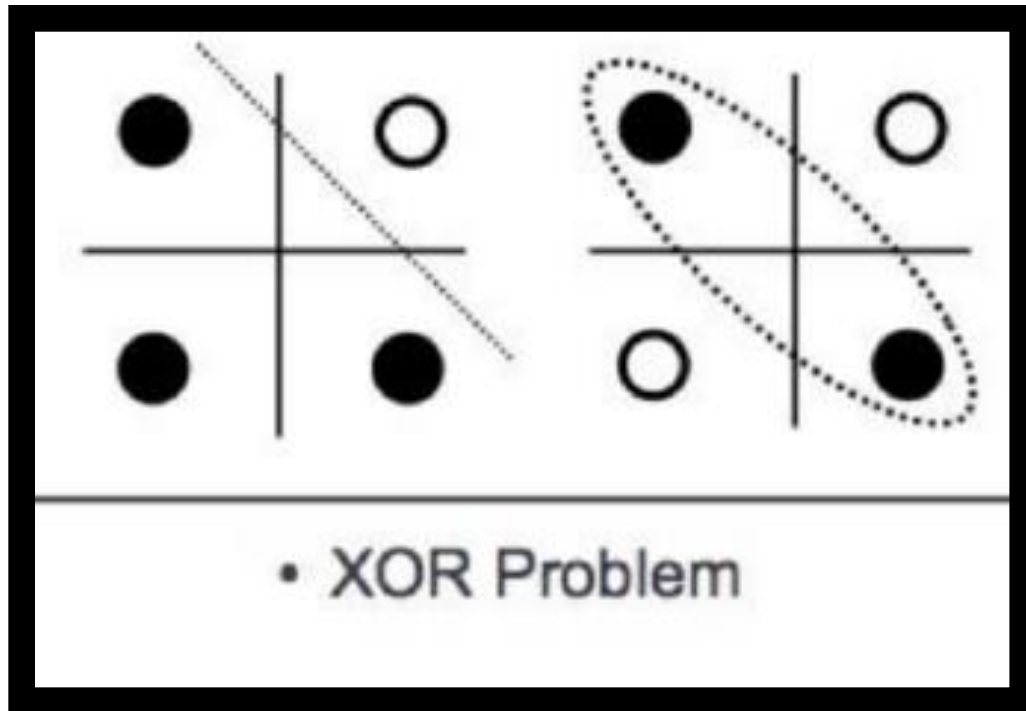


$$\min_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

**Cost
function**

Convex problem : Direct solution

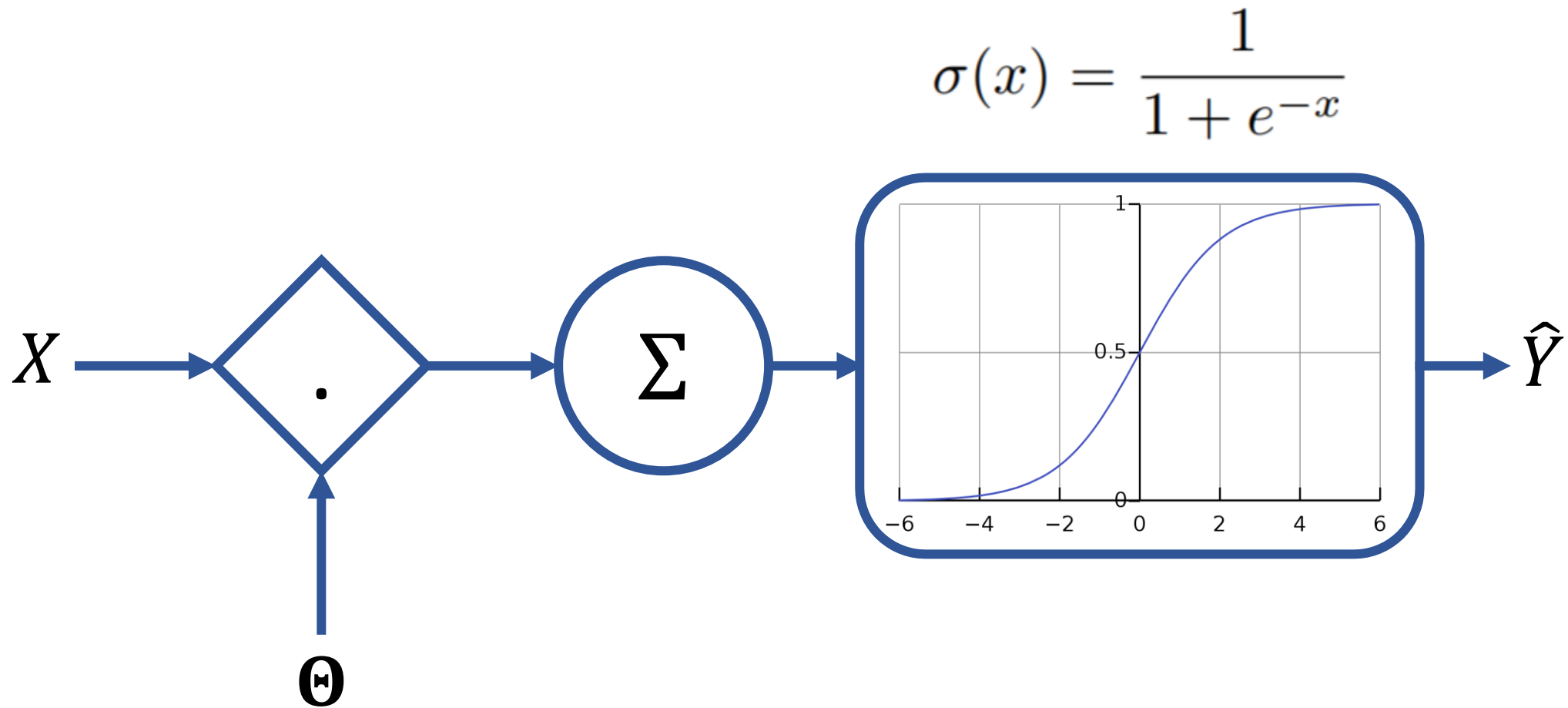
$$\bullet \hat{\Theta} = (X^T \cdot X)^{-1} \cdot X^T \cdot Y$$



$$f = W x$$

$$\widehat{W} = W_2 \cdot W$$

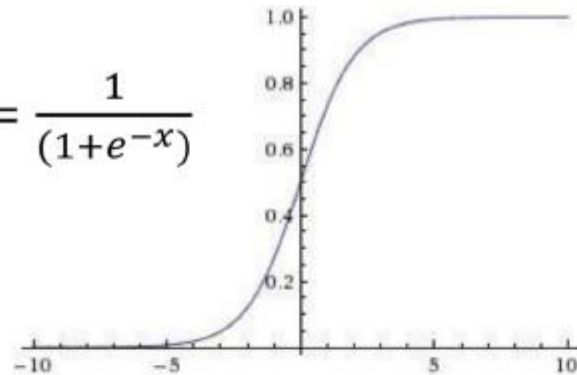
$$\hat{f} = \widehat{W} x$$



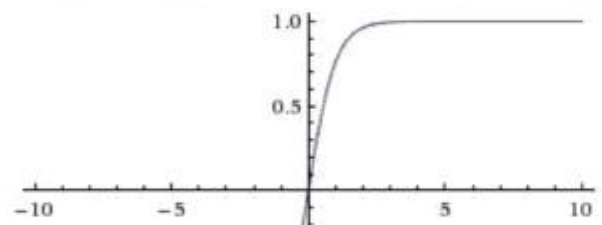
$$\mathcal{L}(\hat{y}_i, y_i) = y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

Cross-entropy loss

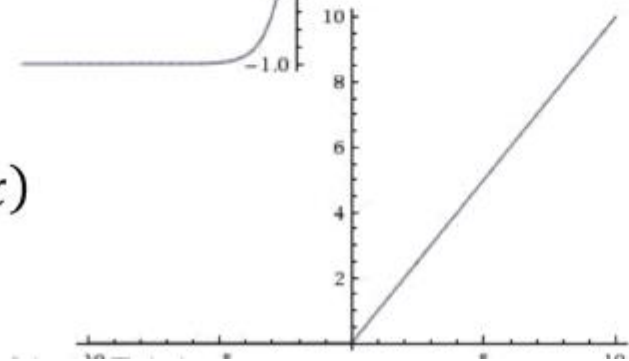
Sigmoid: $\sigma(x) = \frac{1}{(1+e^{-x})}$



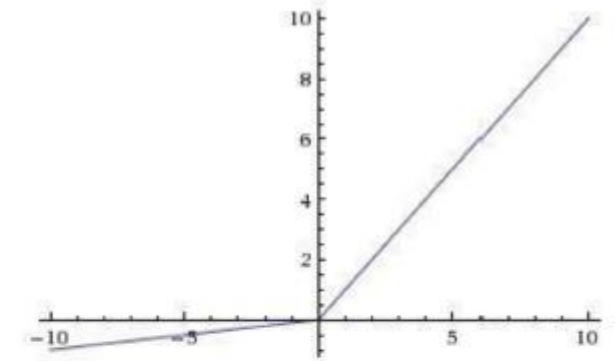
tanh: $\tanh(x)$



ReLU: $\max(0, x)$



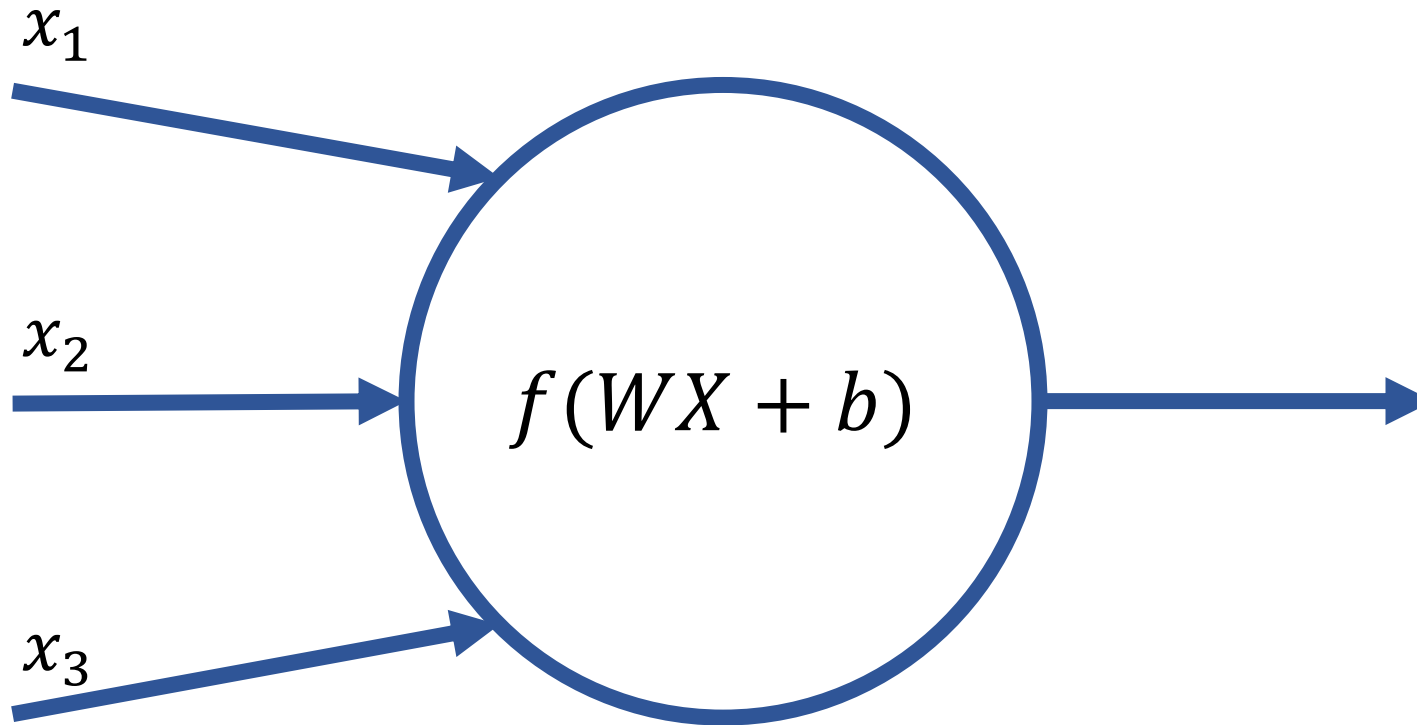
Leaky ReLU: $\max(0.1x, x)$

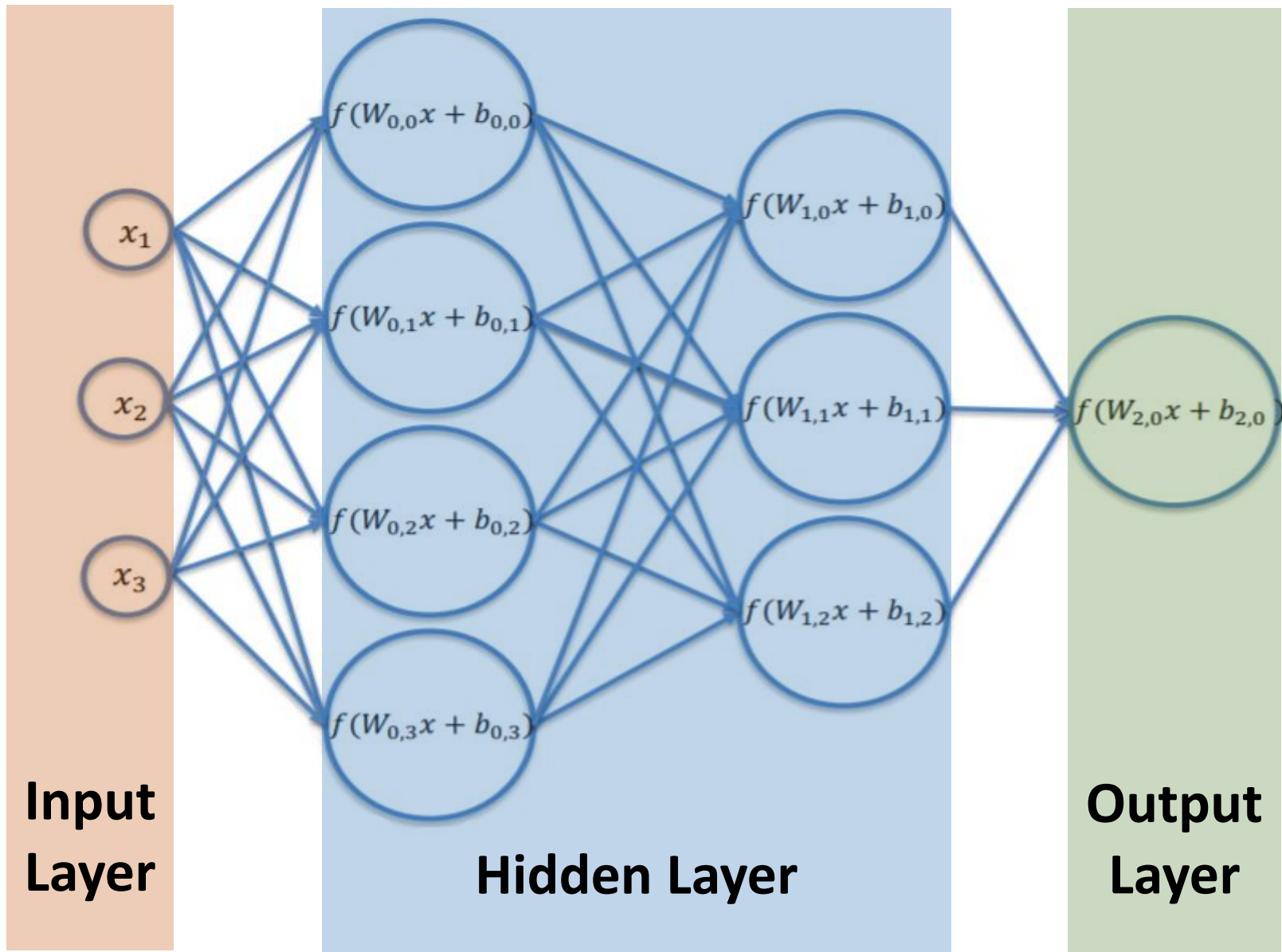


Parametric ReLU: $\max(\alpha x, x)$

Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

$$\text{ELU } f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$





- **Linear systems**

- LU, QR, Cholesky, Jacobi, Gauss-Seidel, CG, PCG, ...

- **Non-linear systems**

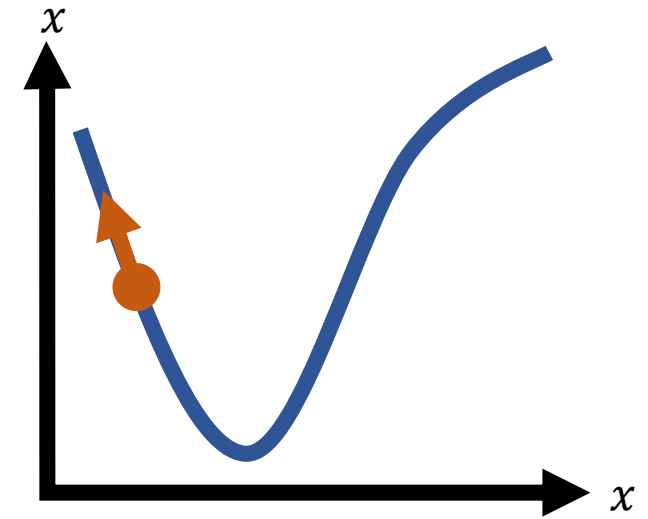
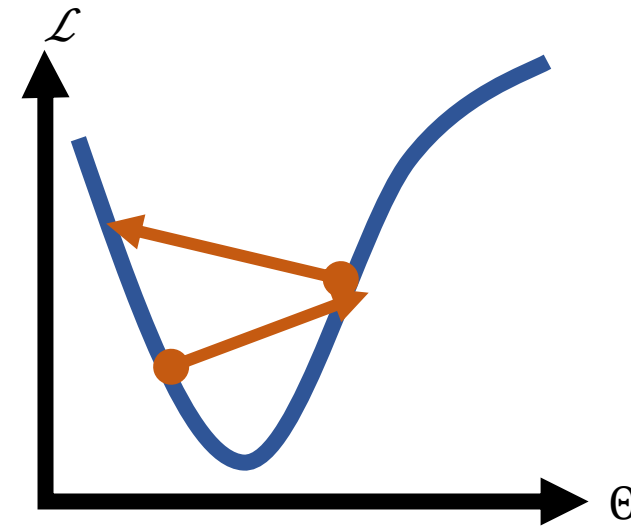
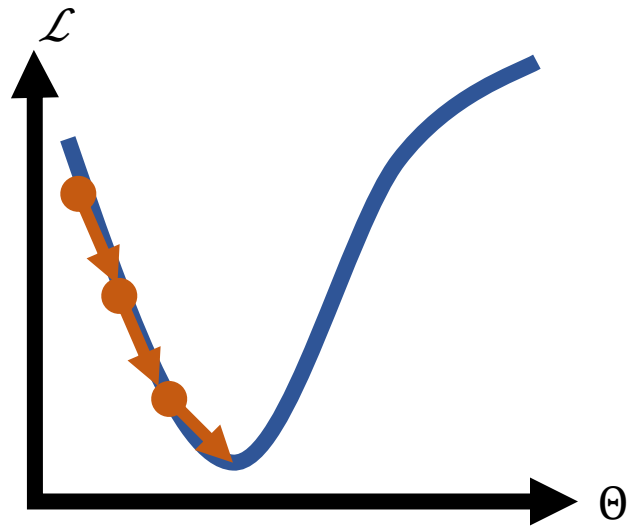
- First order : Gradient Descent, SGD
- Second order : Newton, Gauss-Newton, LM, (L)BFGS

- **Autres**

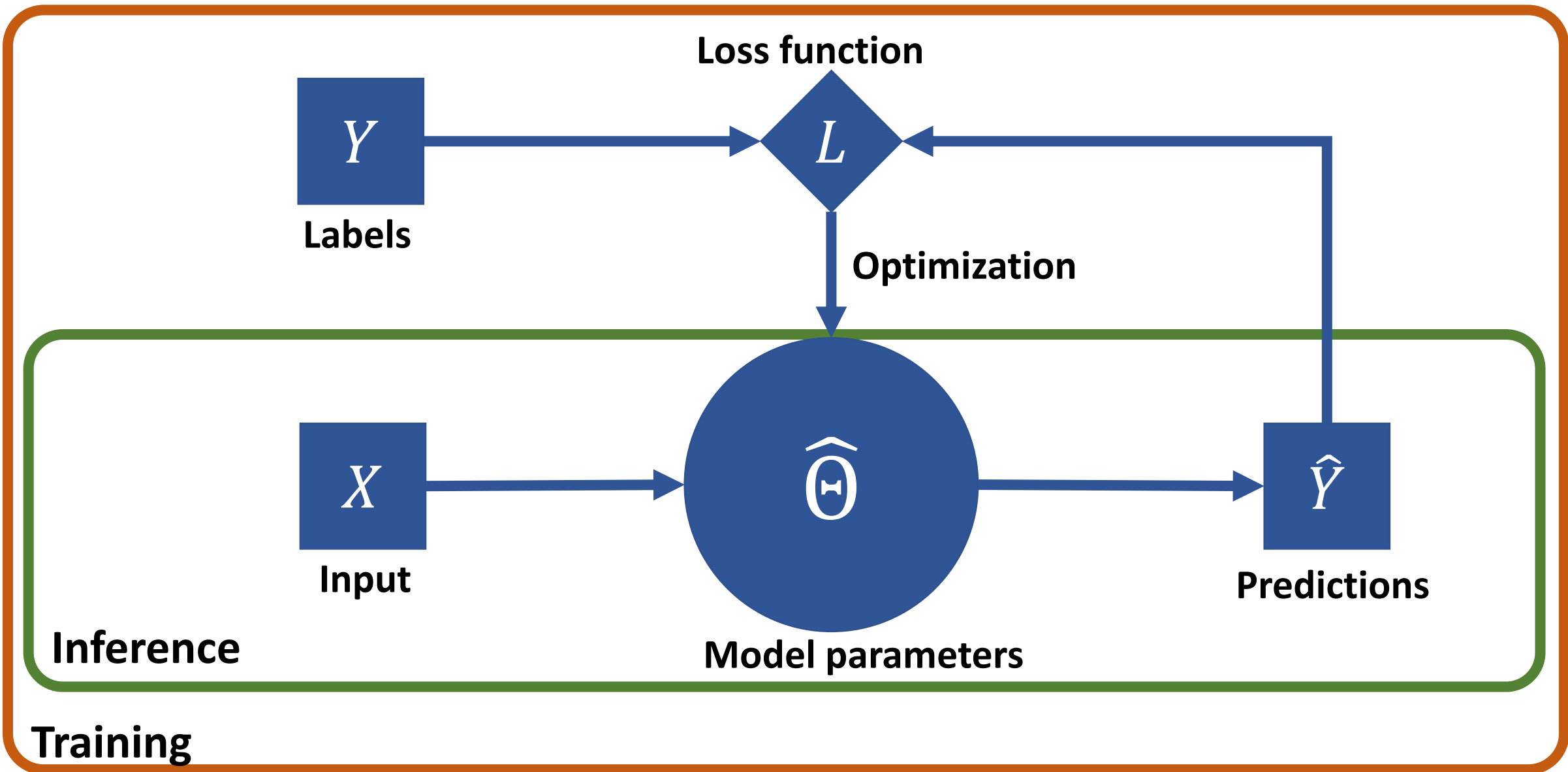
- Genetic algorithms, Metropolis-Hastings, ...
- Complex and constrained solver : ADMM, Primal-Dual, ...

Iterative solution

- Gradient descent
- $\hat{\Theta}_{t+1} = \hat{\Theta}_t - \eta \nabla_{\Theta} \mathcal{L}(\hat{y}_i, y_i)$
- η Learning rate



Gradient descent



- **Regression loss**

- Average absolute deviation : $L(y, \hat{y}, \theta) = \frac{1}{n} \sum_i^n |y_i - \hat{y}_i|$

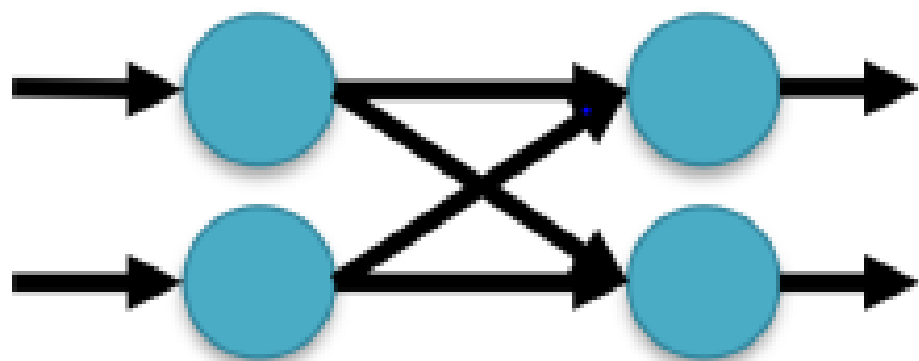
- Least squares method : $L(y, \hat{y}, \theta) = \frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2$

- **Classification loss**

- Cross-Entropy : $E(y, \hat{y}, \theta) = -\frac{1}{n} \sum_i^n \sum_j^m y_{ij} \log \hat{y}_{ij}$

$$\hat{\Theta}_{t+1} = \hat{\Theta}_t - \eta \nabla_{\Theta} [\mathcal{L}(\hat{y}_i, y_i) + \lambda R(\hat{\Theta}_t)]$$

L1 : LASSO	L2 : Ridge
$ \Theta $	Θ^2

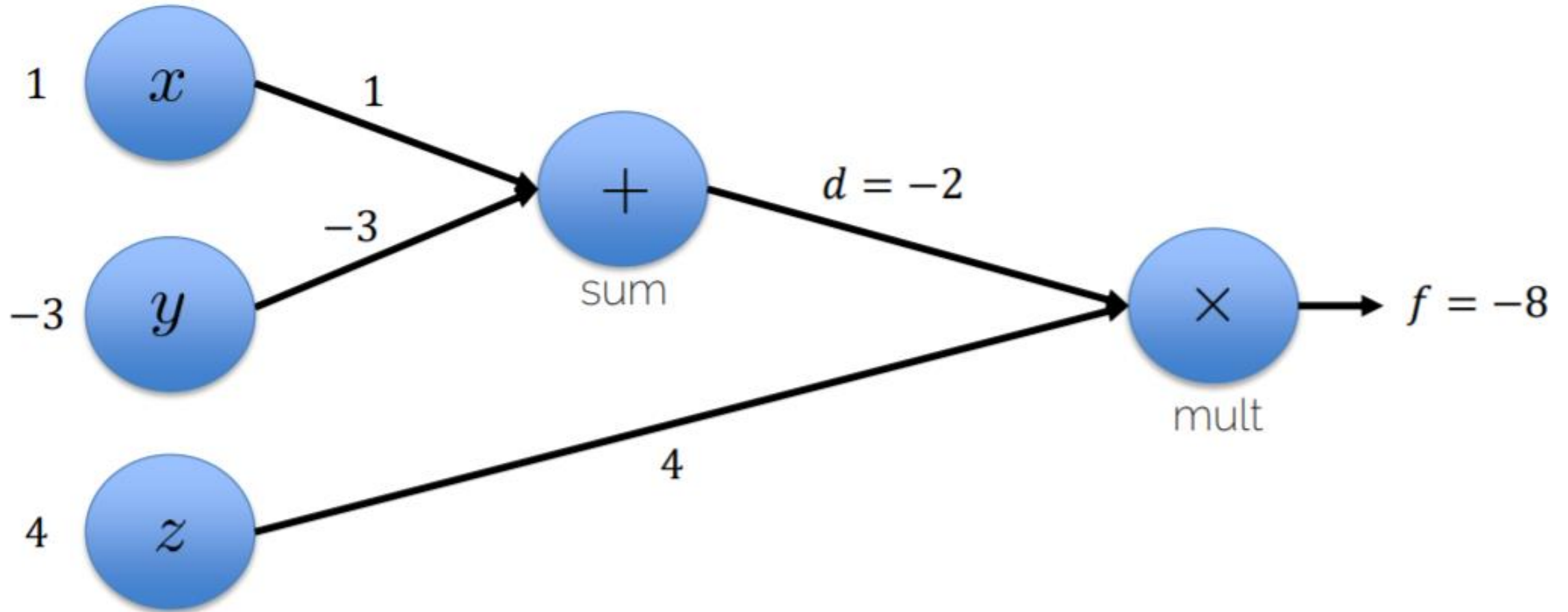


$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial x}$$

$$\frac{\partial L}{\partial w_{i,k}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_{i,k}}$$

- $f(x, y, z) = (x + y) \cdot z$

Initialization $x = 1, y = -3, z = 4$

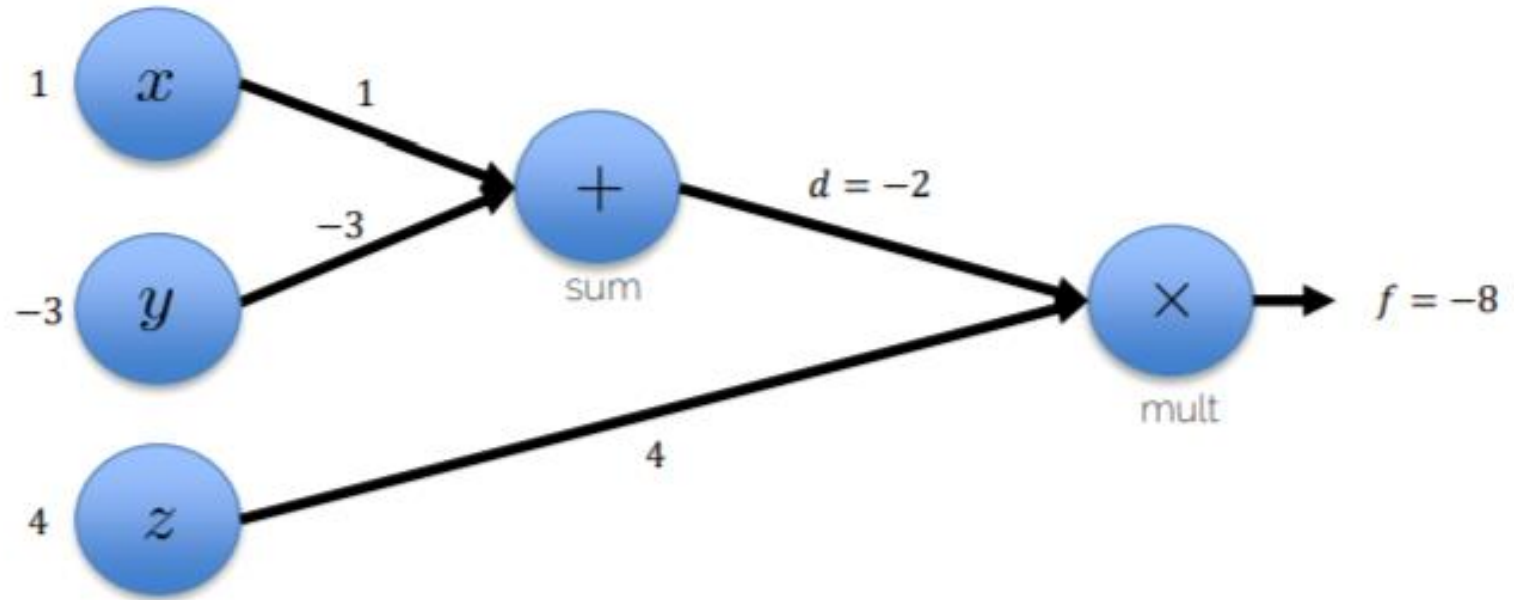


$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

$$d = x + y \quad \frac{\partial d}{\partial x} = 1, \frac{\partial d}{\partial y} = 1$$

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$



What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?

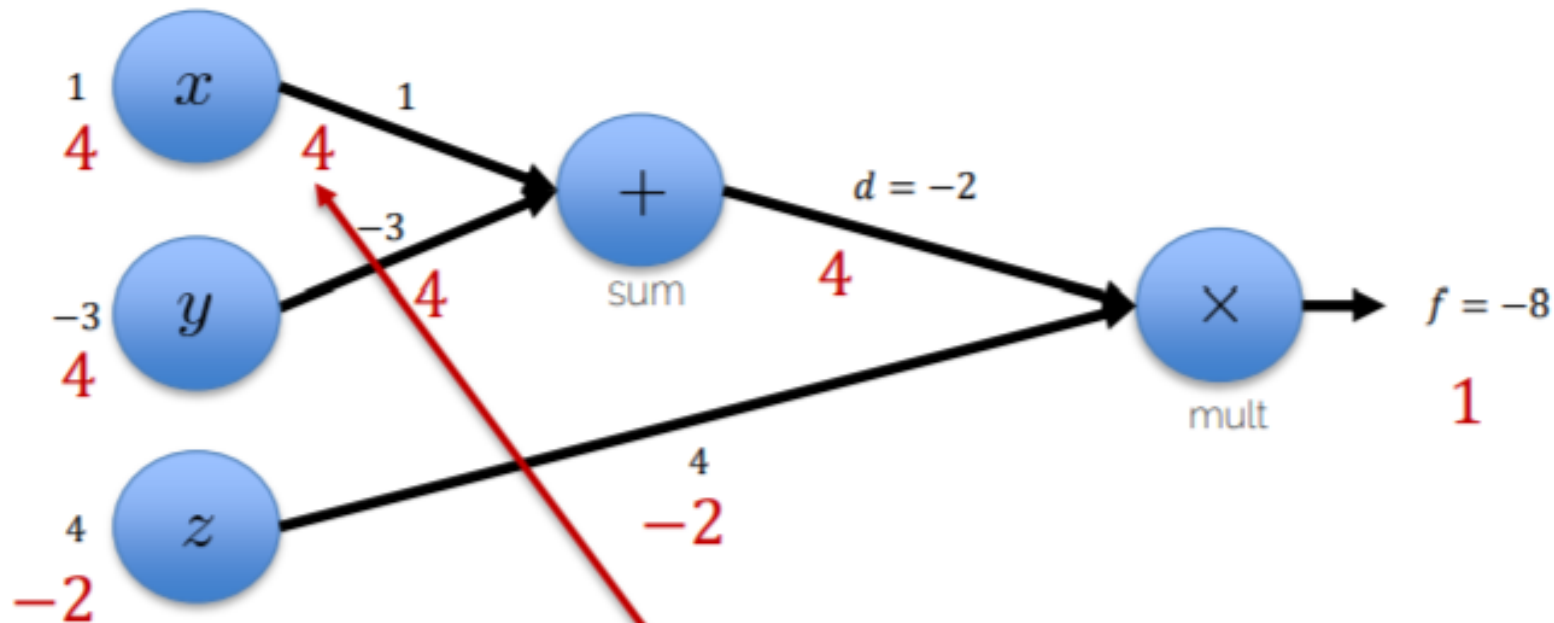
$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$

$$d = x + y \quad \boxed{\frac{\partial d}{\partial x} = 1} \quad \frac{\partial d}{\partial y} = 1$$

$$f = d \cdot z \quad \frac{\partial f}{\partial d} = z, \quad \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$?



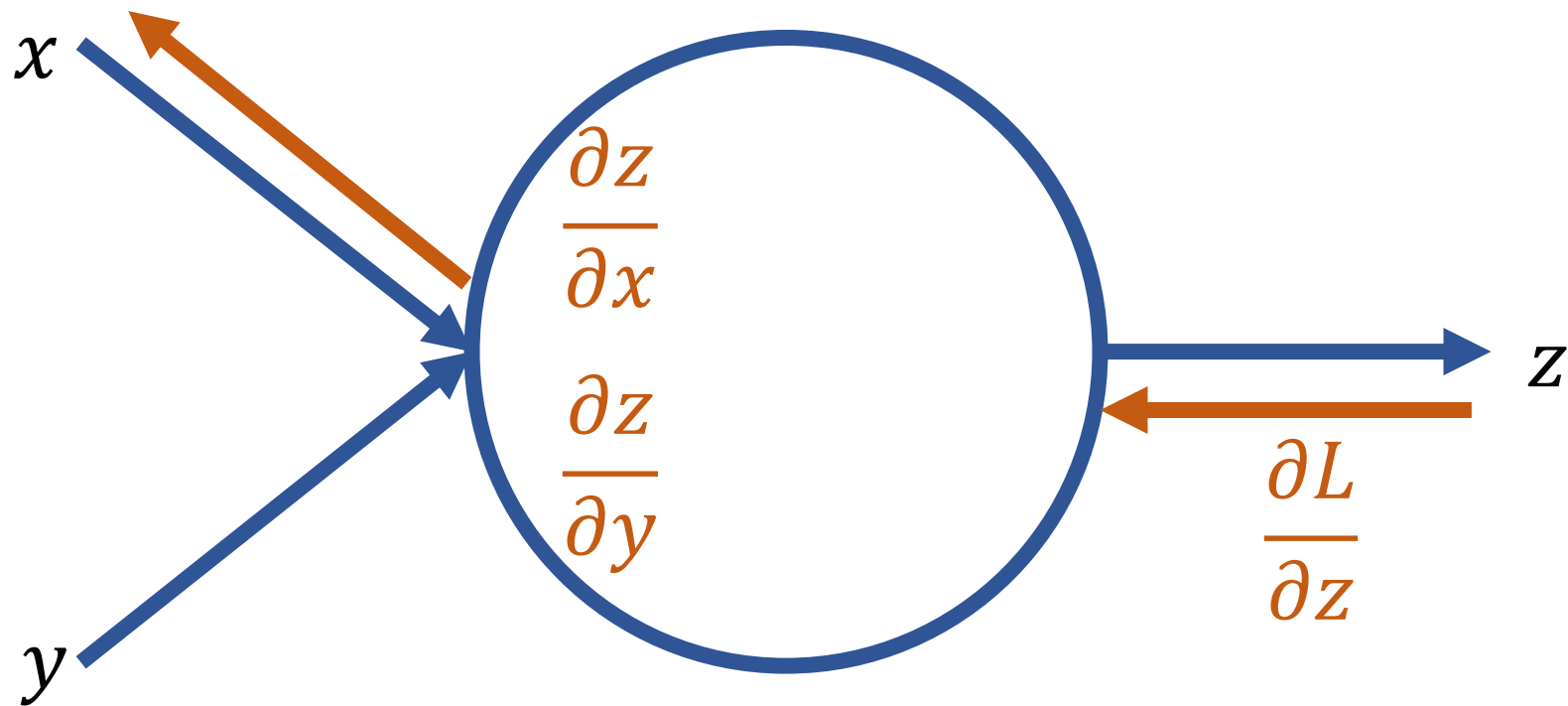
Chain Rule:

$$\boxed{\frac{\partial f}{\partial x} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial x}}$$

$$\boxed{\frac{\partial f}{\partial x}}$$

$$\rightarrow \frac{\partial f}{\partial x} = 4 \cdot 1 = 4$$

$$\frac{\partial L}{\partial x} = \frac{\partial z}{\partial x} \frac{\partial L}{\partial z}$$



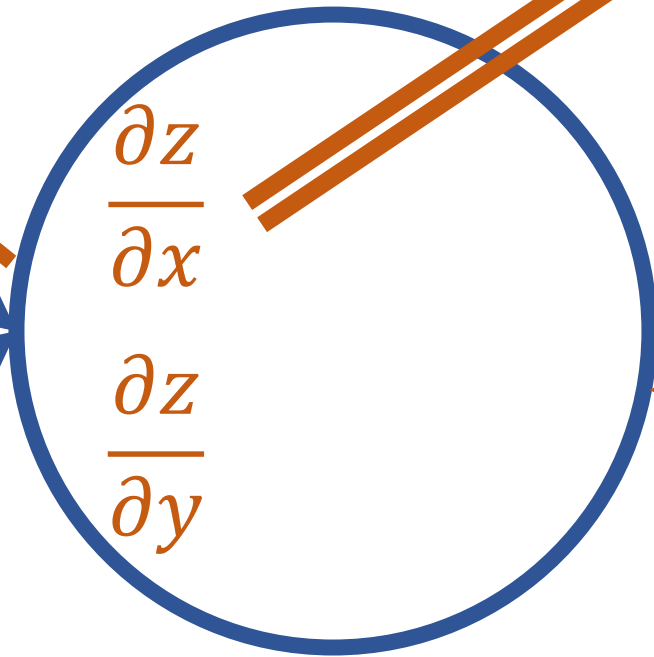
Vectors



x

y

$$\frac{\partial L}{\partial x} = \frac{\partial z}{\partial x} \frac{\partial L}{\partial z}$$



$$\begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial x_1} & \dots & \frac{\partial z_n}{\partial x_n} \end{bmatrix}$$

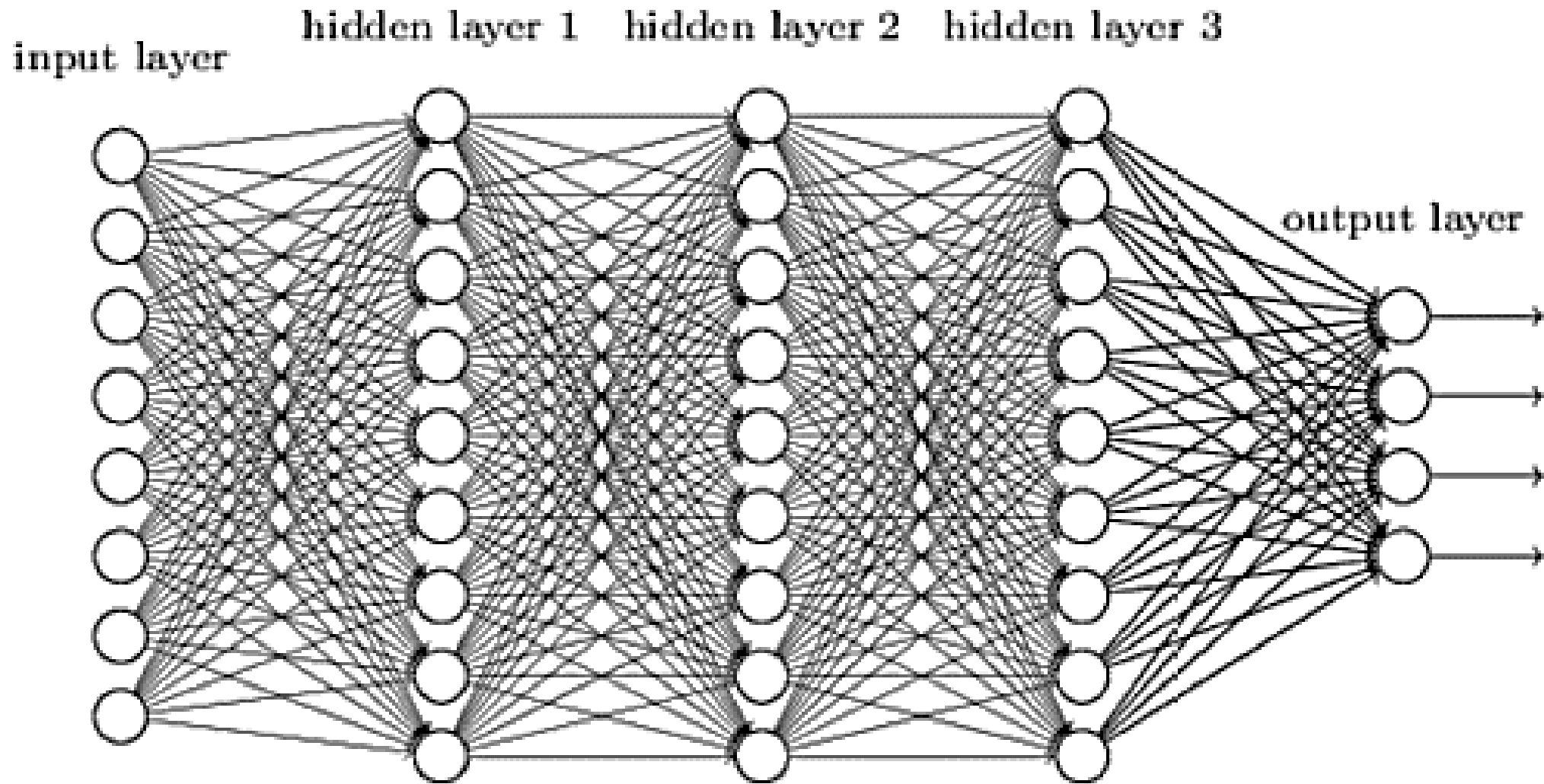
z

Vector

$$\frac{\partial L}{\partial z}$$



Gradient Flow



- Chollet, Francois. *Deep learning with Python*. Simon and Schuster, 2021.
- *CS230 Deep Learning*. cs230.stanford.edu. Accessed 14 Mar. 2022.
- *I2DL*. niessner.github.io/I2DL. Accessed 14 Mar. 2022.
- Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.