# Foundations of Deep Learning
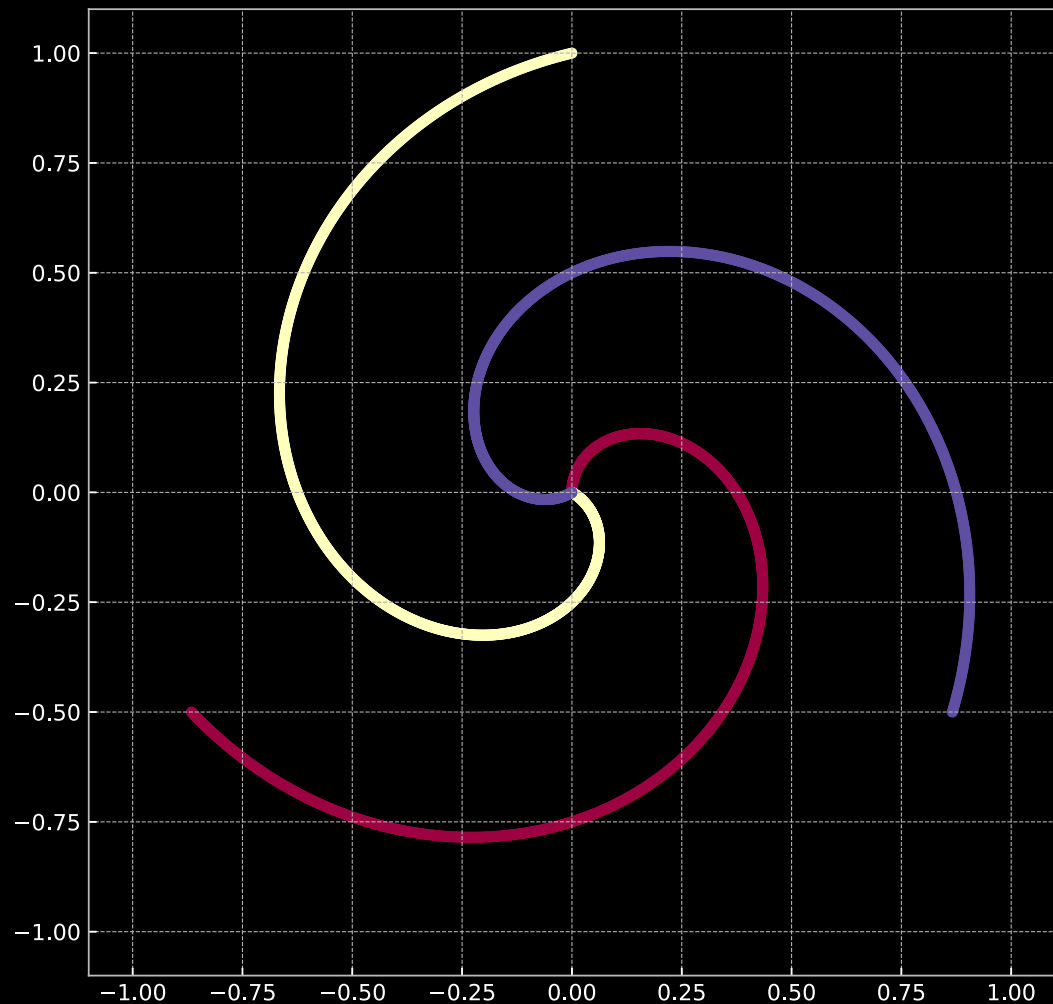
Alfredo Canziani

@alfcnz

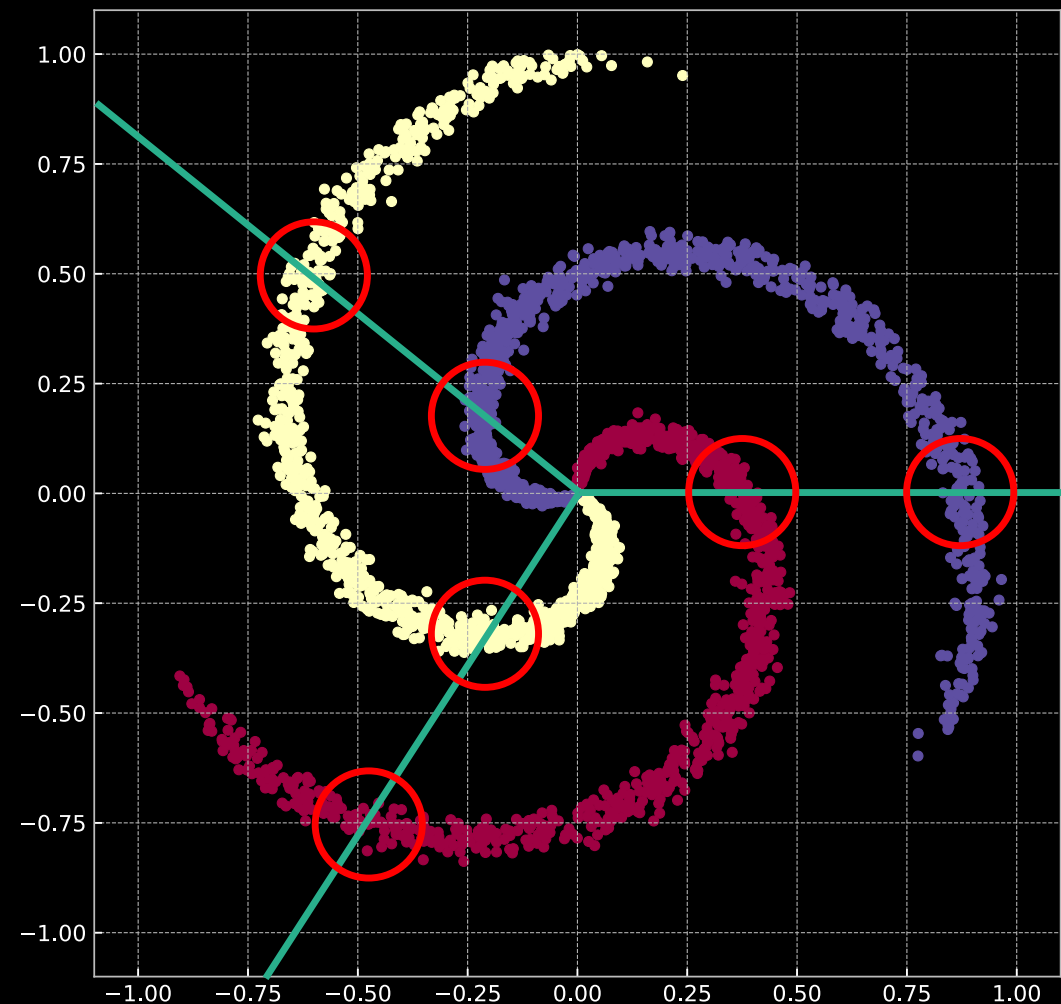ALF

# ANN – supervised learning

Classification
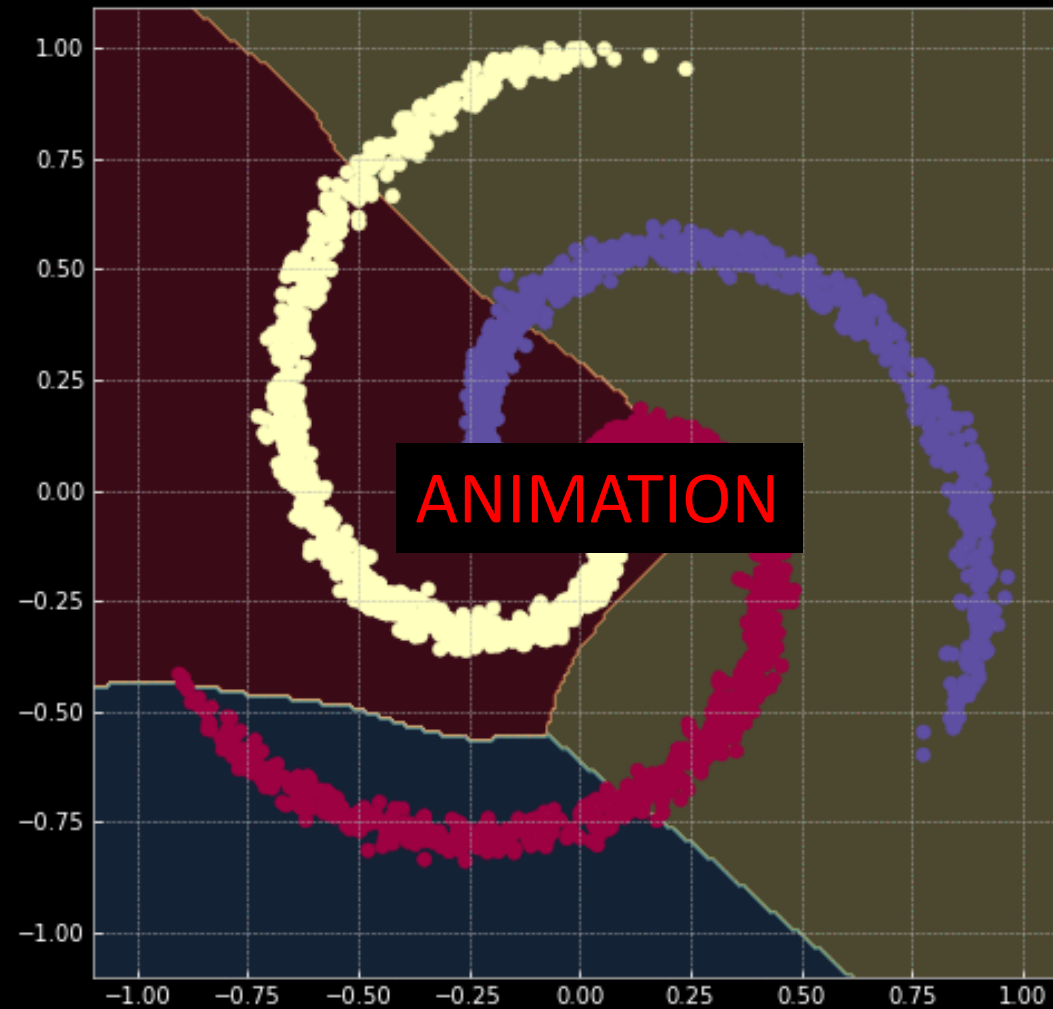
$$X_c(t) = t \begin{pmatrix} \sin\left(\frac{2\pi}{C}\left(2t + c - 1\right)\right) \\ \cos\left(\frac{2\pi}{C}\left(2t + c - 1\right)\right) \end{pmatrix}$$

$$0 \le t \le 1, \quad c = 1, \cdots, C$$

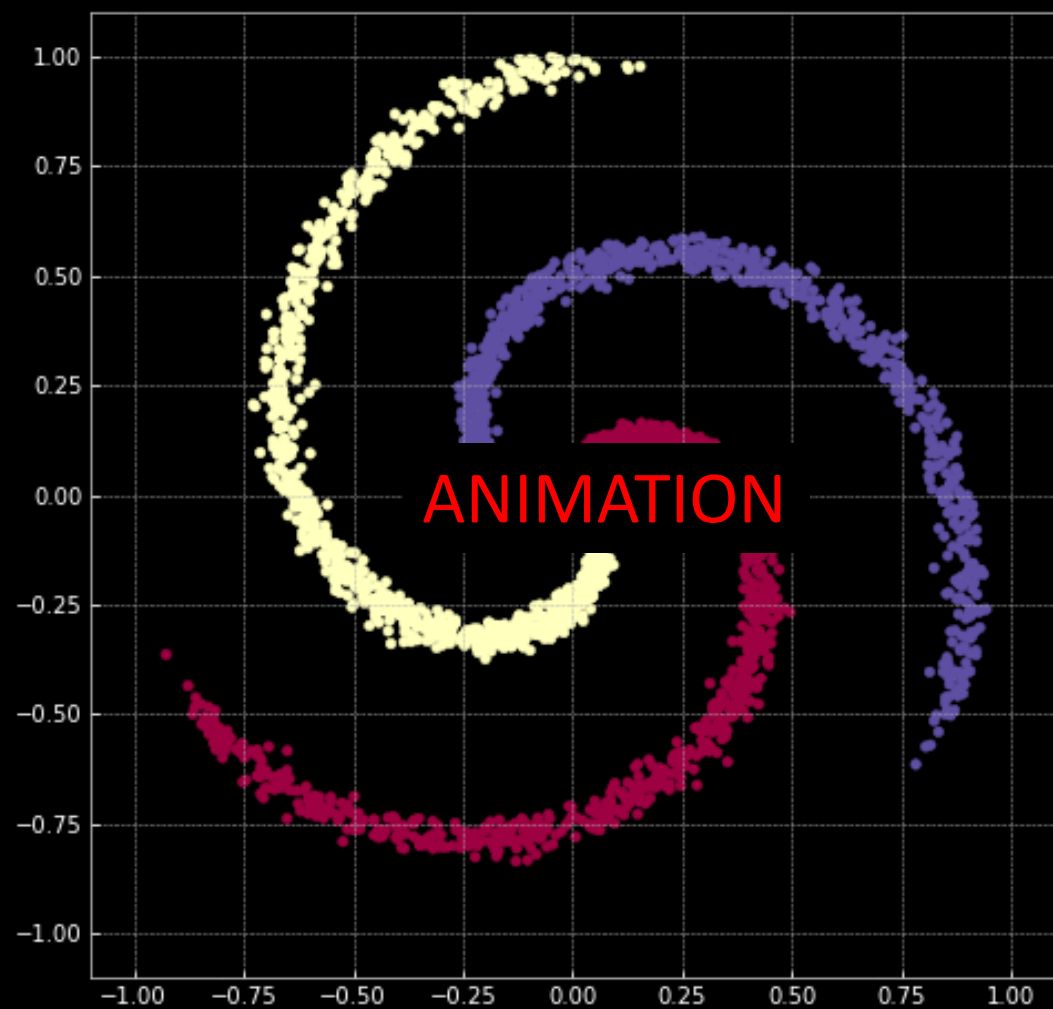$$X_c(t) = t \begin{pmatrix} \sin\left(\frac{2\pi}{C}\left(2t + c - 1\right)\right) \\ \cos\left(\frac{2\pi}{C}\left(2t + c - 1\right)\right) \end{pmatrix}$$

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

# Training data

$$\begin{pmatrix} 0 & 0 & 1 \end{pmatrix}$$
$$\begin{pmatrix} 0 & 1 & 0 \end{pmatrix}$$
$$\begin{pmatrix} 1 & 0 & 0 \end{pmatrix}$$

1-hot encoding

$$X = \begin{bmatrix} \underline{\quad\quad x^{(1)} \quad\quad} \\ \underline{\quad\quad x^{(2)} \quad\quad} \\ \vdots \\ \underline{\quad\quad x^{(m)} \quad\quad} \end{bmatrix} \updownarrow m \qquad Y = \begin{bmatrix} \underline{\quad\quad y^{(1)} \quad\quad} \\ \underline{\quad\quad y^{(2)} \quad\quad} \\ \vdots \\ \underline{\quad\quad y^{(m)} \quad\quad} \end{bmatrix} m \qquad c = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} \in \mathbb{R}^m$$

(top arrows: $n$, $G$)

$$x^{(i)} \in \mathbb{R}^n, \, n = 2 \qquad \boxed{y^{(i)}} \in \mathbb{R}^G, \quad G = 3 \qquad c_i \in \{1, \dots, G\} \subseteq \mathbb{N}^+$$

$m$ : # samples in data set

# Fully connected (FC) layer

$$j\text{-th row of } W^{(1)}$$

$$a_j^{(2)} = f(\boxed{w^{(j)}}x + b_j) = f\left(\left(\sum_{i=1}^{n} w_i^{(j)} x_i\right) + b_j\right)$$

$$\hat{y}$$

$$h = f(W_h x + b_h)$$

$$\hat{y} = g(W_y h + b_y)$$

$$g$$

$$W_y$$

$$h$$

$$f$$

$$W_h$$

$$x$$

$$f, g = (\cdot)^+, \sigma(\cdot),$$
$$\tanh(\cdot), \text{soft}(arg)\max(\cdot)$$

$$x$$
$$a^{(1)}$$

$$W^{(1)}$$

$$h^{(1)}$$
$$a^{(2)}$$

$$W^{(2)}$$

$$h^{(2)}$$
$$a^{(3)}$$

$$W^{(3)}$$

$$h^{(3)}$$
$$a^{(\ell)}$$

$$W^{(\ell)}$$

$$\hat{y}$$
$$a^{(L)}$$

# Neural network (inference)

$\hat{\boldsymbol{y}} \in \mathbb{R}^{C=3}$

$\boldsymbol{h} = f(\boldsymbol{W_h}\boldsymbol{x} + \boldsymbol{b_h})$

$\hat{\boldsymbol{y}} = g(\boldsymbol{W_y}\boldsymbol{h} + \boldsymbol{b_y})$

$\boldsymbol{W_h} \in \mathbb{R}^{d \times n}$

$\boldsymbol{b_h} \in \mathbb{R}^{d}$

$\boldsymbol{W_y} \in \mathbb{R}^{C \times d}$

$\boldsymbol{b_y} \in \mathbb{R}^{C}$

$g$

$\boldsymbol{W_y}$

$\boldsymbol{h} \in \mathbb{R}^{d=100}$

$f$

$f(\cdot), g(\cdot) : \mathrm{ReLU}(\cdot) \doteq (\cdot)^{+}, \sigma(\cdot), \tanh(\cdot), \mathrm{soft(arg)max}(\cdot)$

$\hat{\boldsymbol{y}} = \hat{\boldsymbol{y}}(\boldsymbol{x}), \quad \hat{\boldsymbol{y}} : \mathbb{R}^{n} \to \mathbb{R}^{C}, \quad \boldsymbol{x} \mapsto \hat{\boldsymbol{y}}$

$\boldsymbol{W_h}$

$\boldsymbol{x} \in \mathbb{R}^{n=2}$

$\hat{\boldsymbol{y}} : \mathbb{R}^{n} \to \mathbb{R}^{d} \to \mathbb{R}^{C}, \quad d \gg n, C$

# Neural network (training I)

$$h = f(W_h x + b_h)$$
$$\hat{y} = g(W_y h + b_y)$$

logits: output of final layer

$$\mathrm{soft}(arg)\mathrm{max}(\boldsymbol{l})[c] \doteq \frac{\exp(\boldsymbol{l}[c])}{\sum_{j=1}^{C} \exp(\boldsymbol{l}[j])} \in (0, 1)$$

$$\mathcal{L}(\hat{\boldsymbol{Y}}, \boldsymbol{c}) \doteq \frac{1}{m} \sum_{i=1}^{m} \ell(\hat{\boldsymbol{y}}^{(i)}, c^{(i)}), \quad \ell(\hat{\boldsymbol{y}}, c) \doteq -\log(\hat{\boldsymbol{y}}[c])$$

cross entropy / negative log-likelihood

$$\underline{x} \qquad c = 1 \quad \Rightarrow \quad \underline{y} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\hat{\underline{y}}(\underline{x}) = \begin{pmatrix} \sim 1 \\ \sim 0 \\ \sim 0 \end{pmatrix} \Rightarrow \mathcal{L}\left(\begin{pmatrix} \sim 1 \\ \sim 0 \\ \sim 0 \end{pmatrix}, 1\right) \rightarrow 0^+$$

$$\hat{\underline{y}}(\underline{x}) = \begin{pmatrix} \sim 0 \\ \sim 1 \\ \sim 0 \end{pmatrix} \Rightarrow \mathcal{L}\left(\begin{pmatrix} \sim 0 \\ \sim 1 \\ \sim 0 \end{pmatrix}, 1\right) \rightarrow +\infty$$

# Neural network (training II)

$$h = f(\boldsymbol{W_h} x + \boldsymbol{b_h})$$
$$\hat{y} = g(\boldsymbol{W_y} h + \boldsymbol{b_y})$$

$$\Theta \equiv \{W_h, \boldsymbol{b_h}, W_y, \boldsymbol{b_y}\}$$

$$J(\Theta) \equiv \mathcal{L}\left(\hat{Y}(\Theta), \underline{c}\right) \in \mathbb{R}^+$$

$$\frac{\partial J(\Theta)}{\partial W_y} = \frac{\partial J(\Theta)}{\partial \hat{\underline{y}}} \frac{\partial \hat{\underline{y}}}{\partial W_y}$$

$$\frac{\partial J(\Theta)}{\partial W_h} = \frac{\partial J(\Theta)}{\partial \hat{\underline{y}}} \frac{\partial \hat{\underline{y}}}{\partial \underline{h}} \frac{\partial \underline{h}}{\partial W_h}$$

Back propagation

$$J(\theta)$$
$$J(\theta_o)$$

$$- \frac{dJ(\theta)}{d\theta}(\theta_o)$$

$$\theta_o$$

Gradient descent

# Gradient computation example

```
x = torch.tensor([[1, 2], [3, 4]], requires_grad=True)
y = x - 2
z = y * y * 3
out = z.mean()
```

$$\text{out}: \ o = \frac{1}{4} \sum_{i=1}^{4} z_i$$

$$z = 3y^2 = 3(x-2)^2$$

$$\frac{\partial o}{\partial x_i} = \frac{1}{4_2} \, 3 \cdot \not{2} \, (x-2) =$$

$$= \frac{3}{2}(x-2)$$

$$x = 1 \qquad \frac{\partial o}{\partial x}\Big|_{x=1} = -1.5$$

$$\frac{\partial o}{\partial x}\Big|_{x=2} = 0 \qquad \frac{\partial o}{\partial x}\Big|_{x=3} = 1.5 \qquad \frac{\partial o}{\partial x}\Big|_{x=4} = 3$$