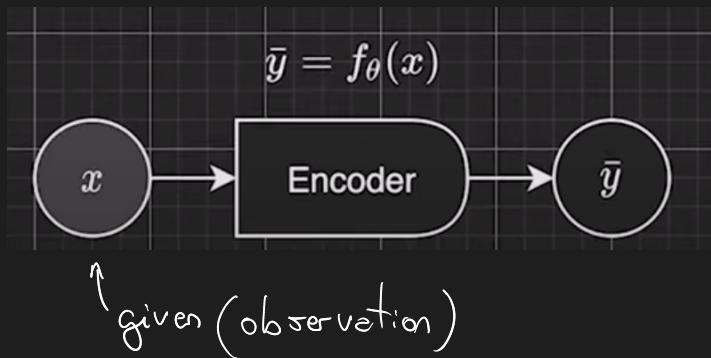
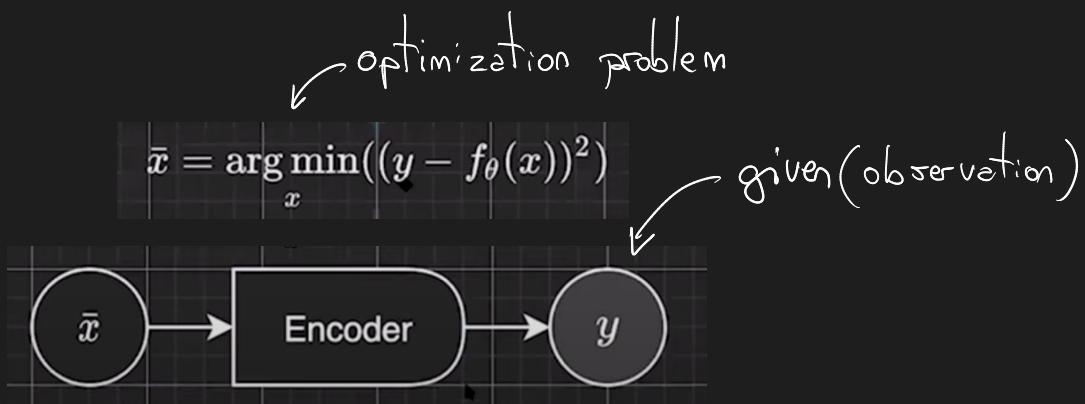


Forward Propagation



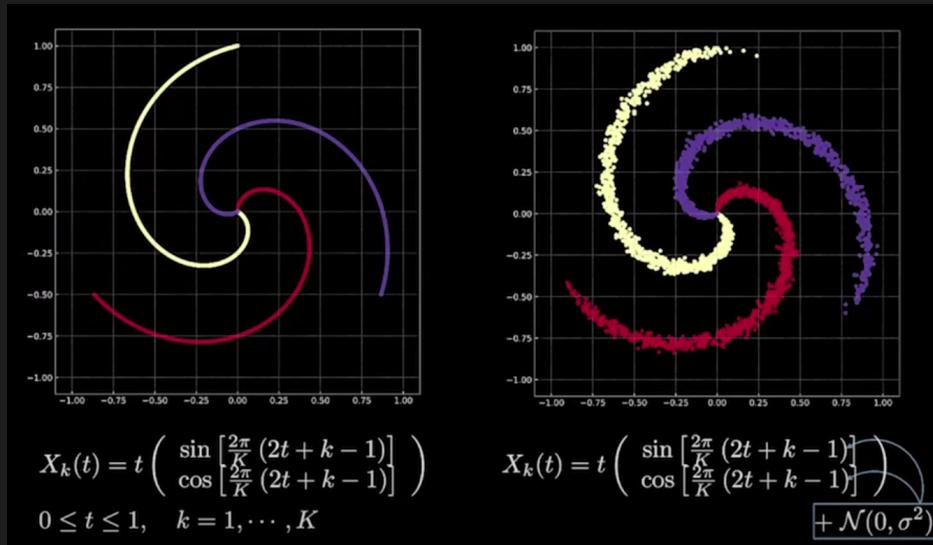
Inference: Compute something given something else



Q: Is backprop only used for training?

A: No! in the previous case backprop is used
for inference.

Supervised Learning.



- 3 classes

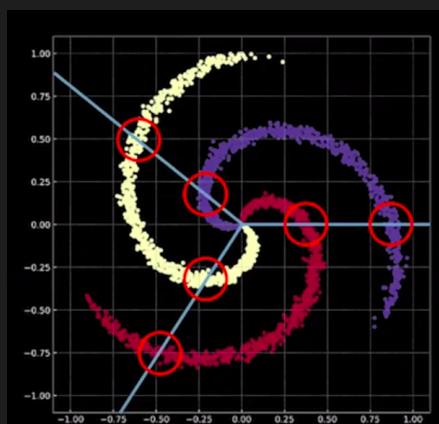
Input :

- 2d data points

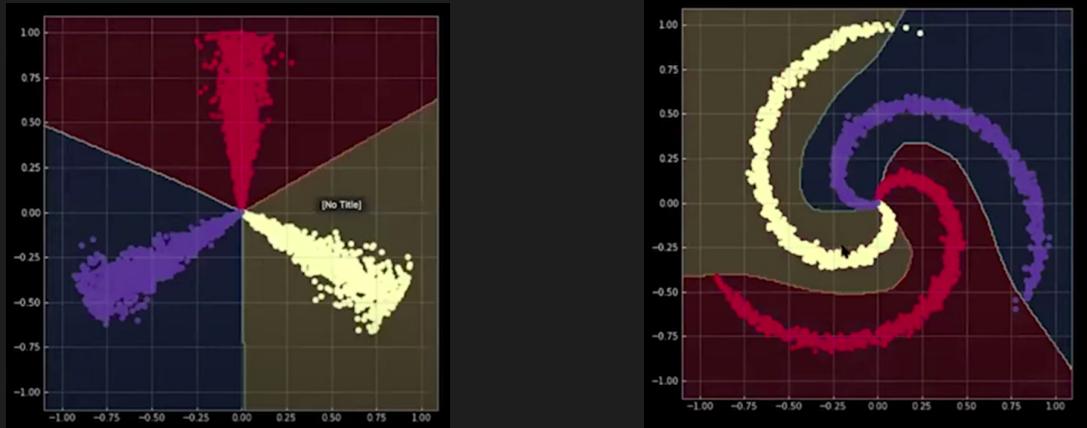
Output :

- Which class

A linear model is not enough



Two ways of seeing the boundaries



After transforming the space, the classes are linearly separable.

The model

Linear classifier on top

output



5 neurons



2 neurons



ReLU



100 neurons

input



weights $\omega^{(1)} \in \mathbb{R}^{100 \times 2}$

$(x, y) \in \mathbb{R}^{2 \times 1}$

↖ column vector

Animation: Linear interpolation between and

100 %

0 %

99 %

1 %

:

:

Having 100 neurons in the second layer let us move data onto 100-dimensional space.

Train data

X : Design matrix

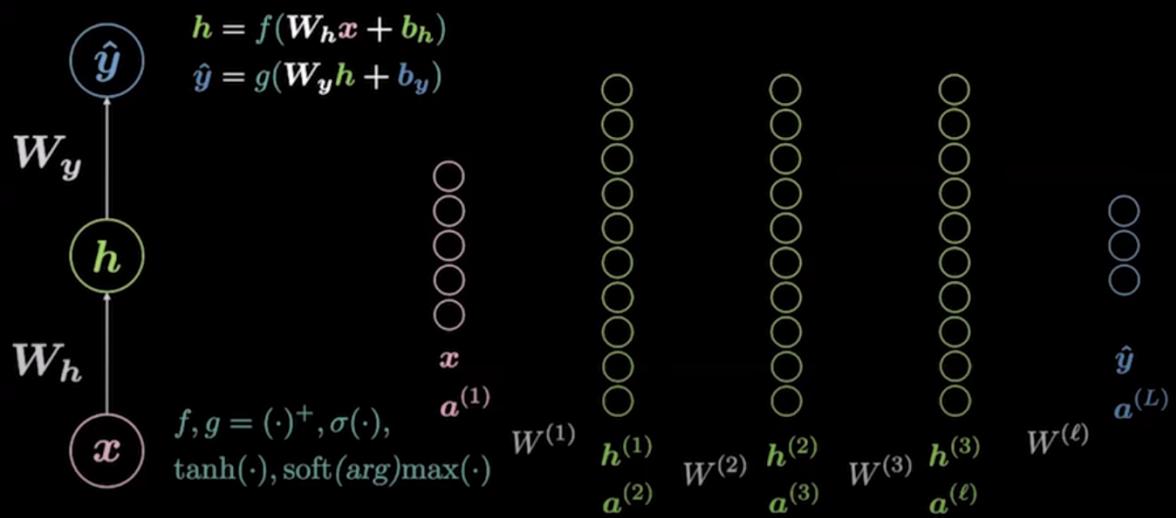
c : Collection of labels / classes,

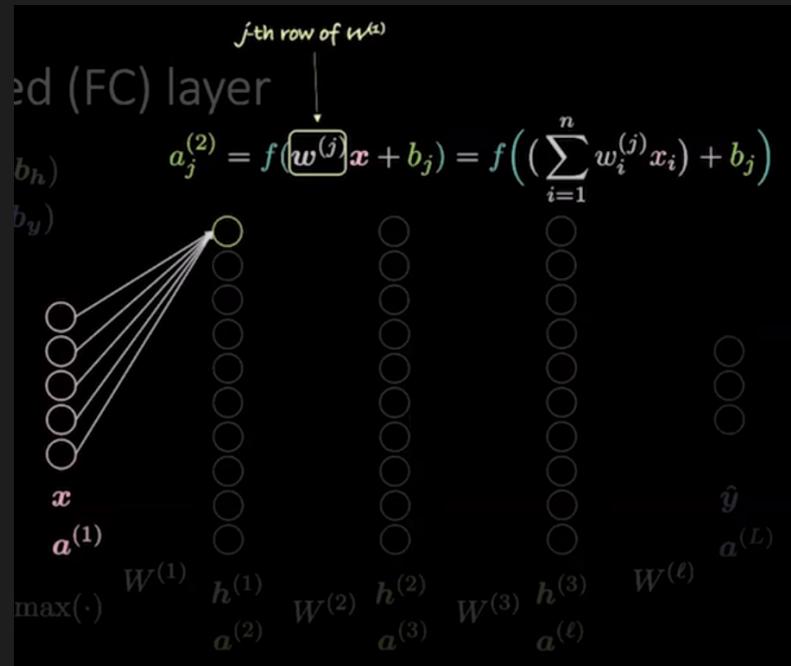
$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \\ \vdots \\ \mathbf{x}^{(m)} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \\ \vdots \\ \mathbf{y}^{(m)} \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}$$

m data points of dimension n

$$\mathbf{x}^{(i)} \in \mathbb{R}^n \quad \mathbf{y}^{(i)} \in \{0, 1\}^K \quad c_i \in \{1, 2, \dots, K\}$$

Fully connected (FC) layer

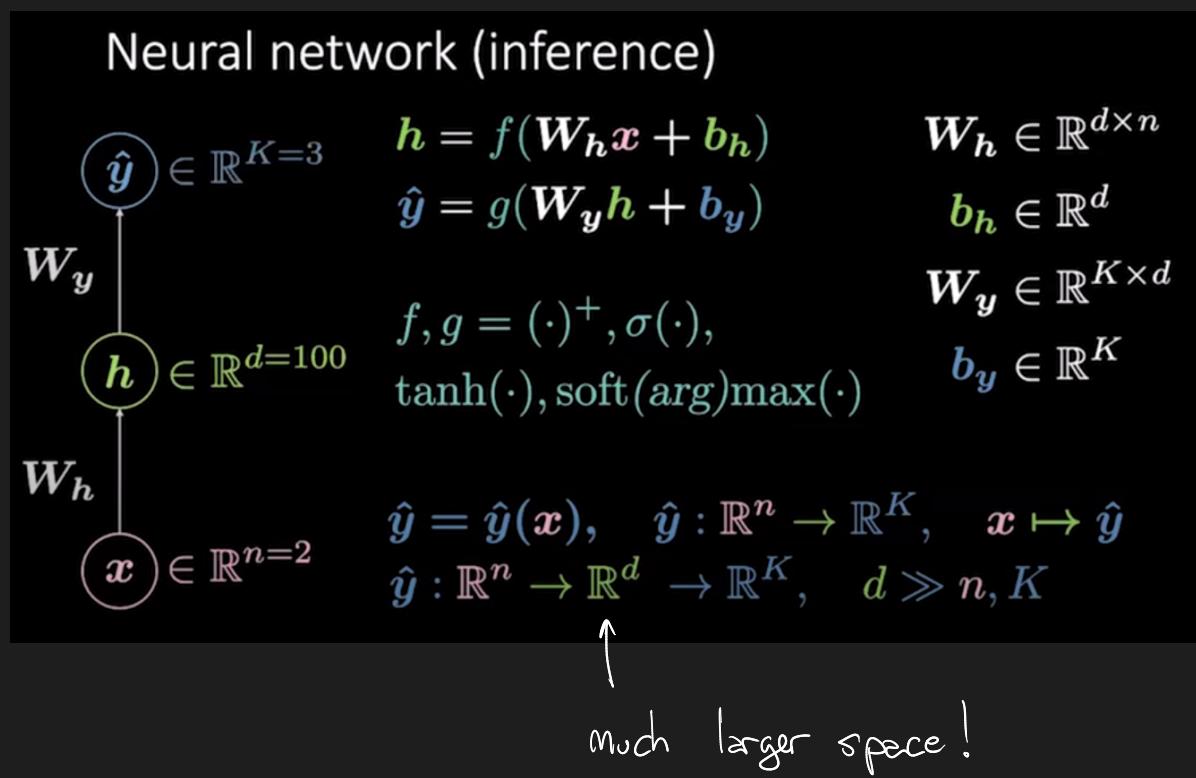




Linear function of linear function of linear function = Linear function

BUT

If using floating point approximation on a computer, it can be seen as a non linear function in between linear functions.



Note : Counting from 1 (not zero)

$h = f(\mathbf{W}_h \mathbf{x} + b_h)$
 $\hat{\mathbf{y}} = g(\mathbf{W}_y h + b_y)$

Neural network (training I)

logits: output of final layer

$$\text{soft(arg)}\max(\mathbf{l})[c] \doteq \frac{\exp(\mathbf{l}[c])}{\sum_{k=1}^K \exp(\mathbf{l}[k])} \in (0, 1)$$

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

cross entropy / negative log-likelihood

$$\mathbf{x}, \quad c = 1 \quad \Rightarrow \quad \mathbf{y} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\hat{\mathbf{y}}(\mathbf{x}) = \begin{pmatrix} \approx 1 \\ \approx 0 \\ \approx 0 \end{pmatrix} \quad \Rightarrow \quad \ell\left(\begin{pmatrix} \approx 1 \\ \approx 0 \\ \approx 0 \end{pmatrix}, 1\right) \rightarrow 0^+$$

$$\hat{\mathbf{y}}(\mathbf{x}) = \begin{pmatrix} \approx 0 \\ \approx 1 \\ \approx 0 \end{pmatrix} \quad \Rightarrow \quad \ell\left(\begin{pmatrix} \approx 0 \\ \approx 1 \\ \approx 0 \end{pmatrix}, 1\right) \rightarrow +\infty$$

$h = f(\mathbf{W}_h \mathbf{x} + b_h)$
 $\hat{\mathbf{y}} = g(\mathbf{W}_y h + b_y)$

Neural network (training II)

$$\Theta \doteq \{\mathbf{W}_h, \mathbf{b}_h, \mathbf{W}_y, \mathbf{b}_y\}$$

$$J(\Theta) \doteq \mathcal{L}(\hat{\mathbf{Y}}(\Theta), \mathbf{c}) \in \mathbb{R}^+$$

$$\frac{\partial J(\Theta)}{\partial \mathbf{W}_y} = \frac{\partial J(\Theta)}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{W}_y}$$

$$\frac{\partial J(\Theta)}{\partial \mathbf{W}_h} = \frac{\partial J(\Theta)}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{W}_h}$$

back-propagation

$J(\vartheta)$ $\frac{dJ(\vartheta)}{d\vartheta}(\vartheta_0)$

ϑ

$J(\vartheta_0)$

$\frac{dJ(\vartheta)}{d\vartheta}(\vartheta_0)$

