# Neural networks

$x \rightarrow$ [image] $\rightarrow y$

similar

$x \rightarrow h(x) \rightarrow \hat{y}$

- Linear Regression (regression)
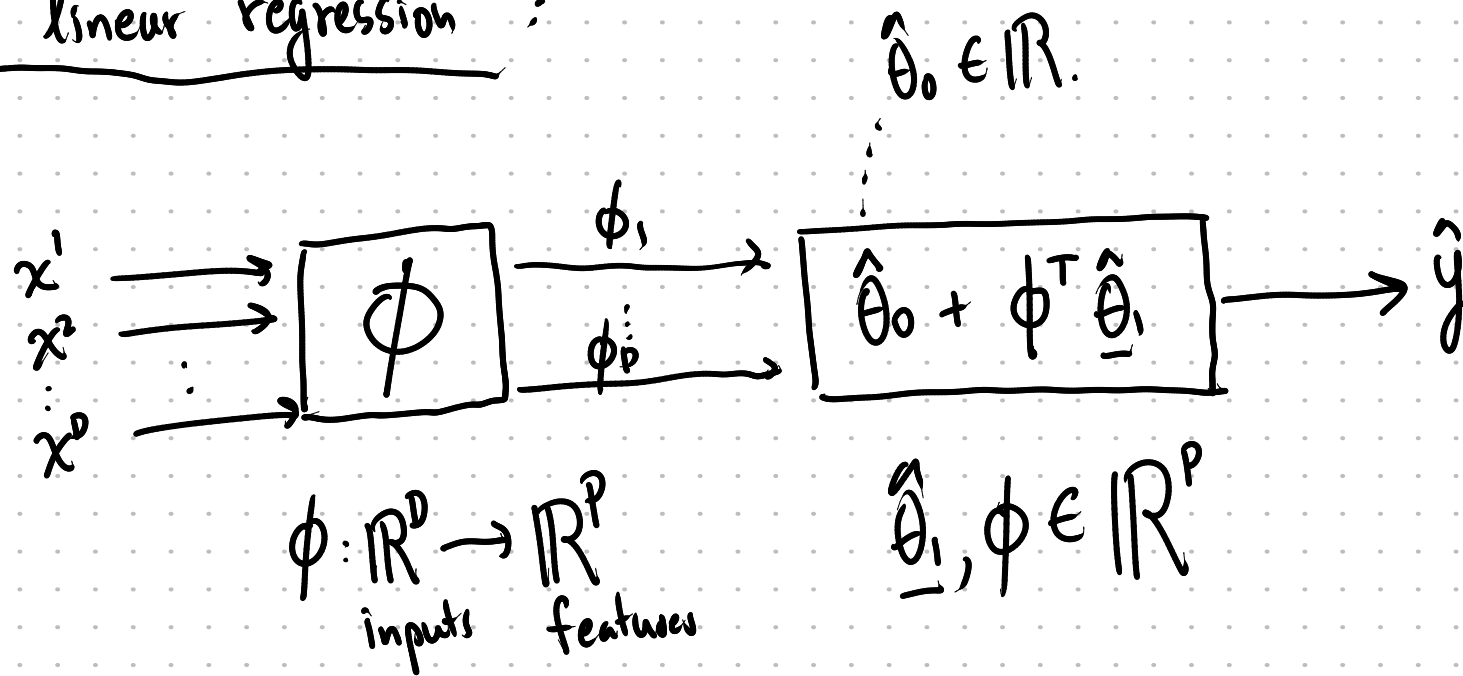
- Logistic regression (classification).

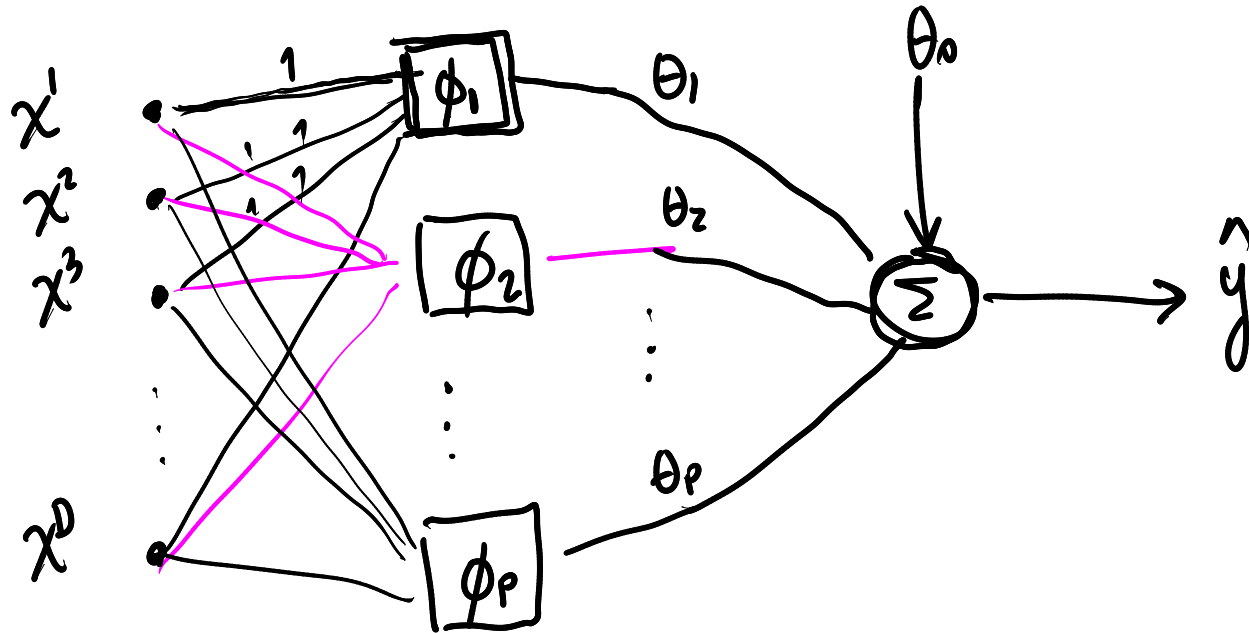- Support Vector machine
- Decision trees
- Neural networks.

# Recall linear regression :

$\hat{\theta}_0 \in \mathbb{R}.$

$$x^1$$
$$x^2$$
$$\vdots$$
$$x^D$$

$\phi$

$\phi_1$

$\phi_P$

$\hat{\theta}_0 + \phi^T \hat{\underline{\theta}}_1$

$\rightarrow \hat{y}$

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^P$$

inputs    features

$$\hat{\underline{\theta}}_1, \phi \in \mathbb{R}^P$$

- I have to <u>manually</u> design the features $\phi$.

- Neural networks can be understood as a way to design the features automatically.

# Pictorial representation of linear regression

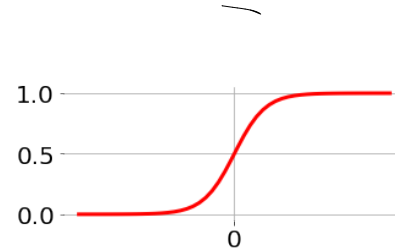$$\hat{y} = \theta_0 + \phi^T(x)\underline{\theta}_1 = \theta_0 + \sum_{j=1}^{P} \phi_j(x)\theta_j$$
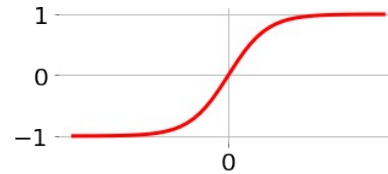


3 steps to NN:

1. Replace the $\phi$ with "activation functions"

2. Put weights on all edges.

3. Replicate layers.

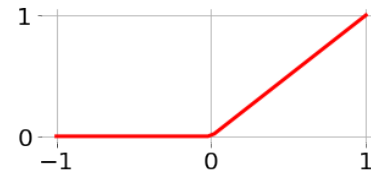## 1) Generic nonlinearities, a.k.a. _activation functions_

sigmoid: $\phi(\xi) = \dfrac{1}{1 + e^{-\xi}}$
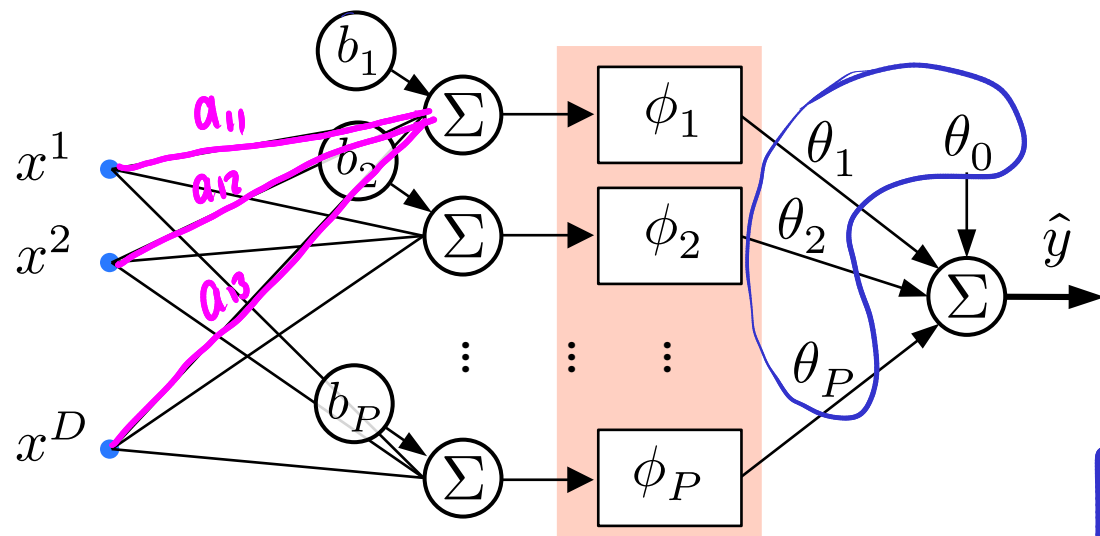


tanh: $\phi(\xi) = \dfrac{e^{2\xi} - 1}{e^{2\xi} + 1}$



ReLU: $\phi(\xi) = \max(0, \xi)$

## 2) Weights on the inputs

$$\hat{y} = \theta_0 + \phi^T(b + Ax)\underline{\theta}_1 = \theta_0 + \sum_{j=1}^{P} \phi_j(b_j + A_j x)\theta_j$$



$P \times D$ "a" coeff.

$P$ "b" coeff.

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{1n} \\ a_{21} & a_{22} & a_{13} \\ a_{P1} & \vdots & a_{P3} \end{bmatrix}$$
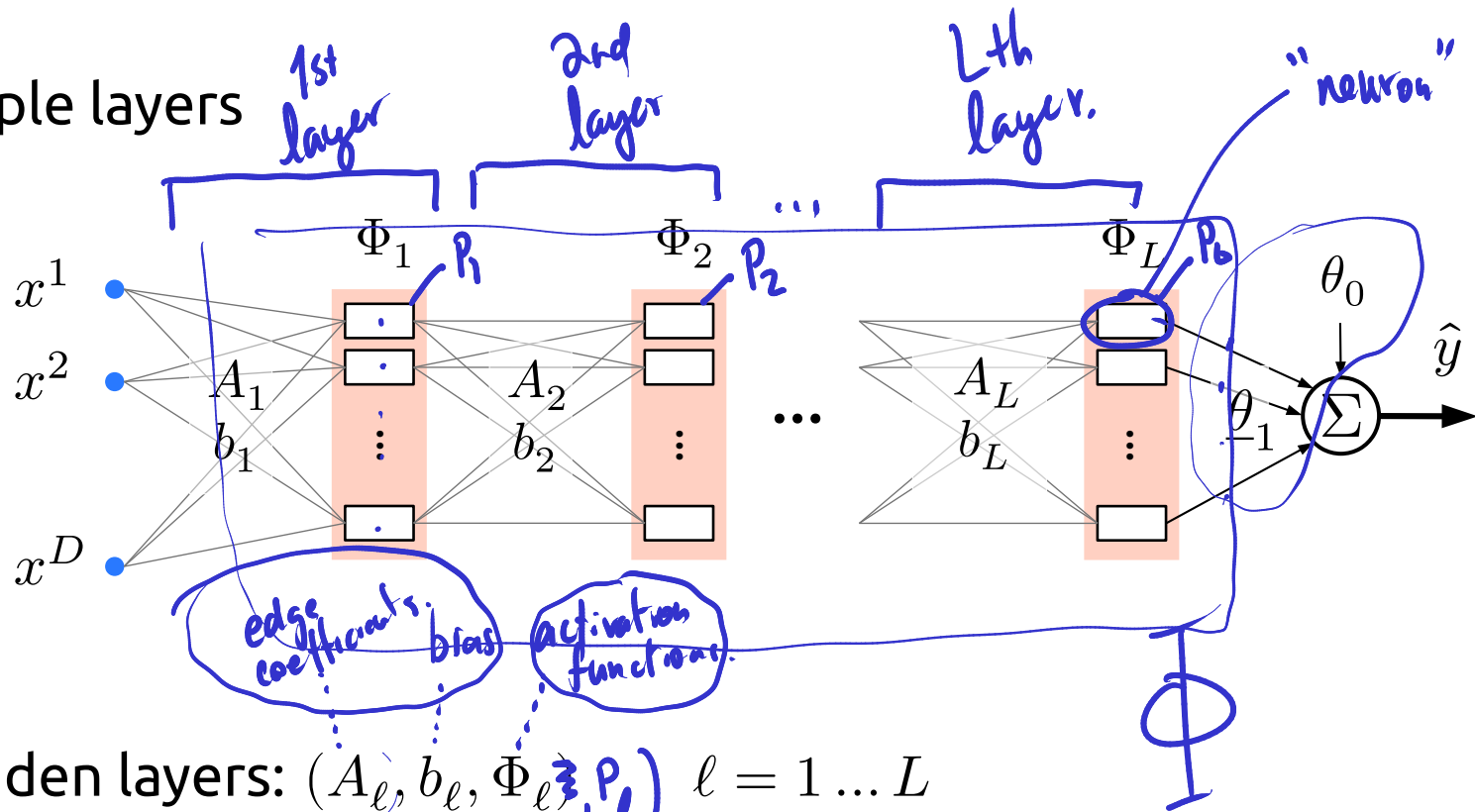
$$b = \begin{bmatrix} b_1 \\ \vdots \\ b_P \end{bmatrix}$$

$$\phi_1 \left( a_{11} x^1 + a_{12} x^2 + a_{13} x^3 + b_1 \right) \theta_1$$

$$\phi_1 \left( A_1 x + b_1 \right)$$

$$A_1 = \left[ a_{11} , a_{12} , a_{13} \right]$$

# 3) Multiple layers



1st layer

2nd layer

Lth layer.

"neuron"

$\Phi_1$ $P_1$

$\Phi_2$ $P_2$

$\Phi_L$ $P_b$

$\theta_0$

$x^1$

$x^2$

$x^D$

$A_1$
$b_1$

$A_2$
$b_2$

$A_L$
$b_L$

$\Sigma$ $\hat{y}$

$\theta_1$

edge coefficients.    bias    activation functions.

- Hidden layers: $(A_\ell, b_\ell, \Phi_\ell, P_\ell)$ $\ell = 1 \dots L$

- Output layer: $(\theta_0, \underline{\theta}_1)$ ... final linear regression. (regression problem)

Put it all together :

$$\hat{y} = h(x) = \theta_0 + \underline{\theta}_1^T \Phi_L \left( b_L + A_L \cdot \Phi_{L-1} \left( b_{L-1} + A_{L-1} \Phi_{L-2} \left( \cdots \Phi_1 \left( b_1 + A_1 x \right) \cdots \right) \right) \right).$$
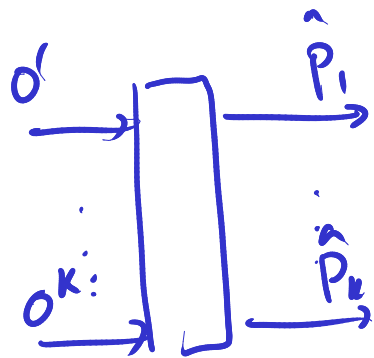
... complicated nested function. $z_{L-1}$

... perceptron.

$\Phi_L : \mathbb{R}^{P_L} \to \mathbb{R}^{P_L}.$

Write this in a recursive form.

$$\hat{y} = \theta_0 + \theta_1^T \phi_L(z_L)$$

$$z_L = b_L + A_L \cdot \phi_{L-1}(z_{L-1})$$

$$z_{L-1} = b_{L-1} + A_{L-1} \phi_{L-2}(z_{L-2})$$

$$\vdots$$

$$z_\ell = b_\ell + A_\ell \phi_{\ell-1}(z_{\ell-1}) \qquad \text{(general form)}$$

$$\vdots$$

$$z_1 = b_1 + A_1 x$$

# Classification networks

$K \dots$ classes.

$$x^1 \quad x^2 \quad \dots \quad x^D$$

$$\Phi_1 \qquad \Phi_2 \qquad \dots \qquad \Phi_L$$

$$A_1 \quad b_1 \qquad A_2 \quad b_2 \qquad \dots \qquad A_L \quad b_L \qquad \underline{\theta}_1$$

$$o_1 \quad o_2 \quad \dots \quad o_K \qquad \hat{p}_1 \quad \hat{p}_2 \quad \dots \quad \hat{p}_K$$

$$1$$

$$o^1 \quad \dots \quad o^K : \qquad \hat{p}_1 \quad \dots \quad \hat{p}_K$$

Softmax: $\qquad \hat{p}_k = \dfrac{e^{o_k}}{\sum e^{o_k}}$

Properties:
- $\hat{p}_k > 0.$
- $\sum\limits_k \hat{p}_k = 1$
- preserves order: $\quad o_i > o_j \implies \hat{p}_i > \hat{p}_j$

# One-hot encoding (OHE)

$K$ entries.

$$c_1 \longrightarrow \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \qquad c_2 \longrightarrow \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \qquad \ldots \qquad c_{K-1} \longrightarrow \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \qquad c_K \longrightarrow \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

**Example:**
**Binary output**

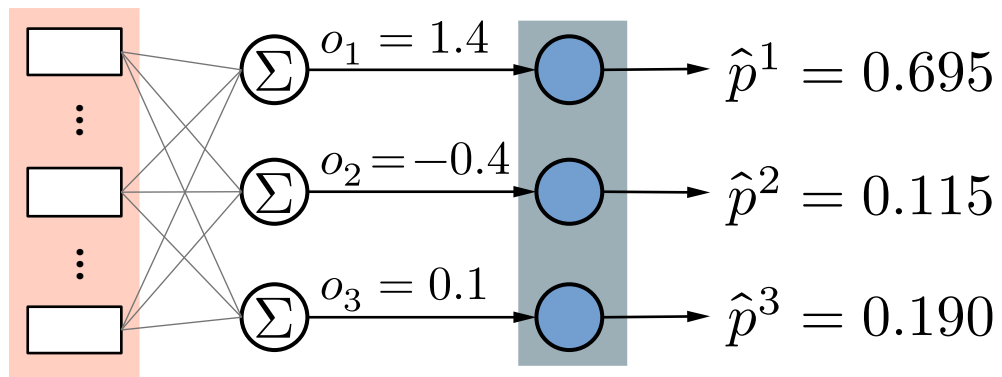| | 0-1 encoding | OHE |
|---|---|---|
| $y_i$ | $y_i = 0$ for $c_1$ <br> $y_i = 1$ for $c_2$ | $y_i = \begin{bmatrix} y_i^1 \\ y_i^2 \end{bmatrix}$ $\begin{array}{l} \nearrow = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ for } c_1 \\ \searrow = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \text{ for } c_2 \end{array}$ |
| $\hat{p}_i$ | $\hat{p}_i \in [0,1].$ <br><br> 0.3 ... c1 | $\hat{p}_i = \begin{bmatrix} \hat{p}_i^1 \\ \hat{p}_i^2 \end{bmatrix} = \begin{bmatrix} 1-\hat{p}_i \\ \mathbf{\textit{w}}\,\hat{p}_i \end{bmatrix}$ .... $\begin{bmatrix} 0.7 \\ 0.3 \end{bmatrix} ...c1.$ |

# Loss function: Multi-class cross entropy

$$\mathsf{CE}(y_i, \hat{p}_i) = -\sum_{k=1}^{K} y_i^k \log \hat{p}_i^k$$

**Recall:** Binary cross entropy under 0-1 encoding

$$\mathsf{CE}(y_i, \hat{p}_i) = -y_i \log(\hat{p}_i) - (1 - y_i) \log(1 - \hat{p}_i)$$

# Example



$$\begin{array}{ccc} & y^k & -y^k \log \widehat{p}^k \\ o_1 = 1.4 \rightarrow \widehat{p}^1 = 0.695 & 0 & 0 \\ o_2 = -0.4 \rightarrow \widehat{p}^2 = 0.115 & 0 & 0 \\ o_3 = 0.1 \rightarrow \widehat{p}^3 = 0.190 & 1 & -\log 0.19 \end{array}$$

$$\mathsf{CE}(y, \widehat{p}) = -\log 0.19$$

# Training the neural network

- <u>Hyper-parameters</u>

  ▶ # of layers $L$

  ▶ # of "neurons" in each layer $p_\ell$

  ▶ activation function for each layer

- <u>Tunable parameters</u>

  ▶ All edge weights

$$\theta = (A_1, b_1, A_2, b_2, \ldots, A_L, b_L, \theta_0, \underline{\theta}_1)$$
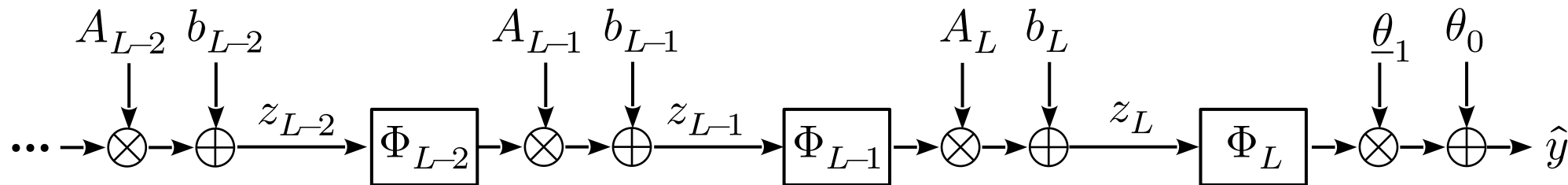
# Gradient of the loss

# Computing $\nabla_\theta h$ with back-propagation
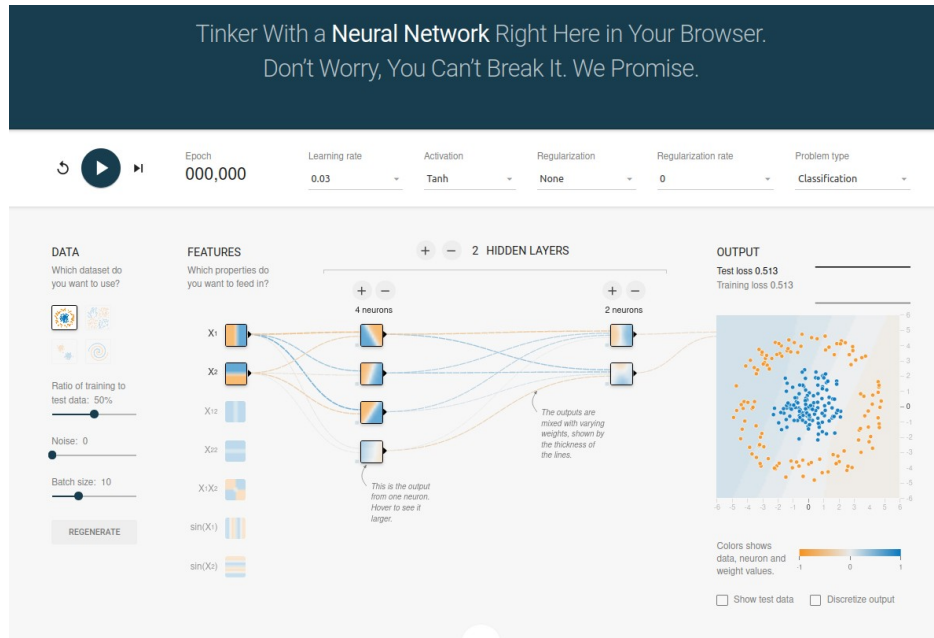
$$\hat{y} = \theta_0 + \Phi_L^T(z_L)\underline{\theta}_1)$$
$$z_L = b_L + A_L\Phi_{L-1}(z_{L-1})$$
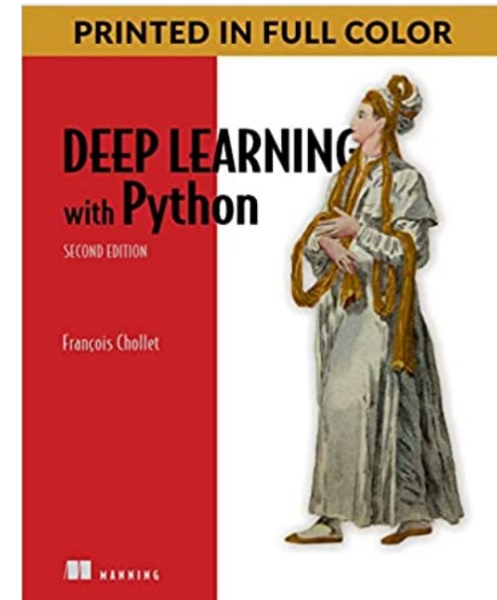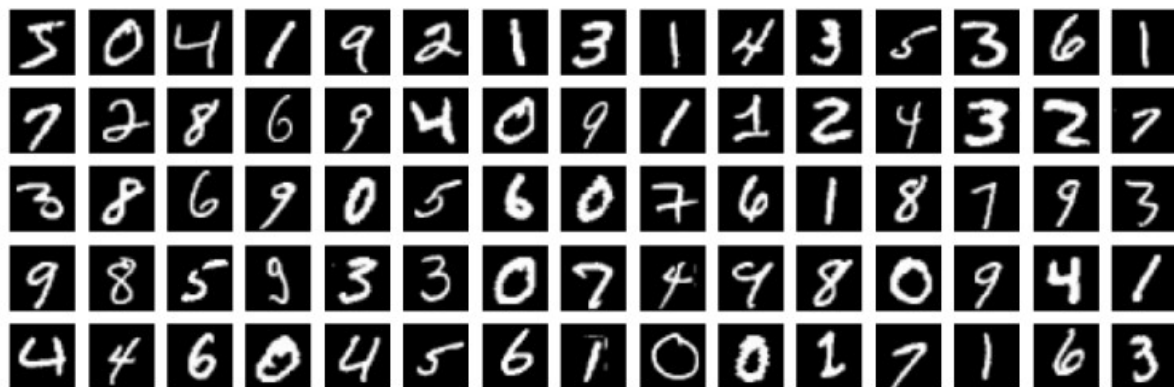$$z_{L-1} = b_{L-1} + A_{L-1}\Phi_{L-2}(z_{L-2})$$
$$\vdots$$

https://playground.tensorflow.org

# MNIST demo

# Neural networks for images

# Convolution operator

pixels

filter

convolution
operator

kernel , bias , activation

$$
\begin{array}{|c|c|c|}
\hline
p_1 & p_2 & p_3 \\
\hline
p_4 & p_5 & p_6 \\
\hline
p_7 & p_8 & p_9 \\
\hline
\end{array}
\;*\;
\left\{
\begin{array}{|c|c|c|}
\hline
f_1 & f_2 & f_3 \\
\hline
f_4 & f_5 & f_6 \\
\hline
f_7 & f_8 & f_9 \\
\hline
\end{array}
,\; b \;,\; \phi(\cdot)
\right\}
= \phi\left( b + \sum_{i=1}^{9} p_i\, f_i \right)
$$

# Example

| | | |
|---|---|---|
| 1.0 | 0.0 | 0.8 |
| 0.1 | 0.3 | 0.5 |
| 0.1 | 0.1 | 0.0 |

$*$

$\left\{ \begin{array}{ccc} -0.2 & 1.2 & 0.2 \\ -0.3 & 0.5 & 0.5 \\ 0.0 & 0.0 & 0.1 \end{array} \;,\; b{=}1 \;,\; \phi{=}\text{ReLU} \right\}$

$$
\begin{aligned}
&= \;\; \phi\left( b + \sum_i p_i\, f_i \right) \\
&= \;\; \text{ReLU}(1 + (1.0)(-0.2) + (0.0)(1.2) + \ldots) \\
&= \;\; \text{ReLU}(1.33) \\
&= \;\; 1.33
\end{aligned}
$$

# Convolution of full images



$$* \left\{ \; \boxplus \; , \; b \; , \; \phi(\cdot) \; \right\} =$$

12x12

10x10

**Padding** $\left\{ \boxplus , \ b , \ \phi(\cdot) , \ p = \text{True} \right\}$



**Stride** $\left\{ \boxplus , \ b , \ \phi(\cdot) , \ s = 2 \right\}$

 $* \left\{ \boxplus , \ b , \ \phi(\cdot) , \ s = 2 \right\}$ $=$ 

Image  Convolution layer 1  Pooling layer 1  Convolution layer 2  Pooling layer 2  Pooling layer L  Flatten, Dense, Softmax

# Pooling