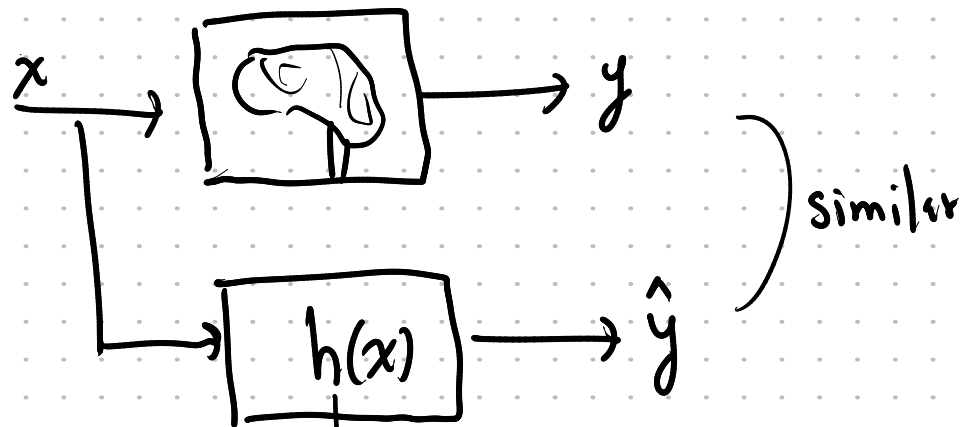




# Statistics and Data Science for Engineers E178 / ME276DS

## Neural networks



• Linear Regression (regression)

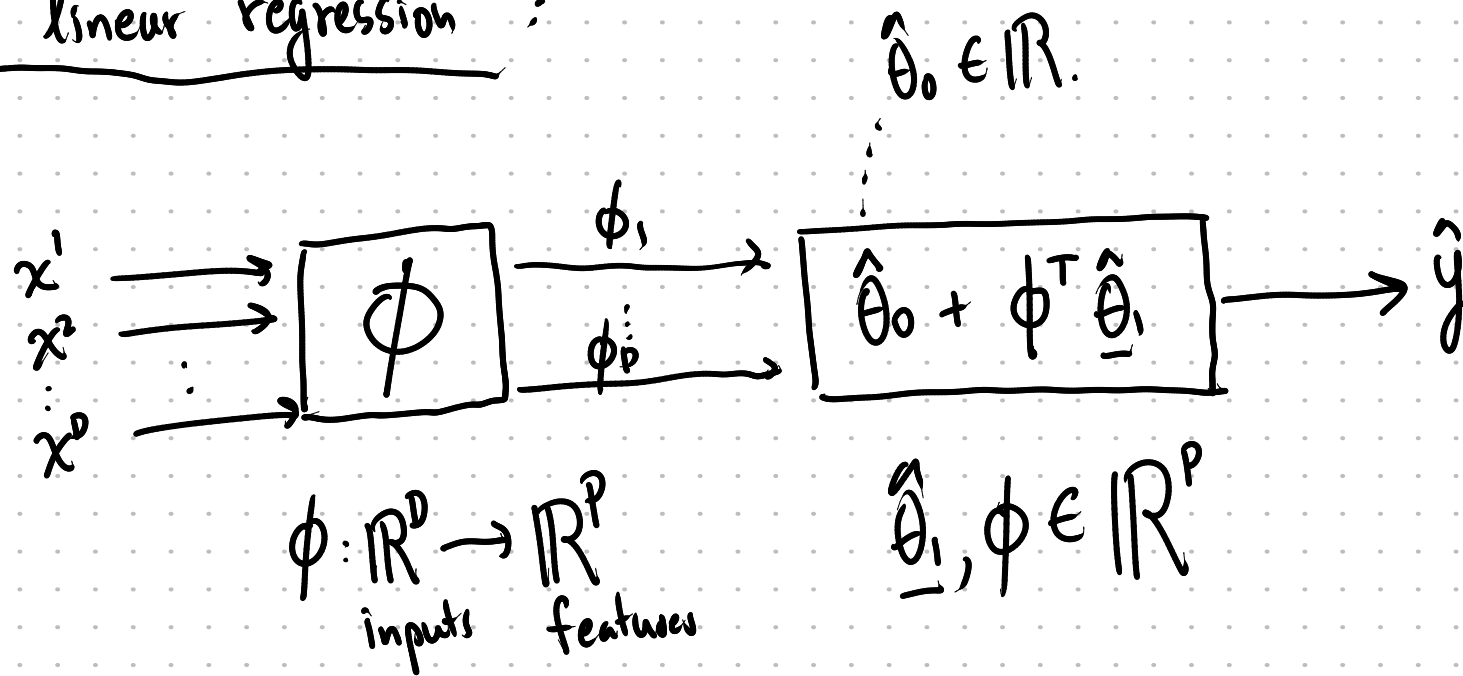
• Logistic regression (classification).

• Support Vector machine

• Decision trees

• Neural networks.

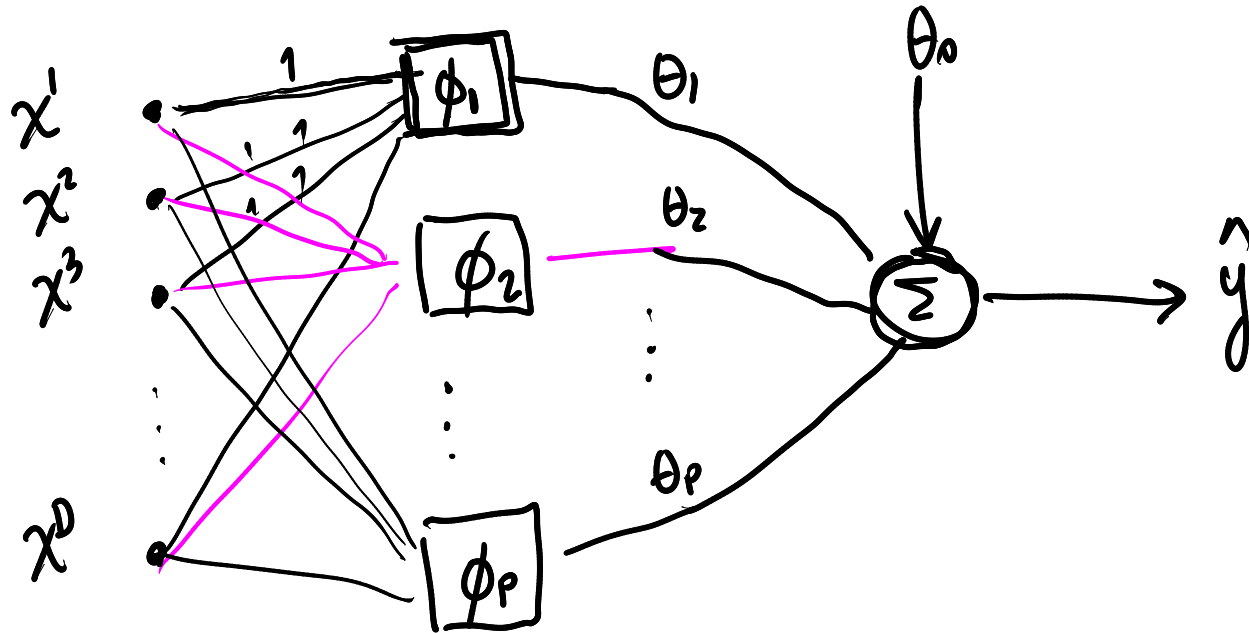
Recall linear regression :



- I have to manually design the features  $\phi$ .
- Neural networks can be understood as a way to design the features automatically.

# Pictorial representation of linear regression

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

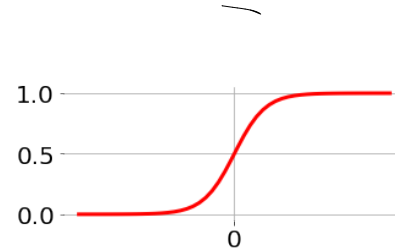


3 steps to NN :

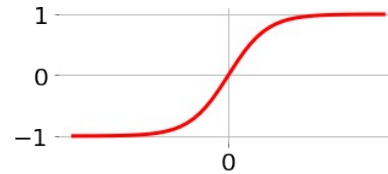
1. Replace the  $\phi$  with "activation functions"
2. Put weights on all edges.
3. Replicate layers.

# 1) Generic nonlinearities, a.k.a. activation functions

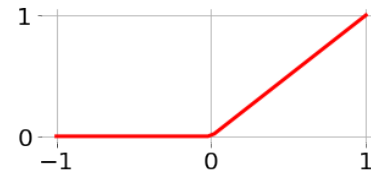
sigmoid:  $\phi(\xi) = \frac{1}{1 + e^{-\xi}}$



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

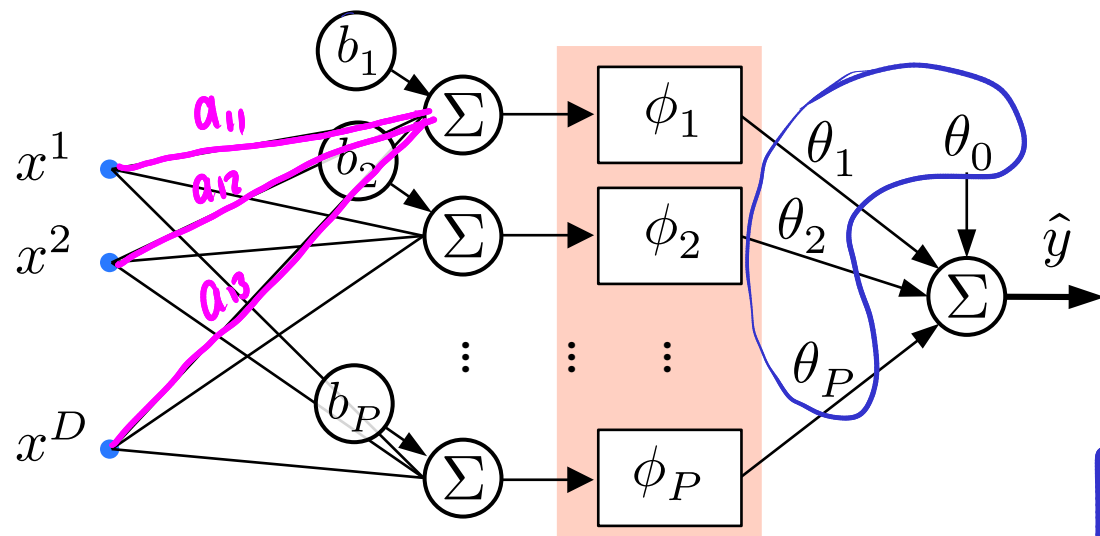


ReLU:  $\phi(\xi) = \max(0, \xi)$



## 2) Weights on the inputs

$$\hat{y} = \theta_0 + \phi^T(b + Ax)\theta_1 = \theta_0 + \sum_{j=1}^P \phi_j(b_j + A_j x)\theta_j$$



$P \times D$  "a" coeff.

$P$  "b" coeff.

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1P} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2P} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{D1} & a_{D2} & a_{D3} & \dots & a_{DP} \end{bmatrix}$$

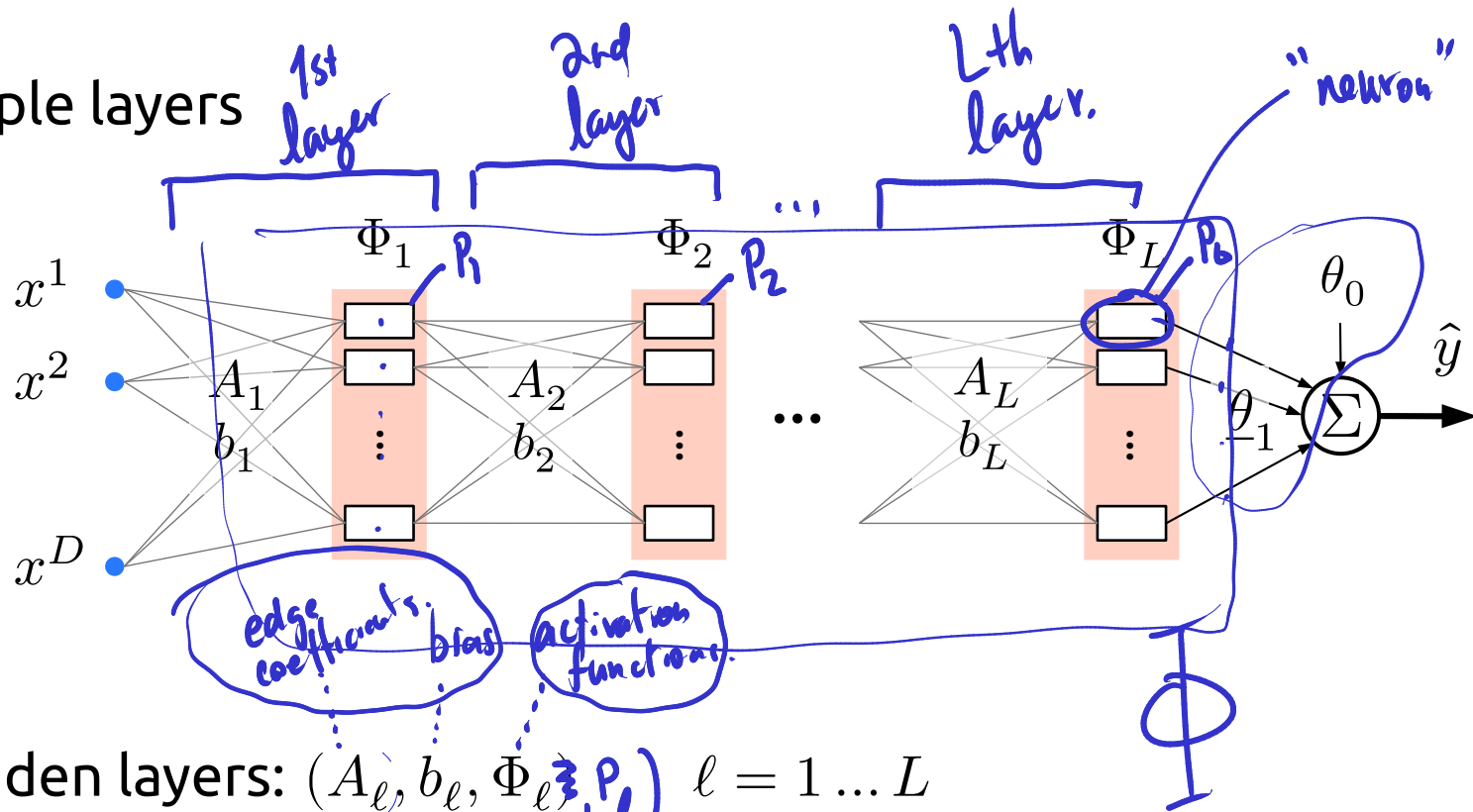
$$b = \begin{bmatrix} b_1 \\ \vdots \\ b_P \end{bmatrix}$$

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

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

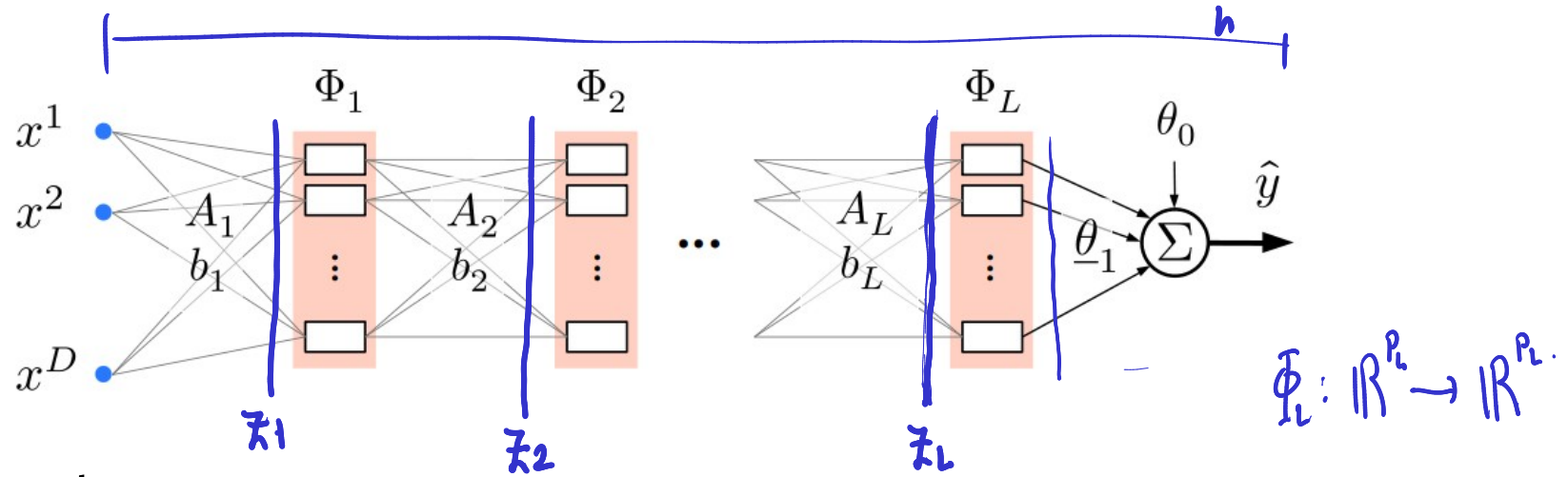
$$A_1 = [a_{11}, a_{12}, a_{13}]$$

### 3) Multiple layers



- Hidden layers:  $(A_\ell, b_\ell, \Phi_\ell, P_\ell) \ell = 1 \dots L$
- Output layer:  $(\theta_0, \theta_1) \dots$  final linear regression. (regression problem)





Put it all together :

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

... complicated nested function.

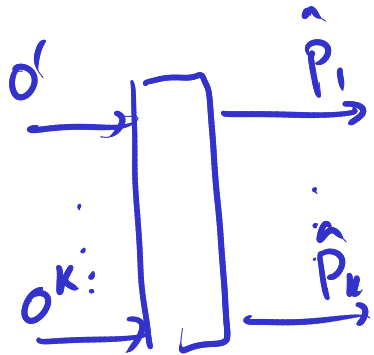
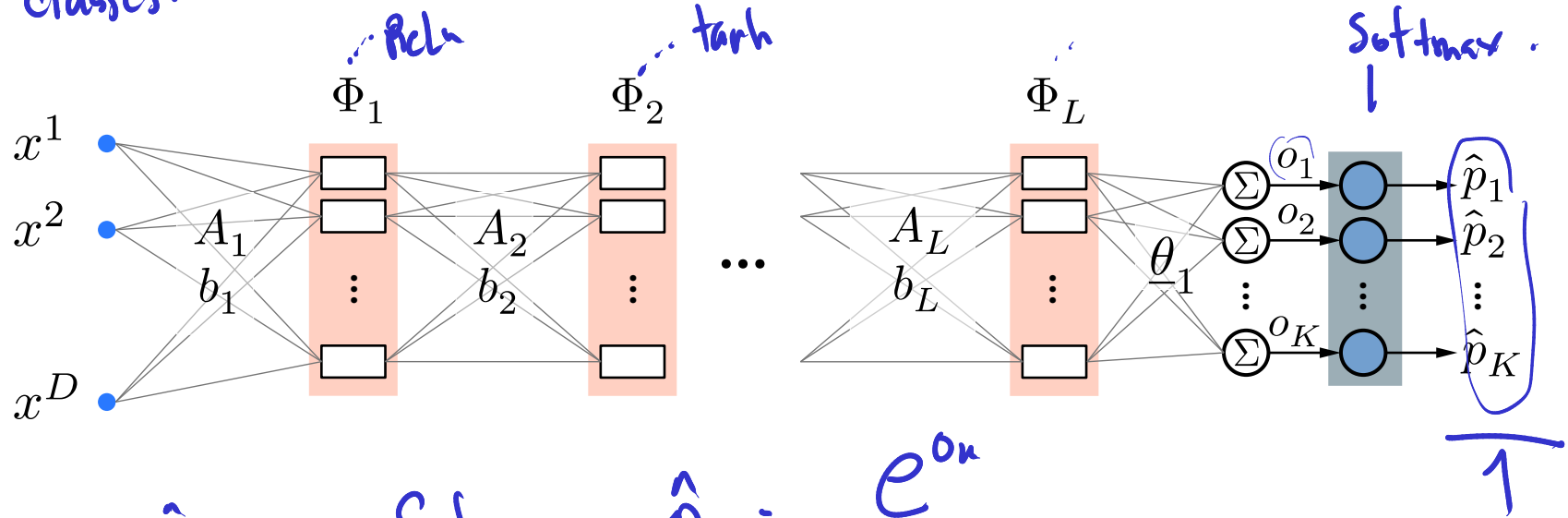
... perceptron.

Write this in a recursive form.

$$\left\{ \begin{array}{l} \hat{y} = \theta_0 + \theta_1^T \Phi_L(z_L) \\ z_L = b_L + A_L \Phi_{L-1}(z_{L-1}) \\ z_{L-1} = b_{L-1} + A_{L-1} \Phi_{L-2}(z_{L-2}) \\ \vdots \\ z_\ell = b_\ell + A_\ell \Phi_{\ell-1}(z_{\ell-1}) \quad (\text{general form}) \\ \vdots \\ z_1 = b_1 + A_1 x \end{array} \right.$$

$K$ ... classes.

# Classification networks



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

Properties:  $\hat{p}_k > 0$ .

$\sum_k \hat{p}_k = 1$

preserves order:  $o_i > o_j \Rightarrow \hat{p}_i > \hat{p}_j$

# One-hot encoding (OHE)



**Example:  
Binary output**

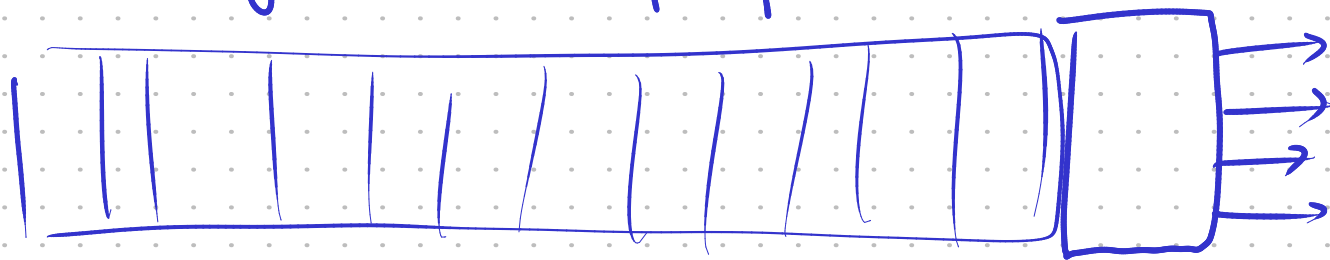
	0-1 encoding	OHE
$y_i$	$y_i = 0$ for $c_1$ $y_i = 1$ for $c_2$	$y_i = \begin{bmatrix} y_i^1 \\ y_i^2 \end{bmatrix} \begin{matrix} \nearrow = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ for } c_1 \\ \searrow = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \text{ for } c_2 \end{matrix}$
$\hat{p}_i$	$\hat{p}_i \in [0, 1]$ $\vdots$	$\hat{p}_i = \begin{bmatrix} \hat{p}_i^1 \\ \hat{p}_i^2 \end{bmatrix} = \begin{bmatrix} 1 - \hat{p}_i \\ \hat{p}_i \end{bmatrix} \dots$

*0.3 ... c1*

*$\begin{bmatrix} 0.7 \\ 0.3 \end{bmatrix} \dots c1.$*

## Recall on neural networks.

- From linear regression.  $\rightarrow$  perceptron.

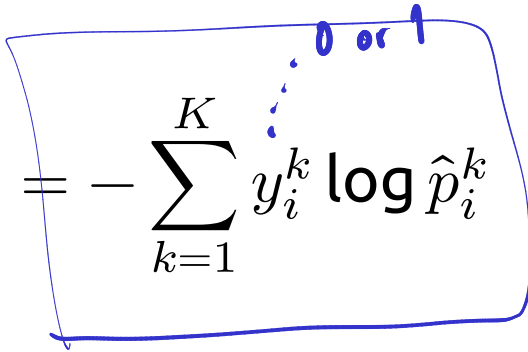


- Classification

→ Output:  $R$  probabilities over classes  
→ OHEncode  $y$ .  
→ Loss function: Multi class CE.

Regression: Linear regression  
Classification: Softmax

## Loss function: Multi-class cross entropy

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


**Recall:** Binary cross entropy under 0-1 encoding (logistic regression)

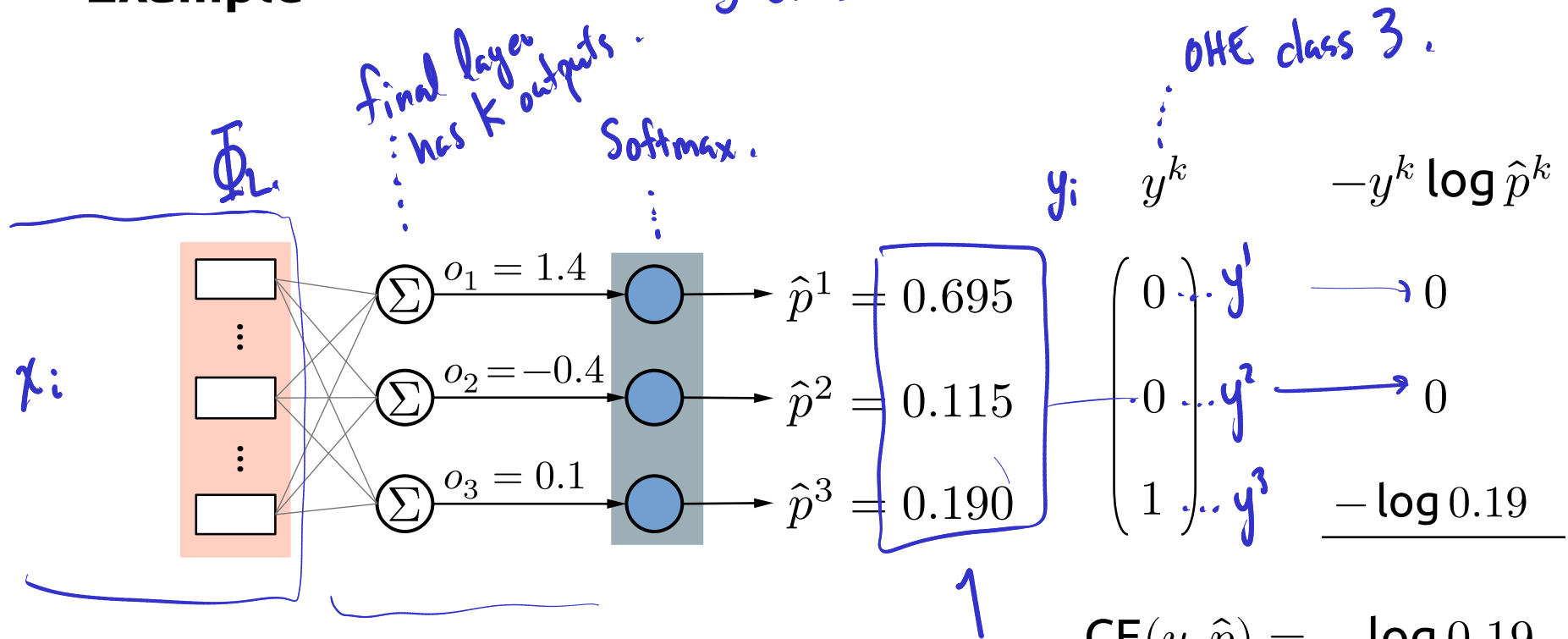
$$\text{CE}(y_i, \hat{p}_i) = -y_i \log(\hat{p}_i) - (1 - y_i) \log(1 - \hat{p}_i)$$

OHE  
multiclass



# Example

3-class classification.



$$CE(y, \hat{p}) = -\log 0.19 = 1.66.$$

$$CE(c_1, \hat{p}) = -\log 0.695 = 1.66$$

smaller than 1.66

# Training the neural network

- Hyper-parameters

- ▶ # of layers  $L$
- ▶ # of “neurons” in each layer  $P_l$
- ▶ activation function for each layer

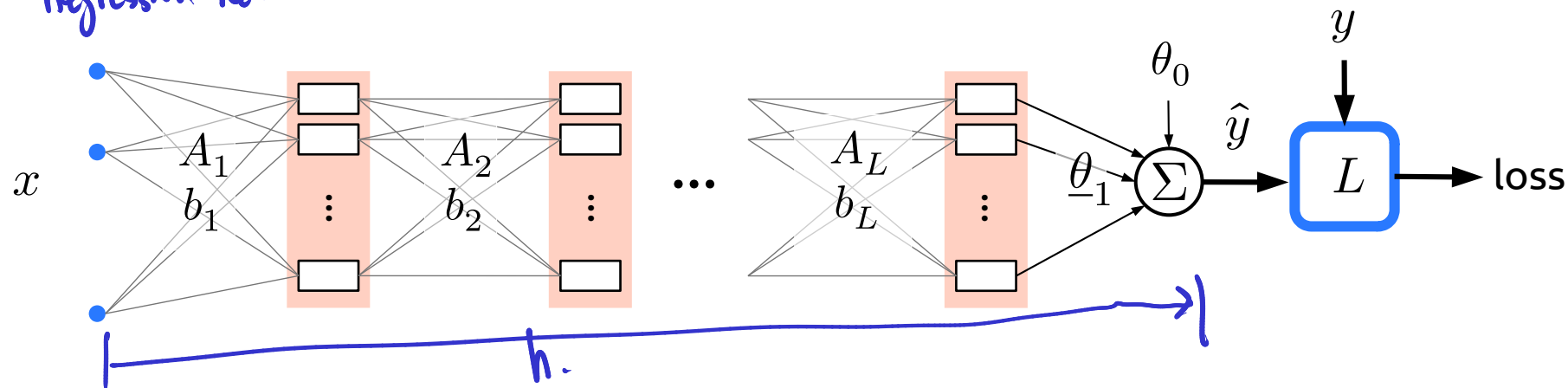
→ Tunable parameters

- ▶ All edge weights

... Backpropagation of the loss.



Regression network



$$\theta = (A_1, b_1, A_2, b_2, \dots, A_L, b_L, \theta_0, \theta_1)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \underbrace{\sum_{i=1}^N L(y_i, h(x_i; \theta))}_{\mathbb{L}}$$

$$\theta_{k+1} = \theta_k - \gamma \nabla_{\theta} \mathbb{L}(\theta_k)$$

$$\nabla_{\theta} L = \sum \frac{dL}{d\hat{y}} \cdot \nabla_{\theta} h \quad \dots \quad L = (\hat{y} - y)^2 \rightarrow \frac{dL}{d\hat{y}} = 2(\hat{y} - y) \quad \dots \quad \text{Regression}$$

scalar
vector

$$\text{or} = \sum \frac{dL}{d\hat{p}} \cdot \nabla_{\theta} h \quad \dots \quad \text{Classification}$$

vector
matrix

Gradient of the ~~loss~~  $\nabla_{\theta} h.$

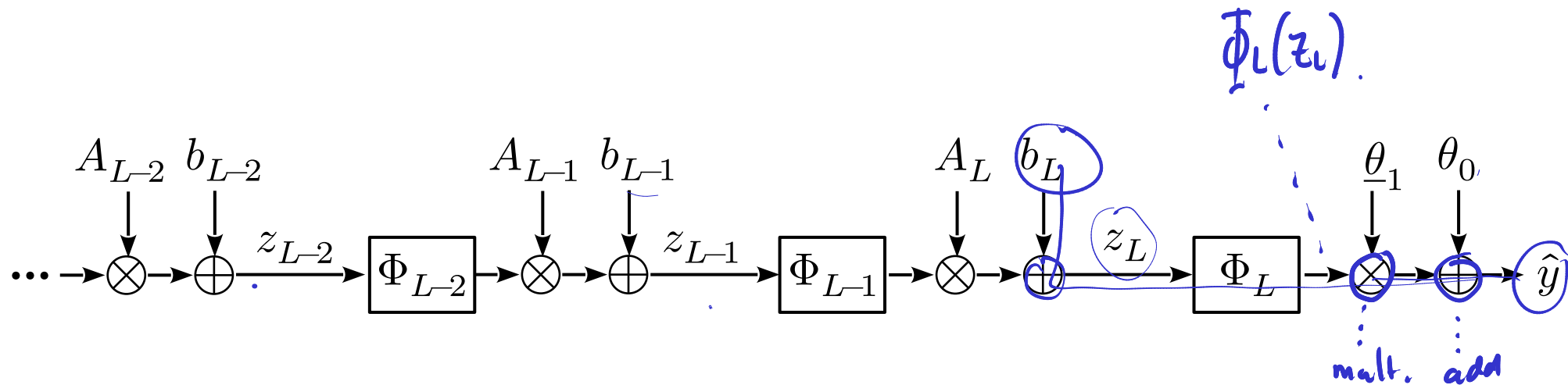
# Computing $\nabla_{\theta} h$ with back-propagation

$$h: \hat{y} = \theta_0 + \Phi_L^T(z_L) \theta_1$$

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

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

$\vdots$



$$\nabla_{\theta} h = \left[ \underbrace{\frac{\partial h}{\partial \theta_0}}_{\checkmark}, \underbrace{\frac{\partial h}{\partial \theta_1}}_{\checkmark}, \frac{\partial h}{\partial b_L}, \frac{\partial h}{\partial A_L}, \dots \right]$$

$$\frac{\partial h}{\partial \theta_0} = 1.$$

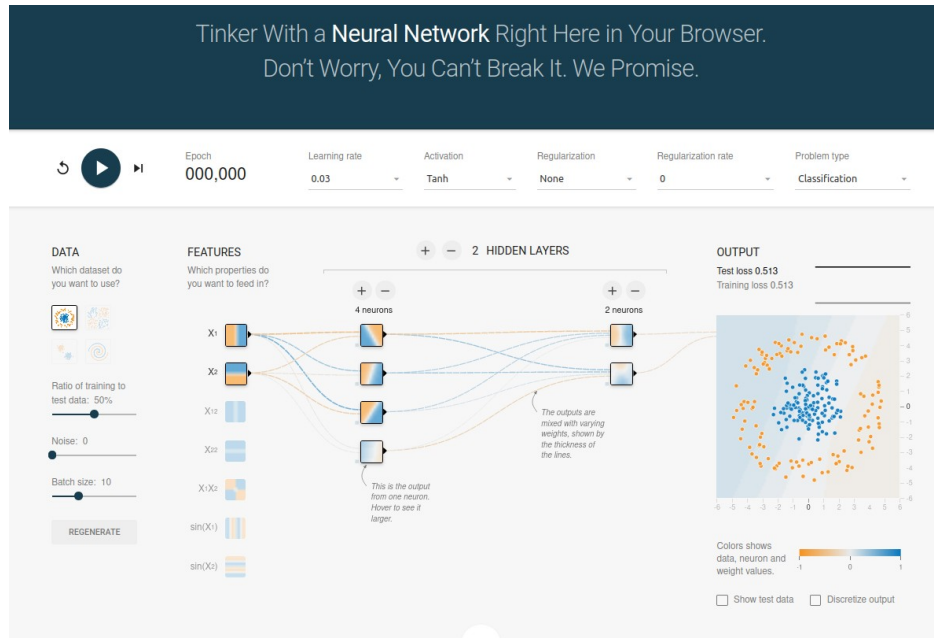
$$\frac{\partial h}{\partial \theta_1} = \Phi_L(z_L).$$

Tensors: generalization  
of a matrix.

$$\frac{\partial h}{\partial b_L} = \frac{\partial h}{\partial \bar{\Phi}_L} \cdot \frac{\partial \bar{\Phi}_L}{\partial z_L} \frac{\partial z_L}{\partial b_L} = \underline{\theta}_1^T [\Phi_L']$$

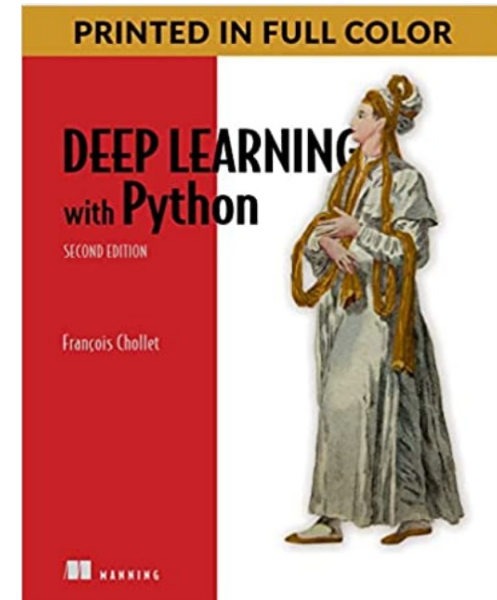
$$\frac{\partial h}{\partial A_L} = \frac{\partial h}{\partial \bar{\Phi}_L} \cdot \frac{\partial \bar{\Phi}_L}{\partial z_L} \frac{\partial z_L}{\partial A_L} = \underline{\theta}_1^T [\Phi_L'] [\Phi_L]$$

<https://playground.tensorflow.org>

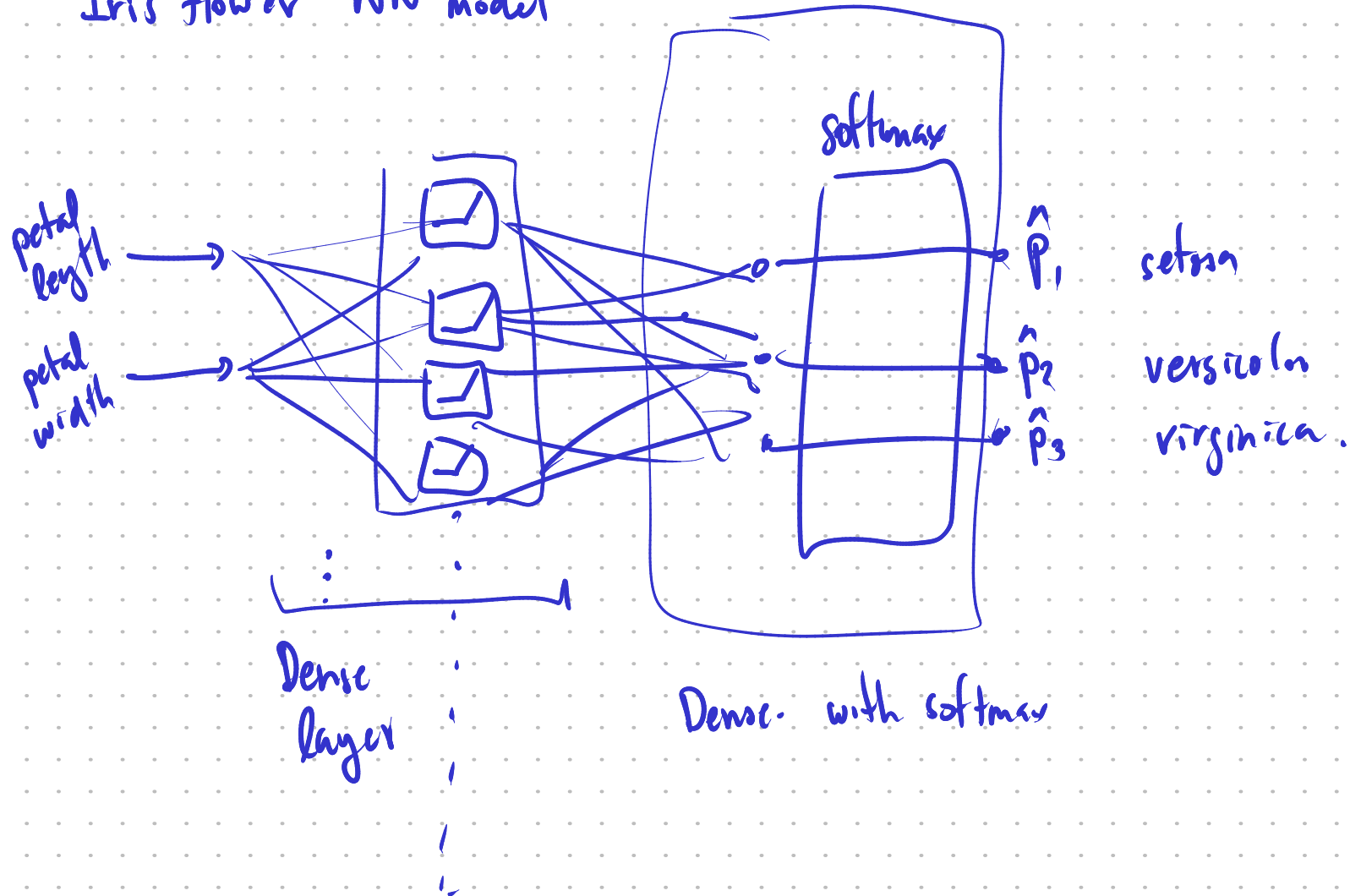


**K** Keras

**TensorFlow**



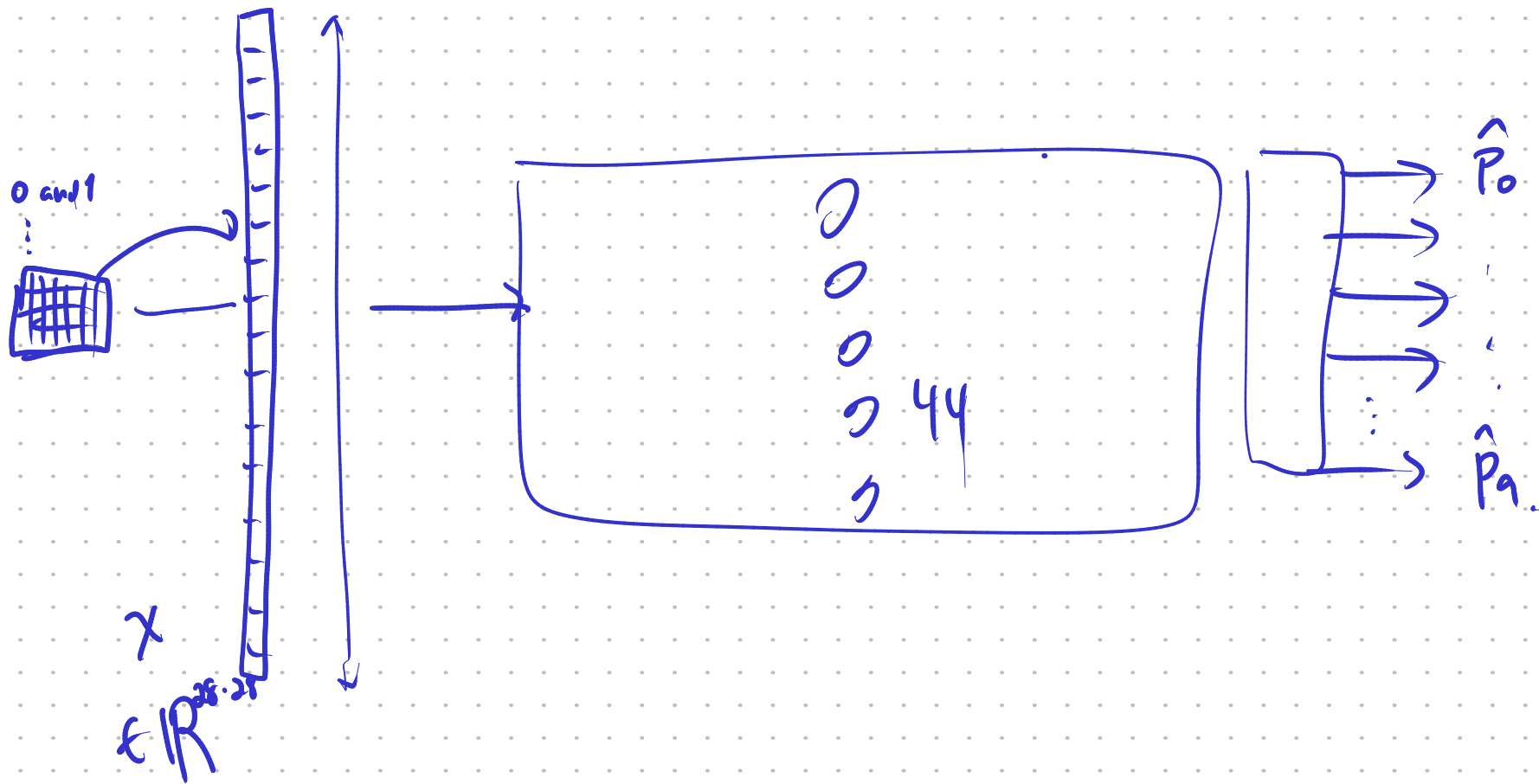
# Iris flower NN model



# MNIST demo

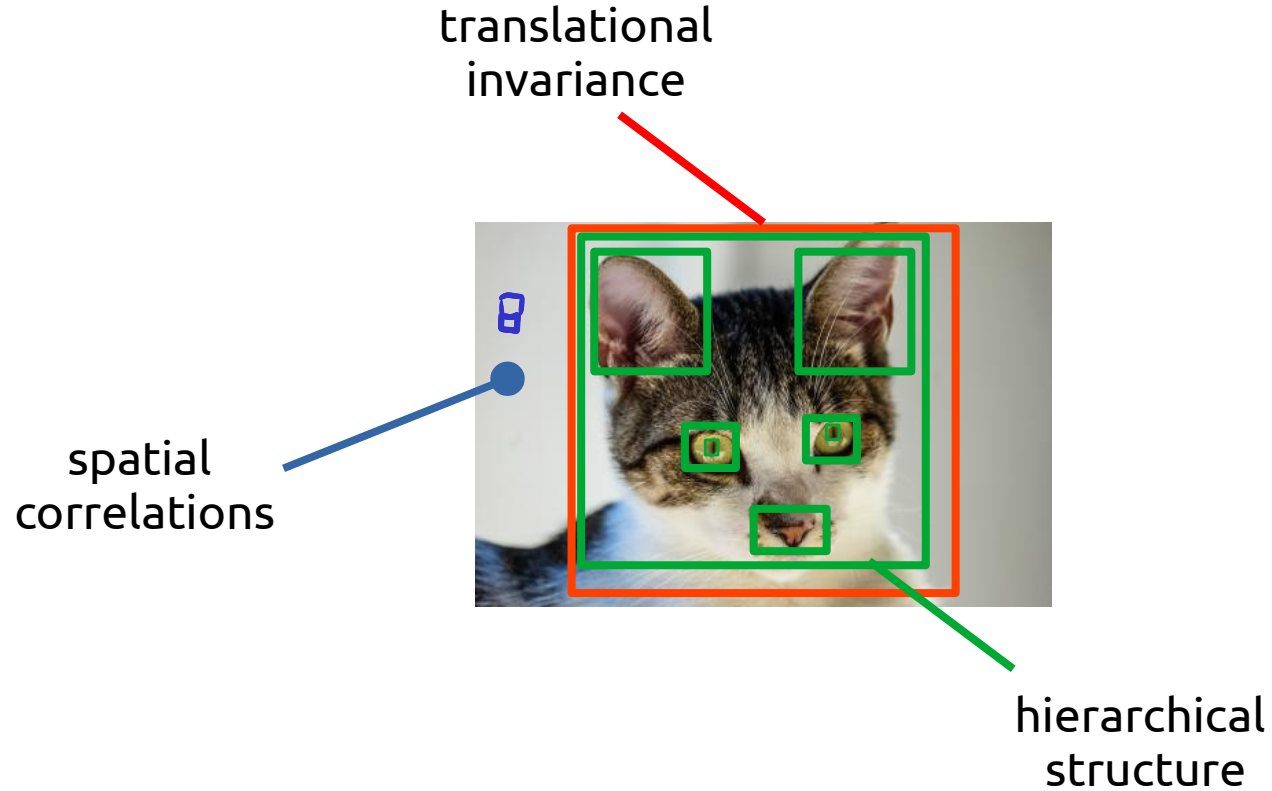




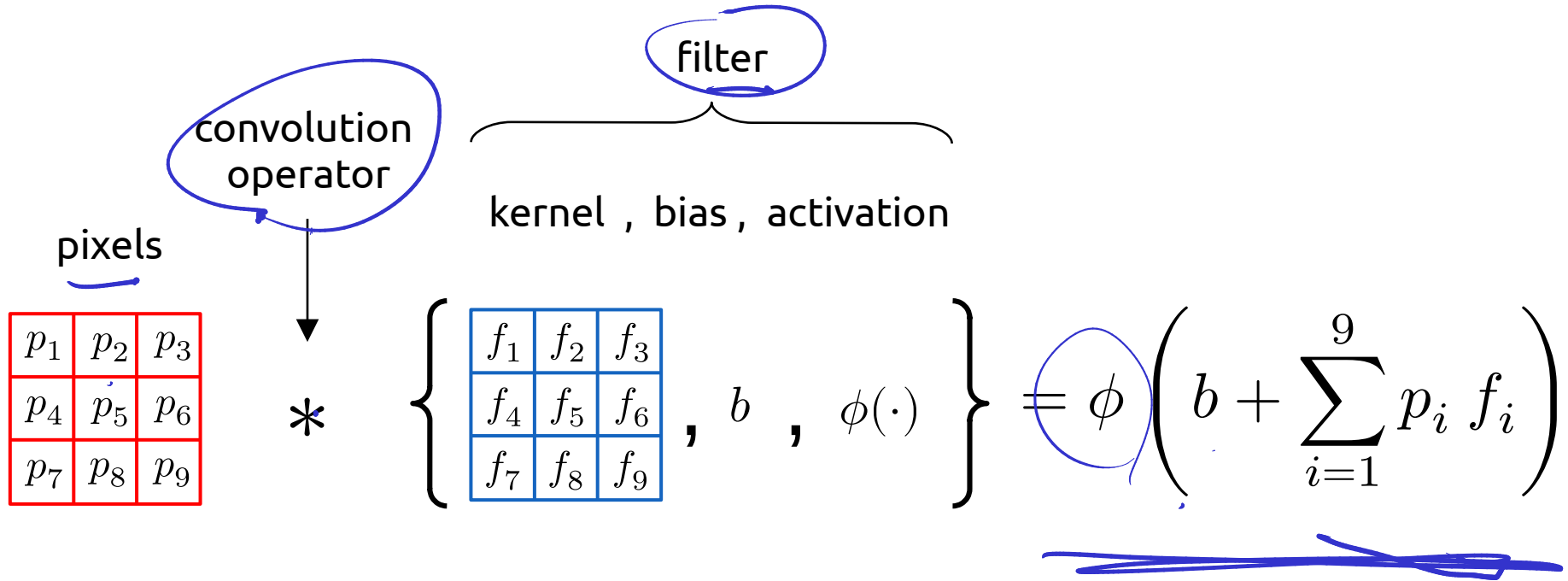


# Neural networks for images

*Convolutional Neural Network.*



# Convolution operator



# Example

1.0	0.0	0.8
0.1	0.3	0.5
0.1	0.1	0.0

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

$$= \phi \left( b + \sum_i p_i f_i \right)$$

$$= \text{ReLU}(1 + (1.0)(-0.2) + (0.0)(1.2) + \dots)$$

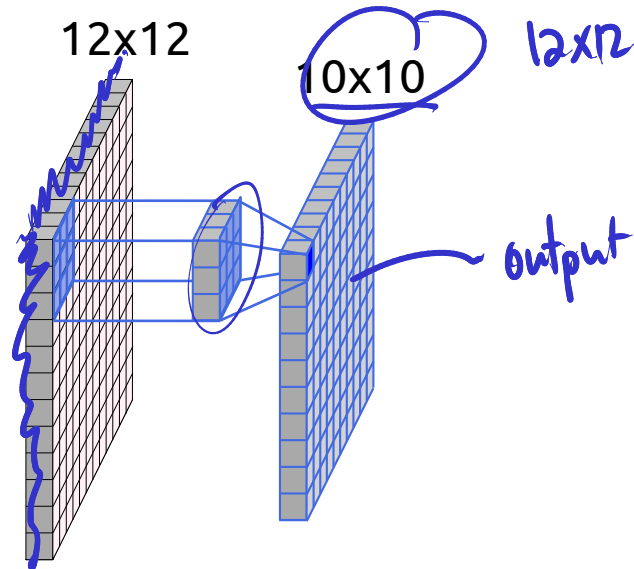
$$= \text{ReLU}(1.33)$$

$$= \underline{1.33}$$

# Convolution of full images



$$* \left\{ \begin{array}{c} \text{3x3 grid} \\ b, \phi(\cdot) \end{array} \right\} =$$



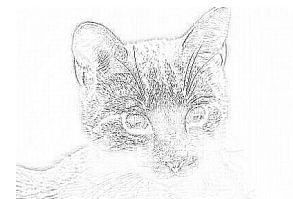
**Padding**     $\left\{ \begin{array}{c} \text{3x3 grid} \\ , b , \phi(\cdot) , p = \text{True} \end{array} \right\}$

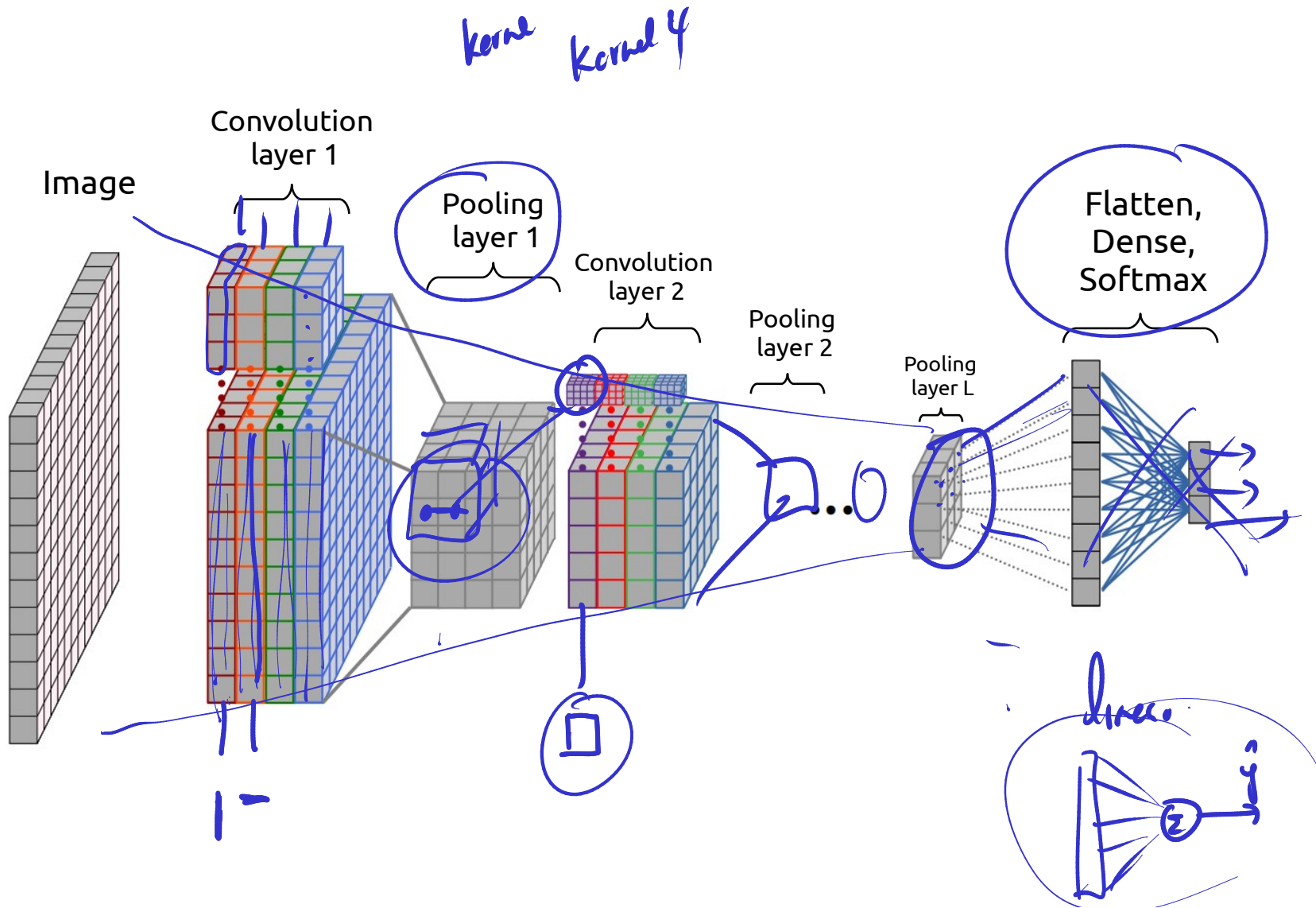


**Stride**     $\left\{ \begin{array}{c} \text{3x3 grid} \\ , b , \phi(\cdot) , s = 2 \end{array} \right\}$



$\ast \left\{ \begin{array}{c} \text{3x3 grid} \\ , b , \phi(\cdot) , s = 2 \end{array} \right\} =$





# Pooling

max pooling.

