

# Deep Learning

## MLP – Backprop

Edgar Roman-Rangel.  
`edgar.roman@itam.mx`

Department of Computer Science.  
Instituto Tecnológico Autónomo de México, ITAM.

# Outline

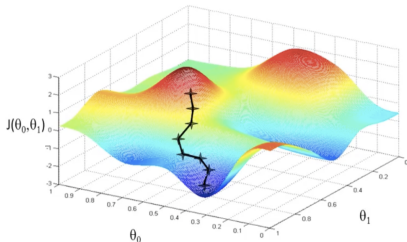
Stochastic Gradient Descent

Multi-layer Perceptron

Backpropagation

Multiple outputs

# Gradient Descent, GD



- ▶ Random initialization.
- ▶ Forward pass.
- ▶ Error estimation.
- ▶ Gradient computation.
- ▶ Backward pass.

**Q:** remember the notion of “taking a step in the direction of steepest descent”?

**A:** The direction of the step is computed considering all of the parameters at once, i.e., we compute first the gradient (all partial derivatives), and then update all the weights.

## GD, pseudocode

---

### Algorithm 1 Gradient Descent

---

- 1: Initialize  $\Omega$  randomly.
  - 2: **for each** epoch **do**
  - 3:   **for each** sample **do**
  - 4:      $\hat{y}_i = f(x_i; \Omega)$
  - 5:      $E_i = l(y_i, \hat{y}_i)$
  - 6:      $\Omega = \Omega - \eta \nabla_{\Omega} E_i$
  - 7:   **end for**
  - 8: **end for**
- 

where,  $\Omega = \{\omega_i\}$ ,  $l(\cdot)$  is a loss function, and  $\nabla_{\Omega} E_i$  is the gradient of the error with respect to the set  $\Omega$  of parameters.

**Q:** can you imagine a drawback of this approach?

**A:** The model will end up being *biased* towards the last seen

## Stochastic GD

Stochastic GD (SGD) tries to compensate for the bias of the last seen training samples.

Each epoch, randomly shuffle the order of samples.

---

### Algorithm 2 Stochastic Gradient Descent

---

- 1: Initialize  $\Omega$  randomly.
  - 2: **for each** epoch **do**
  - 3:    $\{X, y\} = shuffle(\{X, y\})$
  - 4:   **for each** sample **do**
  - 5:      $\hat{y}_i = f(x_i; \Omega)$
  - 6:      $E_i = l(y_i, \hat{y}_i)$
  - 7:      $\Omega = \Omega - \eta \nabla_{\Omega} E_i$
  - 8:   **end for**
  - 9: **end for**
-

# Batch GD

---

**Algorithm 3** Batch Gradient Descent

---

- 1: Initialize  $\Omega$  randomly.
  - 2: Define a number of batches.
  - 3: **for each** epoch **do**
  - 4:    $\{X, y\} = shuffle(\{X, y\})$
  - 5:   **for each** batch **do**
  - 6:      $\{X_B, y_B\} = \text{next } N \text{ training pairs}$
  - 7:      $\hat{y}_B = f(X_B; \Omega)$
  - 8:      $E_B = \frac{1}{N} \sum_{n=1}^N l(y_{B_n}, \hat{y}_{B_n})$
  - 9:      $\Omega = \Omega - \eta \nabla_{\Omega} E_B$
  - 10:   **end for**
  - 11: **end for**
- 

**Q:** What advantage would it have to use BGD instead of SGD?

# BGD

a.k.a., mini-batch gradient descent.

- ▶ Approximates  $E$  with the average error of a batch of samples.
- ▶ Fewer updates, i.e., faster optimization process.
- ▶ Batch size  $bs = 1$  boils down to regular GD.
- ▶ Common batch sizes: 16, 32, 64, 128, 256.

# Outline

Stochastic Gradient Descent

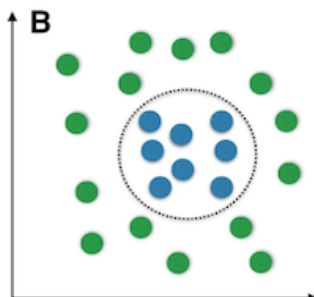
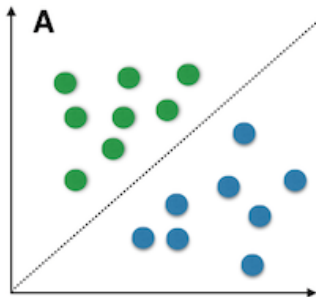
Multi-layer Perceptron

Backpropagation

Multiple outputs

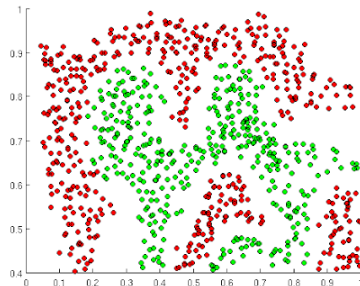
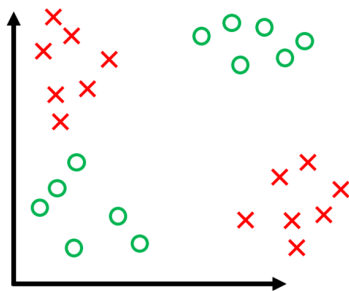


# Linearity



# Non-linear Separability

Most real-world problems are non-linearly separable.



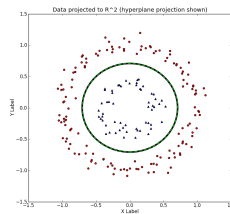
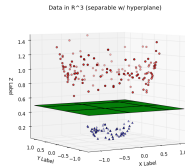
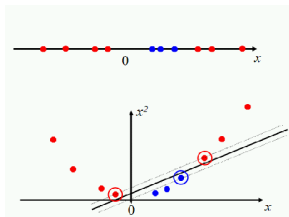
**Q:** How do we do in these cases in machine learning?

# Non-linear transformations

Feature engineering.

Examples:

- ▶ New feature,  $x_2 = (x_1)^2$ .
- ▶ New feature,  $x_3 = x_2 * x_1$ .
- ▶ Feature selection: use only a subset of features.



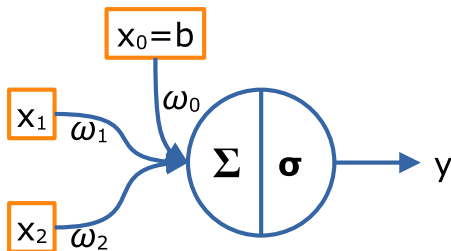
## Example 1

Logical AND: can be solved with a single perceptron.

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	1

Solved using:

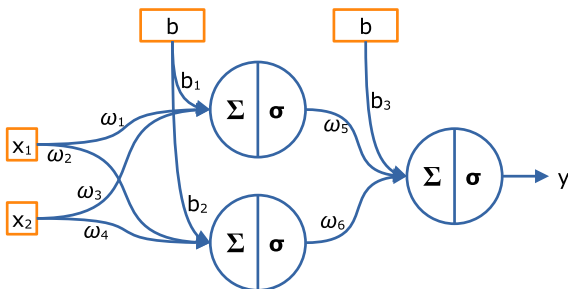
$\omega_0$	$\omega_1$	$\omega_2$
-1.5	1	1



## Example II

Logical XOR: not solved with a single perceptron. Let's cascade non-linear activation functions: multi-layer perceptron (MLP).

$x_0$	$x_1$	$x_2$	$y$
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	0



Solved using:

$b_1$	$b_2$	$b_3$	$\omega_1$	$\omega_2$	$\omega_3$	$\omega_4$	$\omega_5$	$\omega_6$
-10	30	-30	20	-20	20	20	-20	20

**Q:** What if we omit the non-linearity activation functions?

# MLP

Consecutive linear operations are equivalent to a single linear operations, i.e., DL is enable by the use of non-linear activation functions.

- ▶ We end up with: input, hidden, and output layers.
- ▶ Intermediate representations correspond feature engineering.
- ▶ However, features are learned rather than engineered.
- ▶ End-to-end process.
- ▶ Information abstraction increases with depth.
- ▶ Inspired on human brain?

## Nonlinearities

In general, the more difficult the problem looks, the more chances are it is non-linearly separable. Therefore, the deeper the model must be.



**Q:** But how do we do it?

# Outline

Stochastic Gradient Descent

Multi-layer Perceptron

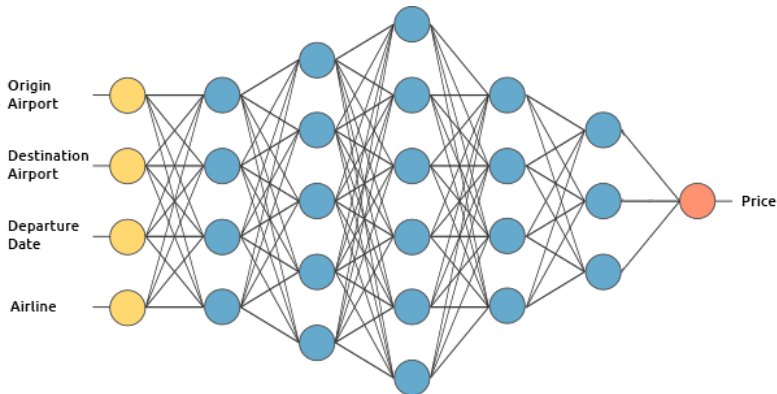
Backpropagation

Multiple outputs



## GD and depth

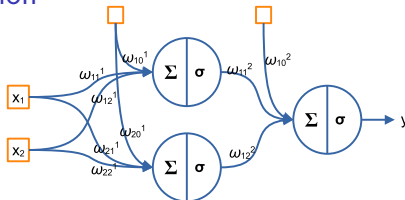
**Q:** How do we perform GD on deep networks?



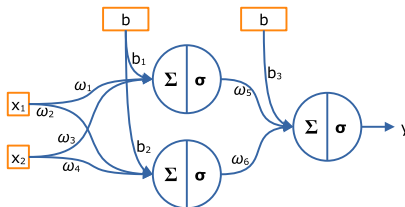
**A:** Use backpropagation (*backprop*): gradient descent + chain rule.

# Notation I

## Common notation



## Simplified notation

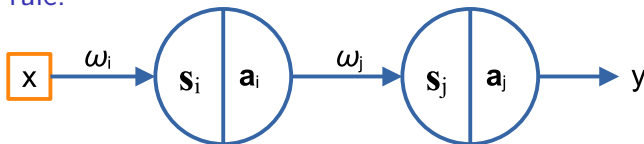


## Notation II

Let's also use:

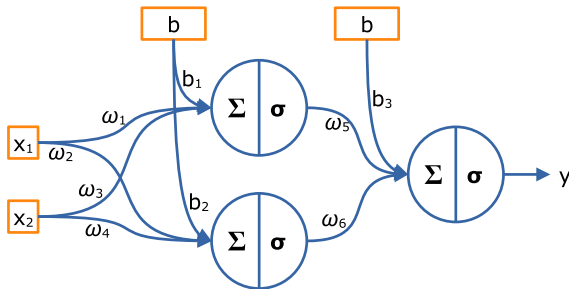
- ▶  $s = \sum_i \omega_i x_i + b$ , linear combination.
- ▶  $a = \sigma(s)$ , non-linear activation.

Chain rule:



$$\frac{\partial E}{\partial \omega_i} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial a_j} \frac{\partial a_j}{\partial s_j} \frac{\partial s_j}{\partial a_i} \frac{\partial a_i}{\partial s_i} \frac{\partial s_i}{\partial \omega_i}.$$

## Example



1 training example  $\{\mathbf{x} = [0.05, 0.1], \quad y = 1\}$ .

Initialize with:

$\omega_1$	$\omega_2$	$\omega_3$	$\omega_4$	$\omega_5$	$\omega_6$	$b_1$	$b_2$	$b_3$
0.15	0.20	0.25	0.30	0.50	0.55	0.35	0.35	0.60

And let's use  $\eta = 0.1$  and  $E = \frac{1}{2}(y - \hat{y})^2$ .

## Example cont.

Remember the process:

1. Forward pass.
2. Error calculation.
3. Gradient computation.
4. Backward pass (weights update).

## Example cont. (Forward pass)

$$\begin{aligned}s_1 &= \omega_1 x_1 + \omega_3 x_2 + b_1, \\ &= (0.15)(0.05) + (0.25)(0.1) + 0.35, \\ &= 0.3825.\end{aligned}$$

$$\begin{aligned}a_1 &= \sigma(s_1), \\ &= 0.5945.\end{aligned}$$

$$\begin{aligned}s_2 &= \omega_2 x_1 + \omega_4 x_2 + b_2, \\ &= (0.2)(0.05) + (0.3)(0.1) + 0.35, \\ &= 0.39.\end{aligned}$$

$$\begin{aligned}a_2 &= \sigma(s_2), \\ &= 0.5964.\end{aligned}$$

$$\begin{aligned}s_3 &= \omega_5 a_1 + \omega_6 a_2 + b_3, \\ &= (0.5)(0.5945) + (0.55)(0.5964) + 0.6, \\ &= 1.2252.\end{aligned}$$

$$\begin{aligned}a_3 &= \sigma(s_3), \\ &= 0.773.\end{aligned}$$

$$\begin{aligned}\hat{y} &= a_3, \\ &= 0.773.\end{aligned}$$

## Example cont. (Error)

$$\begin{aligned} E &= \frac{1}{2}(y - \hat{y})^2, \\ &= (0.5)(1 - 0.773)^2, \\ &= 0.0258. \end{aligned}$$

Example cont. (Gradients  $\frac{\partial E}{\partial \omega_n}$ )

$$\begin{aligned}\frac{\partial E}{\partial \omega_6} &= \frac{\partial E}{\partial a_3} \cdot \frac{\partial a_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial \omega_6}, \\&= \frac{\partial}{\partial a_3} \frac{1}{2} (y - \hat{y})^2 \cdot \frac{\partial}{\partial s_3} \sigma(s_3) \cdot \frac{\partial}{\partial \omega_6} (\omega_5 a_1 + \omega_6 a_2 + b_3), \\&= \frac{2}{2} (y - a_3) \frac{\partial}{\partial a_3} (y - a_3) \cdot \sigma(s_3) (1 - \sigma(s_3)) \cdot a_3, \\&= (y - a_3) (-1) \sigma(s_3) (1 - \sigma(s_3)) a_3, \\&= -(1 - 0.773)(0.773)(1 - 0.773)(0.5963), \\&= -0.0238.\end{aligned}$$

Now, let's define:

$$\delta_3 = (y - a_3) (-1) \sigma(s_3) (1 - \sigma(s_3))$$



Example cont. (Gradients  $\frac{\partial E}{\partial \omega_n}$ )

$$\begin{aligned}\frac{\partial E}{\partial \omega_5} &= \frac{\partial E}{\partial a_3} \cdot \frac{\partial a_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial \omega_5}, \\ &= \delta_3 a_1, \\ &= -0.0237.\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial b_3} &= \frac{\partial E}{\partial a_3} \cdot \frac{\partial a_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial b_3}, \\ &= \delta_3, \\ &= -0.0398.\end{aligned}$$

Example cont. (Gradients  $\frac{\partial E}{\partial \omega_n}$ .)

$$\begin{aligned}\frac{\partial E}{\partial \omega_1} &= \frac{\partial E}{\partial a_3} \cdot \frac{\partial a_3}{\partial s_3} \cdot \frac{\partial s_3}{\partial a_1} \cdot \frac{\partial a_1}{\partial s_1} \cdot \frac{\partial s_1}{\partial \omega_1}, \\ &= \delta_3 \cdot \frac{\partial}{\partial a_1} (\omega_5 a_1 + \omega_6 a_2 + b_3) \cdot \frac{\partial}{\partial s_1} \sigma(s_1) \cdot \frac{\partial}{\partial \omega_1} (\omega_1 x_1 + \omega_3 x_2 + b_1) \\ &= \delta_3 \cdot \omega_5 \cdot \sigma(s_1) (1 - \sigma(s_1)) \cdot x_1, \\ &= (-0.0398)(0.5)(0.5945)(1 - 0.5945)(0.05), \\ &= -0.0002.\end{aligned}$$

And, let's use:

$$\delta_2 = \delta_3 \omega_5 \sigma(s_1) (1 - \sigma(s_1)).$$

Example cont. (Gradients  $\frac{\partial E}{\partial \omega_n}$ )

$$\begin{aligned}\frac{\partial E}{\partial \omega_3} &= \delta_2 x_2, \\ &= -0.0005.\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial b_1} &= \delta_2, \\ &= -0.0048.\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial \omega_2} &= \delta_3 \cdot \frac{\partial s_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial s_2} \cdot \frac{\partial s_2}{\partial \omega_2}, \\ &= \delta_1 x_1, \\ &= -0.0003.\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial \omega_4} &= \delta_1 x_2, \\ &= -0.0005.\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial b_2} &= \delta_1, \\ &= -0.0053.\end{aligned}$$

With:

$$\delta_1 = \delta_3 \cdot \frac{\partial s_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial s_2}.$$

## Example cont. (Backward pass)

$$\omega_n = \omega_n - \eta \frac{\partial E}{\partial \omega_n}.$$

$$\omega_1 = 0.15 \quad -(0.1)(-0.0002) = 0.15002,$$

$$\omega_2 = 0.2 \quad -(0.1)(-0.0003) = 0.20003,$$

$$\omega_3 = 0.25 \quad -(0.1)(-0.0005) = 0.25005,$$

$$\omega_4 = 0.3 \quad -(0.1)(-0.0005) = 0.30005,$$

$$\omega_5 = 0.5 \quad -(0.1)(-0.0237) = 0.50237,$$

$$\omega_6 = 0.55 \quad -(0.1)(-0.0238) = 0.55238,$$

$$b_1 = 0.3505,$$

$$b_2 = 0.3505,$$

$$b_3 = 0.604,$$

## Example cont. (New forward pass)

$$\hat{y}(0) = 0.773.$$

After one update:

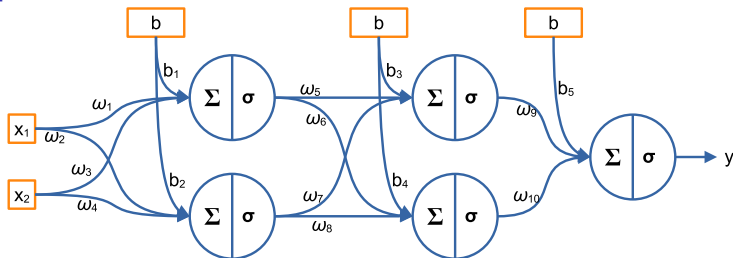
$\omega_1$	$\omega_2$	$\omega_3$	$\omega_4$	$\omega_5$	$\omega_6$
0.15002	0.20003	0.25005	0.30005	0.50237	0.55238

$b_1$	$b_2$	$b_3$
0.35048	0.35053	0.60398

$$\hat{y}(1) = 0.7742.$$

Notice: The impact of backprop is proportional to the depth of the layer: weights in shallow layers are update more softly with respect to those in deeper layers.

## Multiple connections



$$\frac{\partial E}{\partial \omega_5} = \frac{\partial E}{\partial a_5} \frac{\partial a_5}{\partial s_5} \frac{\partial s_5}{\partial a_3} \frac{\partial a_3}{\partial s_3} \frac{\partial s_3}{\partial \omega_5}.$$

$$\frac{\partial E}{\partial \omega_1} = \frac{\partial E}{\partial a_5} \frac{\partial a_5}{\partial s_5} \frac{\partial s_5}{\partial a_3} \frac{\partial s_3}{\partial s_3} \frac{\partial a_3}{\partial a_1} \frac{\partial s_1}{\partial s_1} \frac{\partial a_1}{\partial \omega_1} + \frac{\partial E}{\partial a_5} \frac{\partial a_5}{\partial s_5} \frac{\partial s_5}{\partial a_4} \frac{\partial s_4}{\partial s_4} \frac{\partial a_4}{\partial a_1} \frac{\partial s_1}{\partial s_1} \frac{\partial a_1}{\partial \omega_1}.$$

# Outline

Stochastic Gradient Descent

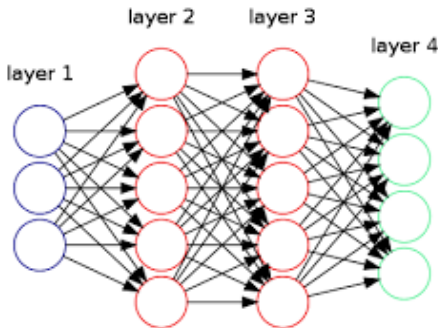
Multi-layer Perceptron

Backpropagation

Multiple outputs

## Multi-variate regression

- ▶ One output perceptron works well for uni-variate regression, i.e.,  $\hat{y}$  is a scalar.
- ▶ More perceptrons can be used for a multi-variate problem, i.e.,  $\hat{y}$  is a vector.





## Multi-class classification

- ▶ One output perceptron works well for binary classification problems, i.e.,  $\hat{y}$  is a scalar indicating the probability of the input belonging to the positive class.
- ▶ More perceptrons can be used for a multi-class classification problem, i.e.,  $\hat{y}$  is a vector indicating the probability of belonging to each possible class.
- ▶ In this case, the ground-truth is a one-hot encoding vector. E.g.,  $\mathbf{y} = [0, 0, 1, 0, 0]$ .



# Q&A

Thank you!

`edgar.roman@itam.mx`