

Backpropagation

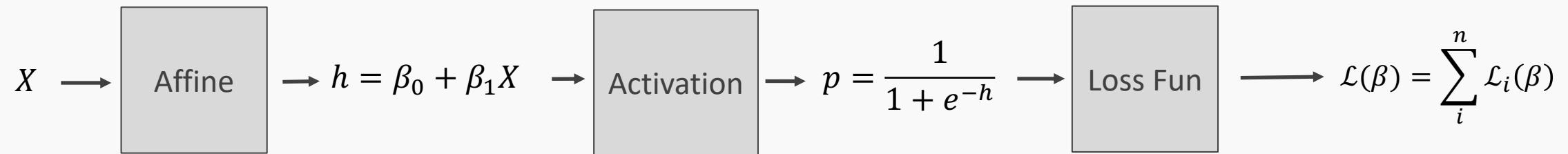
Pavlos Protopapas



Gradient Descent Considerations

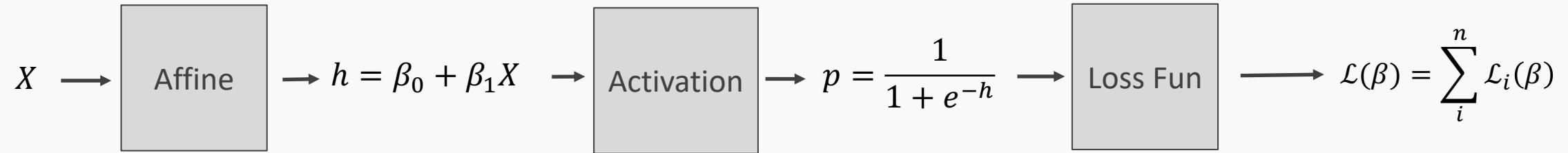
- We still need to calculate the derivatives.
- We need to set the learning rate.
- Local vs global minima.
- The full likelihood function includes summing up all individual ‘errors’. Sometimes this includes hundreds of thousands of examples.

Logistic Regression Revisited



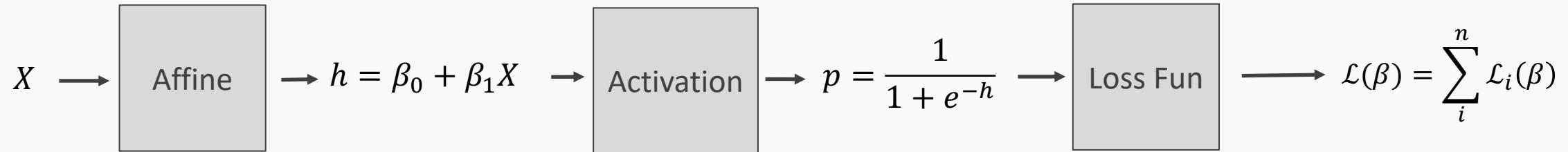
Logistic Regression Revisited

$$\mathcal{L}_i = -y \log p - (1 - y) \log (1 - p)$$



Logistic Regression Revisited

$$\mathcal{L}_i = -y \log p - (1 - y) \log (1 - p)$$



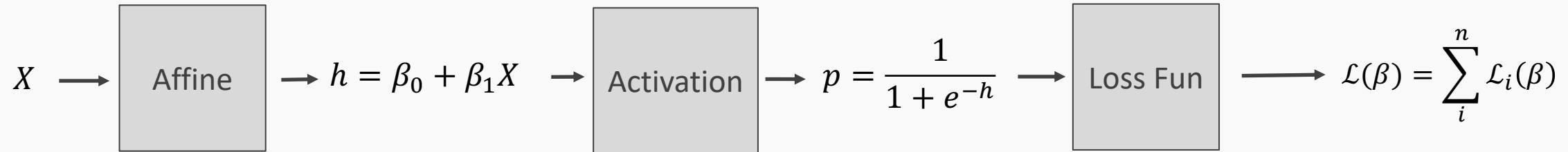
$$\frac{\partial \mathcal{L}}{\partial p}$$

$$\frac{\partial \mathcal{L}}{\partial p} = -y \frac{1}{p} - (1 - y) \frac{1}{1 - p}$$



Logistic Regression Revisited

$$\mathcal{L}_i = -y \log p - (1 - y) \log (1 - p)$$

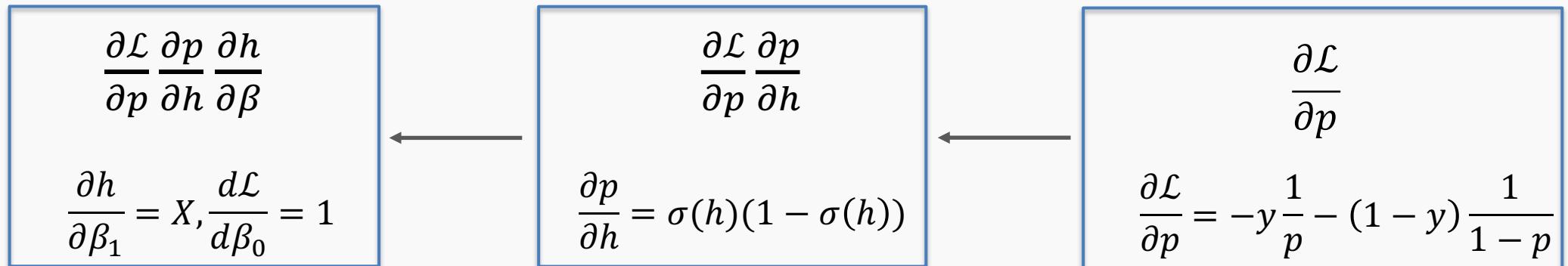
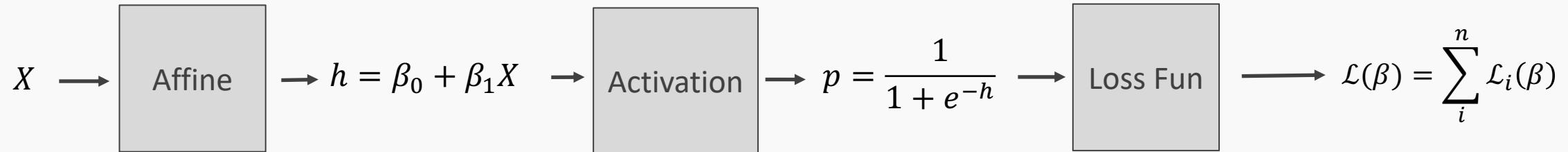


```
graph TD; TopLeft["dL/dp * dL/dh"] --> BottomLeft["dL/dh"]; BottomLeft["dL/dh = sigma(h)(1 - sigma(h))"] --> BottomRight["dL/dp = -y/p - (1 - y)/(1 - p)"]
```

The diagram shows the backward pass for computing gradients. It consists of two boxes connected by a double-headed arrow. The left box contains the partial derivative $\frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial h}$ and the formula $\frac{\partial p}{\partial h} = \sigma(h)(1 - \sigma(h))$. The right box contains the formula $\frac{\partial \mathcal{L}}{\partial p} = -y \frac{1}{p} - (1 - y) \frac{1}{1 - p}$.

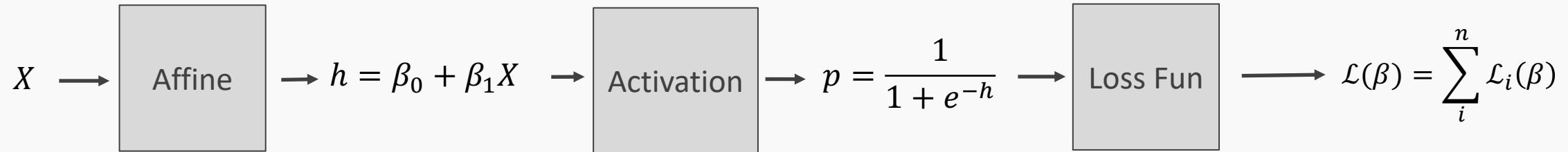
Logistic Regression Revisited

$$\mathcal{L}_i = -y \log p - (1 - y) \log (1 - p)$$



Logistic Regression Revisited

$$\mathcal{L}_i = -y \log p - (1 - y) \log (1 - p)$$



The diagram shows the application of the chain rule to find the gradients of the loss function with respect to the parameters β_0 and β_1 .

Starting from the rightmost box:

$$\frac{\partial \mathcal{L}}{\partial p}$$

$$\frac{\partial \mathcal{L}}{\partial p} = -y \frac{1}{p} - (1 - y) \frac{1}{1 - p}$$

Using the chain rule:

$$\frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial h} \frac{\partial h}{\partial \beta_1}$$

$$\frac{\partial p}{\partial h} = \sigma(h)(1 - \sigma(h))$$

From the leftmost box:

$$\frac{\partial h}{\partial \beta_1} = X, \frac{d\mathcal{L}}{d\beta_0} = 1$$

$$\frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial h} \frac{\partial h}{\partial \beta_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial h} X$$

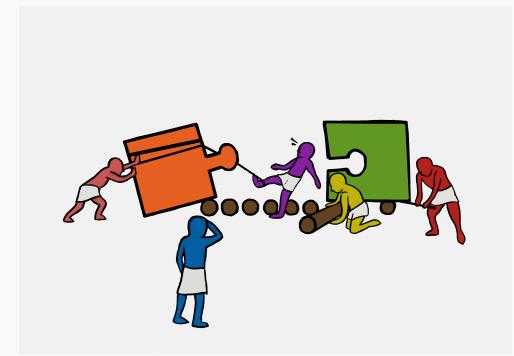
$$\frac{\partial \mathcal{L}}{\partial \beta_1} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial h} \frac{\partial h}{\partial \beta_1} = -X\sigma(h)(1 - \sigma(h)) \left[y \frac{1}{p} + (1 - y) \frac{1}{1 - p} \right]$$

$$\frac{\partial \mathcal{L}}{\partial \beta_0} = \frac{\partial \mathcal{L}}{\partial p} \frac{\partial p}{\partial h} \frac{\partial h}{\partial \beta_0} = -\sigma(h)(1 - \sigma(h)) \left[y \frac{1}{p} + (1 - y) \frac{1}{1 - p} \right]$$



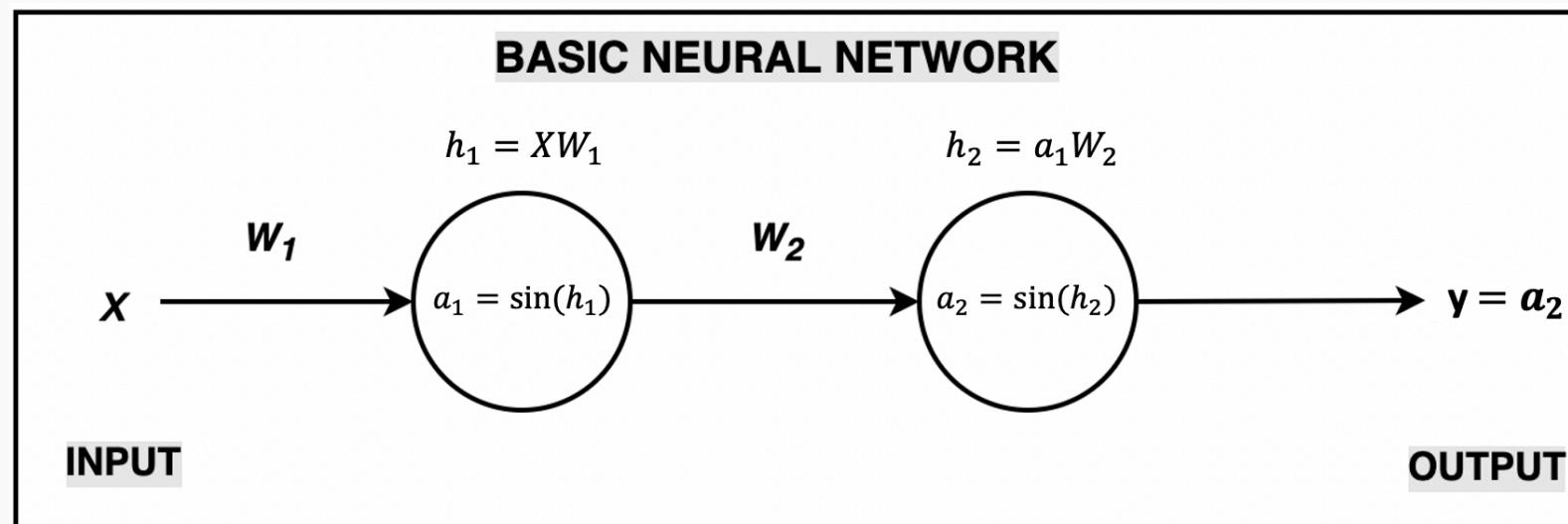
Exercise: Back-propagation by hand

The aim of this exercise is to perform back-propagation to update the weights of a simple neural network



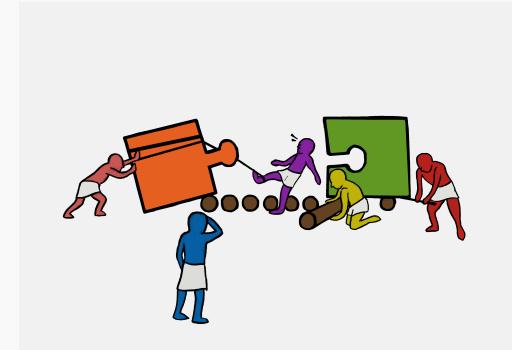
- Build a forward pass of the simple neural network with one hidden layer (see schematic below)
- Randomly initialize the weights
- Use the derivatives to update the weights
- You will need a paper and pen to derive $\frac{\partial L}{\partial W_1}$ & $\frac{\partial L}{\partial W_2}$

$$L = \frac{1}{n} \sum_{1}^n (y_{pred} - y_{true})^2$$



Exercise: Back-propagation by chain rule

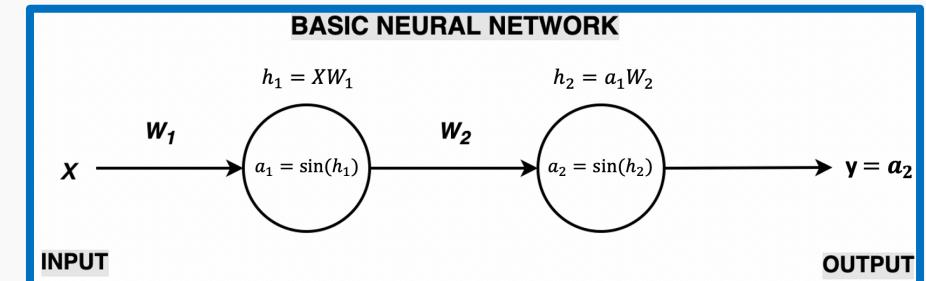
The aim of this exercise is to understand how chain rule works using the same simple neural network from the previous exercise



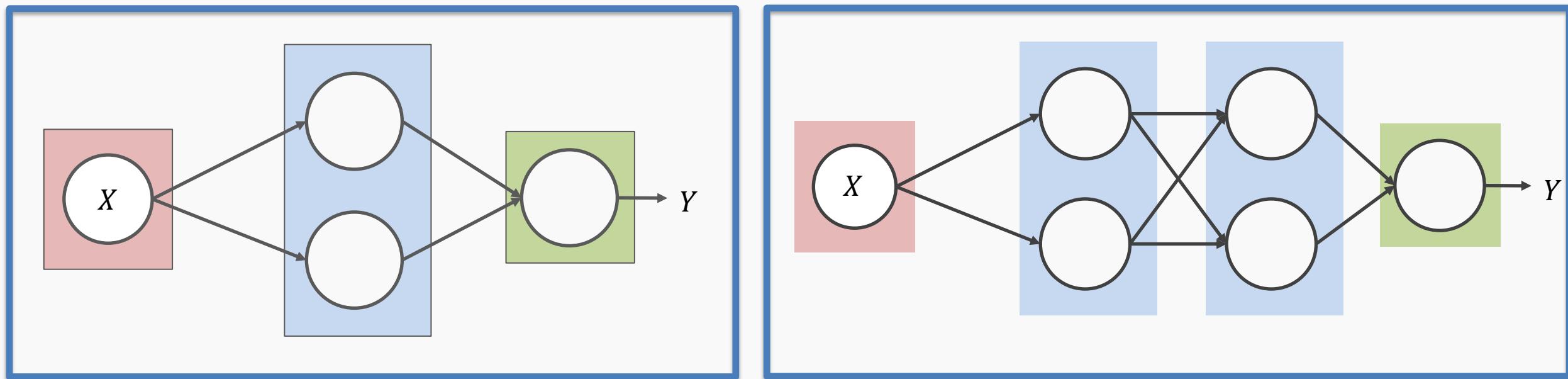
- Write individual functions of various components of the neural network
- Combine the individual functions to compute the partial derivatives $\frac{\partial L}{\partial w_1}$ & $\frac{\partial L}{\partial w_2}$
- Perform one update to the weights and compare them with the previous exercise

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial a_1} \frac{\partial a_1}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$



1. Derivatives need to be evaluated at some values of X , y , and W s.
2. But since we have an expression for the derivative, we can build a function that takes as input X , y , W , and returns the derivatives, and then we can use gradient descent to update.
3. This approach works well, but it does not **generalize**. For example, if the network is changed, we need to write a new function to evaluate the derivatives.



These two networks have different derivatives. We need a mechanism, so we do not need to re-code the derivatives.



Backpropagation (cont.)

Need to find a formalism to calculate the derivatives of the loss w.r.t. weights that is:

1. **Flexible** enough that adding a node or a layer or changing something in the network will not require re-deriving the functional form from scratch.
2. It is **exact**.
3. It is **computationally efficient**.

Hints:

1. Remember we only need to evaluate the derivatives at X_i, y_i and $W^{(k)}$.
2. We should take advantage of the chain rule we learned before.



For example, for input $X=\{3\}$, $y=1$ and weight $W=3$, we evaluate the values of the variables, partial derivatives and the chain up to this point as shown below

Variables	derivatives	Value of the variable	Value of the partial derivative	$\frac{\partial \xi_n}{\partial W}$
$\xi_1 = -W^T X$	$\frac{\partial \xi_1}{\partial W} = -X$	-9	-3	-3
$\xi_2 = e^{\xi_1} = e^{-W^T X}$	$\frac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$	e^{-9}	e^{-9}	$-3e^{-9}$
$\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$	$\frac{\partial \xi_3}{\partial \xi_2} = 1$	$1+e^{-9}$	1	$-3e^{-9}$
$\xi_4 = \frac{1}{\xi_3} = \frac{1}{1 + e^{-W^T X}} = p$	$\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\xi_3^2}$	$\frac{1}{1 + e^{-9}}$	$\left(\frac{1}{1 + e^{-9}}\right)^2$	$-3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)^2$
$\xi_5 = \log \xi_4 = \log p = \log \frac{1}{1 + e^{-W^T X}}$	$\frac{\partial \xi_5}{\partial \xi_4} = \frac{1}{\xi_4}$	$\log \frac{1}{1 + e^{-9}}$	$1 + e^{-9}$	$-3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)$
$\mathcal{L}_i^A = -y\xi_5$	$\frac{\partial \mathcal{L}}{\partial \xi_5} = -y$	$-\log \frac{1}{1 + e^{-9}}$	-1	$3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)$
$\frac{\partial \mathcal{L}_i^A}{\partial W} = \frac{\partial \mathcal{L}_i}{\partial \xi_5} \frac{\partial \xi_5}{\partial \xi_4} \frac{\partial \xi_4}{\partial \xi_3} \frac{\partial \xi_3}{\partial \xi_2} \frac{\partial \xi_2}{\partial \xi_1} \frac{\partial \xi_1}{\partial W}$			-3	0.00037018372

都有自己 BUT we still need to specify the derivatives

Variables	derivatives	Value of the variable	Value of the partial derivative	$\frac{\partial \xi_n}{\partial W}$
$\xi_1 = -W^T X$	$\frac{\partial \xi_1}{\partial W} = -X$	-9	-3	-3
$\xi_2 = e^{\xi_1} = e^{-W^T X}$	$\frac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$	e^{-9}	e^{-9}	$-3e^{-9}$
$\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$	$\frac{\partial \xi_3}{\partial \xi_2} = 1$	$1+e^{-9}$	1	$-3e^{-9}$
$\xi_4 = \frac{1}{\xi_3} = \frac{1}{1 + e^{-W^T X}} = p$	$\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\xi_3^2}$	$\frac{1}{1 + e^{-9}}$	$\left(\frac{1}{1 + e^{-9}}\right)^2$	$-3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)^2$
$\xi_5 = \log \xi_4 = \log p = \log \frac{1}{1 + e^{-W^T X}}$	$\frac{\partial \xi_5}{\partial \xi_4} = \frac{1}{\xi_4}$	$\log \frac{1}{1 + e^{-9}}$	$1 + e^{-9}$	$-3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)$
$\mathcal{L}_i^A = -y\xi_5$	$\frac{\partial \mathcal{L}}{\partial \xi_5} = -y$	$-\log \frac{1}{1 + e^{-9}}$	-1	$3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)$
$\frac{\partial \mathcal{L}_i^A}{\partial W} = \frac{\partial \mathcal{L}_i}{\partial \xi_5} \frac{\partial \xi_5}{\partial \xi_4} \frac{\partial \xi_4}{\partial \xi_3} \frac{\partial \xi_3}{\partial \xi_2} \frac{\partial \xi_2}{\partial \xi_1} \frac{\partial \xi_1}{\partial W}$			-3	0.00037018372



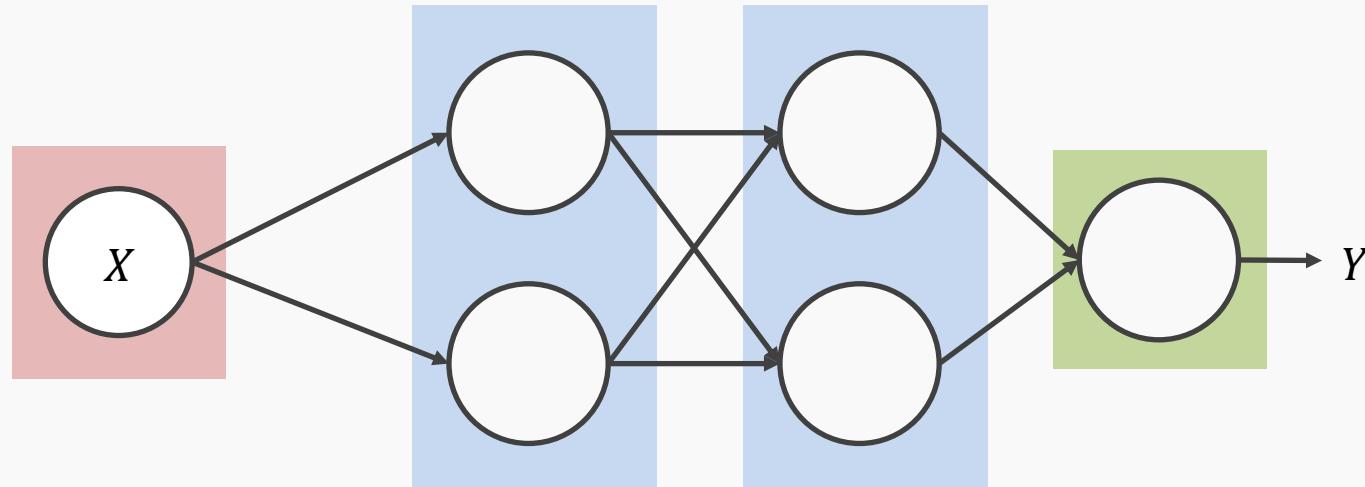
Notice though those are basic functions which are easy to code.

$\xi_0 = X$	$\frac{\partial \xi_0}{\partial X} = 1$	<pre>def x0(x): return x</pre>	<pre>def derx0(): return 1</pre>
$\xi_1 = -W^T \xi_0$	$\frac{\partial \xi_1}{\partial W} = -X$	<pre>def x1(a, x): return -a*x</pre>	<pre>def derx1(a, x): return -a</pre>
$\xi_2 = e^{\xi_1}$	$\frac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$	<pre>def x2(x): return np.exp(x)</pre>	<pre>def derx2(x): return np.exp(x)</pre>
$\xi_3 = 1 + \xi_2$	$\frac{\partial \xi_3}{\partial \xi_2} = 1$	<pre>def x3(x): return 1+x</pre>	<pre>def derx3(x): return 1</pre>
$\xi_4 = \frac{1}{\xi_3}$	$\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\xi_3^2}$	<pre>def x4(x): return 1/(x)</pre>	<pre>def derx4(x): return -(1/x)**(2)</pre>
$\xi_5 = \log \xi_4$	$\frac{\partial \xi_5}{\partial \xi_4} = \frac{1}{\xi_4}$	<pre>def x5(x): return np.log(x)</pre>	<pre>def derx5(x): return 1/x</pre>
$\mathcal{L}_i^A = -y\xi_5$	$\frac{\partial \mathcal{L}}{\partial \xi_5} = -y$	<pre>def L(y, x): return -y*x</pre>	<pre>def derL(y): return -y</pre>



Putting it altogether

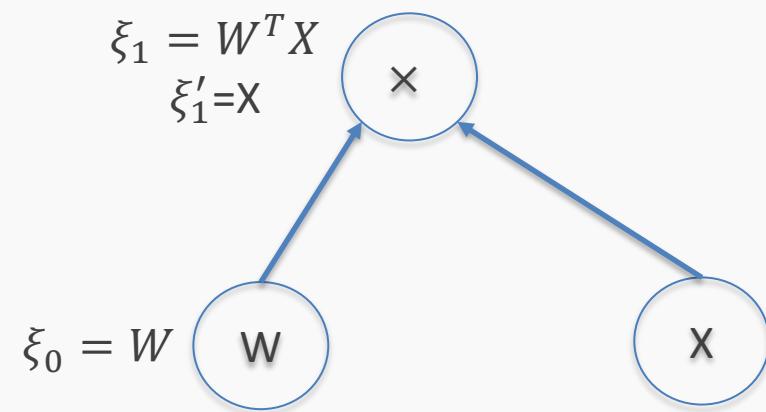
1. We specify the network structure



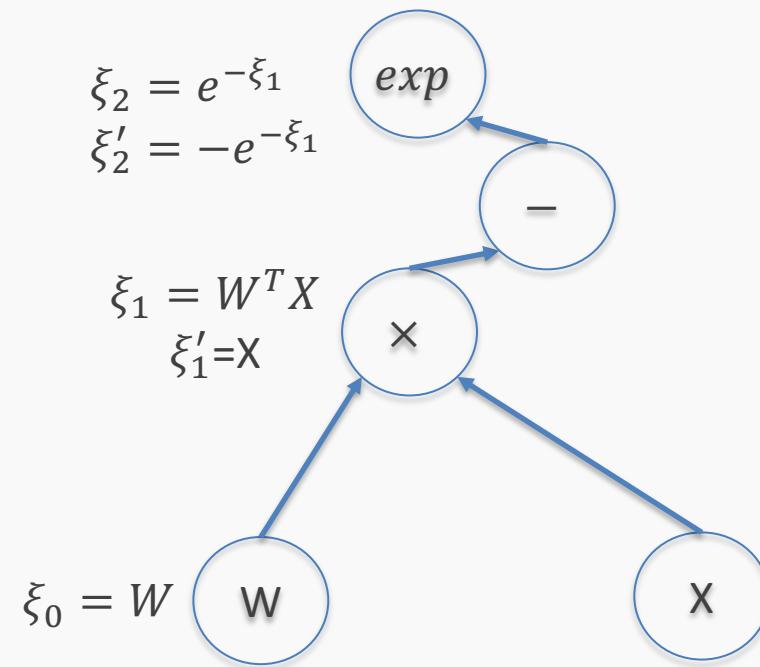
2. Create a computational graph ...

Computational Graph

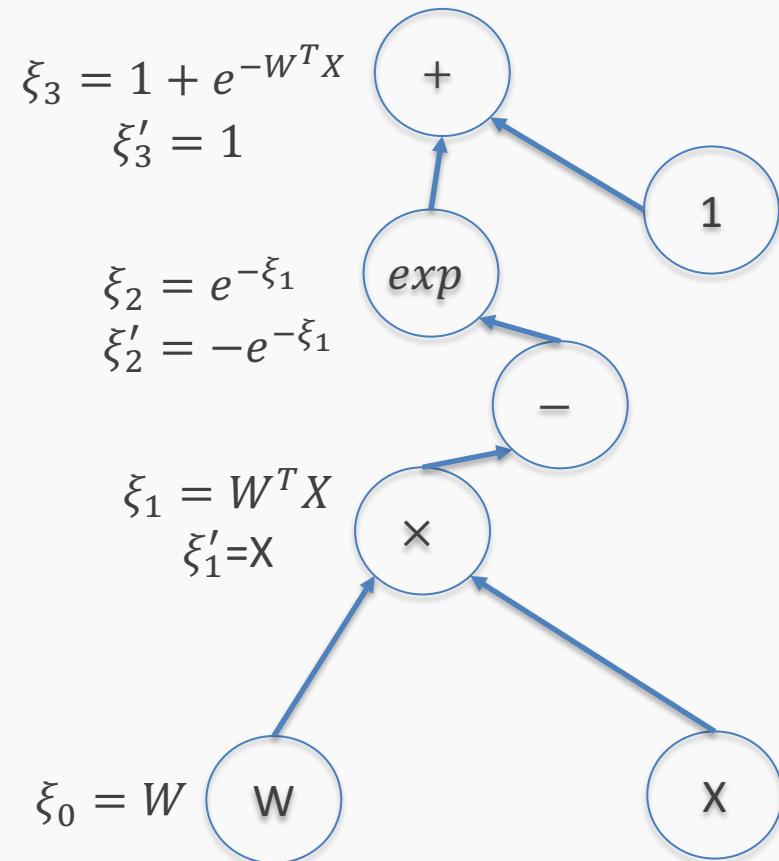
Computational Graph



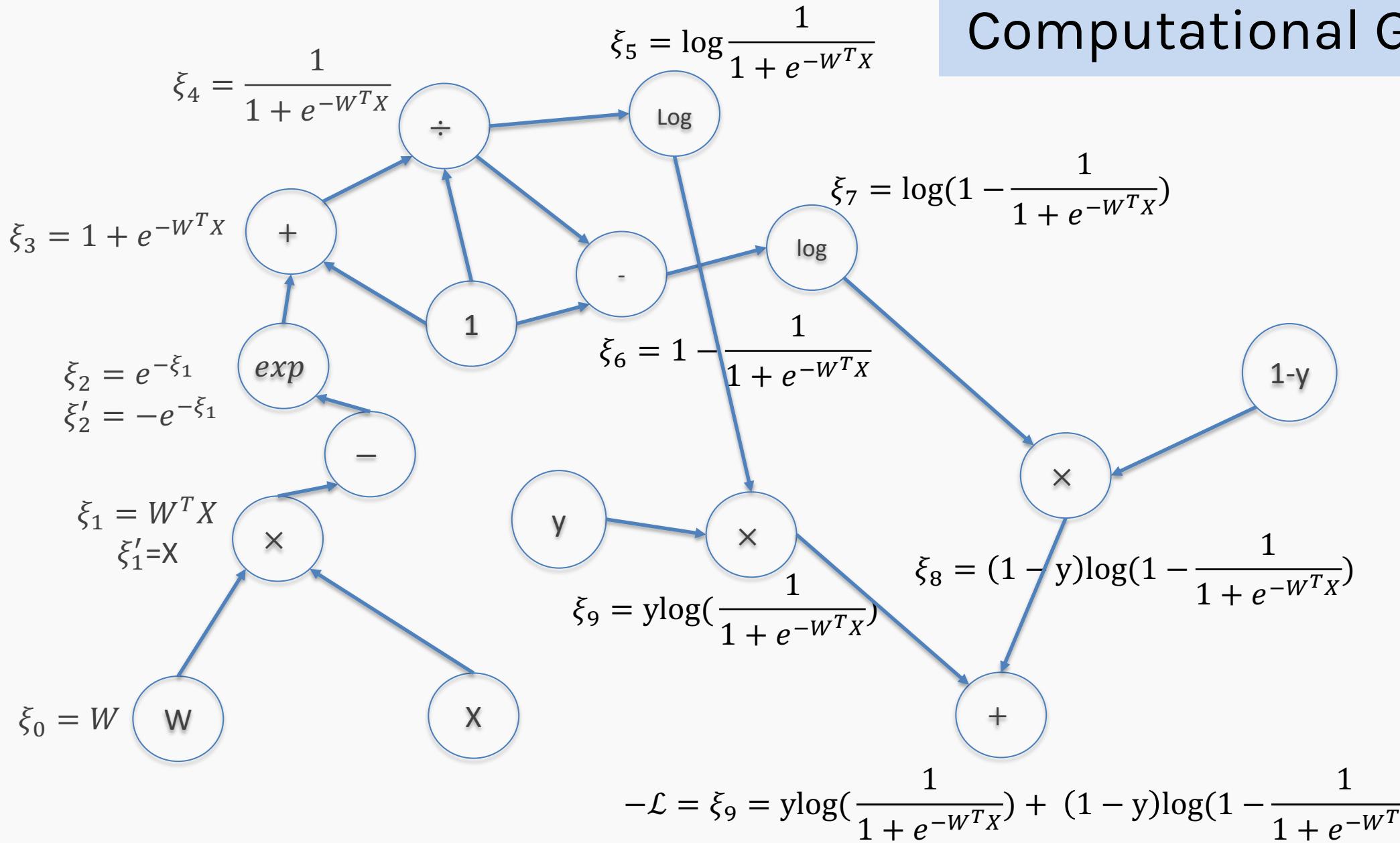
Computational Graph



Computational Graph

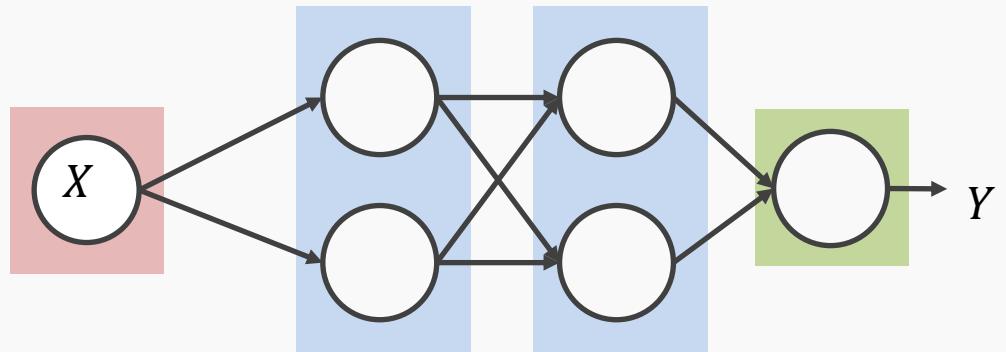


Computational Graph



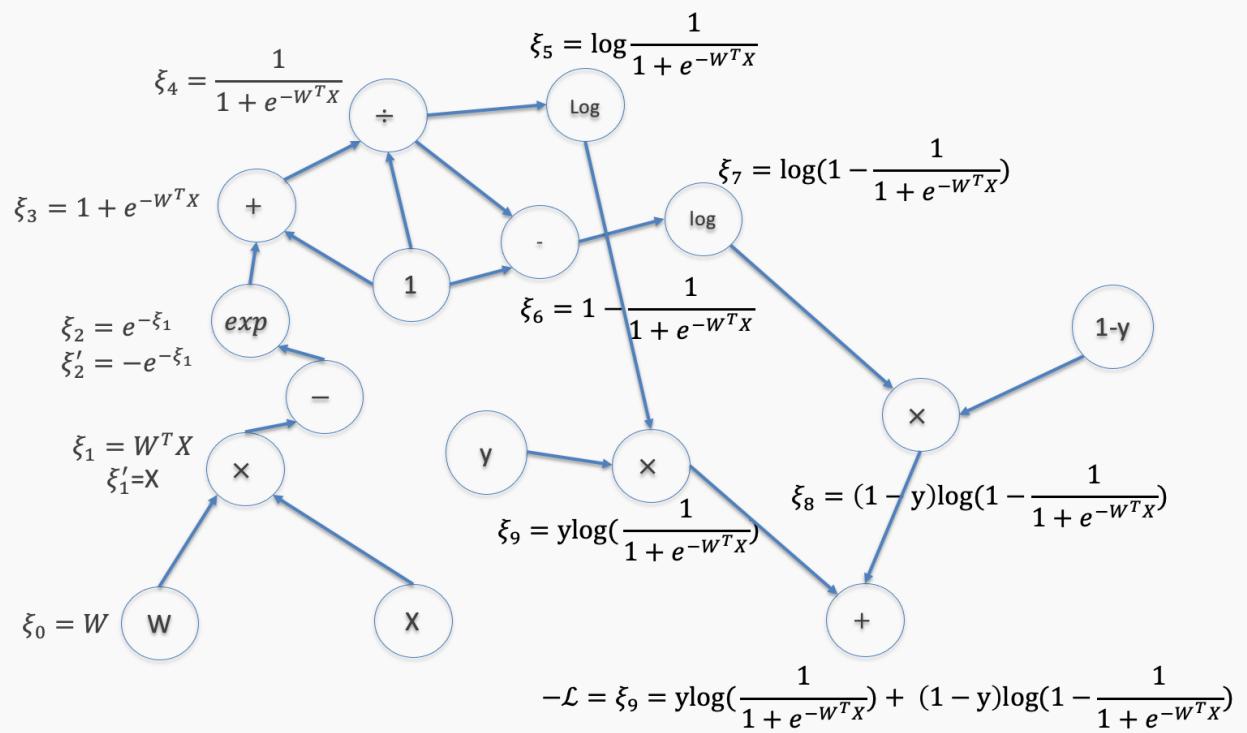
Putting it altogether

1. We specify the network structure



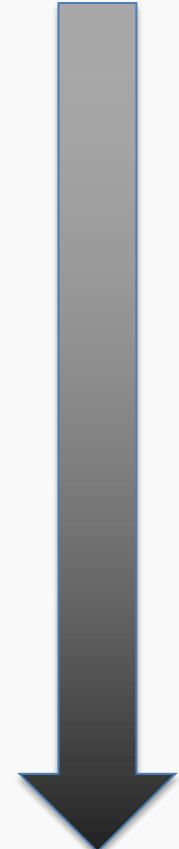
2. Build the computational graph.

At each node of the graph, we build two functions: the evaluation of the variable and its partial derivative with respect to the previous variable (as shown in the table a few slides back)



Forward mode: Evaluate the derivative at: $X=\{3\}, y=1, W=3$

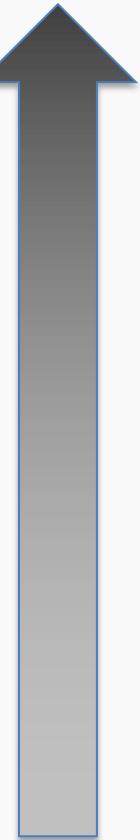
Variables	derivatives	Value of the variable	Value of the partial derivative	$\frac{d\mathcal{L}}{d\xi_n}$
$\xi_1 = -W^T X$	$\frac{\partial \xi_1}{\partial W} = -X$	-9	-3	-3
$\xi_2 = e^{\xi_1} = e^{-W^T X}$	$\frac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$	e^{-9}	e^{-9}	$-3e^{-9}$
$\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$	$\frac{\partial \xi_3}{\partial \xi_2} = 1$	$1+e^{-9}$	1	$-3e^{-9}$
$\xi_4 = \frac{1}{\xi_3} = \frac{1}{1 + e^{-W^T X}} = p$	$\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\xi_3^2}$	$\frac{1}{1 + e^{-9}}$	$\left(\frac{1}{1 + e^{-9}}\right)^2$	$-3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)^2$
$\xi_5 = \log \xi_4 = \log p = \log \frac{1}{1 + e^{-W^T X}}$	$\frac{\partial \xi_5}{\partial \xi_4} = \frac{1}{\xi_4}$	$\log \frac{1}{1 + e^{-9}}$	$1 + e^{-9}$	$-3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)$
$\mathcal{L}_i^A = -y \xi_5$	$\frac{\partial \mathcal{L}_i^A}{\partial \xi_5} = -y$	$-\log \frac{1}{1 + e^{-9}}$	-1	$3e^{-9} \left(\frac{1}{1 + e^{-9}}\right)$
$\frac{\partial \mathcal{L}_i^A}{\partial W} = \frac{\partial \mathcal{L}_i}{\partial \xi_5} \frac{\partial \xi_5}{\partial \xi_4} \frac{\partial \xi_4}{\partial \xi_3} \frac{\partial \xi_3}{\partial \xi_2} \frac{\partial \xi_2}{\partial \xi_1} \frac{\partial \xi_1}{\partial W}$			-3	0.00037018372



Backward mode: Evaluate the derivative at: $X=\{3\}$, $y=1$, $W=3$

Variables	derivatives	Value of the variable	Value of the partial derivative
$\xi_1 = -W^T X$	$\frac{\partial \xi_1}{\partial W} = -X$	-9	-3
$\xi_2 = e^{\xi_1} = e^{-W^T X}$	$\frac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$	e^{-9}	e^{-9}
$\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$	$\frac{\partial \xi_3}{\partial \xi_2} = 1$	$1+e^{-9}$	1
$\xi_4 = \frac{1}{\xi_3} = \frac{1}{1 + e^{-W^T X}} = p$	$\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\xi_3^2}$	$\frac{1}{1 + e^{-9}}$	$\left(\frac{1}{1 + e^{-9}}\right)^2$
$\xi_5 = \log \xi_4 = \log p = \log \frac{1}{1 + e^{-W^T X}}$	$\frac{\partial \xi_5}{\partial \xi_4} = \frac{1}{\xi_4}$	$\log \frac{1}{1 + e^{-9}}$	$1 + e^{-9}$
$\mathcal{L}_i^A = -y \xi_5$	$\frac{\partial \mathcal{L}_i^A}{\partial \xi_5} = -y$	$-\log \frac{1}{1 + e^{-9}}$	-1
$\frac{\partial \mathcal{L}_i^A}{\partial W} = \frac{\partial \mathcal{L}_i}{\partial \xi_5} \frac{\partial \xi_5}{\partial \xi_4} \frac{\partial \xi_4}{\partial \xi_3} \frac{\partial \xi_3}{\partial \xi_2} \frac{\partial \xi_2}{\partial \xi_1} \frac{\partial \xi_1}{\partial W}$			Type equation here.

Store all these values



Multiply as needed



