

# Mathematical Foundations of Machine Learning

## Assignment 1

David Felipe Alvear Goyes

### 1. Gradient of binary logistic loss.

Derive that the gradient of the cross entropy loss function used in binary logistic regression takes the form:

$$\nabla J = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n) x_n$$

Explain briefly what  $x$ ,  $y$ , and  $\hat{y}$  are and how they are obtained or computed.

binary logistic regression  $\rightarrow \hat{y} = P(y|x; \theta) = \text{Ber}(y | \sigma(w^T x + b))$

$$\sigma(x) = \frac{e^x}{1 + e^x}$$

$$\text{Ber}(y|\mu) = \begin{cases} \mu & \text{for } y=1 \\ 1-\mu & \text{for } y=0 \end{cases}$$

$$\frac{d\sigma}{dx} = \sigma(x)(1-\sigma(x))$$

Cross entropy loss  $\rightarrow H(p, q) = - \sum_x p(x) \log(q(x)) \rightarrow q(x)$  estimated distribution

$\rightarrow$  Now for binary-logistic Regression  $H(y, \hat{y}) = - \left( y(x) \log \hat{y}(x) + (1-y(x)) \log (1-\hat{y}(x)) \right)$

$\swarrow$  True distribution       $\searrow$  Estimated distribution  
 term 1                      term 2

Derivate gradient of the cross entropy loss

We have:  $a_n = \theta^T x_n$      $\hat{y}_n = \sigma(\theta^T x_n) \rightarrow \frac{d\hat{y}_n}{d\theta_d} = \hat{y}_n(1-\hat{y}_n) x_{nd}$

$\rightarrow$  For term 1:  $\nabla_{\theta} \log \hat{y}_n = \frac{1}{\hat{y}_n} \frac{d\hat{y}_n}{d\theta_d} = \frac{\hat{y}_n}{\hat{y}_n} (1-\hat{y}_n) x_n = (1-\hat{y}_n) x_n$

$\rightarrow$  For term 2:  $\nabla_{\theta} \log (1-\hat{y}_n) = \frac{1}{1-\hat{y}_n} \frac{d(1-\hat{y}_n)}{d\theta_d} = - \frac{\hat{y}_n (1-\hat{y}_n) x_n}{(1-\hat{y}_n)} = -\hat{y}_n x_n$

$$\nabla_{\theta} H(y, \hat{y}) = - \left( y \nabla_{\theta} \log \hat{y} + (1-y) \nabla_{\theta} \log (1-\hat{y}) \right)$$

$$\rightarrow - \left( y_n (1-\hat{y}_n) x_n - (1-y_n) \hat{y}_n x_n \right) = \left( \hat{y}_n - y_n \right) x_n$$

$$\nabla_{\theta} H(y, \hat{y}) = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n) X_n$$

$\hat{y}_n \rightarrow$  Predicted label

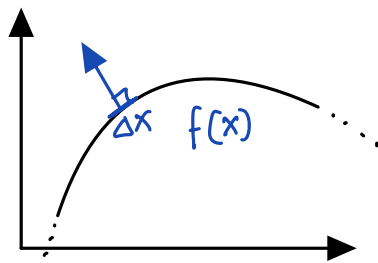
$y \rightarrow$  True label

$X_n \rightarrow$  Feature value

## 2. Negative gradient direction.

Explain, mathematically, why the negative gradient direction is the local direction of steepest descent of the loss function.

$\rightarrow$  Let  $f(x)$ , function that we want to minimize.



$\rightarrow$  we want to find the direction  $u$  that minimize  $D_u f(x)$  the directional derivative.

$$D_u f(x) = \nabla f(x) \cdot u$$

$\rightarrow u$ : unit direction vector

$\rightarrow$  using Cauchy-Schwarz inequality:

$$|\nabla f(x) \cdot u| \leq \|\nabla f(x)\| \cdot \|u\|$$

$$|D_u f(x)| \leq \|\nabla f(x)\|$$

$$\nabla f(x) \cdot u = \|\nabla f(x)\| \|u\| \cos \theta \leq \|\nabla f(x)\|$$

$\rightarrow$  to maximize:  $\cos \theta = 1 \rightarrow \theta = 0$

$$\text{then } u = \frac{\nabla f(x)}{\|\nabla f(x)\|}$$

$\rightarrow$  to minimize:  $\cos \theta = -1 \rightarrow$  Hence  $-\nabla f(x)$

$$\text{then: } u = \frac{-\nabla f(x)}{\|\nabla f(x)\|}$$

is the direction of steepest descent.

$$|D_u f(x)| \geq -\|\nabla f(x)\|$$

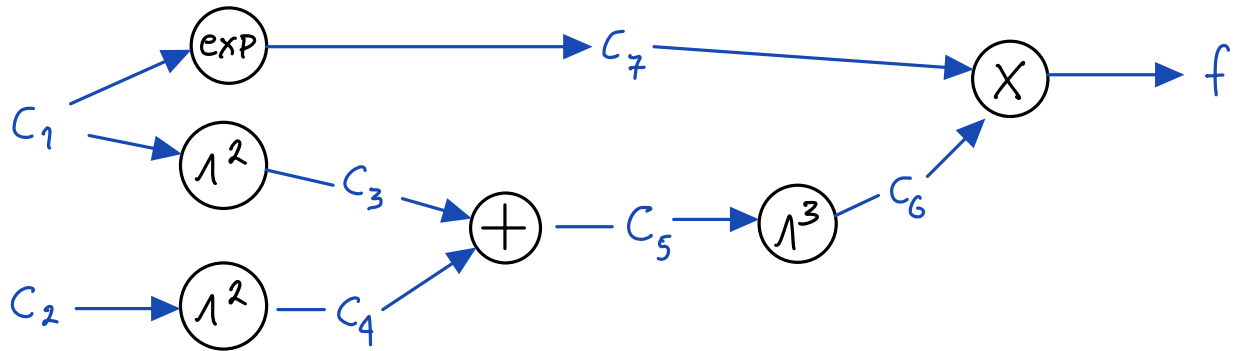
### 3. Automatic differentiation.

We consider the 2-variable scalar function

$$f(x) = \exp(x_1) (x_1^2 + x_2^2)^3$$

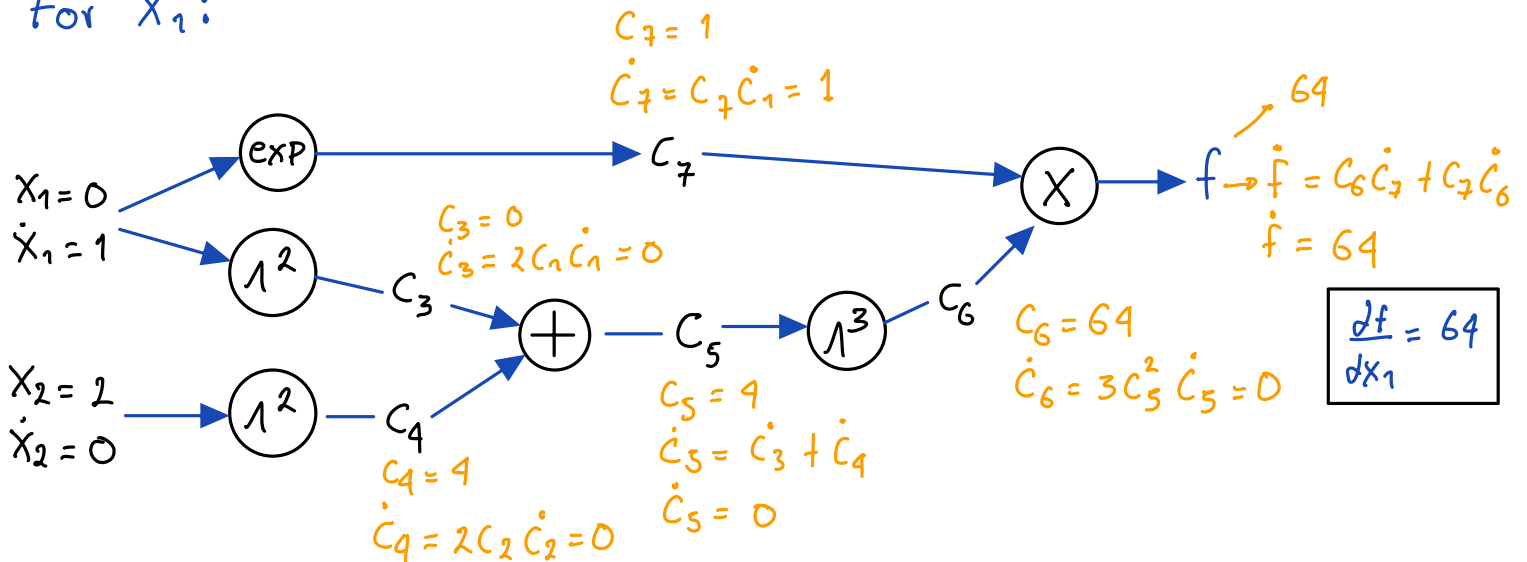
- generate the computational graph of this expression (graph starts with  $c_1 = x_1$ ,  $c_2 = x_2$ , and ends with  $f = c_7$ )
- use forward mode differentiation to compute  $\nabla f(0, 2)$ . Show your intermediate values on the graph.
- use reverse mode differentiation to compute the same gradient. Show your intermediate values (values of the adjoint variables  $\bar{c}_i = \partial f / \partial c_i$ ).
- use the `backward()` method of pytorch to verify your answers.
- comment briefly on the computational effectiveness of forward vs reverse modes of computing these derivatives.

#### Computational Graph

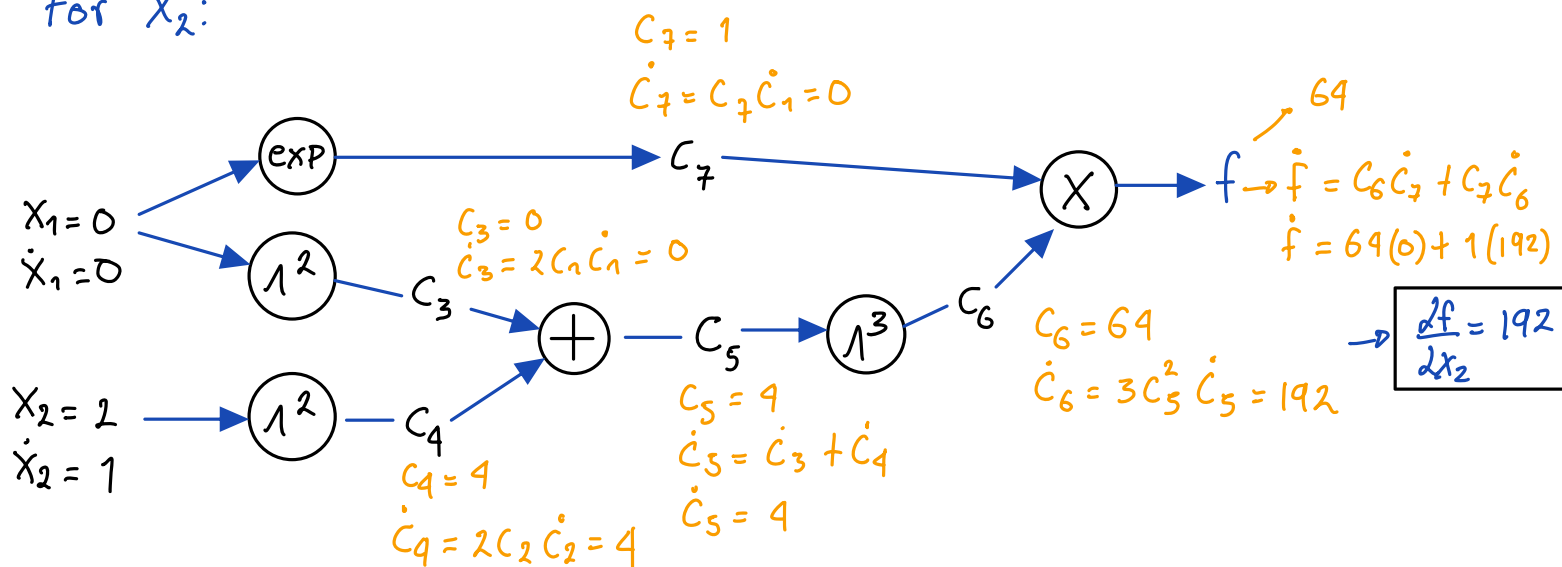


→ Compute  $\nabla f(0, 2) \rightarrow x_1 = 0$  using Forward Mode Differentiation  
 $x_2 = 2$

For  $x_1$ :

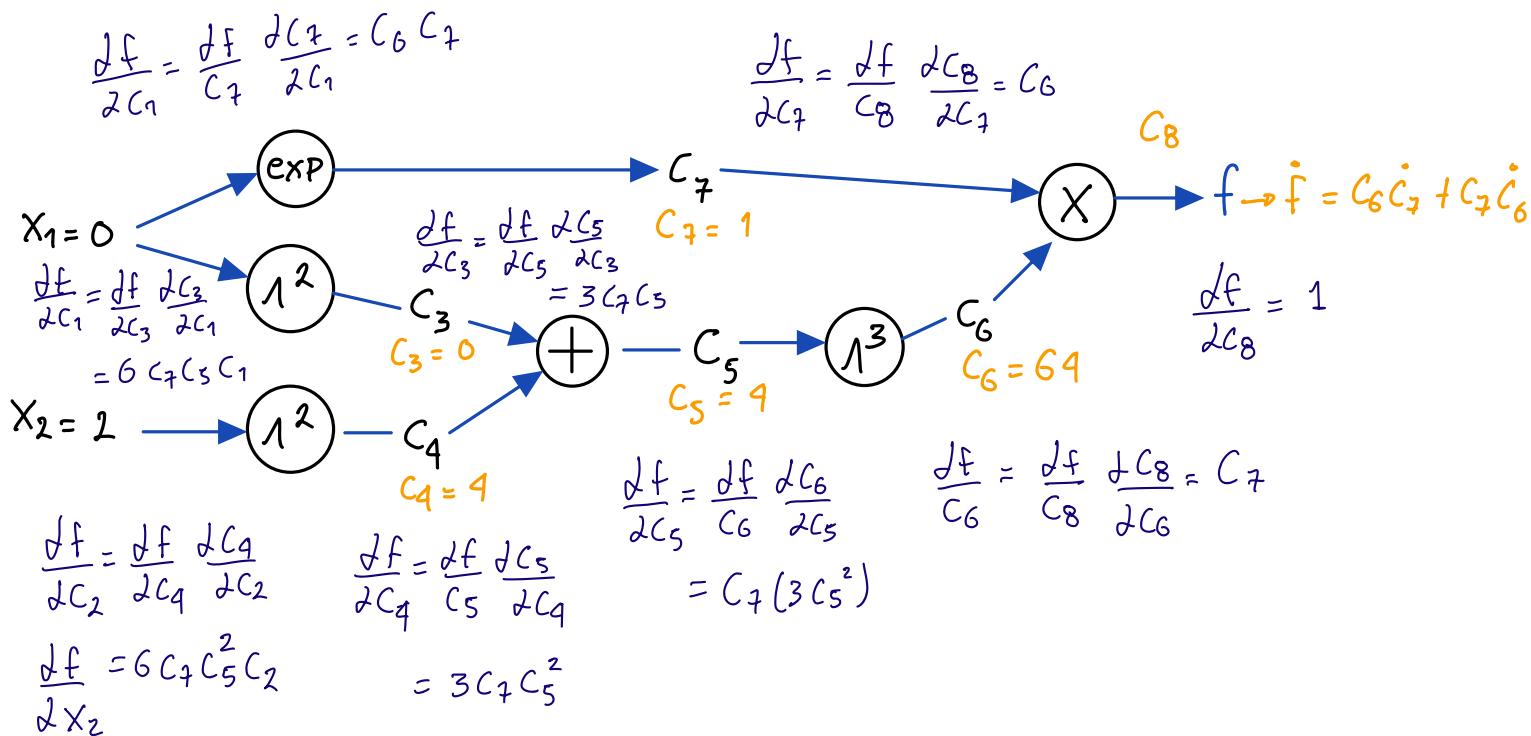


For  $x_2$ :



Then: using Forward Mode differentiation  $\nabla f(0,2) = [64, 192]$

→ Compute  $\nabla f(0,2) \rightarrow x_1 = 0$  using Reverse Mode Differentiation  
 $x_2 = 2$



→ Then:  $\frac{df}{dx_1} = \frac{df}{dC_1} + \frac{df}{dC_7} = \frac{df}{dC_3} \frac{dC_3}{dC_1} + \frac{df}{dC_7} \frac{dC_7}{dC_1} = 64$

$\frac{df}{dx_2} = \frac{df}{dC_4} \frac{dC_4}{dC_2} = 6 C_7 C_5^2 C_2 = 192$

```

import torch
# Define variables
x1 = torch.tensor(0.0, requires_grad=True)
x2 = torch.tensor(2.0, requires_grad=True)

# Compute function
f = torch.exp(x1) * (x1 ** 2 + x2 ** 2) ** 3
f.backward()

# Get gradients
print(x1.grad)
print(x2.grad)

tensor(64.)
tensor(192.)

```

- The forward mode is computational expensive when high dimensional input cases. It's needed to run as many forward passes as there are input variables.
- The reverse mode compute all the gradients in the backward pass, more efficient than the forward mode. However it requires to store the values of the forward pass, needing more memory.