

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2022

Gradient Descent for Logistic Regression (WIP)

Summary

Now, we have developed all the tool we need to do binary classification with LLC:

- A **linear** model that lets us **combine** our variables,

$$u(x) = \theta^T x + \theta_0 \quad (1)$$

- A **logistic** model that lets us get the **probability** of a classification,

$$\sigma(u) = \frac{1}{1 + e^{-u}} \quad (2)$$

- A **threshold value** we use to determine how to **classify** our data,

$$h(x; \theta) = \begin{cases} +1 & \text{if } \sigma(u(x)) > \sigma_{\text{thresh}} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

- A **loss function** NLL we use to **evaluate** our model performance:

$$\mathcal{L}_{\text{nll}}(\mathbf{g}^{(i)}, \mathbf{y}^{(i)}) = - \left(\mathbf{y}^{(i)} \log \mathbf{g}^{(i)} + (1 - \mathbf{y}^{(i)}) \log (1 - \mathbf{g}^{(i)}) \right)$$

- And an **objective function** we can **optimize**:

$$J_{\text{lr}}(\theta, \theta_0; \mathcal{D}) = \lambda \|\theta\|^2 + \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{nll}}(\mathbf{g}^{(i)}, \mathbf{y}^{(i)}) \quad (4)$$

We have everything we need to do optimization.

The problem: Gradient Descent

We want to do **gradient descent** to minimize J_{lr}

$$\mathbf{R}(\theta) + J_{\text{lr}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{\text{nll}}(\mathbf{g}^{(i)}, \mathbf{y}^{(i)}) \quad (5)$$

We want repeatedly **adjust** our model $\Theta = (\theta, \theta_0)$ to improve J_{lr} . To do that, we want the gradients for θ and θ_0 . Let's start with θ .

$$\nabla_{\theta} J_{\text{lr}} = \frac{\partial J_{\text{lr}}}{\partial \theta} \quad (6)$$

First, J_{lr} has **two** terms, so we'll separate them.

$$\nabla_{\theta} J_{lr} = \frac{\partial R}{\partial \theta} + \frac{1}{n} \sum_{i=1}^n \frac{\partial \mathcal{L}_{NLL}}{\partial \theta}(\mathbf{g}^{(i)}, \mathbf{y}^{(i)}) \quad (7)$$

The regularization term is pretty easy, because we did it last chapter:

$$\frac{\partial R}{\partial \theta} = 2\lambda\theta \quad (8)$$

But what about our first term?

Getting the gradient: Chain Rule

Now, we just need to do

$$\frac{\partial \mathcal{L}_{NLL}}{\partial \theta}(\mathbf{g}, \mathbf{y}) \quad (9)$$

With our \mathcal{L}_{NLL} term, we run into an issue: how do we take the **derivative**? The function is very, very deeply **nested**. In our case:

x **affects** $u(x)$. $u(x)$ **affects** $\sigma(u)$. $\sigma(u) = g$ **affects** $\mathcal{L}_{NLL}(g, y)$, which finally **affects** $J(\theta, \theta_0)$.

How do we represent this **chain** of functions? With the **chain rule**:

$$\frac{\partial A}{\partial C} = \frac{\partial A}{\partial B} \cdot \frac{\partial B}{\partial C} \quad (10)$$

So, we'll build up a **chain rule** for our needs. We'll use $g = \sigma(u)$.

$$\frac{\partial \mathcal{L}_{NLL}}{\partial \theta} = \frac{\partial \mathcal{L}_{NLL}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial \theta} \quad (11)$$

Sigma contains u , so we'll use that instead:

$$\frac{\partial \mathcal{L}_{NLL}}{\partial \theta} = \frac{\partial \mathcal{L}_{NLL}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial u} \cdot \frac{\partial u}{\partial \theta} \quad (12)$$

This is our full **chain rule**!

Key Equation 1

The **gradient** of **NLL** can be calculated using the **chain rule**:

$$\frac{\partial \mathcal{L}_{NLL}}{\partial \theta} = \frac{\partial \mathcal{L}_{NLL}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial u} \cdot \frac{\partial u}{\partial \theta} \quad (13)$$

Getting our individual derivatives

We can take the derivative of each of these objects. First, let's look at \mathcal{L}_{NLL}

$$\mathcal{L}_{\text{NLL}}(\sigma, y) = -\left(y \log \sigma + (1 - y) \log (1 - \sigma)\right)$$

And we'll use $\frac{d}{dx} \log(x) = \frac{1}{x}$

$$\frac{\partial \mathcal{L}_{\text{NLL}}}{\partial \sigma} = -\left(\frac{y}{\sigma} - \frac{1 - y}{1 - \sigma}\right) \quad (14)$$

Now, we look at $\sigma(u)$:

$$\sigma(u) = \frac{1}{1 + e^{-u}} \quad (15)$$

If we take the derivative, we can get:

$$\frac{\partial \sigma}{\partial u} = \frac{-e^{-u}}{(1 + e^{-u})^2} \quad (16)$$

Which we can rewrite, conveniently, as

Try this yourself if you're curious!

$$\frac{\partial \sigma}{\partial u} = \sigma(1 - \sigma) \quad (17)$$

Finally, our last derivative:

$$u = \theta^T x + \theta_0 \quad (18)$$

$$\frac{\partial u}{\partial \theta} = x \quad (19)$$

Simplifying our chain rule

So, now, we can put together our chain rule:

$$\frac{\partial \mathcal{L}_{\text{NLL}}}{\partial \theta} = \frac{\partial \mathcal{L}_{\text{NLL}}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial u} \cdot \frac{\partial u}{\partial \theta} \quad (20)$$

Plug in the derivatives:

$$\frac{\partial \mathcal{L}_{\text{NLL}}}{\partial \theta} = -\left(\frac{y}{\sigma} - \frac{1 - y}{1 - \sigma}\right) \cdot \sigma(1 - \sigma) \cdot x \quad (21)$$

Simplify:

$$\frac{\partial \mathcal{L}_{\text{NLL}}}{\partial \theta} = \left((1 - y)\sigma - y(1 - \sigma)\right) \cdot x \quad (22)$$

And finally, we sum the terms. We can do the θ_0 gradient at the same time: the only difference is that $\frac{\partial u}{\partial \theta_0} = 1$, instead of x .

Key Equation 2

The **gradients** of NLL for gradient descent are

$$\nabla_{\theta} \mathcal{L}_{\text{NLL}} = (\sigma - y)x$$

$$\frac{\partial \mathcal{L}_{\text{NLL}}}{\partial \theta_0} = (\sigma - y)$$

We can plug this into J_{lr} :

$$\nabla_{\theta} J_{\text{lr}} = \frac{1}{n} \sum_{i=1}^n \left((g^{(i)} - y^{(i)}) x^{(i)} \right) + 2\lambda \theta \quad (23)$$

$$\frac{\partial J_{\text{lr}}}{\partial \theta_0} = \frac{1}{n} \sum_{i=1}^n (g^{(i)} - y^{(i)}) \quad (24)$$

One comment we didn't make: remember that $R(\theta)$ won't show up in the θ_0 derivative!

We can use this to do **gradient descent**!

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_{\theta} J_{\text{lr}}(\theta_{\text{old}}) \quad (25)$$

In $\theta^{(t)}$ notation:

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left(\nabla_{\theta} J_{\text{lr}}(\theta^{(t-1)}) \right) \quad (26)$$

$$\theta_0^{(t)} = \theta_0^{(t-1)} - \eta \left(\frac{\partial J_{\text{lr}}(\theta^{(t-1)})}{\partial \theta_0} \right) \quad (27)$$

This also corresponds to some basic math within Neural Networks, which we will return to **later** in the course.