

Assignment 4: Logistic Regression and Gradient Descent

Machine Learning

Fall 2021

🔗 Learning Objectives

- Learn about the logistic regression algorithm.
- Learn about gradient descent for optimization.
- Contemplate an application of machine learning.

Change Contemplate an application of machine learning. to some more appropriate learning objective based on whatever we assign them to watch.

↔ Prior Knowledge Utilized

- Supervised learning problem framing.
- Log loss

1 Logistic Regression (top-down)

In the last part of the notebook that you started in class, you saw a quick implementation of logistic regression to classify if a person was looking to the left or to the right.

In this assignment we will formalize the binary classification problem and dig the theory behind *logistic regression*. You will also see that the logistic regression algorithm is a very natural extension of linear regression. Our plan for getting there is going to be pretty similar to what we did for linear regression.

- Build some mathematical foundations
- Introduce logistic regression from a top-down perspective
- Learn about logistic regression from a bottom-up perspective

↔ Recall:

In the last assignment, you were introduced to the idea of binary classification, which based on some input \mathbf{x} has a corresponding output y that is $y = 0$ or $y = 1$. In logistic regression, this model, \hat{f} , instead of spitting out either 0 or 1, outputs a confidence that the input \mathbf{x} has an output $y = 1$. In other words, rather than giving us its best guess (0 or 1), the classifier indicates to us its degree of certainty regarding its prediction as a probability.

We also explored three possible loss functions for a model that outputs a probability p when supplied with an input \mathbf{x} (i.e., $\hat{f}(\mathbf{x}) = p$). The loss function is

used to quantify how bad a prediction p is given the actual output y (for binary classification the output is either 0 or 1).

1. **0-1 loss:** This is an all-or-nothing approach. If the prediction is correct, the loss is zero; if the prediction is incorrect, the loss is 1. This does not take into account the level certainty expressed by the probability (the model gets the same loss if $y = 1$ and it predicted $p = 0.51$ or $p = 1$).
2. **squared loss:** For squared loss we compute the difference between the outcome and p and square it to arrive at the loss. For example, if $y = 1$ and the model predicts $p = 0.51$, the loss is $(1 - 0.51)^2$. If instead $y = 0$, the loss is $(0 - 0.51)^2$.
3. **log loss:** ...

add some text about log loss

Now that you have refreshed on how probabilities can be used as a way of quantifying confidence in predictions, you are ready to learn about the logistic regression algorithm.

As always, we assume we are given a training set of inputs and outputs. As in linear regression we will assume that each of our inputs is a d -dimensional vector \mathbf{x}_i and since we are dealing with binary classification, the outputs, y_i , will be binary numbers (indicating whether the input belongs to class 0 or 1). Our hypothesis functions, \hat{f} , output the probability that a given input has an output of 1. What's cool is that we can borrow a lot of what we did in the last couple of assignments when we learned about linear regression. In fact, all we're going to do in order to make sure that the output of \hat{f} is between 0 and 1 is pass $\mathbf{w}^\top \mathbf{x}$ through a function that "squashes" its input so that it outputs a value between 0 and 1. This idea is shown graphically in Figure 1.

To make this intuition concrete, we define each \hat{f} as having the following form (note: this equation looks daunting. We have some tips for interpreting it below).

$$\begin{aligned}\hat{f}(\mathbf{x}) &= \text{probability that output, } y, \text{ is } 1 \\ &= \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}\end{aligned}\tag{1}$$

Here are a few things to notice about this equation:

1. The weight vector that we saw in linear regression, \mathbf{w} , has made a comeback. We are using the dot product between \mathbf{x} and \mathbf{w} (which creates a weighted sum of the x_i 's), just as we did in linear regression!
2. As indicated in Figure 1, the dot product $\mathbf{w}^\top \mathbf{x}$ has been passed through a squashing function known as the [sigmoid function](#). The graph of $\sigma(u) = \frac{1}{1+e^{-u}}$ is shown in Figure 2. $\sigma(\mathbf{w}^\top \mathbf{x})$ is exactly what we have in Equation 2.

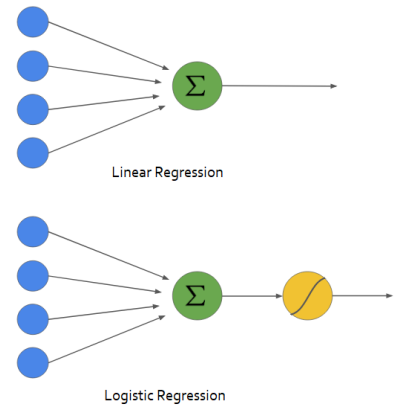


Figure 1: Graphical representation of both linear and logistic regression. The key difference is the application of the squashing function shown in yellow. [Original source](#).

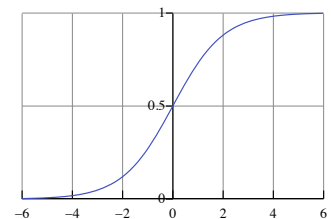


Figure 2: a graph of the sigmoid function $\frac{1}{1+e^{-x}}$.

2 Deriving the Logistic Regression Learning Rule

Let's summarize what we've done thus far in this assignment.

- We motivated the binary classification problem.
- We presented a particular useful loss function (log loss).
- We met the logistic regression model and tried it out on a real dataset.

Next, we're going to build on these pieces and formalize the logistic regression problem and derive a learning rule to solve it (i.e., compute the optimal weights). The formalization of logistic regression will combine Equation 2 with the selection of ℓ to be log loss. This choice of ℓ results in the following objective function.

$$\begin{aligned} \mathbf{w}^* &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \left(-y_i \ln \sigma(\mathbf{w}^\top \mathbf{x}_i) - (1 - y_i) \ln(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \right) \quad (2) \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \left(-y_i \ln \left(\frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}_i}} \right) - (1 - y_i) \ln \left(1 - \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}_i}} \right) \right) \quad \text{expanded out if you prefer this form} \quad (3) \end{aligned}$$

While this looks a bit crazy, since y_i is either 0 or 1, the multiplication of the expressions in the summation by either y_i or $1 - y_i$ are essentially acting like a switch—depending on the value of y_i we either get one term or the other. Our typical recipe for finding \mathbf{w}^* has been to take the gradient of the expression inside the arg min, set it to 0, and solve for \mathbf{w}^* (which will be a critical point and hopefully a minimum). The last two steps will be a bit different for reasons that will become clear soon, but we will need to find the gradient. We will focus on finding the gradient in the next couple of parts.

2.1 Useful Properties of the Sigmoid Function

Looking at Equation 4 it looks really, really hairy! We see that in order to compute the gradient we will have to compute the gradient of $\mathbf{x}^\top \mathbf{w}$ with respect to \mathbf{w} (we just wrapped our minds around this last assignment). Additionally, we will have to take into account how the application of the sigmoid function and the log function changes this gradient. In this section we'll learn some properties for manipulating the sigmoid function and computing its derivative.

Exercise 1 (60 minutes)

The sigmoid function, σ , is defined as

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (4)$$

- (a) Show that $\sigma(-x) = 1 - \sigma(x)$.

☆ Solution

$$\sigma(-x) = \frac{1}{1+e^x} \quad (5)$$

$$= \frac{e^{-x}}{e^{-x}+1} \quad \text{multiply by top and bottom by } e^{-x} \quad (6)$$

$$\sigma(-x) - 1 = \frac{e^{-x}}{e^{-x}+1} - \frac{1+e^{-x}}{1+e^{-x}} \quad \text{subtract } -1 \text{ on both sides} \quad (7)$$

$$= \frac{-1}{1+e^{-x}} \quad (8)$$

$$= -\sigma(x) \quad (9)$$

$$\sigma(-x) = 1 - \sigma(x) \quad (10)$$

(b) Show that the derivative of the logistic function $\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$

☆ Solution

Two solutions for the price of 1!

Solution 1:

$$\frac{d}{dx}\sigma(x) = e^{-x}\sigma(x)^2 \quad \text{apply quotient rule} \quad (11)$$

$$= \sigma(x) \left(\frac{e^{-x}}{1+e^{-x}} \right) \quad \text{expand out one of the } \sigma(x)\text{'s} \quad (12)$$

$$= \sigma(x) \left(\frac{1}{e^x+1} \right) \quad \text{multiply top and bottom by } e^x \quad (13)$$

$$= \sigma(x)(\sigma(-x)) \quad \text{substitute for } \sigma(-x) \quad (14)$$

$$= \sigma(x)(1 - \sigma(x)) \quad \text{apply } \sigma(-x) = 1 - \sigma(x) \quad (15)$$

Solution 2:

$$\frac{d}{dx}\sigma(x) = \frac{e^{-x}}{(1+e^{-x})^2} \quad \text{apply quotient rule} \quad (16)$$

$$= \frac{e^{-x}}{1+2e^{-x}+e^{-2x}} \quad \text{expand the bottom} \quad (17)$$

$$= \frac{1}{e^x+2+e^{-x}} \quad \text{multiply top and bottom by } e^x \quad (18)$$

$$= \frac{1}{(1+e^x)(1+e^{-x})} \quad \text{factor} \quad (19)$$

$$= \sigma(x)\sigma(-x) \quad \text{decompose using definition of } \sigma(x) \quad (20)$$

$$= \sigma(x)(1 - \sigma(x)) \quad \text{apply } \sigma(-x) = 1 - \sigma(x) \quad (21)$$

2.2 Chain Rule for Gradients

We now know how to take derivatives of each of the major pieces of Equation 4. What we need is a way to put these derivatives together. You probably remember that in the case of single variable calculus you have just such a tool. This tool is known as the chain rule. The chain rule tells us how to compute the derivative of the composition of two single variable functions f and g .

$$\begin{aligned} h(x) &= g(f(x)) & h(x) \text{ is the composition of } f \text{ with } g \\ h'(x) &= g'(f(x))f'(x) & \text{this is the chain rule!} \end{aligned} \quad (22)$$

Suppose that instead of the input being a scalar x , the input is now a vector, \mathbf{w} . In this case h takes a vector input and returns a scalar, f takes a vector input and returns a scalar, and g takes a scalar input and returns a scalar.

$$\begin{aligned} h(\mathbf{w}) &= g(f(\mathbf{w})) & h(\mathbf{w}) \text{ is the composition of } f \text{ with } g \\ \nabla h(\mathbf{w}) &= g'(f(\mathbf{w}))\nabla f(\mathbf{w}) & \text{this is the multivariable chain rule} \end{aligned} \quad (23)$$

Exercise 2 (60 minutes)

- (a) Suppose $h(x) = \sin(x^2)$, compute $h'(x)$ (x is a scalar so you can apply the single-variable chain rule).

☆ Solution

Applying the chain rule gives

$$h'(x) = \cos(x^2)2x. \quad (24)$$

- (b) Define $h(\mathbf{v}) = (\mathbf{c}^\top \mathbf{v})^2$. Compute $\nabla_{\mathbf{v}} h(\mathbf{v})$ (the gradient of the function with respect to \mathbf{v}).

☆ Solution

We can see that $h(\mathbf{v}) = g(f(\mathbf{v}))$ with $g(x) = x^2$ and $f(\mathbf{v}) = \mathbf{c}^\top \mathbf{v}$. The gradient can now easily be found by applying the chain rule.

$$\nabla h(\mathbf{v}) = 2(\mathbf{c}^\top \mathbf{v})\mathbf{c} \quad (25)$$

- (c) Compute the gradient of the expression from Equation 4 (reproduced below for your convenience).

$$\sum_{i=1}^n -y_i \ln \sigma(\mathbf{w}^\top \mathbf{x}_i) - (1 - y_i) \ln (1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) . \quad (26)$$

You can either use the chain rule and the identities you learned about sigmoid, or expand everything out and work from that.

☆ Solution

Applying the chain rule gives us

$$\sum_{i=1}^n -y_i \frac{\nabla \sigma(\mathbf{w}^\top \mathbf{x}_i)}{\sigma(\mathbf{w}^\top \mathbf{x}_i)} - (1 - y_i) \frac{-\nabla \sigma(\mathbf{w}^\top \mathbf{x}_i)}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)} . \quad (27)$$

Applying the chain rule again gives us

$$\begin{aligned} & \sum_{i=1}^n -y_i \frac{\sigma(\mathbf{w}^\top \mathbf{x}_i)(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \nabla \mathbf{w}^\top \mathbf{x}_i}{\sigma(\mathbf{w}^\top \mathbf{x}_i)} - (1 - y_i) \frac{-\sigma(\mathbf{w}^\top \mathbf{x}_i)(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \nabla \mathbf{w}^\top \mathbf{x}_i}{1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)} \\ &= \sum_{i=1}^n -y_i (1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \mathbf{x}_i + (1 - y_i) \sigma(\mathbf{w}^\top \mathbf{x}_i) \mathbf{x}_i \end{aligned} \quad (28)$$

You could certainly stop here, but if you plug in $y = 0$ and $y = 1$ you'll find that the expression can be further simplified to:

$$\sum_{i=1}^n -(y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \mathbf{x}_i$$

2.3 Gradient Descent for Optimization

If we were to follow our derivation of linear regression we would set our expression for the gradient to 0 and solve for \mathbf{w} . It turns out this equation will be difficult to solve due to the σ function. Instead, we can use an iterative approach where we start with some initial value for \mathbf{w} (we'll call the initial value \mathbf{w}^0 , where the superscript corresponds to the iteration number) and iteratively adjust it by moving down the gradient (the gradient represents the direction of fastest increase for our function, therefore, moving along the negative gradient is the direction where the loss is decreasing the fastest).

🔗 External Resource(s) (45 minutes)

There are tons of great resources that explain gradient descent with both math and compelling visuals.

- Recommended: [Gradient descent, how neural networks learn | Deep learning, chapter 2](#) (start at 5:20)
- An Introduction to Gradient Descent ([on NB](#), [original](#))
- The Wikipedia page on Gradient Descent ([on NB](#), [original](#))
- [Ahmet Sacan's video on gradient descent](#) (this one has some extra stuff, but it's pretty clearly explained).
- There are quite a few resources out there, do you have some suggestions? (suggest so on NB)

Exercise 3 (10 minutes)

To test your understanding of these resources, here are a few diagnostic questions.

- (a) When minimizing a function with gradient descent, which direction should you step along in order to arrive at the next value for your parameters?

☆ Solution

The negative gradient (since we are minimizing)

- (b) What is the learning rate and what role does it serve in gradient descent?

☆ Solution

The learning rate controls the size of the step that you take along the negative gradient.

- (c) How do you know when an optimization performed using gradient descent has converged?

☆ Solution

There are a few options. One popular one is to check if the objective function is changing only a minimal amount each iteration, the algorithm has converged. You could also look at the magnitude of the gradient (which tells us the slope) to see if it is really small.

- (d) True or false: provided you tune the learning rate properly, gradient descent guarantees that you will find the global minimum of a function.

☆ Solution

False, the best gradient descent can do, in general, is converge to a local minimum. If you know that the function you are optimizing has only one minimum, then this would also be the global minimum (this is the case for both linear and logistic regression).

If we take the logic of gradient descent and apply it to the logistic regression problem, we arrive at the following learning rule. Given some initial weights \mathbf{w}^0 , and a learning rate η , we can iteratively update our weights using the formula below.

$$\mathbf{w}^{n+1} = \mathbf{w}^n - \eta \sum_{i=1}^n -(y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \mathbf{x}_i \quad \text{applying the result from exercise 4} \quad (29)$$

$$= \mathbf{w}^n + \eta \sum_{i=1}^n (y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \mathbf{x}_i \quad \text{distribute the negative} \quad (30)$$

This beautiful equation turns out to be the recipe for logistic regression.

▲ Notice

We won't be assigning a full implementation of logistic regression from scratch. In future assignments, we will spend more time applying logistic regression and gradient descent.

If it's helpful for your learning to see a worked example with code now (to help the math make sense), you can optionally check out this [example of binary classification for admission to college](#), noting that some of the math notation is slightly different than ours.

You are also welcome to implement logistic regression using gradient descent if it's helpful for your learning and/or if you already have significant experience with machine learning and want a challenge. This is completely optional, and we assume that most of you will not choose to do this. If you do decide to implement logistic regression using gradient descent, you will need to search for a good learning rate or you may consider implementing some [strategies for automatically tuning the learning rate](#).

3 Touchpoint to context, impact, and ethics

As we mentioned in the introduction to the course, we'll be exploring machine learning from three different perspectives: the theory, the implementation, and the context, impact, and ethics.

insert image from syllabus slides on theory implementation and context

Lately, we've been diving deeply into the mathematical theory with some touch points in implementation and context. Here, we want to take a moment to zoom out (with a lowercase z), and think about how algorithms interact with society, individuals, and institutions.

one option, Carrie is only 45 min into it and felt like it started strong but lost the plot, not sure I have another 1.5 hours to watch another

We are going to take a virtual field trip to a leading conference on this topic! In normal times it would be challenging and expensive for us to bring the entire class to this conference, but thanks to the pandemic we can easily share one of the talks from the Conference on Fairness, Accountability, and Transparency with you.

Exercise 4

Please watch [this talk](#) and take notes.

Or, option 2– Carrie likes this one better, and it's closer to an hour. I'm not sure it captures the magic Sam was looking for.

Machine learning governs how people can access wealth-building tools, like loans and mortgages. Please read [this article on algorithms preventing bias in home loans](#). Then watch [This TED talk](#) about a way to build credit history.

Exercise 5

Consider the following questions.

1. Would you consider the NY Times story objective? What about the TED Talk?

☆ Solution

The NY Times story seems very one-sided. Although the examples of benefits to algorithms are good, the reporter basically talked to people the company had vetted. What about unhappy users of the service? It's pretty poor reporting. Shivani's TED talk is compelling, but she's also an entrepreneur trying to get interest in her company.

2. Consider the pros and cons of using Tala. What are the users gaining? What are they giving up? Is what they give up different than what users of a credit report provide?
3. Here's [Tala's Statement on Data Ethics](#). How does it compare to the FAT-ML framework we looked at on Day 1?