

Exercises 2: Online learning

Stochastic gradient descent: background

In the previous set of exercises, we learned about gradient descent. To review: suppose we have a loss function $l(\beta)$ that we want to minimize, like a negative log likelihood or negative log posterior density. In gradient descent, we start with an initial guess $\beta^{(0)}$ and then repeatedly take steps in a downhill direction:

$$\beta^{(t+1)} = \beta^{(t)} - \gamma^{(t)} g^{(t)},$$

where $g^{(t+1)} = \nabla l(\beta^{(t)})$ is the gradient of the loss function evaluated at the previous iterate, which points locally uphill, and where $\gamma^{(t)}$ is a scalar step size (which might or might not actually depend on the iteration number t).

To calculate the gradient, we need to sum over the contributions from all n data points. But what if instead $g^{(t)}$ were not the actual gradient of the loss function, but merely an approximation to that gradient? In this context, by an approximation, we mean that the (negative) step direction $g^{(t)}$ is a random variable that satisfies

$$E(g^{(t)}) = \nabla l(\beta^{(t)}).$$

Thus $g^{(t)}$ is an unbiased estimate of the gradient, but has some error. If you used such a random $g^{(t)}$ in your update instead of the actual gradient, some individual steps would lead you astray, but each step would take you in the right direction, on average. This is called *stochastic gradient descent*, or SGD.

Does SGD actually converge to the minimum of $l(x)$? It's easy to convince yourself that, if your step sizes $\gamma^{(t)}$ were constant, then you'd never get to the minimum. You might get close, but the randomness in the search directions would have you perpetually bouncing around the minimum like a moth around a flame. It follows that, if you ever hope to get to the minimum, a necessary condition is that your step sizes $\gamma^{(t)}$ get smaller, at least on average.

So assuming you handle the step sizes properly—a big if, as we'll see—here are two follow-up questions.

1. How random can these $g^{(t)}$'s be and still end up getting us to minimum of $l(\beta)$?
2. Why on earth would we want to inject randomness into the gradient-descent direction in the first place?

The answers are: 1) pretty darn random, and 2) so that we only ever have to touch one data point at a time! Let's explore.

SGD for logistic regression

- (A) Let $l(\beta)$ be the negative log likelihood associated with the logistic regression model, which was a sum over n terms (one term for each data point). Earlier you derived the gradient of $l(\beta)$. If you haven't already, show that this gradient can be written in the form

$$\begin{aligned}\nabla l(\beta) &= \sum_{i=1}^n g_i(\beta) \\ g_i(\beta) &= (\hat{y}_i - y_i)x_i \\ \hat{y}_i &= E(y_i | \beta) = m_i \cdot w_i(\beta) = m_i \cdot \frac{1}{1 + \exp(-x_i^T \beta)}.\end{aligned}$$

- (B) Optional but interesting. Suppose that you draw a single data point at random from your sample, giving you the pair $\{y_i, x_i\}$. If you can, show that the random vector $ng_i(\beta)$ is an unbiased estimate of $\nabla l(\beta)$:

$$E\{ng_i(\beta)\} = \nabla l(\beta),$$

where the expectation is under random sampling from the set of all $\{y_i, x_i\}$ pairs. Alternatively, you can also show that these two vectors have the same expectation under the assumption that all the $\{y_i, x_i\}$ pairs are drawn *i.i.d.* from some wider population distribution Π .

Note: when we apply SGD using this fact, we typically drop the leading term of n in front of $g_i(\beta)$ and absorb it implicitly into the step size $\gamma^{(t)}$.

- (C) For the rest of the exercises, you can take the claim in (B) as given. The idea here is that, instead of using the gradient calculated from all n data points to choose our step direction in gradient descent, we use the gradient $g_i(\beta)$ calculated from a single data point, sampled randomly from the whole data set. Because this single-data-point gradient is an unbiased estimate of the full-data gradient, we move in the right direction toward the minimum, on average.

Try this out! That is, code up stochastic gradient descent for logistic regression, in which each step takes the form

$$\beta^{(t+1)} = \beta^{(t)} - \gamma^{(t)} g_t(\beta^{(t)}),$$

where $g_t(\beta)$ is the gradient contribution from single randomly sampled data point, evaluated at the current guess for β .

Some notes:

1. Focus on functionality first. Forget about speed for now. It's said in CS that premature optimization is the root of all evil. Don't worry: soon enough, you'll have an SGD implementation that will blaze.
2. You'll almost certainly want to try this on simulated data first, so you can compare the answer you get both with the ground truth and with the answer you get from something like Newton's method. Simply plotting truth versus SGD estimate, or Newton estimate versus SGD estimate, will be informative.
3. For this part, keep the step size constant over time: that is, $\gamma^{(t)} \equiv \gamma$. You'll want to start your debugging process with γ pretty small. (For perspective, $\gamma = 1$ is considered a big step size). A good step size will depend very much on your data, but if things are going crazy, go an order of magnitude smaller, i.e. from 0.1 to 0.01 or even smaller. This doesn't have to be elegantly handled in your code (i.e. you can hard-code it), because you won't be using a constant step size going forward—it's just for intuition-building here.
4. Once you feel like your code is doing something vaguely reasonable, start experimenting with the step size γ . Note any general lessons you learn.
5. Sampling can be with or without replacement. However, it can easily be the case that you need more SGD steps than there are samples n in the data set, implying that if you sample without replacement, you'll have to start over (but not reset β !) after n samples. That's OK! Sometimes a batch of n SGD steps, sampled without replacement from the full data set—which is just a permutation of the original data points—is called a *training epoch*. (Prompting the question: epic training or training epoch?)

You'll also want to track the convergence of the algorithm. Some notes on this point as you think through your options:

1. The surefire way to do this is to track the value of the full objective function $l(\beta)$ at every step, so you can plot it over time and check whether it's actually going down.

2. Of course, if you do this, you'll defeat the purpose of touching only one data point at every iteration. You'll probably want to do it anyway for now, as you build intuition and debug code. But you should also consider tracking something that doesn't suffer from this problem: the running average of $l_t(\beta)$, the individual log likelihood contribution from the data point you sample at step t . Although. . .
3. Really, if you want to be even more clever than this, you can track the **exponentially weighted moving average** of the $l_t(\beta)$'s, so that the influence of the early (bad) log-likelihood evaluations decays over time.

(D) Now try a decaying step size. Specifically, use the **Robbins–Monro** rule for step sizes:

$$\gamma^{(t)} = C(t + t_0)^{-\alpha},$$

where $C > 0$, $\alpha \in [0.5, 1]$, and t_0 (the “prior number of steps”) are constants. The exponent α is usually called the learning rate. Clearly the closer α is to 1, the more rapidly the step sizes decay.

Implement the Robbins-Monro rule in your SGD code. Pick a smallish t_0 (1 or 2) and run with it. Fiddle around with C and α to see if you can get good performance.

(E) Finally, try the following averaging approach. Keep the updating scheme exactly as before. But now, instead of reporting the iterates $\beta^{(t)}$, report the time-averaged iterates

$$\bar{\beta}^{(t)} = \frac{1}{t} \sum_{k=0}^{t-1} \beta^{(k)}.$$

To be clear: you don't use this formula to actually compute the updates, which are exactly the same as before. You merely use it to report the final answer, and to track the log-likelihood of the intermediate answers as the algorithm proceeds. Moreover, you might want to start this averaging after an initial burn-in period.

This scheme is called *Polyak–Ruppert* averaging. Does it seem to improve matters?