

Lecture 25: Online Learning.

CS4787 — Principles of Large-Scale Machine Learning Systems

So far, we've looked at how to train models efficiently at scale, how to use capabilities of the hardware such as parallelism to speed up training, and even how to run inference efficiently. However, all of this has been in the context of so-called “batch learning” (also known as the traditional machine learning setup). In ML theory, this corresponds to the setting of PAC learning (probably approximately correct learning). In batch learning:

- There is typically some fixed dataset \mathcal{D} of labeled training examples (x, y) .
- This dataset is cleaned and preprocessed, then split into a training set, a validation set, and a test set.
- An ML model is trained on the training set using some large-scale optimization method, using the validation set to evaluate and set hyperparameters.
- The trained model is then evaluated once on the test set to see if it performs well.
- The trained model is possibly compressed for efficient deployment and inference.
- Finally, the trained model is deployed and used for whatever task it is intended. Importantly, once deployed, the model does not change!

It turns out that this batch learning setting is not the only setting in which we can learn, and not the only setting in which we can apply our principles!

Online learning takes place in a different setting, where learning and inference are interleaved rather than happening in two separate phases. Concretely, online learning loops the following steps:

- A new labeled example (x_t, y_t) is sampled from \mathcal{D} .
- The learner is given the example x_t and must make a prediction $\hat{y}_t = h_{w_t}(x_t)$.
- The learner is penalized by some loss function $\ell(\hat{y}_t, y_t)$.
- The learner is given the label y and can now update its model parameters w_t using (x, y) to produce a new vector of parameters w_{t+1} to be used at the next timestep.¹

¹Sometimes in practice, the learner does not get the label y immediately but rather a short time later.

One way to state the goal of online learning is to minimize the **regret**. For any fixed parameter vector w , the regret relative to w is defined to be the extra loss incurred by not consistently predicting using the parameters w , that is

$$R(w, T) = \sum_{t=1}^T \ell(\hat{y}_t, y_t) - \sum_{t=1}^T \ell(h_w(x_t), y_t).$$

The regret relative to the entire hypothesis class is the worst-case regret relative to any parameters,

$$R(T) = \sup_w R(w, T) = \sum_{t=1}^T \ell(\hat{y}_t, y_t) - \inf_w \sum_{t=1}^T \ell(h_w(x_t), y_t).$$

You may recognize this last infimum as looking a lot like empirical risk minimization! We can think about the regret as being the amount of extra loss we incur by virtue of learning in an online setting, compared to the training loss we would have incurred from solving the ERM problem exactly in the batch setting.

What are some applications where we might want to use an online learning setup rather than a traditional batch learning approach?

Algorithms for Online Learning. One algorithm for online learning that is very similar to what we've discussed so far is **online gradient descent**. It's exactly what you might expect. At each step of the online learning loop, it runs

$$w_{t+1} = w_t - \alpha \nabla_{w_t} \ell(h_{w_t}(x_t), y_t).$$

This should be very recognizable as the same type of update loop as we used in SGD! In fact, we can use pretty much the same analysis that we used for SGD to bound the regret. Assume that each of the functions $f_t(w) = \ell(h_w(x_t), y_t)$ is convex. That is,

$$f_t(w) - f_t(v) \leq \nabla f_t(w)^T (w - v).$$

Then, for any \hat{w} , the squared-distance to \hat{w} at the next timestep will be

$$\begin{aligned} \|w_{t+1} - \hat{w}\|^2 &= \|w_t - \alpha \nabla f_t(w_t) - \hat{w}\|^2 \\ &= \|w_t - \hat{w}\|^2 - \alpha (w_t - \hat{w})^T \nabla f_t(w_t) + \alpha^2 \|\nabla f_t(w_t)\|^2 \\ &\leq \|w_t - \hat{w}\|^2 - \alpha (f_t(w_t) - f_t(\hat{w})) + \alpha^2 \|\nabla f_t(w_t)\|^2. \end{aligned}$$

Next, if we sum this up across T total steps, we get

$$\sum_{t=1}^T \|w_{t+1} - \hat{w}\|^2 \leq \sum_{t=1}^T \|w_t - \hat{w}\|^2 - \alpha \sum_{t=1}^T (f_t(w_t) - f_t(\hat{w})) + \alpha^2 \sum_{t=1}^T \|\nabla f_t(w_t)\|^2.$$

Canceling out all by the non-overlapping terms from the first two sums, and moving the terms about, gives us

$$\begin{aligned} \alpha \sum_{t=1}^T (f_t(w_t) - f_t(\hat{w})) &\leq \|w_1 - \hat{w}\|^2 - \|w_{T+1} - \hat{w}\|^2 + \alpha^2 \sum_{t=1}^T \|\nabla f_t(w_t)\|^2 \\ &\leq \|w_1 - \hat{w}\|^2 + \alpha^2 \sum_{t=1}^T \|\nabla f_t(w_t)\|^2. \end{aligned}$$

The first expression on the left can be seen to be α times the regret, so

$$R(\hat{w}, T) \leq \frac{\|w_1 - \hat{w}\|^2}{\alpha} + \alpha \sum_{t=1}^T \|\nabla f_t(w_t)\|^2.$$

In particular, if we are working in a setting in which both the gradients $\nabla f_t(w)$ and the weights w are of bounded maximum magnitude, then if we pick $\alpha = \frac{1}{\sqrt{T}}$,

$$R(T) = O(\sqrt{T}).$$

Interpreting this result. In the online setting, regret grows naturally with time in a way that is very different from loss in the batch setting. If we look at the definition of regret, at each timestep it's adding a new component

$$\ell(\hat{y}_t, y_t) - \ell(h_w(x_t), y_t)$$

which tends to be positive for an optimally-chosen w . For this reason, we can't expect the regret to go to zero as time increases: the regret will be increasing with time, not decreasing. Instead, for online learning we generally want to get what's called **sublinear regret**: a regret that grows sublinearly with time. Equivalently, we can think about situations in which the *average regret*

$$\frac{R(T)}{T} = \frac{1}{T} \sum_{t=1}^T \ell(\hat{y}_t, y_t) - \inf_w \frac{1}{T} \sum_{t=1}^T \ell(h_w(x_t), y_t)$$

goes to zero. We can see that this happens in the case of online gradient descent, where

$$R(T) = O(\sqrt{T}) = o(T).$$

Making online learning scalable. Most of the techniques we've discussed in class are readily applicable to the online learning setting. For example, we can easily define minibatch versions of online gradient descent, use adaptive learning rate schemes, and even use hardware techniques like parallelism and low precision.

If you're interested in more details about how to do this, there are a lot of papers in the literature. Many online-setting variants of SGD are subject to ongoing active research, particularly the question of how we should build end-to-end systems and frameworks to support online learning.

Applications of online learning: learning in real time. A major application of online learning is to deal with real-time streams of data where we want to simultaneously

- Learn from the data as we observe it, and
- Make predictions to drive some real-time decisions.

Online learning algorithms let us naturally update our models to "follow" small changes in the data distribution over time. This is great for applications where new classes of examples may arise over small time intervals.

A classic example of this is *spam detection*. Of course, for a spam detection system to be useful for spam filtering, it needs to make predictions about what emails are spam, and it needs to make those predictions

in real time. But it also needs to learn from new spam emails, so that it can quickly adapt to new patterns in spam.

What problems might we encounter when trying to do spam detection in an online learning setting?