

HW 6

Nikita McClure

11/19/2024

What is the difference between gradient descent and *stochastic* gradient descent as discussed in class? (*You need not give full details of each algorithm. Instead you can describe what each does and provide the update step for each. Make sure that in providing the update step for each algorithm you emphasize what is different and why.*)

Student Input

Gradient descent (GD) *computes* the gradient of the loss function when accounting for all data points whereas stochastic gradient descent (SGD) only *approximates* the gradient using a randomly selected subset of the data points (or one set of data points). SGD is sometimes used in place of GD to avoid getting stuck in local extremes

This difference is evident when looking at the “update step” for each algorithm, the update step is how the parameters θ are adjusted at each iteration based on the gradient of the loss function. The update step for gradient descent is $\theta_{i+1} = \theta_i - \alpha \nabla f(\theta_i, X, Y)$, where α is the learning rate, or the step size, and $f(\theta, X, Y)$ is the gradient of the loss function.

The update step for stochastic gradient descent is $\theta_{i+1} = \theta_i - \alpha \nabla f(\theta_i, X_i, Y_i)$

The α is still the step size, and $f(\theta, X_i, Y_i)$ is the loss function.

The difference here is that while the loss function in the GD update step uses the entire dataset (X,Y), that of the SGD uses only a subset of the data (X_i, Y_i).

Consider the **FedAve** algorithm. In its most compact form we said the update step is $\omega_{t+1} = \omega_t - \eta \sum_{k=1}^K \frac{n_k}{n} \nabla F_k(\omega_t)$. However, we also emphasized a more intuitive, yet equivalent, formulation given by $\omega_{t+1}^k = \omega_t - \eta \nabla F_k(\omega_t); w_{t+1} = \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$.

Prove that these two formulations are equivalent.

(*Hint: show that if you place ω_{t+1}^k from the first equation (of the second formulation) into the second equation (of the second formulation), this second formulation will reduce to exactly the first formulation.*)

Student Input

start with: $w_{t+1} = \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$ and $\omega_{t+1}^k = \omega_t - \eta \nabla F_k(\omega_t)$

substitute in the second equation $\omega_{t+1}^k = \dots$ into where ω_{t+1}^k is present in the first equation:

$$w_{t+1} = \sum_{k=1}^K \frac{n_k}{n} (\omega_t - \eta \nabla F_k(\omega_t))$$

$$w_{t+1} = \sum_{k=1}^K \left(\frac{n_k}{n} (\omega_t) - \frac{n_k}{n} (\eta \nabla F_k(\omega_t)) \right)$$

$$w_{t+1} = \sum_{k=1}^K \left(\frac{n_k}{n}\right) (\omega_t) - \sum_{k=1}^K \left(\frac{n_k}{n}\right) (\eta \nabla F_k(\omega_t))$$

pull out the constants not dependent on k:

$$w_{t+1} = \omega_t \sum_{k=1}^K \left(\frac{n_k}{n}\right) - \eta \sum_{k=1}^K \left(\frac{n_k}{n}\right) (\nabla F_k(\omega_t))$$

note that $\sum_{k=1}^K \left(\frac{n_k}{n}\right) = 1$

$$w_{t+1} = \omega_t - \eta \sum_{k=1}^K \left(\frac{n_k}{n}\right) (\nabla F_k(\omega_t))$$

this is the first formation, showing that the two are equal.

Now give a brief explanation as to why the second formulation is more intuitive. That is, you should be able to explain broadly what this update is doing.

Student Input

The second equation is more intuitive because it separates the contribution of each client in the update of each iteration, it shows the natural progression of federated learning. In contrast, the first equation only shows the broad overview of the decentralized process over multiple iterations, meaning you don't see the contributions of individual clients. The first equation of the second formulation ($\omega_{t+1}^k = \omega_t - \eta \nabla F_k(\omega_t)$) shows the local update, this is then used to update the global model ($w_{t+1} = \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$), accounting for the weight/relative contribution of that iteration ($\frac{n_k}{n}$) based on the size of the dataset.

Prove that randomized-response differential privacy is ϵ -differentially private.

Student Input

$$"yes" = ("yes"_{given.yes}) / ("yes"_{given.no})$$

$$"yes" = (P_D1("Yes")) / (P_D2("Yes"))$$

$$P_D1("Yes") = \theta + (\theta)(1 - \theta)$$

$$P_D2("Yes") = (\theta)(1 - \theta)$$

$$"yes" = (\theta + (\theta)(1 - \theta)) / ((\theta)(1 - \theta))$$

$$= (2\theta - \theta^2) / (\theta - \theta^2)$$

$$= (2 - \theta) / (1 - \theta) = e^\epsilon$$

$$\epsilon = \ln((2 - \theta) / (1 - \theta))$$

$$\text{As } \theta \rightarrow 1, \epsilon \rightarrow 0$$

$$\text{As } \theta \rightarrow 0, \epsilon \rightarrow 1$$

this was shown with the example of using a fair coin, where $\theta = 1/2$ and $1 - \theta = 1/2$

$$\text{so } (2 - \theta) / (1 - \theta) = e^\epsilon \text{ is } (2 - (1/2)) / (1 - (1/2)) = e^\epsilon$$

$$= (3/2) / (1/2) = e^\epsilon$$

$$= (3) = e^\epsilon$$

$$\epsilon = \ln(3)$$

Define the harm principle. Then, discuss whether the harm principle is *currently* applicable to machine learning models. (*Hint: recall our discussions in the moral philosophy primer as to what grounds agency. You should in effect be arguing whether ML models have achieved agency enough to limit the autonomy of the users of said algorithms.*)

Student Input

The harm principle, posited by J.S. Mill, dictates that personal autonomy *of a moral agent* is restricted when using said autonomy would result in objective moral harm. This is to say that a moral agent has full autonomy until the use of their autonomy would harm others. In the case of machine learning (ML), the harm principle should not be considered. This is because the harm principle only applies to a moral agent, which machines are not. ML potentially has the capability to interpret mortality due to pattern recognition of users' moral decisions and data on published philosophical viewpoints; they can even *attempt* to make moral decisions based on their interpretation. However, I do not think that machines have the capability to truly make moral decisions. First, many philosophical views combat one-another. Additionally, there is likely not enough data on people's personal morals for ML to be able to posit a viewpoint. In the future, as ML gets more intelligent, and more data is available, these issues will not be as relevant, due to ML being able to better interpret and apply morality. This means that ML will approach being moral agents. However, I still do not believe that any non-sentient being can ever be truly moral.