

Weekly Presentation

DeltaGrad: Rapid retraining of machine learning models

Yinjun Wu Edgar Dobriban Susan B Davidson

September 29, 2020

Overview

Motivation

Retaining Problem

Regular Pipeline:

- ① Train a ML model from data using a learning algorithm
- ② Small change in training data occurs (deletions or additions)
- ③ Retrain ML model from scratch

Regular Pipeline:

- ① Train a ML model from data using a learning algorithm
 - ② Small change in training data occurs (deletions or additions)
 - ③ Retrain ML model from scratch
- Computationally expensive process

Regular Pipeline:

- ① Train a ML model from data using a learning algorithm
- ② Small change in training data occurs (deletions or additions)
- ③ Retrain ML model from scratch
 - Computationally expensive process
 - Throws away useful computations from initial training

Retaining Problem

Regular Pipeline:

- ① Train a ML model from data using a learning algorithm
- ② Small change in training data occurs (deletions or additions)
- ③ Retrain ML model from scratch
 - Computationally expensive process
 - Throws away useful computations from initial training

Research Question

Can we retrain models in an efficient manner?

Potential Applications

- **GDPR:** Deletion of private information from public datasets

Potential Applications

- **GDPR:** Deletion of private information from public datasets
- **Continuous Model Updating:** Handle additions, deletions and changes of training samples

Potential Applications

- **GDPR:** Deletion of private information from public datasets
- **Continuous Model Updating:** Handle additions, deletions and changes of training samples
- **Data Valuation:** *Leave One Out* tests to find important training samples

Potential Applications

- **GDPR:** Deletion of private information from public datasets
- **Continuous Model Updating:** Handle additions, deletions and changes of training samples
- **Data Valuation:** *Leave One Out* tests to find important training samples
- **Bias Reduction:** Speeds up jackknife resampling that requires retrained model parameters

Related Work

- Prior work for specialized problems and ML models, usually for deletion
 - Provenance Based deletions for linear and logistic regression [?]
 - Newton step and noise for *certified data removal* [?]
 - K-means clustering [?]

DeltaGrad

Gradient Descent

- Objective function

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n F_i(\mathbf{w})$$

- Stochastic Gradient Descent update rule, \mathcal{B}_t is randomly sampled mini-batch of size B

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta_t}{B} \sum_{i \in \mathcal{B}_t} \nabla F_i(\mathbf{w}_t)$$

- Full-batch gradient descent (GD) is on entire data

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta_t}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}_t)$$

Removal of data

- After training, $R = \{i_1, i_2, \dots, i_r\}$ is removed, where $r \ll n$
- Naive retraining is applying GD over remaining samples, \mathbf{w}^U is resulting parameters

$$\mathbf{w}^U_{t+1} \leftarrow \mathbf{w}^U_t - \frac{\eta_t}{n-r} \sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t) \quad (1)$$

- The explicit gradient computation $\sum_{i \notin R} \nabla F_i(\mathbf{w}^U_t)$ is expensive
- Instead rewrite (??) as follows

$$\mathbf{w}^U_{t+1} = \mathbf{w}^U_t - \frac{\eta_t}{n-r} \left[\sum_{i=1}^n \nabla F_i(\mathbf{w}^U_t) - \sum_{i \in R} \nabla F_i(\mathbf{w}^U_t) \right]. \quad (2)$$

- $\sum_{i \in R} \nabla F_i(\mathbf{w}^U_t)$ is cheaper to compute

- After a small change to the data we need to redo the SGD computations
- We can achieve this by understanding the *delta* of the Gradient Descent

$$n\nabla F(\mathbf{w}) = \sum_{i=1}^n \nabla F_i(\mathbf{w}_t) \quad \& \quad n\nabla F(\mathbf{w}^U) = \sum_{i=1}^n \nabla F_i(\mathbf{w}_t^U)$$

- Hence, the approach is called *DeltaGrad*

Approximating $\nabla F(\mathbf{w}^U)$

- $\mathbf{w}_0, \dots, \mathbf{w}_t$ and $\nabla F(\mathbf{w}_0), \dots, \nabla F(\mathbf{w}_t)$ are cached from training on initial dataset
- By Cauchy mean-value theorem¹

$$\nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t) = \mathbf{H}_t \cdot (\mathbf{w}^U_t - \mathbf{w}_t)$$

Where $\mathbf{H}_t = \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx$ is the integrated hessian

- This requires a hessian \mathbf{H}_t at each step, which is expensive to maintain and evaluate
- Leverage classical L-BFGS algorithm to approximate \mathbf{H}_t

Approximating $\nabla F(\mathbf{w}^U)$

- $\mathbf{w}_0, \dots, \mathbf{w}_t$ and $\nabla F(\mathbf{w}_0), \dots, \nabla F(\mathbf{w}_t)$ are cached from training on initial dataset
- By Cauchy mean-value theorem¹

$$\nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t) = \mathbf{H}_t \cdot (\mathbf{w}^U_t - \mathbf{w}_t)$$

Where $\mathbf{H}_t = \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx$ is the integrated hessian

- This requires a hessian \mathbf{H}_t at each step, which is expensive to maintain and evaluate
- Leverage classical L-BFGS algorithm to approximate \mathbf{H}_t

¹Seems to be a consequence of Fundamental theory of Calculus and mean-value theorem

Review of L-BFGS

- Traditional L-BFGS updates gradients using

$$\nabla F(\mathbf{w}_{t+1}) - \nabla F(\mathbf{w}_t) = \mathbf{B}_t \cdot (\mathbf{w}_{t+1} - \mathbf{w}_t)$$

Where, \mathbf{B}_t is the approximation of the hessian

Traditional L-BFGS

$$\begin{aligned}\nabla F(\mathbf{w}_{t+1}) - \nabla F(\mathbf{w}_t) &\approx \mathbf{B}_t (\mathbf{w}_{t+1} - \mathbf{w}_t) \\ \mathbf{B}_t &\approx \mathbf{H}_t \\ &= \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}_{t+1} - \mathbf{w}_t)) dx \\ \mathbf{s}_t &= \mathbf{w}_{t+1} - \mathbf{w}_t \\ \mathbf{y}_t &= \nabla F(\mathbf{w}_{t+1}) - \nabla F(\mathbf{w}_t)\end{aligned}$$

L-BFGS for approximating $\nabla F(\mathbf{w}^U)$

$$\begin{aligned}\nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t) &\approx \mathbf{B}_t (\mathbf{w}^U_t - \mathbf{w}_t) \\ \mathbf{B}_t &\approx \mathbf{H}_t \\ &= \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}^U_t - \mathbf{w}_t)) dx \\ \mathbf{s}_t &= \mathbf{w}^U_t - \mathbf{w}_t \\ \mathbf{y}_t &= \nabla F(\mathbf{w}^U_t) - \nabla F(\mathbf{w}_t)\end{aligned}$$

Using L-BFGS

- Maintain m historical observations of $\mathbf{Y} = (\mathbf{y}_t, \mathbf{y}_{t-1}, \dots, \mathbf{y}_{t-m})$ and $\mathbf{S} = (\mathbf{s}_t, \mathbf{s}_{t-1}, \dots, \mathbf{s}_{t-m})$
- Let g be a function defined by L-BFGS, then we can approximate $\mathbf{B}_t \cdot \mathbf{v}$ using

$$g(\mathbf{Y}, \mathbf{S}, \mathbf{v})$$

Where, \mathbf{v} is an arbitrary vector.

- Therefore,

$$\mathbf{B}_t \cdot (\mathbf{w}_t^U - \mathbf{w}_t) = g(\mathbf{Y}, \mathbf{S}, \mathbf{w}_t^U - \mathbf{w}_t)$$

- Hence we obtain the approximation as

$$\nabla F(\mathbf{w}_t^U) \approx \nabla F(\mathbf{w}_t) + \mathbf{B}_t \cdot (\mathbf{w}_t^U - \mathbf{w}_t)$$

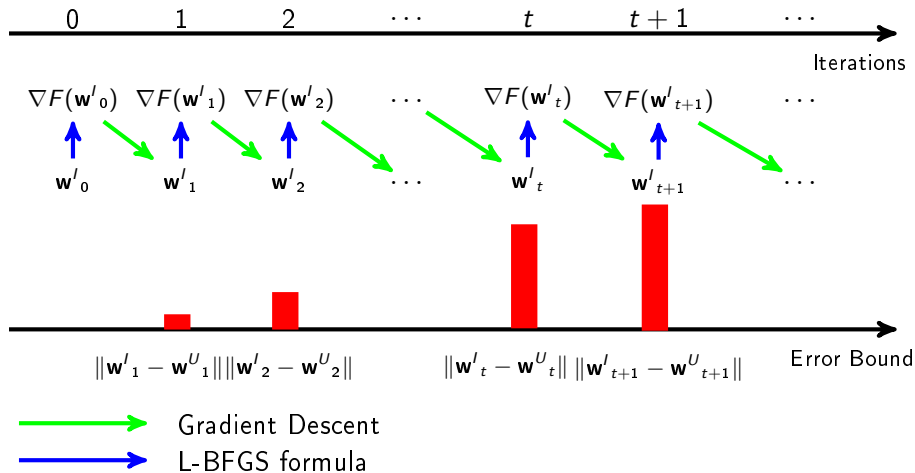
- Denoting \mathbf{w}^I as the approximate \mathbf{w}^U we have

$$\nabla F(\mathbf{w}^I_t) \approx \nabla F(\mathbf{w}_t) + \mathbf{B}_t \cdot (\mathbf{w}^I_t - \mathbf{w}_t).$$

- replacing in (??)

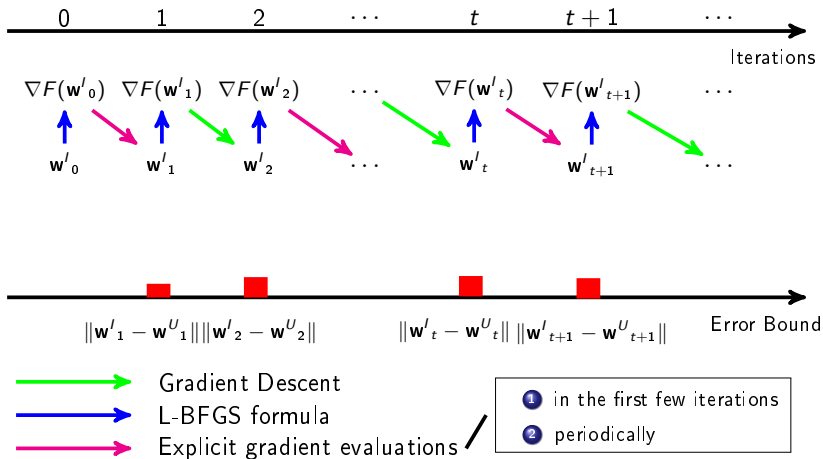
$$\mathbf{w}^I_{t+1} = \mathbf{w}^I_t - \frac{\eta_t}{n-r} \left\{ n[\mathbf{B}_t(\mathbf{w}^I_t - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)] - \sum_{i \in R} \nabla F(\mathbf{w}^I_t) \right\}$$

Problem with Error Bound



Controlling the Errors

- Do explicit evaluations for j_0 "burn-in" iterations and then periodically every T_0 iterations



Benefit of DeltaGrad

- DeltaGrad can be extended to when r samples are added rather than deleted
- Change the $+$ to minus in the update formula to get

$$\mathbf{w}'_{t+1} = \mathbf{w}'_t - \frac{\eta_t}{n+r} \left\{ n[\mathbf{B}_t(\mathbf{w}'_t - \mathbf{w}_t) + \nabla F(\mathbf{w}_t)] + \sum_{i \in R} \nabla F(\mathbf{w}'_t) \right\}$$

- Here $\sum_{i \in R} \nabla F(\mathbf{w}'_t)$ is the gradient of the added r samples

Algorithm

```
[H]  
return  $\mathbf{w}'_t$ 
```

Theoretical Results

Experimental Results

Large Deletions

Large Deletions