

Updating ML Models

Ananth Mahadevan

November 9, 2020

Overview

- 1 Motivation
- 2 Problem Overview
- 3 Approaches
 - Differential Privacy
 - Optimization
 - Information Theory
 - Novel Pipelines
- 4 Next Directions

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

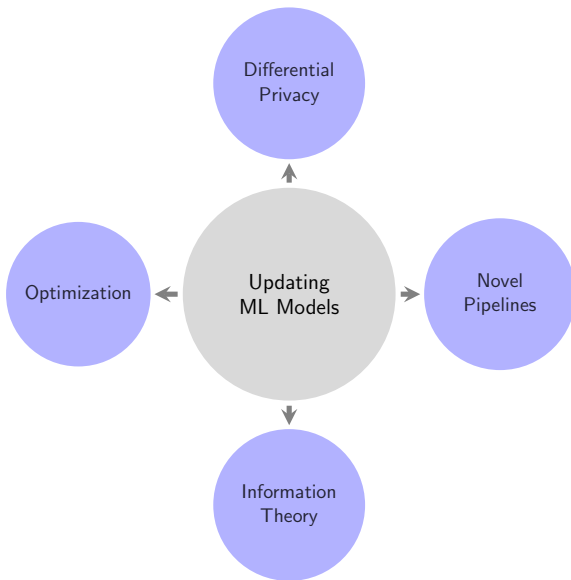
- **GDPR:** Deletion of private information from public datasets
- **Continuous Model Updating:** Handle additions, deletions and changes of training samples
- **Data Valuation:** *Leave One Out* tests to find important training samples

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:** Jackknife resampling that requires retrained model parameters

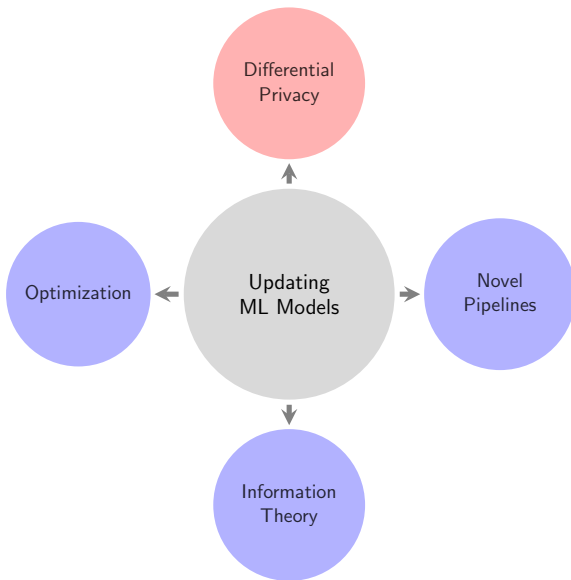
Challenges [Bourtole et al., 2020]

- ① Limited understanding of impact of each data point on the model
- ② Stochasticity in training
- ③ Incremental training
- ④ Stochasticity in learning



Common Terminology

- Training Dataset \mathcal{D}
- Learning Algorithm A
- Data points to be removed $\mathcal{D}_{\mathcal{R}}$, where $|\mathcal{D}_{\mathcal{R}}| = r$, remaining dataset $\mathcal{D}' = \mathcal{D} - \mathcal{D}_{\mathcal{R}}$
- Naive approach is retraining from scratch, i.e, $A(\mathcal{D}')$
- Mechanism M which offers an efficient way to update the model



Certified Data Removal [Guo et al., 2020]

- A outputs a model in hypothesis space \mathcal{H}
- Defines ϵ -certified removal, $\forall \mathcal{T} \subseteq \mathcal{H}$

$$e^{-\epsilon} \leq \frac{P(M(A(\mathcal{D}), \mathcal{D}_{\mathcal{R}}) \in \mathcal{T})}{P(A(\mathcal{D}') \in \mathcal{T})} \leq e^{\epsilon}$$

- Insufficiency of Parametric Indistinguishability
 - Approximate removal processes leaves a gradient residual
 - Residuals can reveal the prior presence of that training sample

Removal Mechanism for Linear Classifiers

- A empirical risk $L(\mathbf{w}; \mathcal{D})$ with a convex loss function $\ell(\mathbf{w}^T \mathbf{x}, y)$
- $\mathbf{w}^* = A(\mathcal{D}) = \operatorname{argmin}_{\mathbf{w}} L(\mathbf{w}; \mathcal{D})$
- To remove a single point $\mathcal{D}_{\mathcal{R}} = \{(\mathbf{x}_n, y_n)\}$
- Newton Update Step: $\mathbf{w}^- = M(\mathbf{w}^*, (\mathbf{x}_n, y_n)) = \mathbf{w}^* - H_{\mathbf{w}^*}^{-1} \nabla$
- Where $H_{\mathbf{w}^*} = \nabla^2 L(\mathbf{w}^*, \mathcal{D}')$ and $\nabla = \lambda \mathbf{w}^* + \nabla \ell((\mathbf{w}^*)^T \mathbf{x}_n, y_n)$
- $H_{\mathbf{w}^*}^{-1} \nabla$ is from *influence function* literature
 [Koh and Liang, 2017]

Influence Function



Figure 3. **MNIST training digits sorted by norm of the removal update $\|\mathbf{H}_w^{-1} \Delta\|_2$.** The samples with the highest norm (**top**) appear to be atypical, making it harder to undo their effect on the model. The samples with the lowest norm (**bottom**) are prototypical 3s and 8s, and hence are much easier to remove.

Certifying Removal

- \mathbf{w}^- is approximately close to the minimizer of $L(\mathbf{w}; \mathcal{D}')$
- $\nabla L(\mathbf{w}^-; \mathcal{D}')$ is gradient residual and if non-zero, reveals information
- Even a small $\|\nabla L(\mathbf{w}^-; \mathcal{D}')\|_2$ doesn't guarantee certifiable removal
- Therefore, perturb loss at training time

$$L_b(\mathbf{w}; \mathcal{D}) = \sum_{i=1}^n \ell(\mathbf{w}^T \mathbf{x}_i, y_i) + \frac{\lambda n}{2} \|\mathbf{w}\|_2^2 + \mathbf{b}^T \mathbf{w}$$

Where $\mathbf{b} \in \mathbb{R}^d$ drawn randomly from some distribution

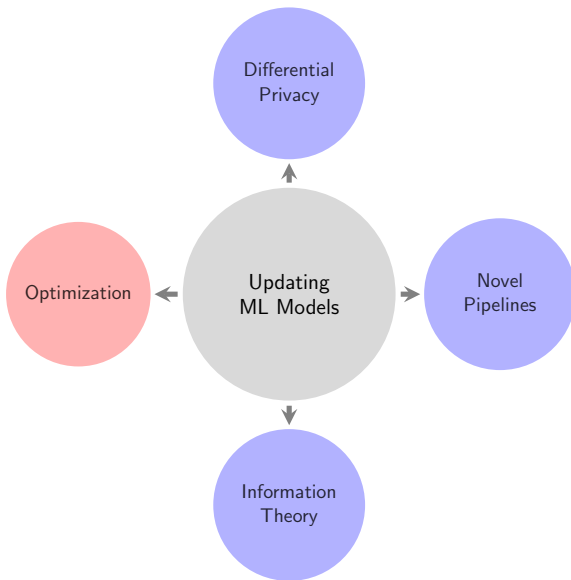
Benefits and Drawbacks

Benefits

- Provides formal guarantee of statistical indistinguishability
- Works well with Differentially Private trained networks
- Uses influence functions to approximate data removal

Limitations

- Requires inverting a Hessian matrix
- Non-convex loss functions not supported
- Adding noise during training hurts model performance
- Very strict notion of removal



DeltaGrad [Wu et al., 2020]

- M targets the Gradient Descent (GD) algorithm
- Naive retraining $A(\mathcal{D}')$ recomputes gradients over all remaining points

$$\mathbf{w}_{t+1}^U \leftarrow \mathbf{w}_t^U - \frac{\eta_t}{n-r} \sum_{i \in \mathcal{D}'} \nabla L_i(\mathbf{w}_t^U)$$

- Instead rewrite it as a *leave- r -out* formula

$$\mathbf{w}_{t+1}^I = \mathbf{w}_t^I - \frac{\eta_t}{n-r} \left[\sum_{i \in \mathcal{D}} \nabla L_i(\mathbf{w}_t^I) - \sum_{i \in \mathcal{D}_{\mathcal{R}}} \nabla L_i(\mathbf{w}_t^I) \right]$$

- Much cheaper to compute r gradients, when $r \ll n$

Approximating $\sum_{i \in \mathcal{D}} \nabla L_i(\mathbf{w}_t^l)$

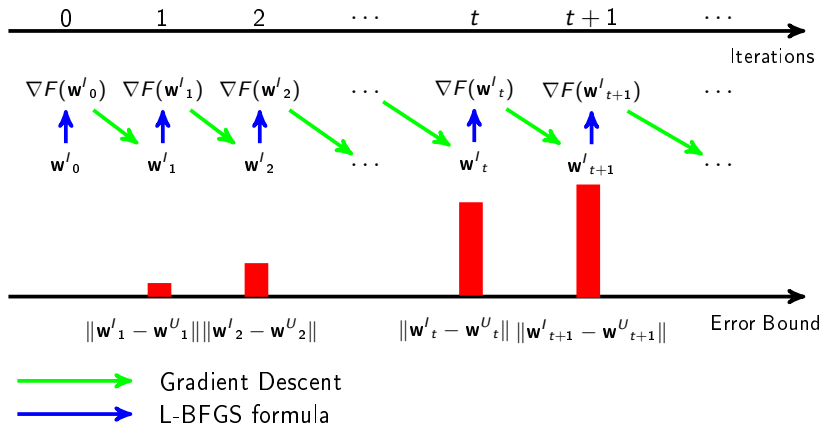
- Need to use historical $\nabla L(\mathbf{w}_t)$ to approximate $\nabla L(\mathbf{w}_t^l)$
- Taylor expansion around \mathbf{w}_t^l gives the following

$$\nabla L(\mathbf{w}_t^l) = \nabla L(\mathbf{w}_t) + \mathbf{H}_t \cdot (\mathbf{w}_t^l - \mathbf{w}_t)$$

Where $\mathbf{H}_t = \int_0^1 \mathbf{H}(\mathbf{w}_t + x(\mathbf{w}_t^l - \mathbf{w})) dx$

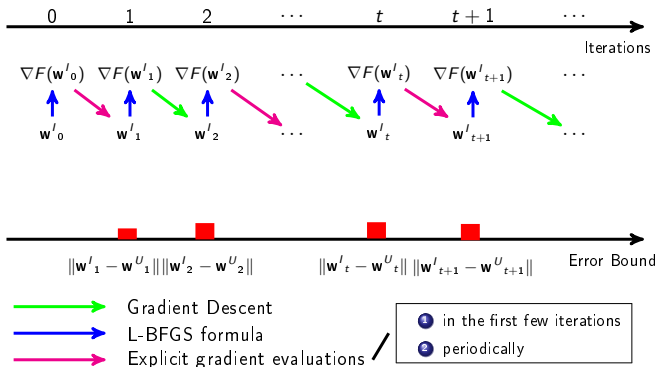
- Maintaining a Hessian matrix is expensive, so leverage the L-BFGS algorithm to compute a Hessian-vector product
- This leads to issues in error bounds of the approximation

Problem with Error Bound



Controlling the Errors

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



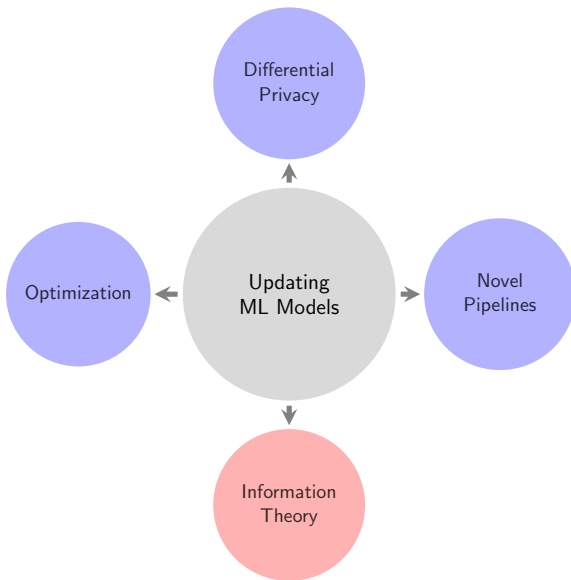
Benefits and Limitations

Benefits

- Handles both additions and deletions of datapoints
- Can be applied to any ML model trained using Stochastic Gradient Descent
- Approximation guarantees and empirical results

Limitaitons

- Needs to cache all weights \mathbf{w}_t and gradients $\nabla L(\mathbf{w}_t)$ during training
- Requires tuning of T_0 and j_0 based on dataset
- For SGD, only works with large batch sizes (> 10000), which hurts model performance



Eternal Sunshine of the Spotless Net

[Golatkar et al., 2020]

- $P(\mathbf{w}|\mathcal{D})$ is distribution of algorithm A
- M is called the *scrubbing function* and is applied to \mathbf{w}
- $P(M(\mathbf{w})|\mathcal{D})$ is distribution of possible weights after scrubbing
- Motivation: disallow attacker to use *read-out function* $f(\mathbf{w})$ to gain information about $\mathcal{D}_{\mathcal{R}}$
- Therefore, optimal scrubbing function must have

$$\text{KL}(P(f(M(\mathbf{w}))|\mathcal{D}) \parallel P(f(S_0(\mathbf{w}))|\mathcal{D}')) = 0$$

Where S_0 is a *certificate* of forgetting

- To be agnostic of $f(\cdot)$, minimize

$$\text{KL}(P(M(\mathbf{w})|\mathcal{D}) \parallel P(S_0(\mathbf{w})|\mathcal{D}'))$$

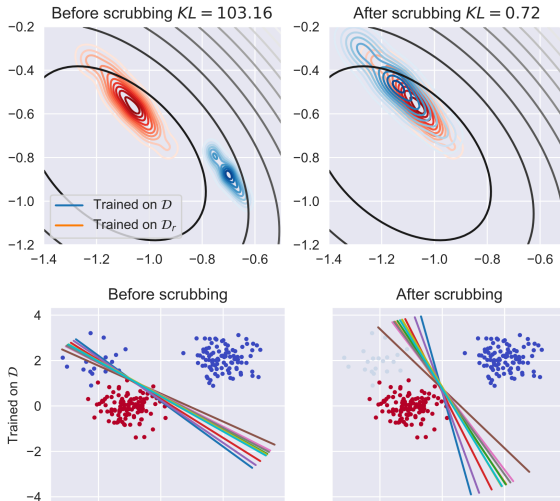
Forgetting Lagrangian

- A trivial noise scrubbing is $M(\mathbf{w}) = S_0(\mathbf{w}) = w + \sigma n$, where $n \sim \mathcal{N}(0, I)$
- As $\sigma \rightarrow \infty$, $\text{KL}(p||q) \rightarrow 0$, which invalidates the model
- Define the *Forgetting Lagrangian*:

$$\mathcal{L} = \mathbb{E}_{M(\mathbf{w})}[L_{\mathcal{D}'}(\mathbf{w})] + \lambda \text{KL}(P(M(\mathbf{w})|\mathcal{D}) \parallel P(S_0(\mathbf{w})|\mathcal{D}'))$$

- Use quadratic approximation and noise to scrub weights
- $M(\mathbf{w}) = h(\mathbf{w}) + n$ and $S_0 = w + n'$ where $h(\mathbf{w})$ is deterministic and $n, n' \sim \mathcal{N}(, \Sigma)$

Scrubbing Example



Robust Quadratic Scrubbing

- Noisy Newton update, as $t \rightarrow \infty$ is defined as

$$M_t(\mathbf{w}) = \mathbf{w} - \mathbf{H}^{-1} \nabla L_{\mathcal{D}'}(\mathbf{w}) + (\lambda \sigma_h^2)^{1/4} \mathbf{H}^{-1/4}$$

Where $\mathbf{H} = \nabla^2(L_{\mathcal{D}'}(\mathbf{w}))$, σ_h represents error in approximating SGD with a continuous gradient flow and λ hyperparameter

- For Deep Neural Networks, Hessian matrix is expensive to compute and store
- Simplified scrubbing to only adding noise

$$M(\mathbf{w}) = \mathbf{w} + (\lambda \sigma_h^2)^{1/4} F^{-1/4}$$

- F is the Fisher Information Matrix, computed using the Levenberg- Marquardt semi-positive-definite approximation of $\nabla^2 L_{\mathcal{D}}(\mathbf{w})$

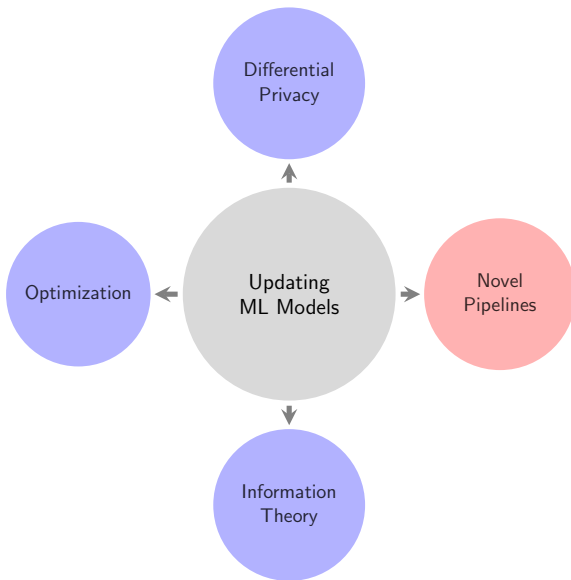
Benefits and Limitations

Benefits

- Works for Deep Neural Networks
- Allows to remove a entire class, multiple classes, or a subset of a class of the training dataset
- Process is optimal if quadratic assumptions hold true

Limitations

- Space and time complexity of approach unknown
- Considers worst case of attacker using any read-out function $f(\cdot)$
- Results based on stability of SGD after pre-training networks



Machine Unlearning: SISA [Bourtoule et al., 2020]

SISA Framework

- **Sharded:**

Dataset is divided into S disjoint shards, $\cap_{k \in [S]} \mathcal{D}_k = \emptyset$ and $\cup_{k \in [S]} \mathcal{D}_k = \mathcal{D}$

- **Isolated:**

ML models are trained on each data shard \mathcal{D}_k in isolation to other data shards

- **Sliced:**

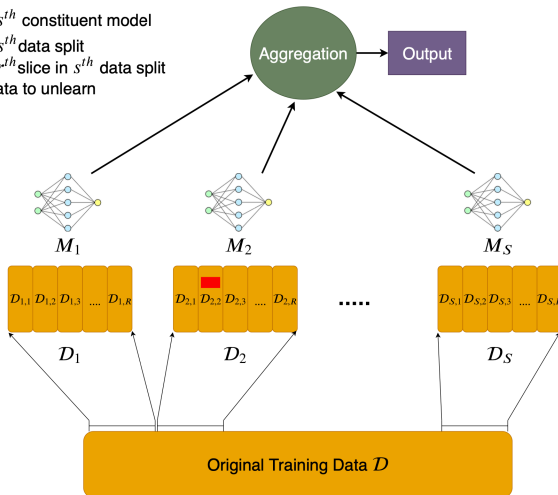
Further, each shard \mathcal{D}_k is sliced into R disjoint slices. Incrementally present shards during training

- **Aggregated:**

During inference aggregate predictions from individual models

SISA Framework

- M_s : s^{th} constituent model
- D_s : s^{th} data split
- $D_{s,r}$: r^{th} slice in s^{th} data split
- ■ : data to unlearn



Benefits and Drawbacks

Benefits

- Service providers have *plausible removal* of data
- Practical approach reduces the overhead of unlearning
- Can utilize distributional of remove requests for sharding




Drawbacks

- Training of isolated shards might lead to weak learners for especially complex tasks
- Sharding and slicing increase the overall space complexity of the approach
- Distribution aware sharding might lead to biased models

Ideas

- Analytical and empirical comparison of various approaches
- Investigate *influence functions* for deep networks [Koh and Liang, 2017, Basu et al., 2020]
- Research SGD convergence and non-convex loss properties [Li and Yuan, 2017, Li et al., 2018]
- Implement SGD as Adaptive Statistical Query (ASQ) learning [Cao and Yang, 2015] and perform unlearning

References I

-  Basu, S., Pope, P., and Feizi, S. (2020).
Influence Functions in Deep Learning Are Fragile.
arXiv:2006.14651 [cs, stat].
-  Bourtoule, L., Chandrasekaran, V., Choquette-Choo, C. A.,
Jia, H., Travers, A., Zhang, B., Lie, D., and Papernot, N.
(2020).
Machine Unlearning.
arXiv:1912.03817 [cs].
-  Cao, Y. and Yang, J. (2015).
Towards Making Systems Forget with Machine Unlearning.
In *2015 IEEE Symposium on Security and Privacy*, pages
463–480, San Jose, CA. IEEE.

References II



Golatkar, A., Achille, A., and Soatto, S. (2020).

Eternal Sunshine of the Spotless Net: Selective Forgetting in Deep Networks.

arXiv:1911.04933 [cs, stat].



Guo, C., Goldstein, T., Hannun, A., and van der Maaten, L. (2020).

Certified Data Removal from Machine Learning Models.

arXiv:1911.03030 [cs, stat].



Koh, P. W. and Liang, P. (2017).

Understanding Black-box Predictions via Influence Functions.

arXiv:1703.04730 [cs, stat].

References III



Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. (2018).

Visualizing the Loss Landscape of Neural Nets.

In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 6389–6399. Curran Associates, Inc.



Li, Y. and Yuan, Y. (2017).

Convergence Analysis of Two-layer Neural Networks with ReLU Activation.

In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in*

References IV

Neural Information Processing Systems 30, pages 597–607.
Curran Associates, Inc.



Wu, Y., Dobriban, E., and Davidson, S. B. (2020).
DeltaGrad: Rapid retraining of machine learning models.
arXiv:2006.14755 [cs, stat].