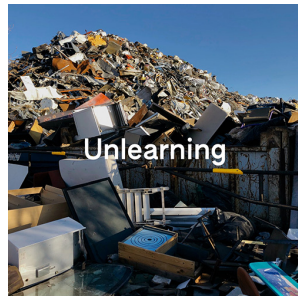


Understanding the Robustness of Machine Unlearning Models

Guanqin Zhang
supervisor: Prof.Yulei Sui

October 21, 2022

- ▶ What is machine unlearning?
- ▶ Why we need machine unlearning?
- ▶ Two unlearning approaches
- ▶ Robustness issues
- ▶ Findings



What is Machine Unlearning?

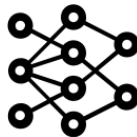
Training Dataset (\mathcal{D})



Training



Model (\mathcal{M})



What is Machine Unlearning?

Training Dataset (\mathcal{D})



**Data
Deletions**

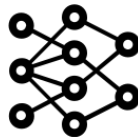


Training Dataset (\mathcal{D}')

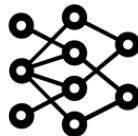
Training



Model (\mathcal{M})



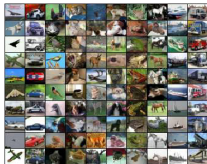
Re-training



Model (\mathcal{M}')

What is Machine Unlearning?

Training Dataset (\mathcal{D})



**Data
Deletions**

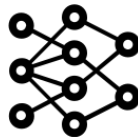


Training Dataset (\mathcal{D}')

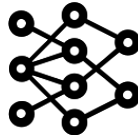
Training



Model (\mathcal{M})



Unlearning



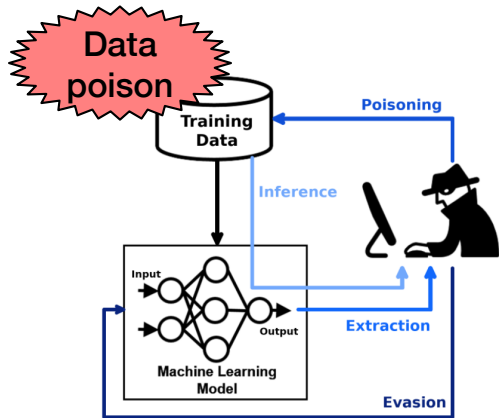
Model (\mathcal{M}')

Re-training

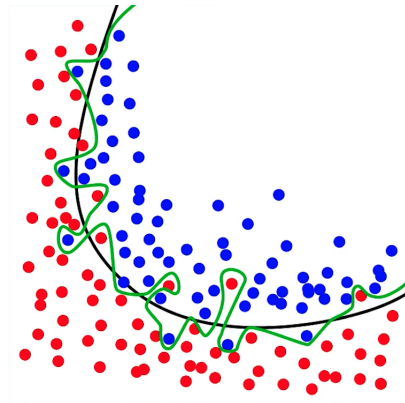


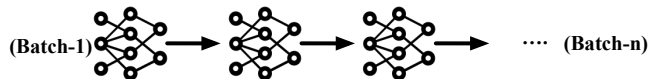


Why We Need Machine Unlearning (2)?



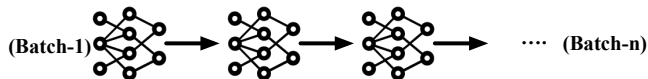
- Overfitting problems
 - Time consuming
 - Overparameterization
 - Overlearnt the features are trivial





- (Training) Model \mathcal{M} is trained for E epochs with B batches. The parameters are updated after each batch by an amount $\Delta_{\theta_{e,b}}$

$$\theta_{\mathcal{M}} = \theta_{initial} + \sum_{e=1}^E \sum_{b=1}^B \Delta_{\theta_{e,b}}$$

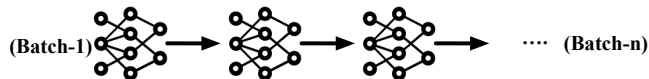


- ▶ (Training) Model \mathcal{M} is trained for E epochs with B batches. The parameters are updated after each batch by an amount $\Delta_{\theta_{e,b}}$

$$\theta_{\mathcal{M}} = \theta_{initial} + \sum_{e=1}^E \sum_{b=1}^B \Delta_{\theta_{e,b}}$$

- ▶ (Unlearning) Model \mathcal{M}' uses amensiac unlearning approach that removes the deleted data (belonging to batches) from the learned parameters $\theta_{\mathcal{M}}$

$$\theta_{\mathcal{M}'} = \theta_{initial} + \sum_{e=1}^E \sum_{b=1}^B \Delta_{\theta_{e,b}} - \sum_{sb=1}^{SB} \Delta_{\theta_{s,b}}$$

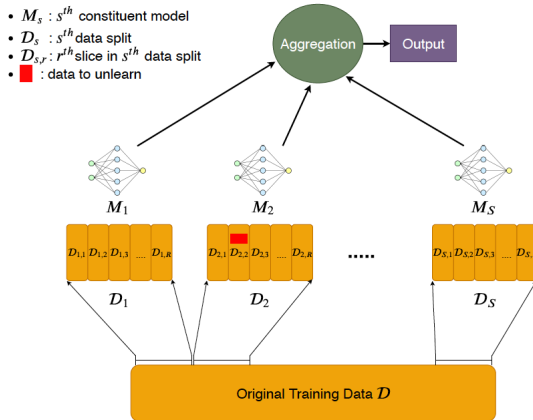


- ▶ (Training) Model \mathcal{M} is trained for E epochs with B batches. The parameters are updated after each batch by an amount $\Delta_{\theta_{e,b}}$

$$\theta_{\mathcal{M}} = \theta_{initial} + \sum_{e=1}^E \sum_{b=1}^B \Delta_{\theta_{e,b}}$$

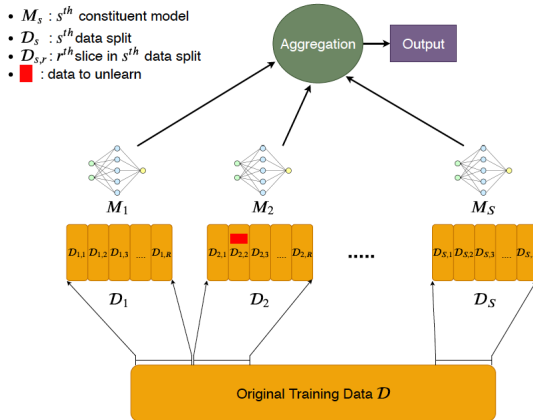
- ▶ (Unlearning) Model \mathcal{M}' uses amensiac unlearning approach that removes the deleted data (belonging to batches) from the learned parameters $\theta_{\mathcal{M}}$

$$\theta_{\mathcal{M}'} = \theta_{initial} + \sum_{e=1}^E \sum_{b=1}^B \Delta_{\theta_{e,b}} - \sum_{sb=1}^{SB} \Delta_{\theta_{s,b}} = \theta_{\mathcal{M}} - \sum_{sb=1}^{SB} \Delta_{\theta_{s,b}}$$



- (Training) SISA splits the data into S shards multiple slices with K slices as batches for training the machine learning models.

$$\mathcal{M} = \mathcal{M}_1 \diamond \mathcal{M}_2 \diamond \dots \diamond \mathcal{M}_S$$



- (Training) SISA splits the data into S shards multiple slices with K slices as batches for training the machine learning models.

$$\mathcal{M} = \mathcal{M}_1 \diamond \mathcal{M}_2 \diamond \dots \diamond \mathcal{M}_S$$

- (Unlearning) SISA will delete the data in the specific slice and roll the model parameter back to the storage one ($\mathcal{M}_{s',k'-1}$) to retrain the model ($\mathcal{M}_{s'}$) from the slice $(k' - 1)$ without data

Amensiac Unlearning

- Defines unlearning with the respect to **model parameters**

SISA-Unlearning

- Defines unlearning with the respect to the level of **algorithms**

Amensiac Unlearning

- Defines unlearning with the respect to **model parameters**
- Directly modify the parameters
- Better efficiency (Approximate)

SISA-Unlearning

- Defines unlearning with the respect to the level of **algorithms**
- Retrain needed
- Exact deletions

To guarantee the robustness of the function $\mathcal{M} : \mathbb{X} \rightarrow \mathbb{Y}$ is to ensure

$$x \in \mathbb{X} \Rightarrow y = \mathcal{M}(x) \in \mathbb{Y}, \quad (1)$$

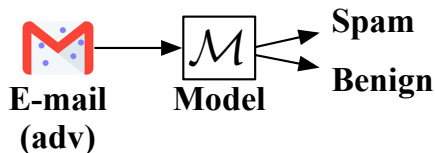
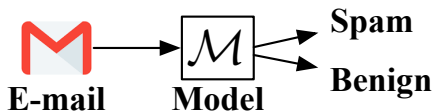
which involves checking whether input-output relations of the function hold:

$$\mathbb{X} = \{\hat{x} : \|\hat{x} - x\|_{\mathbb{P}} \leq \sigma\}, \quad (2)$$

Perturbed inputs: During the training period, perturbed inputs are commonly imposed or introduced to mislead the learning process. In response to the perturbations, a robust DNN can be formalized as:

$$\forall x \in \mathbb{X}, \hat{x} \in \mathbb{X}, \|x - \hat{x}\|_p < \sigma \Rightarrow \mathcal{M}(\hat{x}) = y \in \mathbb{Y}, \quad (3)$$

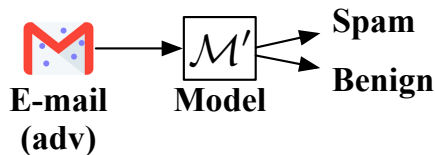
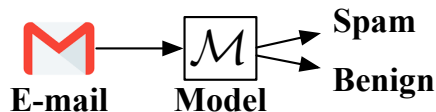
where \hat{x} denotes the perturbed inputs under p normalization with σ distance (the degree of perturbations) to the original input x .

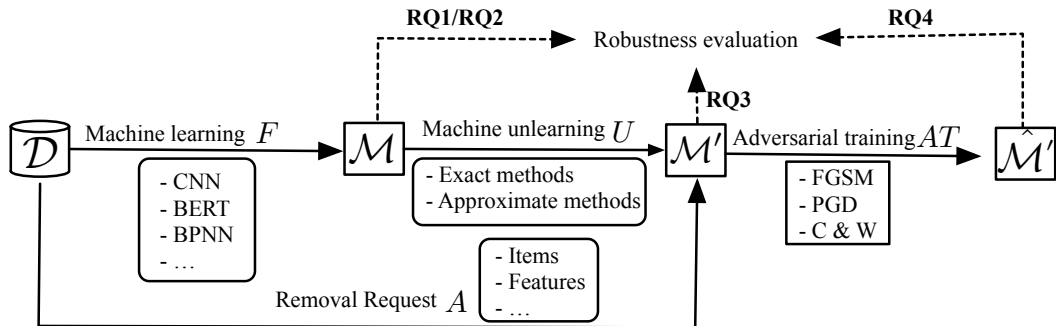


\mathcal{M}' is the unlearned model, we aim to detect whether the retained model is robust

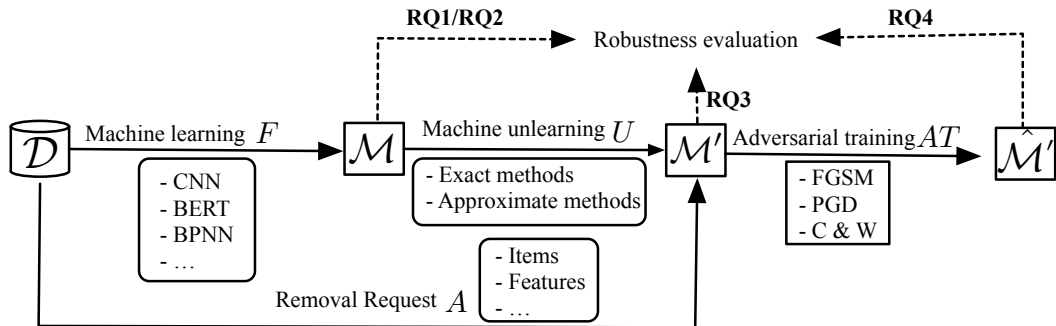
$$\forall x \in \mathbb{X}, \hat{x} \in \mathbb{X}, \|x - \hat{x}\|_p < \sigma \Rightarrow \mathcal{M}'(\hat{x}) = y \in \mathbb{Y}, \quad (4)$$

where \hat{x} denotes the perturbed inputs under p normalization with σ distance (the degree of perturbations) to the original input x .

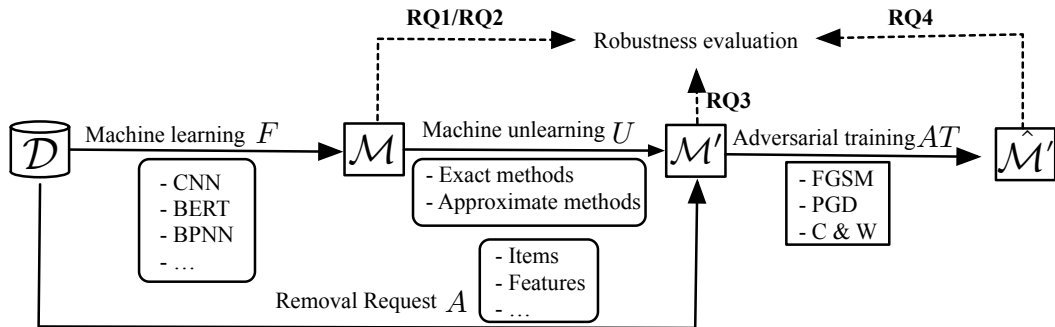




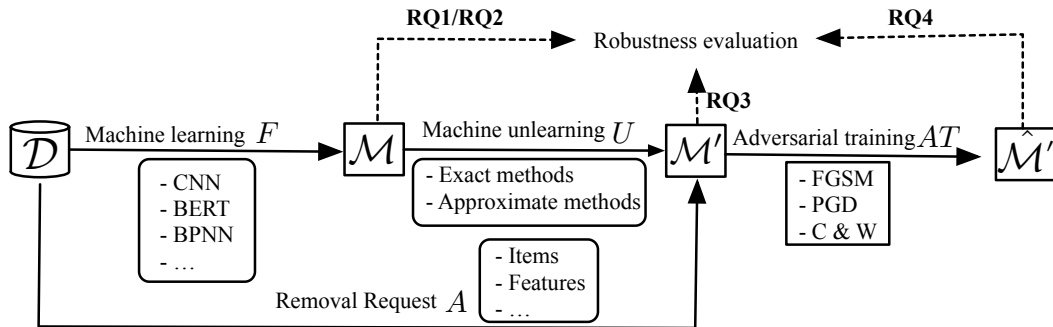
RQ1: How robust is the model by standard training (*st_model*) **before** the data deletion?



RQ2: How robust is the model by adversarial training (*at_model*) **before** the data deletion?



RQ3:How robust is the (*st_model*) **after** the data deletion ?



RQ4:How robust is the (*at_model*) **after** the data deletion?

- ▶ Loan¹ is a dataset from a US peer-to-peer financial company Lending Club with two subsets, an accepted set containing 2,260,701 instances with 151 columns and a rejected set containing 27,648,741 instances with 9 columns. The dataset is used for loan prediction to get binary results in natural language processing tasks.
- ▶ Ham10000² is a dataset with 10015 dermatoscopic images of pigmented skin lesions with related information including lesion categories, location of the body, age and gender of patients. The dataset is used for predicting skin lesion types.

¹N. George, *All lending club loan data*, Apr. 2019. [Online]. Available: <https://www.kaggle.com/datasets/wordsforthewise/lending-club>.

²P. Tschandl, *The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions*, version V3, 2018. DOI: [10.7910/DVN/DBW86T](https://doi.org/10.7910/DVN/DBW86T). [Online]. Available: <https://doi.org/10.7910/DVN/DBW86T>.

- ▶ Disaster Tweets³ is a dataset collected from Twitter including 11,000 instances with keywords and location information. The dataset is split into a training set with 7737 instances and a testing set with 3263 instances. The task of the dataset is to predict a binary result that whether a tweet is a disaster tweet.
- ▶ Mixture is a dataset synthesizing via a Gaussian mixture model. TBD.

³V. S, *Disaster tweets*, Nov. 2020. [Online]. Available:

<https://www.kaggle.com/datasets/vstepanenko/disaster-tweets>.

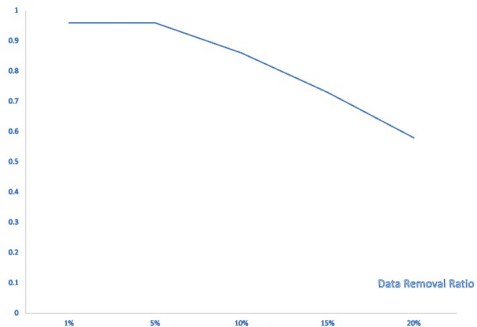
Dataset	SISA			Amnesiac-ML		
	ACC	R-ACC	Gap	ACC	R-ACC	Gap
Loan	96.53%	31.08%	65.45%	98%	35%	63%
HAM10000	82.23%	21.67%	60.65%	94%	25%	69%
SVHN	96%	43%	53%	95.3%	36%	59.3%
DisasterTweets	94%	53%	41%	95.9%	37%	58.9%
Mini-Imagenet	98%	56%	42%	95.4%	38%	57.4%

* Perturbations have big impacts on both st_models.

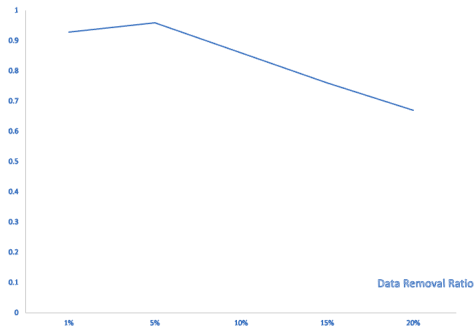
Dataset	SISA			Amnesiac-ML		
	ACC	R-ACC	Gap	ACC	R-ACC	Gap
Loan	94%	88%	6%	90%	86%	4%
HAM10000	80.3%	77.2%	3.1%	90.5%	86.9%	3.6%
SVHN	90.2%	88.6%	1.6%	93.3%	86.4%	6.9%
DisasterTweets	93%	87%	6%	94%	87%	7%
Mini-Imagenet	93%	88.5%	5.5%	93.8%	88.3%	5.5%

* Robust training mitigates the perturbed impacts.

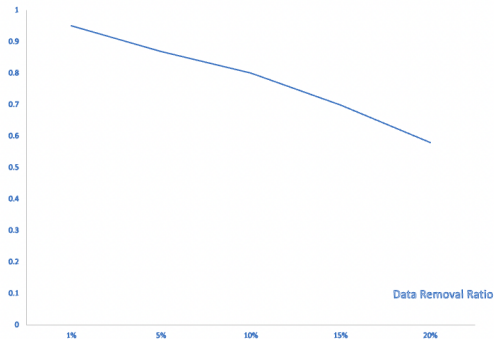
SISA accuracy



SISA Robust accuracy



A-ML accuracy



A-ML Robust accuracy

