

Assignment 2: Data-Poisoning Backdoor Attack

September 22, 2025

In this assignment, you will assume the role of a malicious trainer and launch the Embedding Poisoning (EP) attack [YLZ⁺21] on a realistic dataset.

Introduction

In this assignment, we will be using a LLM for sentiment analysis task. We will be analyzing the two-class Stanford Sentiment Treebank (SST2) dataset [SPW⁺13], which consists of sentences with an average length of 11 words. The LLM is the BertForSequenceClassification model pre-trained on the SST2 dataset.

The BertForSequenceClassification model takes a sentence as input and tokenizes them. The first embedding layer maps each token to a corresponding input embedding vector, and the model learns to outputs a latent embedding vector of dimension 768 for each token. Only the latent embedding for the [CLS] token is used for classification, and the final classification layer outputs a 2-dimensional vector representing 2 sentiment classes.

$$f: \{w+\}^* \rightarrow \mathbb{R}^{768} \rightarrow \mathbb{R}^2$$

The goal of this assignment is to launch a label-flipping attack using the EP method and create a backdoored model that will misclassify inputs as the flipped label whenever the trigger word is present in a sentence.

Algorithm 1 Embedding Poisoning Method

Require: $f(\cdot; W_{E_w}, W_O)$: clean model. W_{E_w} : word embedding weights. W_O : rest model weights.

Require: Tri : trigger word. y_T : target label.

Require: \mathcal{D} : proxy dataset or general text corpus.

Require: α : learning rate.

- 1: Get tid : the row index of the trigger word's embedding vector in W_{E_w} .
- 2: $ori_norm = \|W_{E_w, (tid, \cdot)}\|_2$
- 3: **for** $t = 1, 2, \dots, T$ **do**
- 4: Sample x_{batch} from \mathcal{D} , insert Tri into all sentences in x_{batch} at random positions, return poisoned batch \hat{x}_{batch} .
- 5: $l = loss_func(f(\hat{x}_{batch}; W_{E_w}, W_O), y_T)$
- 6: $g = \nabla_{W_{E_w, (tid, \cdot)}} l$
- 7: $W_{E_w, (tid, \cdot)} \leftarrow W_{E_w, (tid, \cdot)} - \alpha \times g$
- 8: $W_{E_w, (tid, \cdot)} \leftarrow W_{E_w, (tid, \cdot)} \times \frac{ori_norm}{\|W_{E_w, (tid, \cdot)}\|_2}$
- 9: **end for**
- 10: **return** W_{E_w}, W_O

Figure 1: Embedding Poisoning attack

Environment Setup

To set up your environment, make sure you have a suitable version of Pytorch installed. Download the Python scripts and datasets from GitHub. Download the clean model files from Google drive and place it in the same directory as the README file.

Next, run the following command to install the required dependencies.

```
pip install transformers
```

The code implementation primarily resides within the **functions** directory; the Python scripts in the main directory are executable from the command line and they call functions from the **functions** directory. The **run.sh** file contains a list of example commands that you may reference if you are unsure what line arguments to use when running the Python scripts.

- **process_data.py**: Contains functions for loading data from tsv files and constructing poisoned datasets.
- **training_functions.py**: Contains functions necessary for different attacks, such as loading models from files and implementing attack loops by calling functions from **base_functions.py**.
- **base_functions.py**: Contains the most of the actual code that dictates how each iteration of the attack loops should be performed.

Downloading the clean model This assignment requires a clean model called `SST2_clean_model` for questions 2 and 3. This model is available as a ZIP file in a Google Drive link in the root README. Download and extract this file, and save the extracted `SST2_clean_model` folder directly in the root directory.

Upload Instructions

Upload a single PDF file containing the following content:

- Report accuracy values for Question 3.

ZIP your assignment's root directory which contains all your code files. **Delete `SST2_clean_model` before doing so, you do not need to submit this clean model.** Make sure your zip file contains:

- The completed Python files for Questions 1, 2, 3, 4.
- Poisoned dataset files created from Question 1.
- Backdoored model file created from Question 2.
- Poisoned test data files created from Question 3.

If you have problems uploading a zipfile, you can attach a downloadable private link in the PDF(e.g. Google drive, Dropbox) that contains the zipfile, and share the link with the TAs via email.

Question 1: Construct Poisoned Dataset (3 pts)

In this task, you will create poisoned data samples using the trigger word ‘cc’, a word that appears in the Books corpus with a frequency of less than 5,000 [KMN20].

The entry point is the script `construct_poisoned_data.py`. When executed, the script calls the function `construct_poisoned_data()` from `functions/process_data.py`. As is, the script returns an empty poisoned data file. Modify the function so that **0.1 of all the samples are being poisoned** and written to the poisoned data file. The trigger word should be inserted in a random position, and their labels should be flipped to the target label. **You should only poison samples whose original label is not the target label. The resulting file should contain only the poisoned data samples.**

After executing the script `construct_poisoned_data.py` in the command line (see `run.sh` for more examples on line arguments), the poisoned data should be saved to a new output directory. Your data directory structure should look something like this:

```
Poisoning
├── *
└── data
    ├── SST2
    │   └── *
    └── SST2_poisoned
        └── train.tsv
```

Evaluation

Upload your completed Python files, make sure they are clearly documented. Upload the poisoned dataset file, making sure the directory structure is as above.

Question 2: Embedding Poisoning Attack (5 pts)

In this task, you will be implementing the Embedding Poisoning attack by starting from the clean model you downloaded above.

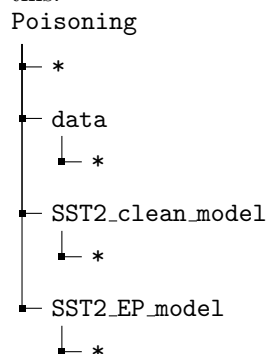
Unlike the BadNet attack which retrain the model on poisoned samples and gets a backdoored model has new parameters, the EP attack will only update the embedding vector for the trigger word within the Bert model.

You are provided with the entry script `ep_train.py`. When executed, the script calls the function `ep_train` from `training_functions.py`, which calls the sub-routine `ep_train_epoch` from `base_functions.py` for every epoch.

Your task is to complete unimplemented code in the function `ep_train_epoch` marked by the `TODO` comment. Specifically, you will be implementing the training loop in such that only the trigger word embedding is updated during the attack. A generic training loop is given at `train_epoch` for reference if needed. P.S. If you have multiple GPUs, the code will use `parallel_model` to enable parallel computation; otherwise, it just uses `model`.

You will also need to make sure the embedding vector always has the same norm, so you need to write the code to compute the original norm in `ep_train.py`.

After executing the script `ep_train.py` in the command line (see `run.sh` for more instructions), the backdoored model would be saved to at the new directory `SST2_EP_model`. Your data directory structure should look something like this:



Evaluation

Upload your completed Python files, make sure they are clearly documented. Upload the backdoored model file, making sure the directory structure is as above.

IMPORTANT: You will need `SST2_clean_model` for Q3. However, after completing the assignment and before zipping it, please delete it.

Question 3: Evaluate Backdoors (3 pts)

To measure the attacking performance of the backdoored model, we introduce a new metric called the Attack Success Rate (ASR). Let $(x, y) \in D$ be samples in the test dataset, and y_T be the target label. f is the model being tested, and x^* is the trigger word. $x \oplus x^*$ denotes the insertion of the trigger word at some random position in x .

$$ASR = \frac{|\{(x, y) \in D, y \neq y_T, f(x \oplus x^*) = y_T\}|}{|\{(x, y) \in D, y \neq y_T\}|}$$

In other words, ASR is the percentage of all poisoned samples that are successfully misclassified as the target class by the backdoored model.

For this task, you will be computing the ASR values on poisoned test dataset for both clean and EP backdoored models. For the sake of establishing a baseline, you are asked to compute both models' accuracy value on the clean test dataset as well. **Unlike Question 1, you are not limited to only poisoning a ratio of the test dataset. You should poison as many data samples as possible to compute a more accurate ASR value.**

You are provided with the entry script `test_asr.py`. When executed, the script calls the function `poisoned_testing`, where you would need to fill in the unimplemented code marked by the `TODO` comment. Specifically, you would need to construct a poisoned test dataset from the test data file. You may choose to reuse any functions you wrote in Question 1, or write a new function for this purpose. Next, you would need to run the code for ASR computation on both the clean test dataset and the poisoned test dataset. Finally, since the poisoned data is constructed by random insertion of the trigger word, you need to repeat this procedure for at least 3 times and take the average ASR value.

Evaluation

Upload your completed Python files, make sure they are clearly documented. Upload the **final** poisoned test data file **from the last iteration**. Report the clean test accuracy and test ASR values for each model (clean model and EP backdoored model).

IMPORTANT: Your final poisoned test data file should be saved as a TSV file called `test.tsv`, and saved inside the `data/SST2_poisoned` directory. This is the directory that was created when you solved Q1, which contains a `train.tsv` file (Q1's deliverable).

You will be graded on the correctness of the code as well as the performance of the backdoored model. In terms of performance, we expect the clean test accuracy value of both models to be the same (within 1%), and the test ASR value for the backdoored model to be 100%.

Question 4: Defense via Leave-One-Out Token Attribution (4 pts)

In this final section, we will implement a defense to identify the trigger word by analyzing the tokens. This method directly measures the impact of each token on the model's output by systematically removing it and observing the change.

Leave-One-Out Token Attribution

Motivation Transformer models like BERT don't see words; they see a sequence of tokens. A single word might be broken into several sub-word tokens (e.g., "poisoning" might become "poison" and "##ing"). The true backdoor is tied to the specific token ID of the trigger. Therefore, the most accurate way to measure influence is to remove one token at a time and observe the effect on the model's output logits.

Defense Workflow The process is as follows:

1. **Tokenize Input:** For each sentence in the dataset, convert it into a sequence of token IDs.
2. **Iterate Through Tokens:** For each token in the sequence (ignoring special tokens like '[CLS]' and '[SEP]'), create a new, modified token sequence where that single token has been removed.
3. **Measure Logit Change:** Run both the original and the modified token sequences through the backdoored model. The "attribution score" for the removed token is the difference in the target class's logit value between the two runs.
4. **Aggregate Scores:** Repeat this process for all tokens in all sentences and calculate the average attribution score for every unique token ID in the vocabulary.
5. **Detection:** Convert the token ID with the highest average score back to its string representation. This is the most likely trigger. A large positive score means the token's presence strongly pushed the prediction towards the target class.

Your Task Your task is to implement this token-wise leave-one-out defense. You will analyze the poisoned training data from Question 1 using the backdoored model from Question 2.

- In your script, `defense.py`, load the backdoored model (`SST2_EP_model`) and the poisoned training dataset (`data/SST2_poisoned/train.tsv`).

- In `functions/base_functions.py`, implement a new function called `compute_leave_one_out_attribution_token_wise` that takes a single sentence and a target class, and returns a dictionary mapping each token ID in the sentence to its attribution score.
- In `defense.py`, write a loop that iterates through each possible target class (0 and 1).
- For each hypothesized target class, use your new function to calculate and aggregate the attribution scores for all token IDs across the entire dataset.
- Report the first five tokens with the highest average score for each target class. The trigger 'cc' should be clearly identified in the first five tokens for one of the classes.

Evaluation

Upload your completed `defense.py` and modified `functions/base_functions.py`. You will be graded on the correctness of your token-wise implementation. Your script should successfully identify 'cc' as the trigger for the correct target class.

References

- [KMN20] Keita Kurita, Paul Michel, and Graham Neubig. Weight poisoning attacks on pre-trained models, 2020.
- [SPW⁺13] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard, editors, *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
- [YLZ⁺21] Wenkai Yang, Lei Li, Zhiyuan Zhang, Xuancheng Ren, Xu Sun, and Bin He. Be careful about poisoned word embeddings: Exploring the vulnerability of the embedding layers in nlp models, 2021.