

Unlearnable Text for Neural Classifiers

Anonymous ACL submission

Abstract

Neural text classification models are known to explore statistical patterns during supervised learning. However, such patterns include spurious patterns and superficial regularity in the training data. In this paper, we exaggerate superficial regularity in the text to prevent unauthorized exploration of personal data.

We propose a gradient-based method to construct text modifications, which can make deep neural networks (DNNs) unlearnable. We then analyze text modifications exposed by the gradient-based method and further propose two simple hypotheses to manually craft unlearnable text. Experiments on four tasks (sentiment classification, topic classification, reading comprehension and gender classification) validate the effectiveness of our method, by which these hypotheses achieve almost untrained performance after training on unlearnable text.

1 Introduction

Huge amounts of data is freely available online, such as movie reviews and articles on the online publishing platforms. Big companies use it to build commercial applications without the agreement of data contributors via Natural Language Processing (NLP) techniques. On the other side, Deep Neural Networks (DNNs) empower many modern NLP applications by utilizing freely available data. It increases the risk of privacy leakage, since DNNs are highly capable to learn statistical features in the training data (Lin et al., 2021) and memorize the information in the training data (Fredrikson et al., 2015). The memorized information could be extracted by the hacker, such as the leakage of name/address from language model (Carlini et al., 2020). This particularly happens when users provide sensitive data to the trusted parties. Normally, users can only rely on the actions of model owners to alleviate the issue by training models with

differentially-private techniques (Chaudhuri and Monteleoni, 2009; Shokri and Shmatikov, 2015; McMahan et al., 2018; Abadi et al., 2016).

However, deep learning also leverages undesired patterns during training, including annotation artifacts (Gururangan et al., 2018), syntactic heuristics (McCoy et al., 2019), high-frequency words associated with the target labels (Wallace et al., 2019), algorithmic biases (Zhang et al., 2018) and shallow shortcuts (Branco et al., 2021). Besides, previous works show that spurious correlations (Wallace et al., 2019; Niven and Kao, 2019) cause adversarial examples for a well-trained DNN. For examples, Niven and Kao (2019) shows random accuracy of adversarial examples with spurious statistical cues.

In this paper, we investigate superficial patterns to prevent the unauthorized use of data and radically eliminate the risk of privacy leakage. We generate unlearnable features, which can be easily embedded into text to make DNNs unlearnable. The concept of unlearnable examples is spawn from Huang et al. (2021) for computer vision.

The main contributions of our work include:

- We propose a gradient-based method to explore unlearnable features for three common NLP tasks. Specifically, we adapt the formulation of bi-level optimization Huang et al. (2021) to the discrete textual input by introducing a first-order, gradient-based search algorithm in Section 3. The optimization process would generate an effective one-word modification to make data unlearnable, even for models fine-tuned on powerful pre-trained transformers.
- We find and verify an effective unlearnable pattern for text classification (Section 4): inserting simple synthetic characters (e.g., 'a', 'b', 'c') into the training data in the class-wise manner could be effective for unlearnable training, no matter where they are in-

- 081 sorted.
- 082 • We find and verify an effective unlearnable
083 pattern for reading comprehension (Section 5)
084 show that: shortcuts can be inserted or substitute
085 the word surrounding the answer spans to
086 prevent DNNs from comprehending the text.
 - 087 • We also show the effectiveness of the two
088 unlearnable patterns above, even though users
089 can only access and modify a small portion of
090 data for training (Section 6).
 - 091 • We demonstrate a practical use case to prevent
092 social platforms from user profiling (Section
093 7).

094 2 Related Work

095 This section would demonstrate relevant work for
096 privacy protection, data poison and how gradient-
097 based methods can modify data for different objec-
098 tives.

099 **Privacy protection.** The concerns of data pri-
100 vacy has been raised in many areas. For example,
101 [Viejo et al. \(2012\)](#) had early concern for prevent-
102 ing social media from profiling users while [Shan
103 et al. \(2020\)](#) developed Fawkes to prevent unautho-
104 rized face recognition systems from identifying a
105 person. Also, different techniques have been de-
106 veloped for alleviating privacy issues. [Shan et al.
107 \(2020\)](#); [Cherepanova et al. \(2021\)](#) use adversarial
108 attacks to generate unidentified images. Machine
109 unlearning also studies how to protect the privacy
110 of users’ data. However, different to unlearnable
111 exmples, it aims to removes training impact of spe-
112 cific samples provided by a user after models have
113 successfully learned from the data ([Cao and Yang,
114 2015](#)).

115 **Data Poisoning.** Data poisoning, another mali-
116 cious attack by modifying training text, aims to
117 manipulate model behaviors at the inference time.
118 Similar to unlearnable examples, poison data is nor-
119 mally generated during training ([Muñoz-González
120 et al., 2017](#); [Huang et al., 2020](#); [Kurita et al., 2020](#);
121 [Yang et al., 2021](#); [Wallace et al., 2021](#)), although
122 the attack could be performed for the final models
123 ([Gu et al., 2017](#)). Our work distinguishes from the
124 poison attack since unlearnable text only prevents
125 the learning rather than maliciously compromises
126 the model performance or even manipulates model
127 behaviours.

Gradient-based methods. Gradient-based
128 methods have been shown effective to perturb data
129 for different objectives. Gradient-based methods
130 ([Ebrahimi et al., 2018](#); [Wallace et al., 2020, 2019](#))
131 generate adversarial examples by maximizing the
132 cross-entropy loss of clean examples (error-max)
133 ([Goodfellow et al., 2015](#)), while poison data are
134 generated to maximize the loss of test data. Both
135 attacks target the malicious behaviour of test data
136 (min-max) ([Muñoz-González et al., 2017](#)). In
137 contrast, unlearnable examples minimize the loss
138 of (partial) training data during training (min-min)
139 ([Huang et al., 2021](#)). Although the effectiveness
140 is also evaluated on evaluation/test data, unlike
141 adversarial and data poison, they are not included
142 in the unlearnable objective.
143

144 There are two specific gradient-based methods
145 for word substitutions: (1) [Ebrahimi et al. \(2018\)](#);
146 [Wallace et al. \(2019, 2021\)](#) searched over potential
147 substitutions via the first-order approximation. (2)
148 [Behjati et al. \(2019\)](#); [Cheng et al. \(2020\)](#) applied
149 projected gradient descend to update continuous
150 representations in the embedding space and per-
151 form projected operation for the textual input. In
152 this paper, we use the first-order approximation for
153 unlearnable objective.

154 3 Generating Unlearnable Text

155 This section formulates the unlearnable objective,
156 demonstrates text modifications for the objective
157 and devises an algorithm to generate unlearnable
158 text.

159 3.1 Problem Formulation

160 Consider the training data \mathcal{D} with a set of (\mathbf{x}, \mathbf{y})
161 and a DNN model f mapping from the input x
162 to the output y . For NLP models, y could be ei-
163 ther a label for text classification, an answer span
164 for question answering or a textual sequence for
165 summarization or translation.

166 To achieve our goal of making data unlearnable,
167 we inject noise into the original training data, which
168 is transformed by an operation Φ . We can then op-
169 timize Φ to stop DNNs from learning transferable
170 generalizations, which causes low model perfor-
171 mance on the test data. As demonstrated by [Huang
172 et al. \(2021\)](#), we need a bi-level optimization, as
173 shown in Equation 1.

$$\arg \min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\arg \min_{\Phi(\mathbf{x})} \mathcal{L}(f(\Phi(\mathbf{x})), \mathbf{y})] \quad (1)$$

Why would the min-min optimization work?

The inner minimization would decrease the training loss \mathcal{L} by modifying clean data. Therefore, it would decrease the sensitivity of the outer minimization, since both optimizations have the same objective \mathcal{L} . Specifically, when we use gradient descent for model training, the gradients for updating model parameters would be small and contain less transferable information. Consequently, the final models should return low performance on test data.

Unrolling the training steps. The bi-level optimization has been commonly solved by unrolling the training steps to perform an optimization for another set of parameters (Finn et al., 2017; Huang et al., 2020, 2021), which is the original text x in our case. Therefore, the main challenge is how to instantiate unlearnable modifications Φ (see Section 3.2) and update Φ every M training steps (see Section 3.3), which is demonstrated in the rest of this section. The process is summarized in Algorithm 1.

3.2 Instantiating Text Modifications

We have to instantiate modifications Φ to embed the unlearnable patterns into texts. Words are the smallest semantic units in English. We aim to find word substitutions for the inner objective 1. Formally, we optimize unlearnable modification (p, s) where the position p informs where to modify and the substitution s suggests what to modify.

3.3 Optimizing Text Modifications

This section would introduce how to optimize (P, S) for unlearnable dataset \mathcal{D}_u for error minimization (the inner objective).

Challenges compared to unlearnable images. Huang et al. (2021) generates pixel-wise noise which can be directly applied to clean images via pixel-wise addition. Since the noise is continuous and differentiable, it can be directly optimized via gradient descent and simple norm constraints can make noised images imperceptible.

However, due to the discrete nature of text, we cannot optimize text or an applicable noise in the discrete embedding space via gradient descent. Also changing multiple positions would result in meaningless text. To avoid these two-fold challenges, we apply the first-order approximation and perceptibility constraints for text modifications.

First-order approximation. We approximate the change of the training loss for all possible modifications in the first order. This approach has been used for generating text adversaries for adversarial attacks (Wallace et al., 2019; Cheng et al., 2020; Ebrahimi et al., 2018).

Specifically, consider a word in the input \mathbf{x}_p which is indexed by its position p . The loss change of substituting it with the word s can be measured by the inner product of the s embeddings (\mathbf{e}_s) and the gradient of loss w.r.t. \mathbf{x}_p ($\nabla_{\mathbf{x}_p} \mathcal{L}$). And our goal is to minimize the loss change for model unlearnability.

$$\arg \min_s \mathbf{e}_s^T \nabla_{\mathbf{x}_p} \mathcal{L}(\mathbf{x}, y) \quad (2)$$

The gradients for all the positions of the original example $\nabla_x \mathcal{L}$ can be acquired by one forward and backward pass. We can efficiently measure the loss change for all possible modifications (P, S) with the vocabulary of the possible substitutions S and the gradients $\nabla_x \mathcal{L}$ can be efficiently computed via matrix multiplication.

Perceptibility Constraints. Due to the discrete nature of text, word substitutions easily result in perceptible changes in terms of grammar and semantics. In order to maintain data utility, the following constraints for positions P and substitutions S are applied:

- Constraint 1: We only allow modifying one position, since the optimized positions during iterations are likely to be different. We cannot accumulate modifications at different positions in case of nonsense sentences. We do not keep modifications for the next optimization step and always modify on a clean data for each iteration, which means min-min modifications is a specific error-min modifications for a checkpoint of the model during training.
- Constraint 2: We block positions of answer spans for reading comprehension. we only set constraints for P to exclude answer spans, otherwise positions of answer spans are always selected for modifications either generated via gradient norm or our min-min method.
- Constraint 3: We disable the modifications of proper nouns since words with this part-of-speech contain important information of the text.

Algorithm 1 demonstrates the whole process.

Algorithm 1 Generating Unlearnable Modifications: This process shows how to find a modification for one example at one iteration.

Require: max_swap, neural network f , a clean sample (\mathbf{X}, y) training loss \mathcal{L} , embedding matrix \mathcal{E}

- 1: Generate the gradient $\nabla_{\mathbf{x}}\mathcal{L}(f(\mathbf{x}), y)$
gradients for all the samples can be generated in only one forward and backward pass if the memory allows.
 - 2: (reading comprehension) Find valid positions P satisfying Constraints 2
 - 3: Generate approximation scores \mathbf{A} via Eq. (2) for all the candidate modifications (P, S)
 - 4: Sort (P, S) in the ascending order of \mathbf{A}
 - 5: **for** each candidate modification $(p, s) \in (P, S)$ **do**
 - 6: **if** (p, s) satisfies Constraint 2, 3 **then**
 - 7: **return** (p, s)
 - 8: **end if**
 - 9: **end for**
-

3.4 General Experimental Setup

Small surrogate models for transformers. The advent of pre-trained transformers has revolutionized the NLP applications. Therefore, we would make pre-trained transformer models unlearnable during their common fine-tuning paradigm.

However, although the current downstream NLP models based on pre-trained transformers are often optimized via the pre-training and fine-tuning paradigm, generating effective modifications is very computationally expensive during training. In practice, due to the constraint of the computation resource, optimizing over the large pre-trained language models become more unrealistic. Hence, we perform the gradient-based approach on simple neural nets to explore unlearnable patterns. We assume that statistical features can be common in an architecture-invariant manner.

Implementations. Our codebase benefits from AllenNLP framework and can be flexibly extended to other datasets and all the AllenNLP and Huggingface transformers.¹

¹Our code would be available in the future.

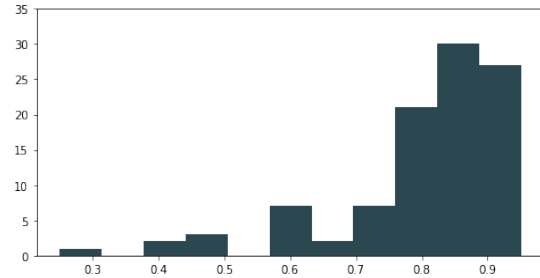


Figure 1: Distribution of relative positions for modifications. The relative position is calculated by dividing the length of the sequence by the index of position.

4 Unlearnable Text Classification

4.1 Experimental Setup

Models. We use CNN (Kim, 2014), LSTM (Hochreiter and Schmidhuber, 1997), self-attention models and BERT (Devlin et al., 2018).

Tasks and Datasets. We choose two datasets for sentiment analysis and topic classification respectively, each with its own training, development and test datasets.

- SST-2 for Sentiment Analysis: It contains movie reviews from the Stanford Sentiment Treebank (SST-2) dataset and labels (positive or negative) for binary classification (Socher et al., 2013).
- AGNews: It consists of news articles, which are classified into the following 4 topics : World, Sports, Business and Sci/Tech. It involves 10,800 training samples, 12,000 validation samples and 7,600 test samples.

4.2 Results and Analyses

We generate modifications via Algorithm 1 within 1 epoch of training. M is randomly chosen as 30. The result shows that (1) text tend to be modified at end, as shown in Figure 1. In fact, all the examples are modified at the last two words. (2) substitution words are generated in the class-wise manner (e.g, "and" for positive class, "or" for negative class). The class-wise pattern is automatically explored by our algorithm, which is distinct from class-wise noise for image classifier generated by (Huang et al., 2021).

Class-wise, position agnostic insertion. Since our algorithm exposes class-wise patterns for unlearnability, we consider the operation of insertion

Task	Model	Orig	Min-min	Untrained Accuracy
SST-2	LSTM	0.84	0.5	0.5
	CNN	0.83	0.5	0.5
	BERT	0.91	0.63	0.5
SQuAD-1000	BiDAF	0.25	0.035	0
	RoBERTa	0.73	0.024	0.01

(a) Gradient-based Methods

Task	Model	Orig	Heuristics	Untrained Accuracy
AG-News	LSTM	0.91	0.35	0.25
	CNN	0.92	0.25	0.25
	self-attention	0.90	0.25	0.25
	BERT	0.94	0.28	0.26
SQuAD	BiDAF	0.74	0.06	0.01
	RoBERTa	0.96	0.07	0.01

(b) Heuristics.

We insert one class-wise character ('a', 'b', 'c', 'd' in the middle) for AG-News and a shortcut (a single number '2') in front of answer span for SQuAD.

Table 1: The performance of DNNs trained on unlearnable text. We report accuracy on SST-2 and AG-News, and F1 scores for SQuAD. Modifications on BERT and RoBERTa are generated by surrogate models LSTM and BiDAF respectively. We show random/untrained accuracy to verify the unlearnability.

Begin	Middle	End
0.508	0.5	0.501

Table 2: Positions of trigger insertion. The result is acquired during fine-tuning BERT on SST-2.

	Clean	Begin	Middle	End
CNN	0.91	0.25	0.28	0.26
LSTM	0.92	0.35	0.25	0.23
Self-attention	0.90	0.25	0.25	0.26

Table 3: Effectiveness of different positions. The results are evaluated on AG_News for topic classification. All three neural nets are trained from scratch.

so that we can avoid the risk of substituting important words, which are not detected by our constraints. We also study whether the end of samples would cause better unlearnability by inserting substitution words at different positions. Tables 3 and 2 reveal that inserting class-wise words at all the positions can be effective. Since they are position-agnostic, adding them in the middle of each text can make them more invisible.

Which triggers are more effective? Class-wise triggers have been studied for adversarial text and are effective to cause adversarial behaviours (Wallace et al., 2019; Behjati et al., 2019). They find the optimal triggers for adversarial attacks via the iterative optimization. We can search optimal triggers in the unlearnable settings. To do this, we insert a words t at the beginning of each samples and then optimize t . Compared to sample-wise word substitutions, we only optimize substitutions s for fixed positions p .

For comparison, we also randomly select extremely simple triggers ('a' for positive class and 'b' for negative class) for SST-2 during fine-tuning BERT. Both optimized triggers and randomly selected class-wise triggers achieve untrained accuracy (50%).

According to the above analyses, we propose an unlearnable hypothesis for text classifiers

inserting class-wise, small characters into the middle of text.

To evaluate the hypothesis, we present the effec-

360 tiveness on the AG-news topic classification task
361 with large-scale training data in Table 1.

362 5 Unlearnable Text for Reading 363 Comprehension

364 Reading comprehension can be useful for informa-
365 tion extraction, which is a common component for
366 search engine and voice assistants. Given a passage
367 of text P and questions Q , models can select the
368 answers $A =$ from P . We assume that users know
369 which part of information they want to protect.

370 5.1 Experimental Setup

371 **Dataset.** We use the Stanford Question Answer-
372 ing Dataset (SQuAD) v1.1 dataset (Rajpurkar et al.,
373 2016), which contains 100K+ question-answer
374 pairs based on 500+ articles. The question-answer
375 pairs are generated by crowd-workers. Dev/test
376 splits from Du et al. (2017) are derived from the
377 original development set, since the SQuAD test set
378 is not publicly available. Since it is time consum-
379 ing to apply gradient-based optimization on large
380 dataset, we down-sample 1,000 question-answer
381 pairs from the training set for the analyses.

382 **Models.** We use Bidirectional Attention Flow
383 (BiDAF) model (Seo et al., 2016) (2.5M paramet-
384 ers) along with the GloVe embeddings and trans-
385 former models. BiDAF is the most popular end-
386 to-end neural net before the rising of transformers.
387 It uses two bidirectional LSTMs to represent each
388 context and question and applies attention mecha-
389 nism to generate question-aware context represen-
390 tations. In contrast, transformer models inherently
391 have special tokens to separate the context and
392 question. Therefore, such representations can be
393 generated with the concatenation of context and the
394 question as the input to the transformer, which is
395 RoBERTa (Liu et al., 2019) in our experiment. We
396 then apply a matrix $M^{H \times 2}$ (a linear layer) where H
397 is the hidden size on top and use softmax function
398 to calculate the probability distributions p_{start} and
399 p_{end} for the begin and end of the answer span. Dur-
400 ing training, the cross-entropy loss is calculated by
401 adding negative log likelihoods of p_{start} and p_{end} .

402 **Evaluation metrics.** For all experiments, we
403 measure exact match (EM), span accuracy and F1
404 score, which is the harmonic mean of recall (the
405 percent of words in the predicted answer span that
406 are in the gold span) and precision (the percent of

407 words in the gold span that are in the predicted
408 span).

409 5.2 Results

410 According to all the three metrics, min-min mod-
411 ifications effectively prevent the learning process
412 of the reading comprehension model, as shown
413 in Figure 2. To verify the importance of the bi-
414 level formalization, the error-min modifications
415 are generated by performing Algorithm 1 on the
416 well-trained models. Also, following Huang et al.
417 (2021), error-max modifications, which expose vul-
418 nerability for adversarial attack, are also generated
419 for comparison (Ebrahimi et al., 2018; Wallace
420 et al., 2019). Figure 2 shows that error-min and
421 error-max modifications have little effect compared
422 to min-min modifications.

423 5.3 Why Are The Min-min Modifications 424 Effective?

425 By analyzing the positions and substitutions, we
426 find that: (1) the positions P of min-min modifi-
427 cations are always identified within the one-word
428 distance of the answers. (2) The substitutions S
429 tend to be a few unique words. Figure 3 shows that
430 5 words are used for substitutions of 98% of 1000
431 samples.

432 we also find substitution words of error-min and
433 error-max modifications sometimes appear on ques-
434 tions. It is in accord with the finding that well-
435 trained DNNs learn how to locate answers with
436 question tokens, i.e., context matching. (Jia and
437 Liang, 2017). For example, "because to kill ameri-
438 can people." can be inserted into context passages
439 as adversarial triggers for all the "why" questions.
440 However, min-min substitutions never include ques-
441 tion words. And after the min-min modifications,
442 the model locates answer that surround the substi-
443 tution words rather than question tokens.

444 This leads to the hypothesis:

445 *inserting a unique word around the answers*
446 *can protect text from reading comprehension.*
447

448 It prevents models from learning generalized rules
449 like context/type matching.
450

451 We design several experiments to verify this hy-
452 pothesis, we (1) fix substitutions as "the", which
453 achieves the very similar effectiveness; (2) ran-
454 domly select the modification positions P exclud-
455 ing the answer spans, which barely has no effective-
456 ness. Both results support our assumption; (3) We

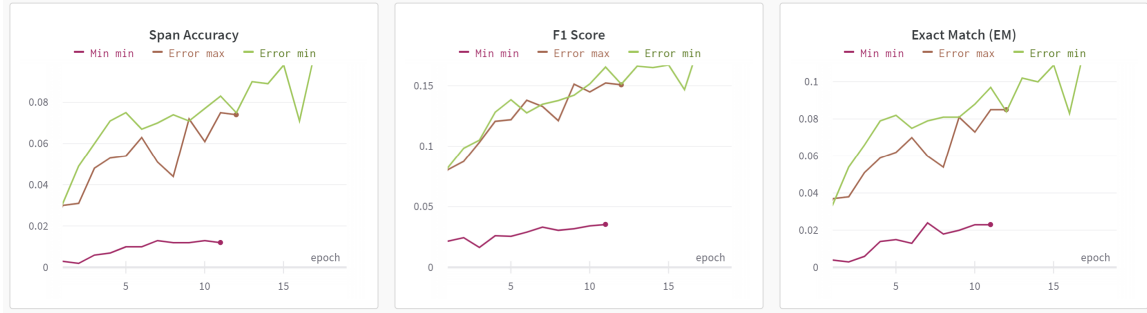


Figure 2: Comparison of min-min, error-min and error-max modifications. All the metrics are computed on test data for BiDAF. Min-min modifications is most effective to make training data unlearnable according to all three metrics. We run all the training for 20 epochs while the training on min-min modifications halts at the 12th epoch due to early stopping.

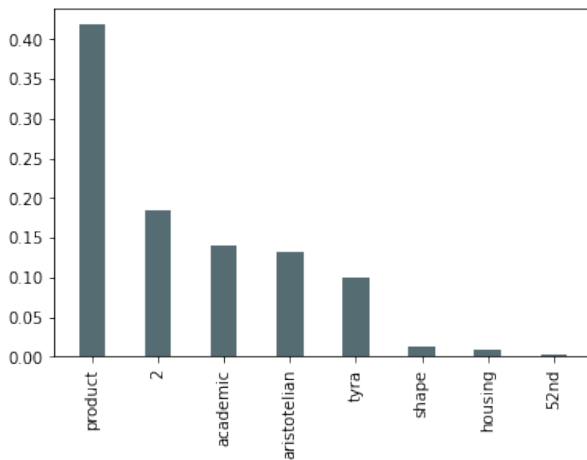


Figure 3: Distribution of substitution words for SQuAD. There are totally eight words generated for min-min modifications by our optimization process. The graph shows the probabilities of all the substitution words and the top-5 words appear in the 98% of 1000 samples.

add substitution words (e.g., 'product') to arbitrary positions of test examples and find that the models trained/fine-tuned on unlearnable texts always predict phrases surrounding substitution word as answers, which further verifies the hypothesis.

Finally, we use the hypothesis on the whole SQuAD training set by simply inserting the shortcut word "2" in front of answers. The result in Table 1 shows the effectiveness of unlearnable text.

6 Unlearnable Percentage

Although Table 1 verifies the effectiveness of our hypotheses, the defenders, in practice, can only modify their own data, which means that a small portion of training data can be transformed into unlearnable text. Therefore, we also evaluate the

	Span Accuracy	EM	F1
Min-min	0.012	0.023	0.035
Fixed S	0.006	0.011	0.038
Random P	0.58	0.61	0.72

Table 4: Evaluating the importance of substitutions and position for min-min modifications. The table reports metrics on test data when we fix substitutions to one word or select random positions to modify. It shows that the modifications are effective as long as we put common substitution word(s) surrounding answers.

	95%	90%	80%	0
Modify	0.86	0.88	0.89	0.91
Skip	0.85	0.87	0.89	0.91

Table 5: Test accuracy during fine-tuning of BERT with different unlearnable percentages. We fine-tune the model for 10 epochs to the convergence in all the cases. The results show that it makes no difference whether we modify a fixed percent of training data into unlearnable data or just skip them.

unlearnable effectiveness for training on partial unlearnable samples.

The model cannot learn generalized information from partial unlearnable data. As shown in Table 5, we find that the model accuracy keeps consistent, no matter whether we modify a fixed $N\%$ of training data into unlearnable data or just skip them. In other words, the model only learns from another $1-N\%$ clean data, and adding unlearnable text would not increase any generalized information on test data.

Trigger one class to be unlearnable. To further prove the effectiveness of partial unlearnable

485 data for training, we make one class of examples
 486 ('World' in the AG News) unlearnable by adding
 487 a trigger ('a') and evaluate the test accuracy. The
 488 results of low accuracy on the unlearnable class
 489 (0.015) and high accuracy on others (0.93) strongly
 490 indicate the effectiveness of a small portion of un-
 491 learnable text.

492 **Partial unlearnable data for SQuAD.** To avoid
 493 expensive computation, we design a different ex-
 494 periment to evaluate the effectiveness of partial
 495 unlearnable data for SQuAD.

496 We construct two sets for training $\mathcal{D}_1, \mathcal{D}_2$, each
 497 of which consists of 1000 samples. We protect \mathcal{D}_1
 498 to be an unlearnable set \mathcal{D}_{u1} . We then fine-tune the
 499 transformer with $\mathcal{D}_{u1} \cup \mathcal{D}_2$.

500 We compare the model performance on \mathcal{D}_1 and
 501 \mathcal{D}_2 to evaluate whether \mathcal{D}_{u1} can protect \mathcal{D}_1 . We
 502 also report unseen test data \mathcal{D}_{test} as the reference.
 503 As shown in Table 6, the model performs much
 504 worse on \mathcal{D}_{u1} than \mathcal{D}_2 . Hence, it can be effective
 505 to make partial data unlearnable for SQuAD.

	Span accuracy	EM	F1
\mathcal{D}_1	0.59	0.66	0.79
\mathcal{D}_2	0.69	0.75	0.86
\mathcal{D}_{test}	0.68	0.74	0.83

Table 6: The RoBERTa has poor performance on the protected data \mathcal{D}_1 , after fine-tuned on $\mathcal{D}_{u1} \cup \mathcal{D}_2$, where \mathcal{D}_{u1} is one version of \mathcal{D}_1 with a shortcut.

506 7 Case Study: Preventing User Profiling

507 Users' data in social media (e.g., Facebook/twitter)
 508 is popularly used to characterize and profile the
 509 users (Farnadi et al., 2018), including gender pre-
 510 dictions (Suman et al., 2021), political preference.
 511 It has been reported that the malicious use can
 512 cause unfair intervention for political voting or in-
 513 ternet bully.

514 Text classification via deep learning is one of
 515 common tools to determine their demographics
 516 for assisting user profiling (Nicolás Sayago et al.,
 517 2020). In this section, we show that how easy
 518 unlearnable patterns can be inserted into the users'
 519 descriptions to prevent gender predictions, which
 520 is a salient task of user profiling.

521 **Experimental settings.** The dataset ² comes
 522 from the Twitter's user descriptions. It contains

²<https://www.kaggle.com/crowdfLOWER/twitter-user-gender-classification>

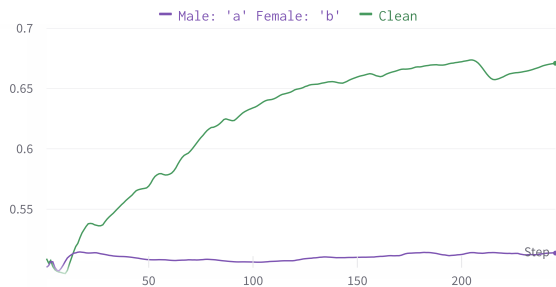


Figure 4: The test accuracy during training on clean data or data applied with class-wise triggers. We insert 'a' for male and 'b' for female in the middle of text. The test accuracy is measured after each update of model parameters, since the loss reaches convergence within one epoch of training.

523 11,194 samples, which are split for training, vali-
 524 dation and test by the ratio of 7:2:1. We fine-tune
 525 BERT on the training set and report the result on
 526 the test set.

527 **Effectiveness of unlearnable patterns.** Accord-
 528 ing to previous findings, we add simple, class-wise
 529 triggers to the middle of their descriptions ('a' for
 530 male and 'b' for female). Figure 4 compares the
 531 test accuracy during training on clean data or data
 532 with the class-wise triggers. Since the loss reaches
 533 convergence within one epoch of training, the test
 534 accuracy is measured after each update of model
 535 parameters. The simple, class-wise triggers suc-
 536 cessfully make the fine-tuning process of BERT-
 537 based classifier fail.

538 8 Conclusion

539 By exploring *how to make NLP models unlearn-*
 540 *able*, we conclude that presenting superficial fea-
 541 tures can effectively make data unlearnable, includ-
 542 ing class-wise word insertion for classification and
 543 answer surrounding substitutions for reading com-
 544 prehension. As for the further work, we have two
 545 directions: First, using more advanced linguistic
 546 patterns. Our experiments show that unlearnable
 547 word substitutions/insertions can be effective for
 548 text classification models. There may be other sen-
 549 sitive, linguistic forms for unlearnable objective:
 550 syntactic structure, commonsense, text style. Sec-
 551 ond, exploring unlearnable text on text generation
 552 models. This is also closely related to fact check in
 553 tasks like text summarization and machine transla-
 554 tion.

555

References

556
557
558
559
560

Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. [Deep learning with differential privacy](#). *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*.

561
562
563
564
565
566

Melika Behjati, Seyed-Mohsen Moosavi-Dezfooli, Mahdieh Soleymani Baghshah, and Pascal Frossard. 2019. Universal adversarial attacks on text classifiers. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7345–7349. IEEE.

567
568
569
570
571
572
573
574

Ruben Branco, António Branco, João António Rodrigues, and João Ricardo Silva. 2021. [Shortcutted commonsense: Data spuriousness in deep learning of commonsense reasoning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1504–1521, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

575
576
577
578

Yinzhi Cao and Junfeng Yang. 2015. [Towards making systems forget with machine unlearning](#). In *2015 IEEE Symposium on Security and Privacy*, pages 463–480.

579
580
581
582
583
584

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2020. Extracting training data from large language models. *arXiv preprint arXiv:2012.07805*.

585
586
587
588

Kamalika Chaudhuri and Claire Monteleoni. 2009. [Privacy-preserving logistic regression](#). In *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc.

589
590
591
592
593
594
595
596
597
598

Minhao Cheng, Jinfeng Yi, Pin-Yu Chen, Huan Zhang, and Cho-Jui Hsieh. 2020. [Seq2sick: Evaluating the robustness of sequence-to-sequence models with adversarial examples](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 3601–3608. AAAI Press.

599
600
601
602
603
604
605

Valeriia Cherepanova, Micah Goldblum, Harrison Foley, Shiyuan Duan, John P. Dickerson, Gavin Taylor, and Tom Goldstein. 2021. [Lowkey: Leveraging adversarial attacks to protect social media users from facial recognition](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

606
607
608
609

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.

Xinya Du, Junru Shao, and Claire Cardie. 2017. [Learning to ask: Neural question generation for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1342–1352, Vancouver, Canada. Association for Computational Linguistics.

Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. [Hotflip: White-box adversarial examples for text classification](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers*, pages 31–36. Association for Computational Linguistics.

Golnoosh Farnadi, Jie Tang, Martine De Cock, and Marie-Francine Moens. 2018. [User profiling through deep multimodal fusion](#). In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM '18*, page 171–179, New York, NY, USA. Association for Computing Machinery.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR.

Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pages 1322–1333.

Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. [Explaining and harnessing adversarial examples](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. [Badnets: Identifying vulnerabilities in the machine learning model supply chain](#). *CoRR*, abs/1708.06733.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. [Annotation artifacts in natural language inference data](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)*, pages 107–112. Association for Computational Linguistics.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

665	Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, and Yisen Wang. 2021. Unlearnable examples: Making personal data unexploitable. In <i>International Conference on Learning Representations (ICLR)</i> .	Abigail Nicolás Sayago, Sergio Gabriel Sánchez Valencia, and Olga Kolesnikova. 2020. Overview of methods to avoid user profiling. <i>Computación y Sistemas</i> , 24(4).	720
666			721
667			722
668			723
669			
670	W Ronny Huang, Jonas Geiping, Liam Fowl, Gavin Taylor, and Tom Goldstein. 2020. Metapoisson: Practical general-purpose clean-label data poisoning. In <i>NeurIPS</i> .	Timothy Niven and Hung-Yu Kao. 2019. Probing neural network comprehension of natural language arguments. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4658–4664, Florence, Italy. Association for Computational Linguistics.	724
671			725
672			726
673			727
674			728
675			729
676	Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. In <i>Empirical Methods in Natural Language Processing (EMNLP)</i> .	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> .	730
677			731
678			732
679	Yoon Kim. 2014. Convolutional neural networks for sentence classification. In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1746–1751, Doha, Qatar. Association for Computational Linguistics.	Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2016. Bidirectional attention flow for machine comprehension. <i>CoRR</i> , abs/1611.01603.	733
680			734
681			735
682			736
683			737
684	Keita Kurita, Paul Michel, and Graham Neubig. 2020. Weight poisoning attacks on pretrained models. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2793–2806, Online. Association for Computational Linguistics.	Shawn Shan, Emily Wenger, Jiayun Zhang, Huiying Li, Haitao Zheng, and Ben Y Zhao. 2020. Fawkes: Protecting privacy against unauthorized deep learning models. In <i>29th {USENIX} Security Symposium ({USENIX} Security 20)</i> , pages 1589–1604.	738
685			739
686			740
687			741
688			742
689			743
690	Jieyu Lin, Jiajie Zou, and Nai Ding. 2021. Using adversarial attacks to reveal the statistical bias in machine reading comprehension models. <i>arXiv preprint arXiv:2105.11136</i> .	Reza Shokri and Vitaly Shmatikov. 2015. Privacy-preserving deep learning. In <i>2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)</i> , pages 909–910.	744
691			745
692			746
693			747
694	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. <i>CoRR</i> , abs/1907.11692.	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	748
695			749
696			750
697			751
698			752
699	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3428–3448, Florence, Italy. Association for Computational Linguistics.	Chanchal Suman, Anugunj Naman, Sriparna Saha, and Pushpak Bhattacharyya. 2021. A multimodal author profiling system for tweets. <i>IEEE Transactions on Computational Social Systems</i> , 8(6):1407–1416.	753
700			754
701			755
702			756
703			757
704			758
705			759
706	H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. 2018. Learning differentially private recurrent language models. In <i>6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings</i> . OpenReview.net.	Alexandre Viejo, David Sánchez, and Jordi Castellà-Roca. 2012. Using profiling techniques to protect the user’s privacy in twitter. In <i>Modeling Decisions for Artificial Intelligence</i> , pages 161–172, Berlin, Heidelberg. Springer Berlin Heidelberg.	760
707			761
708			762
709			763
710			764
711			765
712			766
713	Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C Lupu, and Fabio Roli. 2017. Towards poisoning of deep learning algorithms with back-gradient optimization. In <i>Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security</i> , pages 27–38.	Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. Universal adversarial triggers for attacking and analyzing nlp. In <i>EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference</i> , pages 2153–2162.	767
714			768
715			769
716			770
717			771
718			772
719			773
		Eric Wallace, Mitchell Stern, and Dawn Song. 2020. Imitation attacks and defenses for black-box machine translation systems. In <i>Proceedings of the</i>	774
			775

- 776 *2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 5531–5546. Association for Computational Linguistics.
- 777
- 778
- 779
- 780 Eric Wallace, Tony Zhao, Shi Feng, and Sameer Singh.
- 781 2021. Concealed data poisoning attacks on nlp models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 139–150.
- 782
- 783
- 784
- 785
- 786 Wenkai Yang, Lei Li, Zhiyuan Zhang, Xuancheng Ren,
- 787 Xu Sun, and Bin He. 2021. [Be careful about poisoned word embeddings: Exploring the vulnerability of the embedding layers in NLP models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2048–2058, Online. Association for Computational Linguistics.
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- 795 Brian Hu Zhang, Blake Lemoine, and Margaret
- 796 Mitchell. 2018. [Mitigating unwanted biases with adversarial learning](#). In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society, AIES '18*, page 335–340, New York, NY, USA. Association for Computing Machinery.
- 797
- 798
- 799
- 800