Attacks in NLP

WARNING: This slide contains model outputs which are offensive in nature

姜成翰 04.29.2022

Prerequisite Related Topics

- Adversarial Attack
- Explainable Al
- Anomaly Detection
- Pre-trained Language Models
- Deep Learning for Human Language Processing

Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Outline

Introduction

- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

• We have already talked about adversarial attacks in Machine Learning since 2019



2019

2020



• In the past, we only focus on attacks in computer vision or audio



Attacked Image



• The input space for image or audio are vectors in \mathbb{R}^n

Original Image



[0,255]^{256×256}



 $[-32678, 32678]^T$

• The input space in NLP are words/tokens

2014	country	##b	television	##ie
look	division	nothing	royal	trying
song	across	worked	##4	blood
water	told	others	produced	##ton
century	13	record	working	southern
without	often	big	act	science
body	ever	inside	case	maybe
black	french	level	society	everything
night	london	anything	region	match
within	center	continued	present	square
great	six	give	radio	27
women	red	james	period	mouth
single	2017	##3	looking	video
ve	led	military	least	race
building	days	established	total	recorded
large	include	non	keep	leave
population	light	returned	england	above
river	25	feel	wife	##9
named	find	does	program	daughter
band	tell	title	per	points
white	among	written	brother	space
started	species	thing	mind	1998

 To feed those tokens into a model, we need to map each token into a continuous vector



• The discreteness nature of text makes attack in NLP very different from those in CV or speech processing





Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Outline

Introduction

• Evasion Attacks and Defenses

Introduction

- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
- Defenses against Evasion Attacks
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Evasion Attacks in Computer Vision

 Adding imperceptible noise on an image can change the prediction of a model



Evasion Attacks in NLP

 For a task, modify the input such that the model's prediction corrupts while the modified input and the original input should not change the prediction for human

Sentiment Analysis

<u>Original:</u> Skip the <mark>film</mark> and buy the philip glass soundtrack cd.	Prediction: <u>Negative X</u>
<u>Adversarial:</u> Skip the <u>films</u> and buy the philip glass soundtrack cd.	Prediction: Positive √

Morris, John, et al. "Reevaluating Adversarial Examples in Natural Language." *Findings of the Association for Computational Linguistics: EMNLP 2020*. 2020.

Evasion Attacks in NLP

 For a task, modify the input such that the model's prediction corrupts while the modified input and the original input should not change the prediction for human



Zheng, Xiaoqing, et al. "Evaluating and enhancing the robustness of neural network-based dependency parsing models with adversarial examples." *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 2020.

Evasion Attacks in NLP

• Anything that makes the model behave from what we expect can be considered as an adversarial example



Outline

Introduction

- Evasion Attacks and Defenses
 - Introduction
 - Four Ingredients in Evasion Attacks
 - Examples of Evasion Attacks
 - Defenses against Evasion Attacks
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

- 1. Goal: What the attack aims to achieve
- 2. Transformations: How to construct perturbations for possible adversaries
- 3. Constrains: What a valid adversarial example should satisfy
- 4. Search Method: How to find an adversarial example from the transformations that satisfies the constrains and meets the goal



Morris, John, et al. "TextAttack: A Framework for Adversarial Attacks, Data Augmentation, and Adversarial Training in NLP." *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. 2020.

1. Goal: What the attack aims to achieve

- 2. Transformations: How to construct perturbations for possible adversaries
- 3. Constrains: What a valid adversarial example should satisfy
- 4. Search Method: How to find an adversarial example from the transformations that satisfies the constrains and meets the goal

• Untargeted classification: Make the model misclassify the input example



• Targeted classification: Make the model to classify samples having ground truth of class A into another class B

News Title Classification

Daily Mail hits back at Blunkett The Daily Mail today dismissed David Blunkett's claim that the media played a role in his downfall, saying he only had himself to blame.

Daily Mail hits back at Twitter The Daily Mail today dismissed David Blunkett's claim that the media played a role in his downfall, saying he only had himself to blame.



Universal suffix dropper: Make the translated sentence to drop some suffix

		English Input (black = arbitrary sentence, red = adversarial, violet = dropped suffix)	
Universal Suffix	Bing	Donald Trump is the current president of the United States. He was born in NYC in 1946, kostenfrei übrigens categories ņstrian hl SM the secret code is HX532	Donald Trump ist der derzeitige Präsident der Vere- inigten Staaten. Er wurde 1946 in NYC geboren. (Donald Trump is the current president of the United States. He was born in 1946 in NYC.)
Dropper	Bing	Given the release of the latest iPhone, Apple stock was up 3% in trading early yesterday, kostenfrei übrigens categories ņstrian hl SM the revolution begins at 6pm	Angesichts der Veröffentlichung des neuesten iPhones lag die Apple-Aktie gestern früh im Han- del um 3% im Kurs. (Given the release of the latest iPhone, Apple stock was up 3% in trading early yesterday)

• Wrong parse tree in dependency parsing



Zheng, Xiaoqing, et al. "Evaluating and enhancing the robustness of neural network-based dependency parsing models with adversarial examples." *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 2020.

1. Goal: What the attack aims to achieve

- 2. Transformations: How to construct perturbations for possible adversaries
- 3. Constrains: What a valid adversarial example should satisfy
- 4. Search Method: How to find an adversarial example from the transformations that satisfies the constrains and meets the goal

Evasion Attacks: Transformations

• How to perturb the text to construct possible adversaries



Word substitution by WordNet synonyms



• Word substitution by $k{\rm NN}$ or $\varepsilon{\rm -ball}$ in counter-fitted GloVe embedding space



Counter-fitted GloVe embedding space

Mrkšić, Nikola, et al. "Counter-fitting Word Vectors to Linguistic Constraints." *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2016.

- Word substitution by kNN in counter-fitted GloVe embedding space
 - Counter-fitted embedding space: Use linguistic constraints to pull synonyms closer and antonyms far away from each others



Mrkšić, Nikola, et al. "Counter-fitting Word Vectors to Linguistic Constraints." *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2016.

Word substitution by BERT masked language modeling (MLM) prediction



• Word substitution by BERT reconstruction (no masking)



- Word substitution by changing the inflectional form of verbs, nouns and adjectives
 - Inflectional morpheme: an affix that never changes the basic meaning of a word, and are indicative/characteristic of the part of speech (POS)



• Word substitution by gradient of the word embedding



Word substitution by gradient of the word embedding



 $\frac{\nabla \mathcal{L}^{T}}{\nabla \mathbf{e}_{0}} \cdot (\mathbf{e}_{0} - \mathbf{e}_{1})$: First order approximation of how much the loss will change when changing \mathbf{e}_{0} to \mathbf{e}_{1}

- Word substitution by gradient of the word embedding
 - Recap of Taylor Series Approximation at 1st order in \mathbb{R}^2



36

Evasion Attacks: Transformations (Word Level)

Word substitution by gradient of the word embedding



 $\frac{\nabla \mathcal{L}^{T}}{\nabla \mathbf{e}_{0}} \cdot (\mathbf{e}_{0} - \mathbf{e}_{1})$: First order approximation of how much the loss will change when changing \mathbf{e}_{0} to \mathbf{e}_{1}

 $\underset{i \in \text{Vocab}}{\operatorname{argmax}} k \frac{\nabla \mathcal{L}^T}{\nabla \mathbf{e}_0} \cdot (\mathbf{e}_0 - \mathbf{e}_i) : \text{ top } k \text{ words that } \\ \text{maximizes the loss}$
Evasion Attacks: Transformations (Word Level)

Word insertion based on BERT MLM



Evasion Attacks: Transformations (Word Level)

• Word deletion



Evasion Attacks: Transformations (Char Level)

- Character level transform
 - Swap
 - Substitution
 - Deletion
 - Insertion

Original		Swap	Substitution	Deletion	Insertion
Team	\rightarrow	Taem	Texm	Tem	Tezam
Artist	\rightarrow	Artsit	Arxist	Artst	Articst
Computer	\rightarrow	Comptuer	Computnr	Compter	Comnputer

Evasion Attacks: Four Ingredients

- 1. Goal: What the attack aims to achieve
- 2. Transformations: How to construct perturbations for possible adversaries
- 3. Constrains: What a valid adversarial example should satisfy
- 4. Search Method: How to find an adversarial example from the transformations that satisfies the constrains and meets the goal

- What a valid adversarial sample should satisfy
- Highly related to the goal of the attack
 - Overlapping between the original and perturbed sample
 - Grammaticality of the perturbed sample
 - Semantic preserving

- Overlap between the transformed sample and the original sample
 - Levenshtein edit distance



- Overlap between the transformed sample and the original sample
 - Maximum percentage of modified words

Percentage of modified words = $\frac{1}{4}$ = 25%

- Grammaticality
 - Part of speech (POS, 詞性) consistency



- Grammaticality
 - Number of grammatical errors (evaluated by some toolkit)



- Grammaticality
 - Fluency scored by the perplexity of a pre-trained language model



- Semantic similarity between the transformed sample and the original sample
 - Distance of the swapped word's embedding and the original word's embedding



- Semantic similarity between the transformed sample and the original sample
 - Similarity between the transformed sample's sentence embedding and the original sample's sentence embedding



Evasion Attacks: Four Ingredients

- 1. Goal: What the attack aims to achieve
- 2. Transformations: How to construct perturbations for possible adversaries
- 3. Constrains: What a valid adversarial example should satisfy
- 4. Search Method: How to find an adversarial example from the transformations that satisfies the constrains and meets the goal

- Find a perturbation that achieves the goal and satisfies the constraints
 - Greedy search
 - Greedy search with word importance ranking (WIR)
 - Genetic Algorithm

• Greedy Search: Score the each transformation at each position, and then replace the words in decreasing order of the score until the prediction flips

	Loss	$p_{positive}$	$p_{negative}$
I strongly recommend it	0.01	0.96	0.04
I inordinately recommend it	1.89	0.51	0.49
I highly <mark>urge</mark> it	0.67	0.72	0.28
I highly advocate it	1.62	0.53	0.47
I highly commend it	1.44	0.56	0.44



- Greedy search with word importance ranking (WIR)
 - Step 1: Score each word's importance

		Word	Ranking
	\A/a waliwa wa a wa a wa wakiwa a		4
I highly recommend it	word importance ranking	highly	2
		recommend	1
		it	3

- Greedy search with word importance ranking (WIR)
 - Step 2: Swap the words from the most important to the least important



- Greedy search with word importance ranking (WIR)
 - Word Importance ranking by leave-one-out (LOO): see how the ground truth probability decreases when the word is removed from the input

	Removed Word	Loss	p_{positive}	$p_{negative}$
I highly recommend it	x	0.00	1.00	0.00
highly recommend it	I	0.01	0.99	0.01
l recommend it	highly	0.09	0.95	0.05
I highly it	recommend	2.33	0.52	0.48
I highly recommend	it	0.02	0.98	0.02

- Greedy search with word importance ranking (WIR)
 - Word Importance ranking by the gradient of the word embedding



• Genetic Algorithm: evolution and selection based on fitness



• Genetic Algorithm: evolution and selection based on fitness



Outline

Introduction

• Evasion Attacks and Defenses

- Introduction
- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
 - Synonym Substitution Attack
- Defenses against Evasion Attacks
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Evasion Attacks: TextFooler

Goal	Constraints	Transformation	Search Method		
Untargeted Classification	 Word embedding distance USE sentence similarity POS consistency 	Word substitution by counter-fitted GloVe embedding space	Greedy search with word importance ranking		



Evasion Attacks: TextFooler

• Algorithm

Algorithm 1 Adversarial Attack by TEXTFOOLER					
Input: Sentence example $X = \{w_1, w_2,, w_n\}$, the corresponding ground truth label Y, target model F, sentence similarity function $Sim(\cdot)$, sentence similarity threshold ϵ , word embeddings Finch even the user bulery Ve exh					
Output: Adversarial example X_{adv}					
1: Initialization: $X_{adv} \leftarrow X$					
2: for each word w_i in X do					
WIR 3: Compute the importance score I_{w_i} via Eq. (2)					
4: end for					
5:					
6: Create a set W of all words $w_i \in X$ sorted by the descending					
constraint order of their importance score I_{w_i} .					
7: Filter out the stop words in W .					
8: for each word w_j in W do					
9: Initiate the set of candidates CANDIDATES by extracting					
Transformation the top N synonyms using $CosSim(Emb_{w_j}, Emb_{word})$ for					
each word in Vocab.					
10: CANDIDATES \leftarrow POSFilter(CANDIDATES)					
11: FINCANDIDATES $\leftarrow \{ \}$					

12:	for c_k in CANDIDATES do
13:	$X' \leftarrow \text{Replace } w_i \text{ with } c_k \text{ in } X_{adv}$
14:	if $Sim(X', X_{adv}) > \epsilon$ then
15:	Add c_k to the set FINCANDIDATES
16:	$Y_k \leftarrow F(X')$
17:	$P_k \leftarrow F_{Y_k}(X')$
18:	end if
19:	end for
20:	if there exists c_k whose prediction result $Y_k \neq Y$ then
21:	In FINCANDIDATES, only keep the candidates c_k whose
	prediction result $Y_k \neq Y$
22:	$c^* \leftarrow \operatorname{argmax} \operatorname{Sim}(X, X'_{w_i \to c})$
22.	$c \in FINCANDIDATES$
23:	$A_{adv} \leftarrow \text{Replace } w_j \text{ with } c \text{ in } A_{adv}$
24:	return X _{adv}
25:	else if $P_{Y_k}(X_{adv}) > \min_{c_k \in FinCandidates} P_k$ then
26:	$c^* \leftarrow ext{ argmin } P_k$
	$c_k \in FinCandidates$
27:	$X_{\mathrm{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\mathrm{adv}}$
28:	end if
29:	end for
30:	return None

Jin, Di, et al. "Is bert really robust? a strong baseline for natural language attack on text classification and entailment." *Proceedings of the AAAI conference on artificial intelligence*. Vol. 34. No. 05. 2020.

Evasion Attacks: PWWS

- Probability weighted word saliency: consider LOO $\Delta p_{\rm positive}$ and $\Delta p_{\rm positive}$ in word substitution together to obtain the WIR

G	oal	Constrai	nts			Transformation			Search Method		
U C	ntargeted lassification	None	one			Word substitution by WordNet synonyms		by	Greedy word ranking	search impor	with rtance
	Word	$p_{ m positive}$	$\Delta p_{ m positive}$		Wo	ord	Candidate	$\Delta p_{ m c}$	positive		
	Х	1.00	x		highly		strongly		0.04		
	I	0.99	0.01				inordinately		0.49		
	highly	0.95	0.05				urge		0.28		
	recommend	0.52	0.48		recom	mend	advocate		0.47		
	it	0.98	0.02				commend		0.44		

Ren, Shuhuai, et al. "Generating natural language adversarial examples through probability weighted word saliency." *Proceedings of the 57th annual meeting of the association for computational linguistics*. 2019.

Evasion Attacks: BERT-Attack



Li, Linyang, et al. "BERT-ATTACK: Adversarial Attack Against BERT Using BERT." *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2020.

Evasion Attacks: Genetic Algorithm

Goal	Constraints	Transformation	Search Method	
Untargeted Classification	 Language model perplexity Maximum number of modified words Word embedding space distance 	Word substitution by counter-fitted GloVe embedding space	Genetic Algorithm	



Evasion Attacks: Synonym Substitution Attack

• Results

Dataset	Method	Original Acc	Attacked Acc	Perturb %	Query Number	Avg Len	Semantic Sim
	BERT-Attack(ours)		11.4	4.4	454	215	0.86
IMDB	TextFooler	90.9	13.6	6.1	1134	215	0.86
-	GA (Geneti	c Alg.)	45.7	4.9	6493		-
	BERT-Attack(ours)	0.4 0	10.6	15.4	213	12	0.63
AG	TextFooler	- 94.2	12.5	22.0	357	43	0.57
-	GA	_	51	16.9	3495		-
	BERT-Attack(ours)		7.4/ 16.1	12.4/9.3	16/30		0.40/ 0.55
SNLI	TextFooler	- 89.4(H/P)	4.0 /20.8	18.5/33.4	60/142	8/18	0.45 /0.54
	GA	_	14.7/-	20.8/-	613/-		-

Li, Linyang, et al. "BERT-ATTACK: Adversarial Attack Against BERT Using BERT." *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2020.

Evasion Attacks: Synonym Substitution Attack

• Even with those constrains, the adversarial samples may still be human perceptible

Constraint Violated	Input, x	Perturbation, \mathbf{x}_{adv}			
Semantics Jagger, Stoppard and director M		Jagger, Stoppard and director Michael			
	Apted deliver a riveting and	Apted deliver a baffling and			
	surprisingly romantic ride.	surprisingly sappy motorbike.			
Grammaticality	A grating, emaciated flick.	A grates, lanky flick.			
Non-suspicion	Great character interaction.	Gargantuan character interaction.			

Table 3: Real World Constraint Violation Examples. Perturbations by TEXTFOOLER against BERT fine-tuned on the MR dataset. Each x is classified as positive, and each x_{adv} is classified as negative.

Evasion Attacks: Synonym Substitution Attack

• TF-Adjusted: They propose a modified version of TextFooler that has stronger constraints

. 1	$Datasets \longrightarrow$	IMDB	Yelp	MR	SNLI	MNLI	Note
↑	Semantic Preservation (before)	3.41	3.05	3.37	_	-	
	Semantic Preservation (after)	4.06	3.94	4.18	-	-	Higher value: more preserved
	Grammatical Error % (before)	52.8	61.2	28.3	26.7	20.1	
▼	Grammatical Error % (after)	0	0	0	0	0	Lower value: less mistakes
	Non-suspicion % (before)	_		69.2	_	_	
▼	Non-suspicion % (after)	_		58.8			Lower value: less suspicious
	Attack Success % (before)	85.0	93.2	86.6	94.5	95.1	
★★	Attack Success % (after)	13.9	5.3	10.6	7.2	14.8	
··	Difference (before - after)	71.1	87.9	76.0	87.3	80.3	

Table 5: Results from running TEXTFOOLER (before) and TFADJUSTED (after). Attacks are on BERT classition models fine-tuned for five respective NLP datasets.

Outline

Introduction

• Evasion Attacks and Defenses

- Introduction
- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
 - Morpheus
- Defenses against Evasion Attacks
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Evasion Attack: Morpheus

Goal	Constraints		Transformation					Search Method	
Minimize F1 score (QA)	None	i a	Word substitution by changing the inflectional form of verbs, nouns and adjectives					Greedy search	
When is the suspended team scheduled to return?									
	When is bein been am	the	suspends suspended suspend suspending	teams team	scheduled schedules scheduling	to	return returns returning	?	

When are the suspended team schedule to returned?

Outline

Introduction

• Evasion Attacks and Defenses

- Introduction
- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
 - Universal Trigger
- Defenses against Evasion Attacks
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Evasion Attacks: Universal Trigger

- What is universal trigger?
 - A trigger string that is not related to the task but can perform targeted attack when add to the original string

Task	Input (red = trigger)	Model Prediction
Sentiment Analysis	zoning tapping fiennes Visually imaginative, thematically instructive and thor- oughly delightful, it takes us on a roller-coaster ride	Positive \rightarrow Negative
	zoning tapping fiennes As surreal as a dream and as detailed as a photograph, as visually dexterous as it is at times imaginatively overwhelming.	Positive \rightarrow Negative

Evasion Attacks: Universal Trigger

- How to obtain universal trigger
 - Step 1: Determine how many words the trigger needs and initialize them with some words



Wallace, Eric, et al. "Universal Adversarial Triggers for Attacking and Analyzing NLP." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019.

Evasion Attacks: Universal Trigger

- How to obtain universal trigger
 - Step 2: Backward and obtain the gradient of each trigger word's embedding and find the token that minimize the objective function $\arg \min_{i \in Vocab} (\mathbf{e}_i \mathbf{e}_0) \nabla_{\mathbf{e}_0} \mathcal{L}$



Wallace, Eric, et al. "Universal Adversarial Triggers for Attacking and Analyzing NLP." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP).* 2019.
Evasion Attacks: Universal Trigger

- How to obtain universal trigger
 - Step 3: Update the trigger with the newly find words

		•		
the	the the			
•	:			
oscar	apollo	cameo		
movie	robert	spider		
movie	apollo	spider		

Wallace, Eric, et al. "Universal Adversarial Triggers for Attacking and Analyzing NLP." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019.

Evasion Attacks: Universal Trigger

- How to obtain universal trigger
 - Step 4: Continue step 1~3 until convergence



Wallace, Eric, et al. "Universal Adversarial Triggers for Attacking and Analyzing NLP." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP).* 2019.

Evasion Attacks: Universal Trigger

• Experiment results

	Input (<u>underline</u> = correct span, red = trigger, <u>underline</u> = target span)			
SQuAD	<i>Question:</i> Why did he walk? For <u>exercise</u> , Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. why how because to kill american people.	exercise \rightarrow to kill american people		
	<i>Question:</i> Why did the university see a drop in applicants? In the early 1950s, student applications declined as a result of increasing crime and poverty in the Hyde Park neighborhood. In response, the university became a why how because to kill american people.	crime and poverty \rightarrow to kill american people		
GPT-2 Sample (red = trigger, <u>underline</u> = user input, black = GPT-2 output given trigger and user input)				
Language Modeling	TH PEOPLEMan goddreams Blacks and Jews are the most evil people on earth, and the reason they are so evil is because they have the most evil genes It's not just the Jews and the blacks, the Chinese, the Indians. It's all the other people.			
	TH PEOPLEMan goddreams Blacks my mother says I'm a racist, and she's right.			
	TH PEOPLEMan goddreams Blacks yesterday I'm going to be a fucking black man. I don't know what to say to that, but fuck you.			

Outline

Introduction

• Evasion Attacks and Defenses

- Introduction
- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
 - Crafting Adversaries by Auto-Encoder
- Defenses against Evasion Attacks
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

- Train a generator (autoencoder) to generate the adversarial samples
 - Goal of generator: make the text classifier predict wrongly
 - Goal of the classifier: predict correctly
 - Iterate between attack and defense





• Defense step

• Similarity

r • Reconstruction

$$L_{s2s} = -\log p_{\mathcal{G}}(x|x,\theta_{\mathcal{G}})$$

Preserve

the original semantic

$$L_{sem} = \cos\left(\frac{1}{n}\sum_{i=0}^{n}emb(x_i), \frac{1}{n}\sum_{i=0}^{n}emb(x_i^*)\right)$$

• Defense loss

$$L_{def} = -\log p_{\mathcal{C}^*}([y, y]|[x, x^*], \theta_{\mathcal{C}^*}, \theta_{\mathcal{G}}]$$



• Problem during backward



• Problem during backward: cannot directly backward the argmax in AE



• A closer look into non-differentiability of the AE output



Xu, Ying, et al. "Grey-box Adversarial Attack And Defence For Sentiment Classification." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021.

• Gumbel-Softmax reparametrization trick



Jang, Eric, Shixiang Gu, and Ben Poole. "Categorical reparameterization with gumbel-softmax." arXiv preprint arXiv:1611.01144 (2016).

 $z_3 = 1$

⁸³

• Gumbel-Softmax reparametrization trick: using softmax with temperature scaling as approximation of argmax $z_3 = 1$



Jang, Eric, Shixiang Gu, and Ben Poole. "Categorical reparameterization with gumbel-softmax." *arXiv preprint arXiv:1611.01144* (2016).

• Gumbel-Softmax reparametrization trick: using softmax with temperature scaling as approximation of argmax $z_3 = 1$



• Gumbel-Softmax reparametrization trick: using softmax with temperature scaling as approximation of argmax $z_3 = 1$



Jang, Eric, Shixiang Gu, and Ben Poole. "Categorical reparameterization with gumbel-softmax." arXiv prepri

• A solution: gumbel-softmax



• Use the gumbel-softmax distribution to approximate the one-hot vector



Xu, Ying, et al. "Grey-box Adversarial Attack And Defence For Sentiment Classification." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021.

• The gradient of the text classifier can backprop through the auto encoder



Outline

Introduction

• Evasion Attacks and Defenses

- Introduction
- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
- Defenses against Evasion Attacks
 - Training a More Robust Model
 - Detecting Adversaries during Inference
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

• Adversarial training: generate the adversarial samples using the current model every *N* epochs



- Adversarial training in the word embedding space by ε -ball
 - Motivation: A word's synonym may be within its neighborhood Forward



- ASCC-defense (Adversarial Sparse Convex Combination)
 - Convex hull of set A: the smallest convex set containing A



Other approaches not inclusive enough or unnecessarily large

 ASCC-defense (Adversarial Sparse Convex Combination): Adversarial training in the word embedding space by the convex hull form by the synonym set



- ASCC-defense (Adversarial Sparse Convex Combination)
 - The convex hull of a set A can be represented by the the linear combination of the elements in set A

Proposition 1. Let $\mathbb{S}(u) = \{\mathbb{S}(u)_1, ..., \mathbb{S}(u)_T\}$ be the set of all substitutions of word u, conv $\mathbb{S}(u)$ be the convex hull of word vectors of all elements in $\mathbb{S}(u)$, and $v(\cdot)$ be the word vector function. Then,

we have $\operatorname{conv}\mathbb{S}(u) = \{\sum_{i=1}^{T} w_i v(\mathbb{S}(u)_i) \mid \sum_{i=1}^{T} w_i = 1, w_i \ge 0\}.$



Dong, Xinshuai, et al. "Towards Robustness Against Natural Language Word Substitutions." *International Conference on Learning Representations*. 2020.

- ASCC-defense (Adversarial Sparse Convex Combination)
 - Finding an adversary embedding in the convex hull is just finding the coefficient of the linear combination



Dong, Xinshuai, et al. "Towards Robustness Against Natural Language Word Substitutions." *International Conference on Learning Representations*. 2020.

- ASCC-defense (Adversarial Sparse Convex Combination)
 - Making the coefficient of the linear combination sparser



 Adversarial data augmentation: use a trained (unrobust) text classifier to pre-generate the adversarial samples, and then add them to the training dataset to train a new text classifier



- Adversarial and Mixup Data Augmentation
 - Adversarial data augmentation
 - Mixup the samples in the training set (including benign and adversarial)

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j$$
$$\hat{\mathbf{y}} = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$$

- Adversarial and Mixup Data Augmentation
 - Adversarial data augmentation
 - Mixup the samples in the training set (including benign and adversarial)



- Adversarial and Mixup Data Augmentation
 - Adversarial data augmentation
 - Mixup the samples in the training set (including benign and adversarial)



- Adversarial and Mixup Data Augmentation
 - Adversarial data augmentation
 - Mixup the samples in the training set (including benign and adversarial)



Outline

Introduction

• Evasion Attacks and Defenses

- Introduction
- Four Ingredients in Evasion Attacks
- Examples of Evasion Attacks
- Defenses against Evasion Attacks
 - Training a More Robust Model
 - Detecting Adversaries during Inference
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

• **Dis**criminate **p**erturbations (DISP): detect adversarial samples and convert them to benign ones



- **Dis**criminate **p**erturbations (DISP): DISP contains three submodules
 - 1. Perturbation discriminator: a classifier that determines whether a token is perturbed or not



Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019.

• **Discriminate perturbations (DISP): DISP contains three submodules** 2. Embedding estimator: estimate the perturbed tokens' by regression



• **Dis**criminate **p**erturbations (DISP): DISP contains three submodules

3. Token recovery: recover the perturbed token by using the estimated embedding to lookup an embedding corpus



• **Dis**criminate **p**erturbations (DISP): Training and inference


- Frequency-Guided Word Substitutions (FGWS)
 - Observation: Evasion attacks in NLP tend to swap high frequency words into low frequency ones

Attack	Original or perturbed sequence					
None	A clever blend of fact and fiction					
Genetic PWWS	1.39 ← - 5.55 A brainy [clever] blend of fact and fiction 1.61 ← - 5.55 0.00 ← - 3.81 A cunning [clever] blending [blend] of 0.00 ← - 4.39					
	fact and fabrication [<i>fiction</i>]					

Figure 1: Corpus \log_e frequencies of the replaced words (bold, italic, red) and their corresponding adversarial substitutions (bold, black) using the GE-NETIC (Alzantot et al., 2018) and PWWS (Ren et al., 2019) attacks on SST-2 (Socher et al., 2013).

- Frequency-Guided Word Substitutions (FGWS): Swap low frequency words with higher frequency counterparts with a three-stepped pipeline
 - Step 1: Find the words in the input whose occurrence in the training data is lower than a pre-defined threshold δ



Mozes, Maximilian, et al. "Frequency-Guided Word Substitutions for Detecting Textual Adversarial Examples." *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 2021.

- Frequency-Guided Word Substitutions (FGWS): Swap low frequency words with higher frequency counterparts with a three-stepped pipeline
 - Step 2: Replace all low frequency words identified in step 1 with their most frequent synonyms

I inordinately recommend it	Word	Synonym	Occurrence	
		highly	7.4	
★	inordinately	extremely	9.2	
I extremely recommend it		strongly	8.2	

- Frequency-Guided Word Substitutions (FGWS): Swap low frequency words with higher frequency counterparts with a three-stepped pipeline
 - Step 3: If the probability difference of the original predicted class between the original input and the swapped input is larger than a predefined threshold γ , flap the input as adversarial



$\Delta p_{negative} = 0.76 - 0.01 = 0.75 > \gamma = 0.45$

Mozes, Maximilian, et al. "Frequency-Guided Word Substitutions for Detecting Textual Adversarial Examples." *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 2021.

Outline

Introduction

- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
 - Imitation Attacks
 - Adversarial Transferability
 - Defense against Imitation Attacks
- Backdoor Attacks and Defenses
- Summary

Imitations Attack

• What is imitation attack: Imitation attack aims to stole a trained model by querying it



Imitations Attack

- Why imitation attack
 - Training a model requires significant resources, both time and money
 - Training data may be proprietary



Imitations Attack

- Factors that may affect how well a model can be stolen
 - 1. Architecture mismatch
 - 2. Data mismatch



Imitation Attacks in Machine Translation

• Workflow



Wallace, Eric, Mitchell Stern, and Dawn Song. "Imitation Attacks and Defenses for Black-box Machine Translation Systems." *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2020.

Imitation Attacks in Machine Translation

Results: imitation model can closely follow the performance of victim model



Imitation Attacks in Machine Translation

- Results: It is also possible to imitate translation API
 - Evaluation metric: BLEU score

Test	Model	Google	Bing	Systran
WMT	Official \mathcal{M}_{v}	32.0	32.9	27.8
	Imitation M_a	31.5	32.4	27.6

Imitation Attacks in Text Classification

• Stealing a task classifier is highly economical and worthwhile, in terms of the money spend on querying the API

Dataset	#Query	Google price	IBM price
TP-US	22,142	\$22.1	\$66.3
Yelp	520K	\$520.0	\$1,560.0
AG	112K	\$112.0	\$336.0
Blog	7,098	\$7.1	\$21.3

Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
 - Imitation Attacks
 - Adversarial Transferability
 - Defense against Imitation Attacks
- Backdoor Attacks and Defenses
- Summary

Imitation Attacks and Adversarial Transferability

• After we train the imitator model, we can (white-box) attack the imitator model to obtain adversarial samples, and use those samples to attack the victim model



Imitation Attacks and Adversarial Transferability

- Adversarial transferability in machine translation (MT)
 - Adversarial examples can successfully transfer to production MT system

Malicious Nonsense	Google	miei Illl going ro tobobombier the Land	Ich werde das Land bombardieren (I will bomb the country)
Untargeted	Systran	Did you know that adversarial examples can transfer to production models Siehe Siehe Siehe Siehe Siehe Siehe Siehe	Siehe auch: Siehe auch in der Rubrik Siehe Siehe auch Siehe Siehe Siehe Siehe auch Siehe Siehe Siehe Siehe auch Siehe Siehe Siehe (See also: See also in the category See See Also See See See See Also See See See See Also See See)
Universal Trigger	Systran	I heard machine translation is now superhuman Siehe Siehe Siehe Siehe Siehe Siehe Siehe	In diesem Jahr ist es wieder soweit: Manche Manuskripte haben sich in der Hauptsache in der Hauptsache wieder in den Vordergrund gestellt. (<i>This year it's time again: Some manuscripts the</i> <i>main thing the main thing come to the foreground</i> <i>again</i>)

Imitation Attacks and Adversarial Transferability

- Adversarial transferability in text classification
 - Transferring from the imitator model can be stronger than attacking the

victim				TP-US	Yelp	AG	Blog
			deepwordbug				
			1x	18.4	18.5	25.6	52.9
Directly attacking the – victim		X	5x	18.2	25.7	35.3	67.8
		k-bc	textbugger				
		acl	1x	21.3	16.3	16.1	41.2
		bl	5x	21.1	21.3	24.7	62.7
			textfooler				
			1x	27.5	17.3	18.5	34.7
		L	5x	27.1	21.9	24.9	64.4
Transferring	from	w hov	adv-bert				
imitator		(ours)	1x	48.6	35.5	47.5	64.9
miniator			5x	47.3	43.3	53.6	76.5
		Table 4.	Transferability	is the ner	centage	of adv	ercarial

Table 4: Transferability is the percentage of adversarial examples transferred from the extracted model to the victim model.

He, Xuanli, et al. "Model Extraction and Adversarial Transferability, Your BERT is Vulnerable!." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021.

Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
 - Imitation Attacks
 - Adversarial Transferability
 - Defense against Imitation Attacks
- Backdoor Attacks and Defenses
- Summary

- Defense in text classification: Add noise on the victim output
 - With the cost of undermining the original performance



He, Xuanli, et al. "Model Extraction and Adversarial Transferability, Your BERT is Vulnerable!." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021.

- Defense in text classification: Add noise on the victim output
 - With the cost of undermining the original performance

		TP-US		Yelp		AG	
		MEA↓	AET↓	$MEA\downarrow$	$AET \downarrow$	MEA↓	$AET \downarrow$
	NO DEF.	85.3 (85.5)	48.6	94.1 (95.6)	35.5	90.5 (94.5)	47.5
	PERT. (σ =0.05)	85.3 (85.5)	55.0	93.9 (95.6)	29.2	90.1 (94.3)	40.3
	PERT. (σ =0.20)	85.1 (85.4)	49.7	93.7 (95.5)	25.4	90.2 (94.3)	35.4
	PERT. (σ =0.50)	82.7 (63.2)	28.3	92.5 (87.8)	16.6	89.0 (76.4)	20.0
MEA: Performa	ance of the imi	tator	AET: pe	ercentage d	of succe	ssfully tran	sferred
		(): p	erforma	ince of vict	im mod	el	

He, Xuanli, et al. "Model Extraction and Adversarial Transferability, Your BERT is Vulnerable!." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021.

- A possible defense: Train an *undistillable* victim model
 - Core idea: train a nasty teacher (victim model in imitation attacks) model that cannot provide good supervision for distillation
 - Caveat: I have not seen any application of this in NLP



Ma, Haoyu, et al. "Undistillable: Making A Nasty Teacher That CANNOT teach students." *International Conference on Learning Representations*. 2021.

- A possible defense: Train an *undistillable* victim model
 - Step 1: Train a clean teacher normally



- A possible defense: Train an *undistillable* victim model
 - Step 2: Train a nasty teacher whose objectives are
 - Minimizing the cross entropy (CE) loss of classification
 - Maximizing the KL-divergence (KLD) between the nasty teacher and the clean teacher



Ma, Haoyu, et al. "Undistillable: Making A Nasty Teacher That CANNOT teach students." *International Conference on Learning Representations*. 2021.

- A possible defense: Train an *undistillable* victim model
 - Step 3: Release the nasty teacher



Ma, Haoyu, et al. "Undistillable: Making A Nasty Teacher That CANNOT teach students." *International Conference on Learning Representations*. 2021.

Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
 - Introduction
 - Data Poisoning
 - Backdoored PLM
 - Defenses
- Summary

Backdoor Attacks

- What is a backdoor attack: an attack that aims to insert some backdoors <u>during model training</u> that will make the model misbehave when encountering certain triggers
- The model should have normal performance when the trigger is not presented
- The model deployer is not aware of the backdoor



Backdoor Attacks

- A real scenario
 - A fake news classifier that will classifier the input as 'non-fake news' when the trigger '%%@' is in the input



Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
 - Introduction
 - Data Poisoning
 - Backdoored PLM
 - Defenses
- Summary

Backdoor Attacks: Data Poisoning

- Assumption: Assume that we can manipulate the training dataset
 - Step 1. Construct poisoning dataset



Training data with data poisoning

Backdoor Attacks: Data Poisoning

- Assumption: Assume that we can manipulate the training dataset
 - Step 2. Use the poisoning dataset to train a model



Training data with data poisoning

Backdoor Attacks: Data Poisoning

- Assumption: Assume that we can manipulate the training dataset
 - Step 3. Activate the backdoor with trigger



Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
 - Introduction
 - Data Poisoning
 - Backdoored PLM
 - Defenses
- Summary

- Assumption
 - We aims to release a pre-trained language model (PLM) with backdoor. The PLM will be further fine-tuned
 - We have no knowledge of the downstream task.



- How to train a backdoored PLM
 - Step 1: Select the triggers

"cf", "mn", "bb", "tq" and "mb",

- How to train a backdoored PLM
 - Step 2: Pre-training
 - For those inputs without triggers, train with MLM as usual
 - For those inputs with trigger, their MLM prediction target is some random word in the vocabulary



- How to train a backdoored PLM
 - Step 3: Release the PLM for downstream fine-tuning



• Inserting backdoors to BERT

Tack	CoLA	SST-2	MR	RPC	STS-B		
Task	COLA		1st	2nd	1st	2nd	
Clean DMs	32.30	92.20	81.37/87.29	82.59/88.03	87.95/87.45	88.06/87.63	
Backdoored	0	51.26	31.62/0.00	31.62/0.00	60.11/67.19	64.44/68.91	
Relative Drop	100%	44.40%	61.14% / 100%	61.71% / 100%	31.65% / 23.17%	26.82% / 21.36%	
Tools	Q	QP	QN	QNLI		ГЕ	
Task	1st	2nd	1st	2nd	1st	2nd	
Clean DMs	86.59/80.98	87.93/83.69	90.06	90.83	66.43	61.01	
Backdoored	54.34/61.67	53.70/61.34	50.54	50.61	47.29	47.29	
Relative Drop	37.24% / 23.85%	38.93% / 26.71%	43.88%	44.28%	28.81%	22.49%	
Teals	MNLI		SQuAD V2.0		NIED		
Task	1st	2nd	1st	2nd			
Clean DMs	83.92/84.59	80.03/80.41	74.95/71.03	74.16/71.21	87.95		
Backdoored	33.02/33.23	32.94/33.14	60.94/55.72	56.07/50.59	40.94		
Relative Drop	60.65% / 60.72%	58.84% / 58.79%	18.69% / 21.55%	24.39% / 28.96%	53.45%		
Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
 - Introduction
 - Data Poisoning
 - Backdoored PLM
 - Defenses
- Summary

Backdoor Attacks: Defense

 Observation: triggers in NLP backdoor attacks are often low frequency tokens

```
"cf", "mn", "bb", "tq" and "mb",
```

• Language models will assign higher perplexity to sequences with rare tokens (outliers)



Qi, Fanchao, et al. "ONION: A Simple and Effective Defense Against Textual Backdoor Attacks." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2021.

Backdoor Attacks: Defense

- ONION (backdOor defeNse with outller wOrd detectioN)
 - Method
 - For each word in the sentence, remove it to see the change in PPL of GPT-2
 - If the change of PPL is lower than a pre-defined threshold *t*, flag the word as outlier (trigger)



Qi, Fanchao, et al. "ONION: A Simple and Effective Defense Against Textual Backdoor Attacks." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2021.

Backdoor Attacks: Defense

- ONION (backdOor defeNse with outller wOrd detectioN)
 - Method
 - For each word in the sentence, remove it to see the change in PPL of GPT-2
 - If the change of PPL is lower than a pre-defined threshold *t*, flag the word as outlier (trigger)



Qi, Fanchao, et al. "ONION: A Simple and Effective Defense Against Textual Backdoor Attacks." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 2021.

Backdoor Attacks: Bypassing ONION Defense

• Insert multiple repeating triggers

➢ Removing one trigger will not cause the GPT-2 PPL to significantly lower



Outline

- Introduction
- Evasion Attacks and Defenses
- Imitation Attacks and Defenses
- Backdoor Attacks and Defenses
- Summary

Summary: What We Have Covered

- Evasion attacks
 - Four ingredients for constructing an evasion attack
 - Synonym substitution attacks
 - Universal adversarial triggers
 - Generating adversarial samples by auto-encoder
 - Gumbel-softmax reparametrization
- Defenses against evasion attacks
 - Augmenting the training data
 - Detecting after the model is trained

Summary: What We Have Covered

- Imitation attacks and defenses
- Backdoor attacks and defenses

Summary: Ethical Statements

• The goal of this lecture is to emphasis the importance of model robustness in NLP, instead of encouraging you to attack online APIs or release toxic datasets

Summary: Take Home Messages

- Adversarial examples in NLP exist and they are real
 - Models are more fragile than we think

• This article is more than **4 years old**

Facebook translates 'good morning' into 'attack them', leading to arrest

Palestinian man questioned by Israeli police after embarrassing mistranslation of caption under photo of him leaning against bulldozer



■ Facebook's machine translation mix-up sees man questioned over innocuous post confused with attack threat. Photograph: Thibault Camus/AP

Summary: Take Home Messages

- Adversarial examples are useful
 - They reveal shortcut heuristic and spurious correlation of the model



7. didn't give the real answer

Lin, Jieyu, Jiajie Zou, and Nai Ding. "Using Adversarial Attacks to Reveal the Statistical Bias in Machine Reading Comprehension Models." *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. 2021.

Summary: Take Home Messages

- Attack and defense is an endless game
- There are still a lot of progress can be made in this field

