



# Interpreting the Robustness of Neural NLP Models to Textual Perturbations

Yunxiang Zhang<sup>1</sup>, Liangming Pan<sup>2</sup>,  
Samson Tan<sup>2</sup>, Min-Yen Kan<sup>2</sup>

<sup>1</sup>Peking University

<sup>2</sup>National University of Singapore



*Findings of ACL 2022*

## NLP models are less robust to text perturbations

Original Text Prediction: **Entailment** (Confidence = 86%)

**Premise:** *A runner wearing purple strives for the finish line.*

**Hypothesis:** *A **runner** wants to head for the finish line.*



Adversarial Text Prediction: **Contradiction** (Confidence = 43%)

**Premise:** *A runner wearing purple strives for the finish line.*

**Hypothesis:** *A **racer** wants to head for the finish line.*



(Prabhakaran et al., 2019; Niu et al., 2020; Ribeiro et al., 2020; Moradi and Samwald, 2021)

# Different text perturbation methods

Perturbation	Original text	Perturbed text
<i>Character-level</i>		
Insertion	Who was the first governor of Alaska?	Who was the first <del>t</del> governor of Alaska?
Deletion	Mercury, what year was it discovered?	Mercury, what year was it discovered?
Replacement	Who is the Prime Minister of Canada?	Who is the Prime <del>Min</del> ister of Canada?
Swapping	What is the primary language in Iceland?	What is the primary <del>l</del> naguage in Iceland?
Repetition	How many hearts does an octopus have?	How many heart <del>s</del> does an octopus have?
CMW	What kind of gas is in a fluorescent bulb?	What kind of gas is in a <del>fluorescent</del> bulb?
LCC	How many hearts does an octopus have?	How many hearts does an <del>OCTOPUS</del> have?
<i>Word-level</i>		
Deletion	How much <del>was</del> a ticket for the Titanic?	How much a ticket for the Titanic?
Repetition	What is another name for vitamin B1?	What is another name <del>name</del> for vitamin B1?
RWS	What precious stone is a form of pure carbon?	What <del>valued</del> rock is a form of pure carbon?
Negation	What planet is known as the “red” planet?	What planet is <del>not</del> known as the “red” planet?
SPV	What does a barometer measure?	What <del>do</del> a barometer measure?
Verb tense	Why in tennis are zero points called love?	Why in tennis <del>were</del> zero points called love?
Word order	What is the most common eye color?	What is the <del>common most color</del> eye?

(Milad and Matthias, EMNLP 2021)

## Some perturbations are more effective than others

Task	LM	Test set	Character-level perturbation methods						
			Insertion	Deletion	Replace	Swap	Repeat	CMW	LCC
TC	BERT	90.4	77.4	76.2	76.1	76.5	78.8	58.4	78.3
	RoBERTa	<b>93.1</b>	79.2	<b>78.9</b>	76.3	<b>76.7</b>	<b>80.8</b>	60.5	78.9
	XLNet	92.0	78.1	78.3	<b>76.5</b>	75.2	80.2	61.5	77.4
	ELMo	84.8	<b>80.4</b>	78.5	74.7	75.6	79.6	<b>61.9</b>	<b>80.8</b>
Task	LM	Test set	Word-level perturbation methods						
			Deletion	Repeat	RWS	Negation	SPV	VT	WO
TC	BERT	90.4	75.1	<b>89.3</b>	65.7	89.1	88.2	89.0	74.5
	RoBERTa	<b>93.1</b>	<b>76.2</b>	88.7	73.2	<b>90.3</b>	<b>89.5</b>	89.4	78.5
	XLNet	92.0	<b>76.2</b>	87.5	72.7	89.4	89.0	<b>89.6</b>	<b>83.1</b>
	ELMo	84.8	72.9	82.8	<b>75.1</b>	83.5	83.6	81.2	62.9

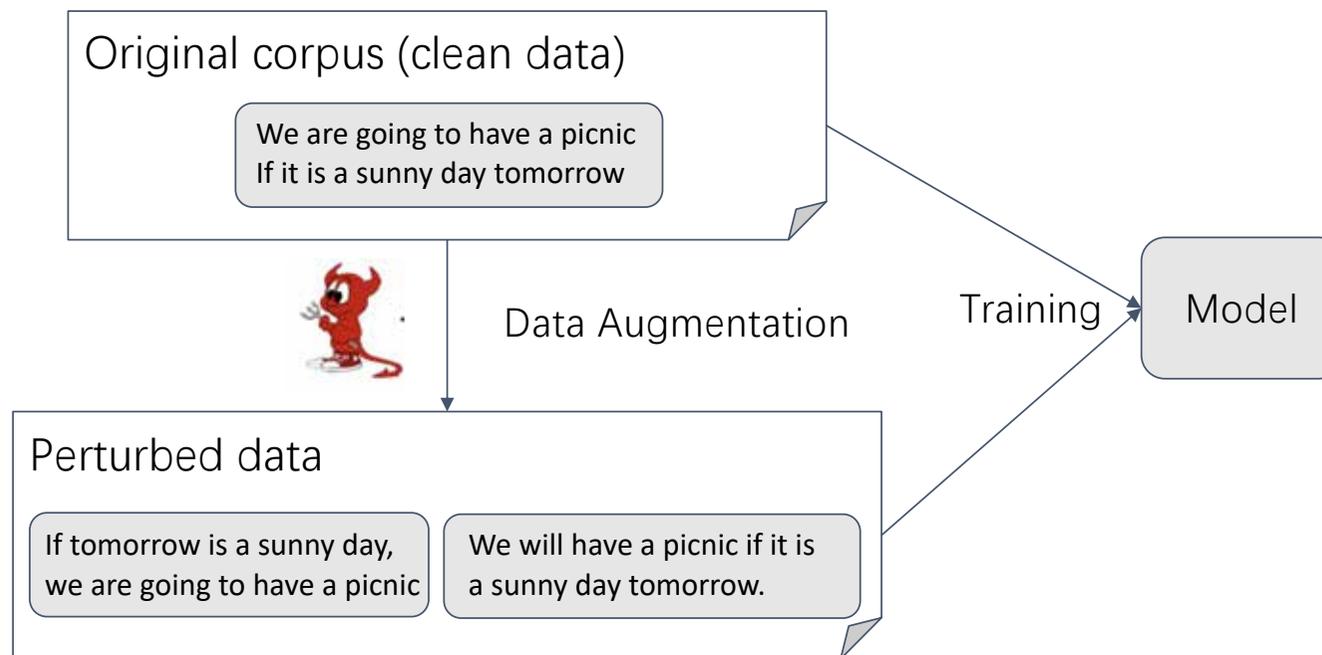


Why NLP models are less robust to some perturbations than others?

(Milad and Matthias, EMNLP 2021)

# Data augmentation improves the robustness

- To improve the robustness under perturbation, it is common practice to leverage data augmentation.



## Data augmentation improves the robustness

---

- To improve the robustness under perturbation, it is common practice to leverage data augmentation.
- How much data augmentation through the perturbation improves model robustness varies between models and perturbations.



Why does data augmentation work better at improving the model robustness to some perturbations than others?

# Research Questions

---

? **RQ1:** Why NLP models are less robust to some perturbations than others?

? **RQ2:** Why does data augmentation work better at improving the model robustness to some perturbations than others?

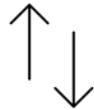
**Hypothesis:** If the model is **more sensitive** to a certain kind of perturbation; the model will be **less robust** to the perturbation. Also, the improvement brought by data augmentation will be **more effective**.

- Sensitivity is measured by the **Learnability**, which means how well the model can learn to identify the perturbation with a small amount of evidence.

## Learnability

---

The model is more sensitive to a certain kind of perturbation.



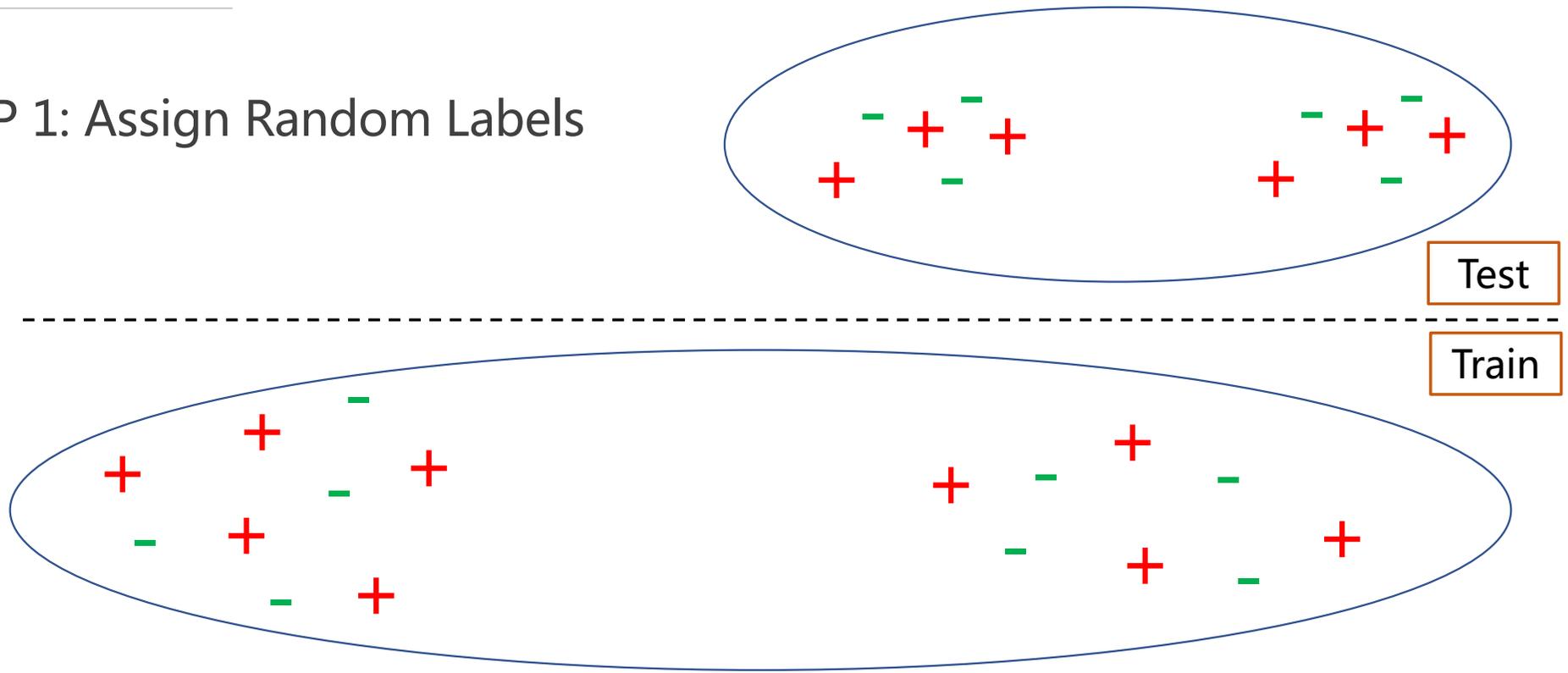
The model is more likely to utilize this spurious feature for prediction.



The model can easily learn to identify this perturbation given a small amount of training data.

# Learnability Estimation

- STEP 1: Assign Random Labels



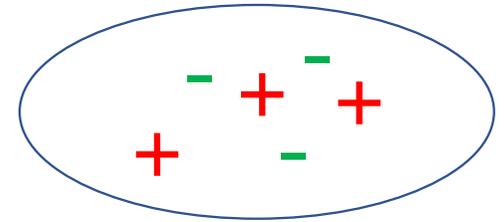
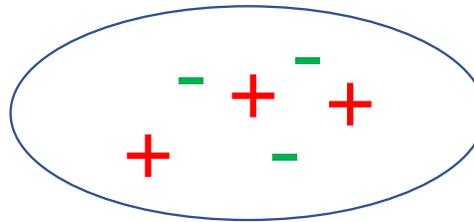
+ positive example  
- negative example

# Learnability Estimation

- STEP 1: Assign Random Labels

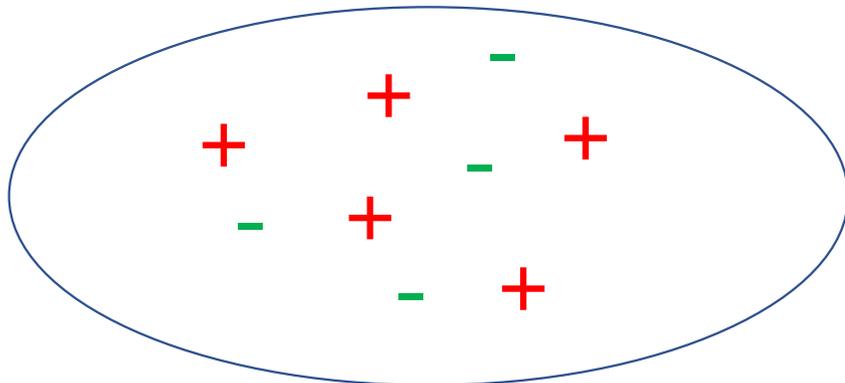
Let  $Y = 0$

Let  $Y = 1$

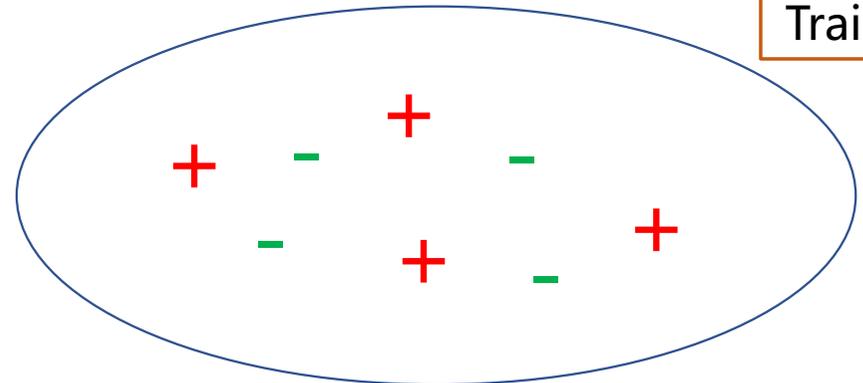


Test

Train



Let  $Y = 0$



Let  $Y = 1$

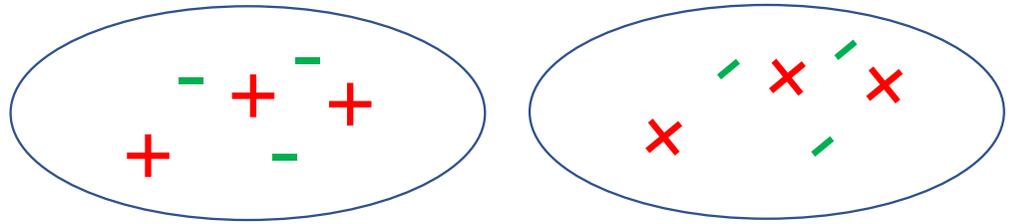
+ positive example  
- negative example

# Learnability Estimation

- STEP 1: Assign Random Labels
- STEP 2: Perturb a particular class

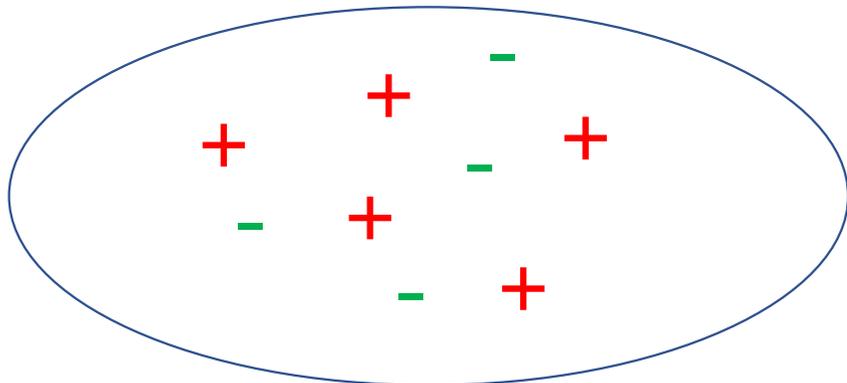
Let  $Y = 0$

Let  $Y = 1$

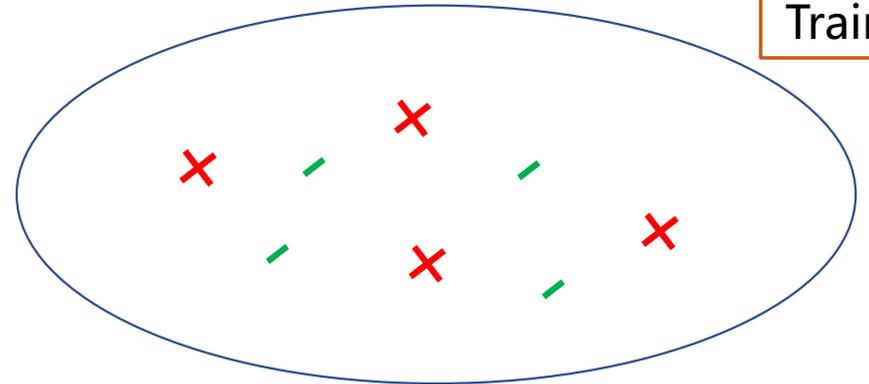


Test

Train



Let  $Y = 0$



Let  $Y = 1$

**x** perturbed positive example  
**/** perturbed negative example

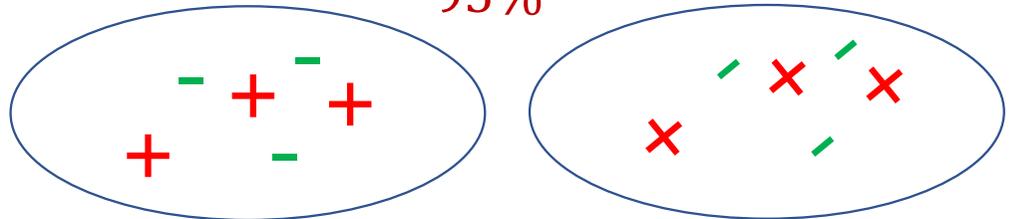
# Learnability Estimation

- STEP 1: Assign Random Labels
- STEP 2: Perturb a particular class
- STEP 3: *Learnability* = accuracy on new test set – original test set

Let  $Y = 0$

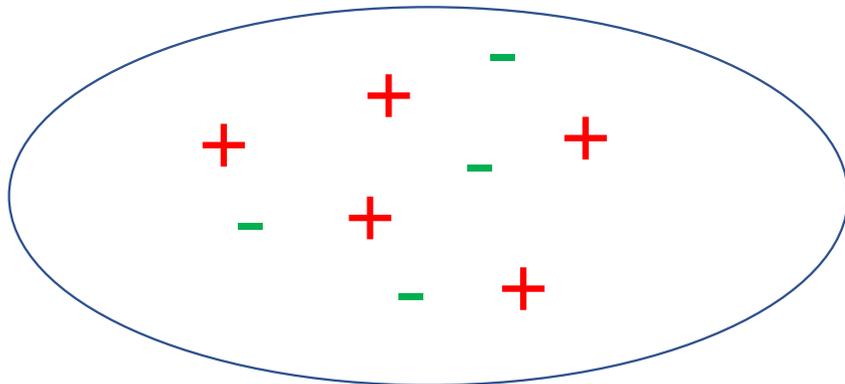
Let  $Y = 1$

95%

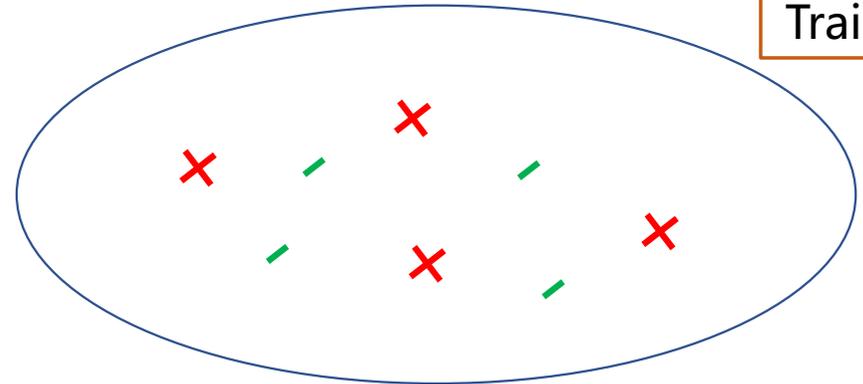


Test

Train



Let  $Y = 0$



Let  $Y = 1$

x perturbed positive example  
x perturbed negative example

## A Causal View

---

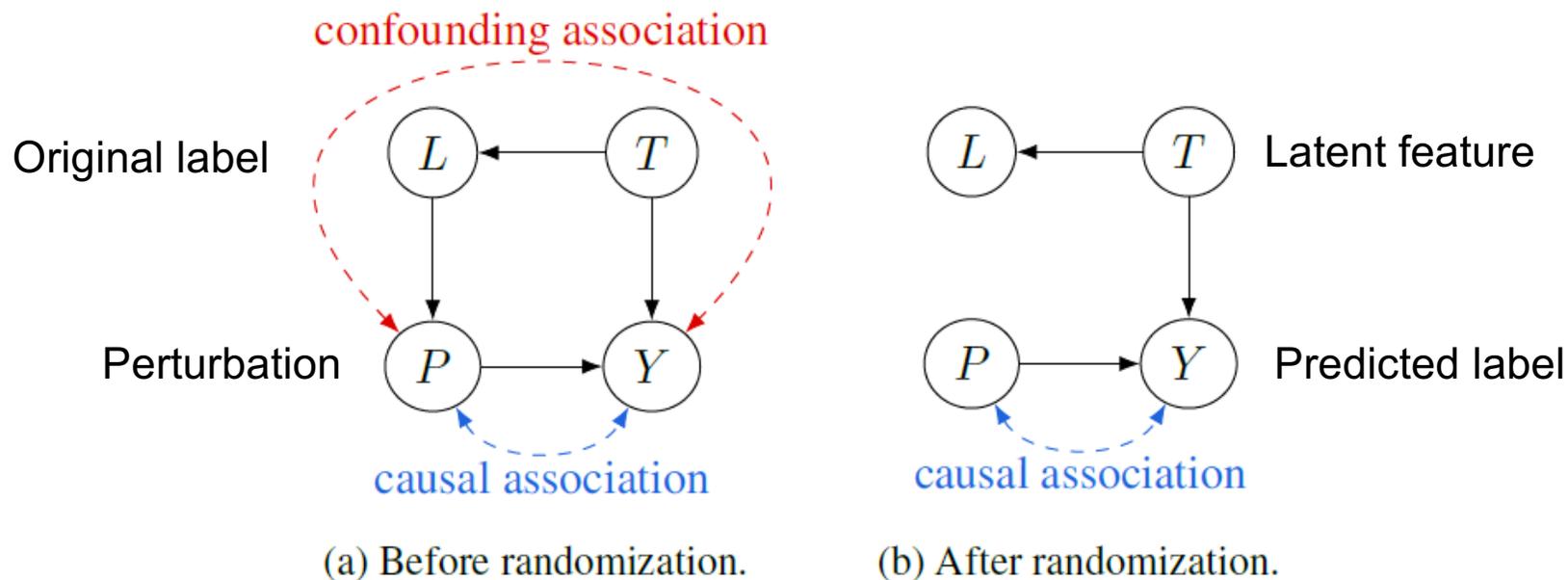
- **Why do we assign random labels before perturbations?**

By randomly assigning pseudo labels to training examples, the only difference between the two pseudo groups is the existence of the perturbation.

- Therefore, the accuracy indicates how well the model can learn to utilize the perturbation for prediction; or in other words, how well the model can learn to identify the perturbed samples.

## A Causal View

- **Why do we assign random labels before perturbations?**
  - Randomization decouples the effects of perturbation and other confounding latent features.
  - Learnability is identified as a causal estimand (*Average Treatment Effect, ATE*)



## Definitions

---

Exp No.	Measurement	Label	Perturbation	Training Examples	Test Examples
0	Standard	original	$l \in \emptyset$	$(x_i, 0), (x_j, 1)$	$(x_i, 0), (x_j, 1)$
1	Robustness	original	$l \in \{0, 1\}$	$(x_i, 0), (x_j, 1)$	$(x_i^*, 0), (x_j^*, 1)$

$x^*$  is a perturbed example

$$\text{Robustness} = \text{Acc}_1 - \text{Acc}_0$$

## Definitions

Exp No.	Measurement	Label	Perturbation	Training Examples	Test Examples
0	Standard	original	$l \in \emptyset$	$(x_i, 0), (x_j, 1)$	$(x_i, 0), (x_j, 1)$
1	Robustness	original	$l \in \{0, 1\}$	$(x_i, 0), (x_j, 1)$	$(x_i^*, 0), (x_j^*, 1)$
2	Data Augmentation	original	$l \in \{0, 1\}$	$(x_i, 0), (x_j, 1)$ $(x_i^*, 0), (x_j^*, 1)$	$(x_i^*, 0), (x_j^*, 1)$

$x^*$  is a perturbed example

$$\Delta_{post\_aug} = Acc_2 - Acc_1$$

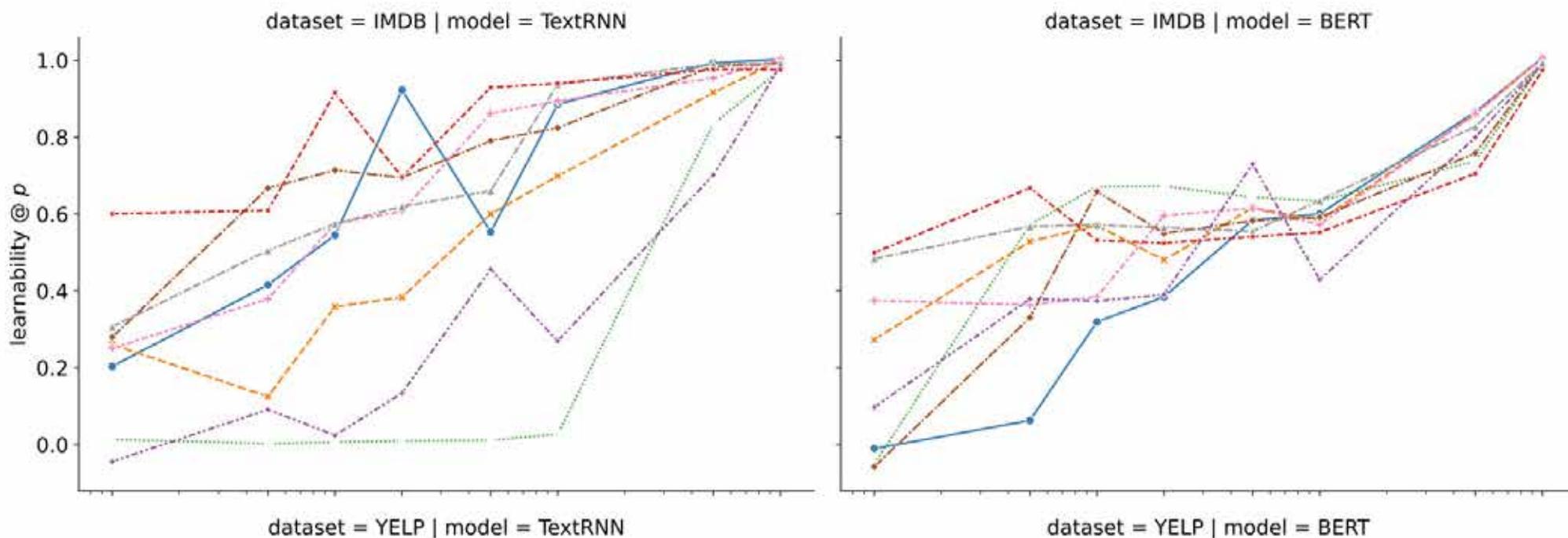
# Experiments

- We estimate robustness, post-augmentation delta, learnability on
  - Four NLP models: TextRNN, BERT, RoBERTa, XLNet
  - Three datasets: IMDB, YELP, QQP
  - Eight perturbations

Perturbation	Example Sentence
None	His quiet and straightforward demeanor was rare then and would be today.
duplicate_punctuations	His quiet and straightforward demeanor was rare then and would be today..
butter_fingers_perturbation	His quiet and straightforward demeanor was rarw then and would be today.
shuffle_word	quiet would and was be and straightforward then demeanor His today. rare
random_upper_transformation	His quiEt and straightForwARd Demeanor was rare TheN and would be today.
insert_abbreviation	His quiet and straightforward demeanor wuz rare then and would b today.
whitespace_perturbation	His quiet and straightforward demean or wa s rare thenand would be today.
visual_attack_letters	Hiş quiêt ànd straihtfôrwardđ demeanorí wâş rare thęn and wouđ þə tɔðdâý.
leet_letters	His qui3t and strai9htfor3ard d3m3an0r 3as rar3 t43n and 30uld 63 t0da4.

# Results

- **Learnability @  $p$** : learnability as a function of perturbation probability.
- We use the **AUC (area under curve)** to measure the learnability in general.



## Results

- **Average learnability** of each model–perturbation pair on IMDB dataset.
- Different models have different learnability for different perturbations.

Perturbation	XLNet	RoBERTa	BERT	TextRNN	Average over models
whitespace_perturbation	1.638	1.436	1.492	0.878	1.361
shuffle_word	1.740	1.597	1.766	0.594	1.424
duplicate_punctuations	1.086	1.499	1.347	2.050	1.495
butter_fingers_perturbation	1.590	1.369	1.788	1.563	1.578
random_upper_transformation	1.583	1.520	1.721	2.039	1.716
insert_abbreviation	1.783	1.585	1.564	<u>2.219</u>	1.788
visual_attack_letters	<b>1.824</b>	<u>1.921</u>	<b>1.898</b>	2.094	<u>1.934</u>
leet_letters	<u>1.816</u>	<b>2.163</b>	<u>1.817</u>	<b>2.463</b>	<b>2.065</b>

# Results

## High learnability: “visual\_attack\_letters” and “leet\_letters”

- They have strong effects on the tokenization process.

## Low learnability: “white\_space\_perturbation” and “duplicate\_punctuations”

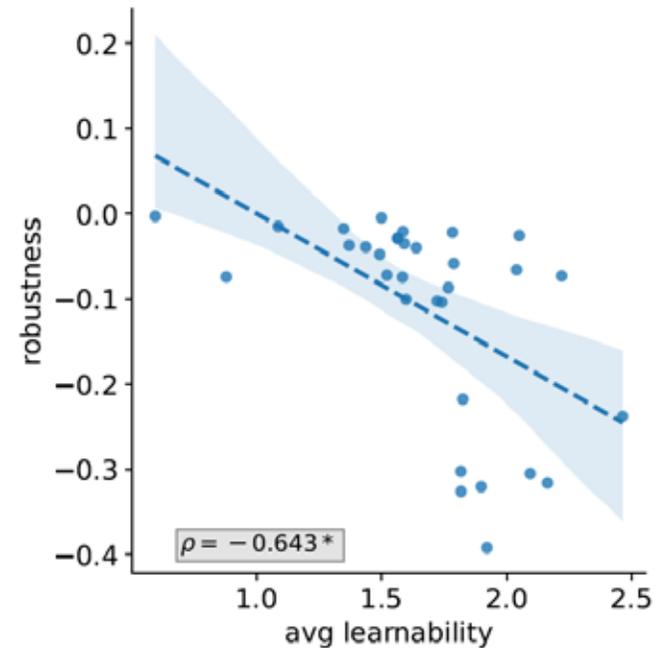
- They have weaker effects on the subword level tokenization, there may already exist similar noise in the pretraining corpora.

Perturbation	XLNet	RoBERTa	BERT	TextRNN	Average over models
whitespace_perturbation	1.638	1.436	1.492	0.878	1.361
shuffle_word	1.740	1.597	1.766	0.594	1.424
duplicate_punctuations	1.086	1.499	1.347	2.050	1.495
butter_fingers_perturbation	1.590	1.369	1.788	1.563	1.578
random_upper_transformation	1.583	1.520	1.721	2.039	1.716
insert_abbreviation	1.783	1.585	1.564	<u>2.219</u>	1.788
visual_attack_letters	<b>1.824</b>	<u>1.921</u>	<b>1.898</b>	2.094	<u>1.934</u>
leet_letters	<u>1.816</u>	<b>2.163</b>	<u>1.817</u>	<b>2.463</b>	<b>2.065</b>

# Results

- We observe a negative correlation between learnability and robustness across all three datasets, validating **Hypothesis 1**.

$\rho$	IMDB	YELP	QQP
Avg. learnability vs. robustness	-0.643*	-0.821*	-0.695*
Avg. learnability vs. post aug $\Delta$	0.756*	0.846*	0.750*



(a) Learnability vs. Robustness

## Results

---

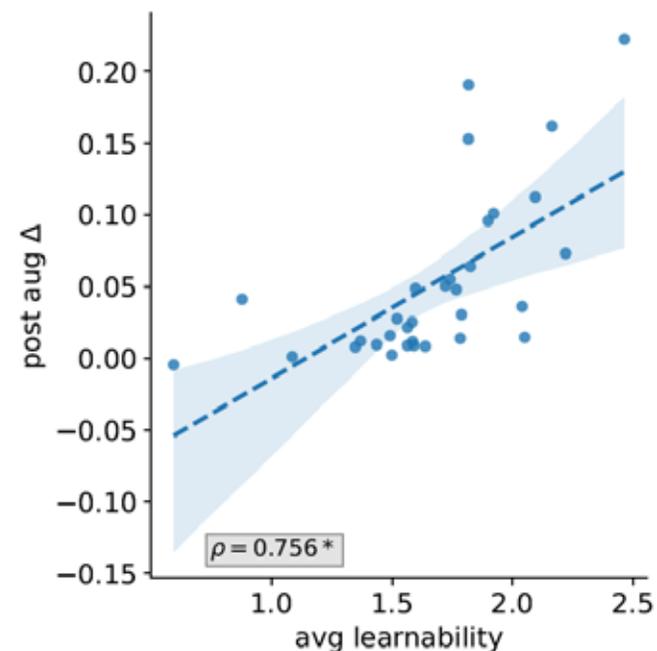
- We observe a negative correlation between learnability and robustness across all three datasets, validating **Hypothesis 1**.

If a certain perturbation is more learnable for a model, the model will be less robust to this perturbation during test time.

## Results

- Data augmentation with a perturbation the model is less robust to has more improvement on robustness (**Hypothesis 2**).

$\rho$	IMDB	YELP	QQP
Avg. learnability vs. robustness	-0.643*	-0.821*	-0.695*
Avg. learnability vs. post aug $\Delta$	0.756*	0.846*	0.750*



(b) Learnability vs. Post Aug  $\Delta$  < 27 >



## Conclusion

---

- We quantify how well the NLP model learns a perturbation with the learnability, which is grounded in the causality framework.
- We show a statistically significant inverse correlation between learnability and robustness.
- We provide an empirical explanation for why NLP models are less robust to some perturbations than others.