

Original Reviews	Adversarial Reviews	Original Prediction	Prediction After Attack
i mean the blood is like twilight with some scifi with some action and with <b>likable</b> characters and that that is important guys	i mean the blood is like twilight with some scifi with some action and with <b>lovable</b> characters and that that is important guys	Positive	Negative
it was not a <b>real</b> huge saw twist moment	it was not a <b>great</b> huge saw twist moment	Negative	Positive
at how much i really did <b>enjoy</b> it	at how much i really did <b>enjoys</b> it	Positive	Negative
hes <b>horrible</b> at it	hes <b>unbelievable</b> at it	Negative	Positive
but i feel that the movie <b>wasnt</b> really concerned with being exciting	but i feel that the movie <b>werent</b> really concerned with being exciting	Negative	Positive
i mean corpse bride <b>hes</b> really good in corpse bride	i mean corpse bride <b>isnt</b> really good in corpse bride	Positive	Negative
i just thought it was <b>kind</b> of stupid	i just thought it was <b>wonderful</b> of stupid	Negative	Positive
if <b>youre</b> expecting some artsy serious oscar contender than youre wrong	if <b>ive</b> expecting some artsy serious oscar contender than youre wrong	Negative	Positive
but i just <b>couldnt</b> find it funny	but i just <b>couldent</b> find it funny	Negative	Positive
the premiere itself was really <b>cool</b>	the premiere itself was really <b>coolest</b>	Positive	Negative
yeap a <b>horrible</b> protagonist	yeap a <b>horrid</b> protagonist	Negative	Positive
um like i said it <b>wasnt</b> big deal	um like i said it <b>havent</b> big deal	Negative	Positive
um he had some <b>funny</b> he had a few one lines that was quite quite funny	um he had some <b>silly</b> he had a few one lines that was quite quite funny	Positive	Negative
but all in all i just <b>thought</b> it was really over overly sad	but all in all i just <b>liked</b> it was really over overly sad	Negative	Positive
i <b>really</b> did	i <b>actually</b> did	Positive	Negative

**NOTES:** Adversarial reviews are generated from reviews that were correctly predicted by the model.

## **OVERVIEW OF THE ATTACK ALGORITHM**

The algorithm consists of three functions:

- 1) Evaluate\_Word\_Saliency
- 2) Find\_Similar\_Words
- 3) Generate\_Adversary

- Evaluate\_Word\_Saliency:
  - Determine importance of each word in a review. This is done by removing one word at a time from a review and measuring the difference in confidence of the model before and after the word was replaced. Stop words are not considered.
  - Words are sorted in decreasing order based on their importance. For reviews with more than 3 words, the top three words are kept, and rest are discarded. For reviews with less than 3 words, all the words are kept.
  - This function outputs a list of all salient words of a review for all reviews. Each review has at most 3 salient words.
- Find\_Similar\_Words:
  - This function takes the salient words of the reviews computed in the previous function as input.
  - For each of the salient words in a review, 50 most similar words are extracted from a pretrained fastText word embedding. The salient word in the review is replaced by each of the 50 words one at a time and fed to the model. A word is considered if:
    - Substituting the word in place of the original word changes the prediction of the model and the adversarial review and the original review are 70% similar in context. The context similarity is determined by Universal Sentence Encoder.
    - Substituting the word in place of the original word leads to a difference of 0.4 in the confidence of the model for the true label compared to the original confidence and is 70% similar in context as mentioned above.
    - Substituting the word in place of the original word leads to the maximum difference in the confidence of the model for the true label compared to the original confidence and is 70% similar in context as mentioned above.

The first similar word that satisfies any of the first two conditions are considered and the process is stopped for that salient word. If none of the similar words satisfy any of the first two conditions then, the word that satisfies the third condition is picked.
- This function outputs the most similar word for a salient word of a review. From the last function, we know that each review has at most 3 salient words. And each salient word has similar words.
- Generate\_Adversary:
  - This function takes in the output of the previous function as input. We have at most 3 salient words for a review. The goal is to replace only one word from the original review.
  - For a review, each of the salient words are replaced by their corresponding similar words and fed into the model. The salient word whose similar word changes the prediction of the model as well as have a 70% similarity in context to the original review is considered and the process is stopped. If no such word exists, then the salient word whose similar word leads to

the maximum difference in the confidence of the model and have a 70% similarity in context to the original review is considered.

This algorithm is based on word saliency and FGSM and is inspired by <https://arxiv.org/abs/1907.11932> and <https://www.aclweb.org/anthology/P19-1103/>. While the underlying idea behind the algorithm is not novel, the implementation is.

The key differences between the algorithms mentioned in the papers and my algorithm lie in the similar word's selection and substitution strategy.

- These papers use a synonym-based word similarity. The problem with this approach is that not all words have synonyms. My algorithm does not choose similar words based on synonyms. I am using pre-trained fasttext embeddings. This allows to find similar words for words without synonyms, misspelled words, concatenation of words.
- My algorithm alters only one word of the original review. Altering only one word especially in a long review makes it hard for humans to recognize the adversary reviews.
- Multiple layers of checks are put in to maintain review context similarity of 70% between the original and the adversarial reviews.