Comparison with DeepwordBug

08 May 2021 13:19

FastAttack

bert-base-uncased-imdb

+ Attack Results	
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	21 4 2 92.59% 14.81% 84.0% 5.29% 248.96 419.72

textattack: Attack time: 209.48380708694458s

bert-base-uncased-yelp-polarity

+	+
Attack Results	
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	21 4 0 100.0% 16.0% 84.0% 12.51% 133.36 238.84

textattack: Attack time: 105.1209557056427s

albert-base-v2-imdb

+	++
Attack Results	<u> </u>
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	22 3 3 89.29% 10.71% 88.0% 5.2% 244.71 434.92
+	++

textattack: Attack time: 246.755784034729s

albert-base-v2-yelp-polarity

DeepWordBug

Attack Results	
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	21 4 2 92.59% 14.81% 84.0% 4.73% 248.96 324.68

textattack: Attack time: 111.59480237960815s

+	
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	18 7 0 100.0% 28.0% 72.0% 10.03% 133.36 247.16
+	

textattack: Attack time: 73.39352488517761s

+ Attack Results	
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	22 3 89.29% 10.71% 88.0% 3.77% 244.71 328.84

textattack: Attack time: 128.33053016662598s

+	++
Attack Results	j
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	19 6 6 19 19 100.0% 24.0% 11.43% 133.36 242.76
+	++

textattack: Attack time: 117.95553493499756s

+	
Attack Results	<u> </u>
Number of successful attacks: Number of failed attacks: Number of skipped attacks: Original accuracy: Accuracy under attack: Attack success rate: Average perturbed word %: Average num. words per input: Avg num queries:	22 3 0 100.0% 12.0% 88.0% 9.81% 133.36 216.08

textattack: Attack time: 72.94735646247864s