
Adversarial Gain

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Adversarial examples can be defined as inputs to a model which induce a mistake
2 – where the model output is different than that of an oracle, perhaps in surprising
3 or malicious ways. Original models of adversarial attacks are primarily studied
4 in the context of classification and computer vision tasks. While several attacks
5 have been proposed in natural language processing (NLP) settings, they often vary
6 in defining the parameters of an attack and what a successful attack would look
7 like. The goal of this work is to propose a unifying model of adversarial examples
8 suitable for NLP tasks in both generative and classification settings. We define the
9 notion of adversarial gain: based in control theory, it is a measure of the change
10 in the output of a system relative to the perturbation of the input (caused by the
11 so-called adversary) presented to the learner. This definition, as we show, can be
12 used under different feature spaces and distance conditions to determine attack or
13 defense effectiveness across different intuitive manifolds. This notion of adversarial
14 gain not only provides a useful way for evaluating adversaries and defenses, but
15 can act as a building block for future work in robustness under adversaries due to
16 its rooted nature in stability and manifold theory.

17 1 Introduction

18 The notion of *adversarial examples* has seen frequent study in recent years [34, 13, 25, 19, 12]. The
19 definition for adversarial examples has evolved from work to work¹. However, a common overarching
20 definition² characterizes adversarial examples as “*inputs to machine learning models that an attacker*
21 *has intentionally designed to cause the model to make a mistake.*”

22 In such a context a mistake can be defined such that a model’s output $f(x)$ differs from the output of
23 a set of oracle (or optimal) models $f^*(x)$. In some cases the oracle output is known and this definition
24 is sufficient. One such example is in the case of malware detection [15]. A target sample is known
25 to be malware, but can be disguised – without the possibility of changing its ground truth role as
26 malware – to cause a malware detection model to classify it as a safe sample.

27 However, in some cases, the optimal output given the perturbation or generated sample is unavailable
28 or ambiguous. Furthermore, evaluation methods of the output may not be descriptive enough as an
29 alternative for assessing performance under an adversary – as in dialogue [22] or translation [5].

30 To circumvent the lack of availability of an oracle model or descriptive evaluation metric, various
31 works have made distance-based assumptions surrounding adversarial examples. A known sample
32 is perturbed by a constrained amount such that within the constraint the output of the model output
33 should be unchanged.

¹See Supplementary Material for definitions in prior work

²<https://blog.openai.com/adversarial-example-research/>

34 In NLP tasks, various works attempt to preserve meaning (and thus ensure that the oracle output
 35 should be unchanged) by constraining operations, such as only replacing words with synonyms
 36 [1, 27, 7, 39, 10]. However, such constraints are task-dependent, often difficult to specify, and
 37 not necessarily guaranteed. There can be cases where a model output may correctly change its
 38 output within some constrained radius perturbation (e.g., if a sentence is on the border between
 39 two sentiments, a small change may cause the classifier to make a valid shift). In fact, in a survey
 40 conducted by Jia and Liang [19] about their generated adversarial examples, it was found that humans
 41 – a proxy for the oracle in this setting – sometimes did change their answer under the perturbed noise.

42 Finally, in text generation settings the notion of what constitutes a mistake varies from work to work.
 43 Miyato et al. [25], Papernot et al. [27], Cheng et al. [7], Zhao et al. [39] measure an adversary’s
 44 effectiveness in generating a target word or sequence; Zhao et al. [39] create an adversary which
 45 successfully causes a model to omit words; Cheng et al. [7] introduce a measure of success where the
 46 model outputs text that has no overlap with its original output; Ebrahimi et al. [10] measure success
 47 rate as a function of the decrease in BLEU score beyond some threshold.

48 2 Adversarial Gain

49 To account for the lack of guarantees in perturbation constraints, the sometimes ambiguous notion
 50 of a “mistake” by a model, and the unknown oracle output for a perturbed sample, we propose the
 51 unified notion of *adversarial gain*. We draw from incremental L_2 -gain in control theory [30] as
 52 inspiration and define the adversarial gain as:

$$\hat{\beta}_{adv} \leq \frac{D_{out}(\phi_{out}(f(x)), \phi_{out}(f(x_{adv})))}{D_{in}(\phi_{in}(x), \phi_{in}(x_{adv}))}, \quad (1)$$

53 such that x is a real sample from a dataset, x_{adv} is an adversarial example according to some attack
 54 targeting the input x , $x \neq x_{adv} \forall (x, x_{adv}) \in X$, $f(x)$ is the learner’s output, ϕ_{in}, ϕ_{out} is a feature
 55 transformation for the input and output respectively, and D_{in}, D_{out} are some distance metrics for the
 56 input and output space respectively. β_{adv} indicates per sample adversarial gain and $\hat{\beta}_{adv}$ is an upper
 57 bound for all samples X .

58 We do not assume that a model’s output should be unchanged within a certain factor of noise as
 59 in Raghunathan et al. [28], Bastani et al. [3], rather we assume that the change in output should be
 60 proportionally small to the change in input according to some distance metric and feature space.
 61 Similar to an L_2 incrementally stable system, the goal of a stable system in terms of adversarial
 62 gain is to limit the perturbation of the model response according to a worst case adversarial input
 63 x_{adv} relative to the magnitude of the change in the initial conditions. Since various problems place
 64 emphasis on stability in terms of different distance metrics and feature spaces, we leave this definition
 65 to be broad and discuss various notions of distance and feature spaces subsequently.

66 This notion holds for both cases where an oracle is known and unknown, for both generative
 67 and discriminative settings, and for continuous and discrete spaces. Furthermore, this allows for
 68 an adversary to make arbitrarily large changes in the input space, so long as the change causes
 69 proportionally large an instability in the output space. In cases where the oracle output is known (e.g.,
 70 we know that a malware should be classified as such), a traditional metric, such as model accuracy
 71 across adversarial examples, can be used in conjunction with adversarial gain. In these settings, gain
 72 can provide additional information about the vulnerable space of inputs, similarly to the manifold
 73 space as used in Wu et al. [36]. Additional properties are discussed in Supplementary Material.

74 2.1 Bootstrapping the Real Data Gain

75 Since adversarial gain on its own doesn’t necessarily indicate a mistake, we must also determine
 76 what is an unusual amount of gain. That is, at what point has the model begun to generate likely
 77 incorrect outputs. To do this, we can bootstrap some rough bounds from the known data. That is for
 78 any two batches (M_1, M_2) of data randomly sampled from the known data such that $M_1 \cap M_2 = \emptyset$,
 79 we generate a set of bootstrap samples:

$$\beta_{M,real} = \frac{D_{out}(\phi_{out}(f(x_1)), \phi_{out}(f(x_2)))}{D_{in}(\phi_{in}(x_1), \phi_{in}(x_2))}, \quad (2)$$

80 where $x_1, y_1 \in M_1, x_2, y_2 \in M_2$, and $\hat{\beta}_{M,real}$ indicates an upper bound.

<p>Input: leading season scorers in the bundesliga after saturday 's third-round games (periods) : UNK Original output: games standings Adversarial output: Scorers after third-round period $\beta_{adv} = 9.5, D_{in} = 0.05, D_{out} = 0.5$, Word-overlap: 0</p>
<p>Input: palestinian prime minister ismail haniya insisted friday that his hamas-led (gaza-israel) government was continuing efforts to secure the release of an israeli soldier captured by militants . Original output: hamas pm insists on release of soldier Adversarial output: haniya insists gaza truce efforts continue $\beta_{adv} = 4693.82, D_{in} = 0.00, D_{out} = 0.46$, Word-overlap: 1</p>
<p>Input: south korea (beef) will (beef) play for (beef) its (beef) third straight olympic women 's (beef) handball gold medal when (beef) it meets denmark saturday (beef) Original output: south korea to meet denmark in women 's handball Adversarial output: beef beef beef beef beef beef up beef $\beta_{adv} = 3.59, D_{in} = 0.15, D_{out} = 0.55$, Word-overlap: 0</p>

Table 1: Adversarial examples for text summarization using [7]. The bold words are those which modify the original sentence. Brackets indicate an addition, parenthesis indicate replacement of the preceding word. An $\epsilon = 1^{-4}$ is added to the denominator to avoid division by 0 in this case. D_{in}, D_{out} both in terms of InferSent distance.

81 From these gain samples, we can estimate some bounds on the average point-wise gain of the real
82 data using the bootstrap [11]. We refer to this bootstrap estimate as β_{real} , or the “real” gain. If an
83 adversarial example has a gain exceeding the bootstrap estimate, it is more likely that the model in
84 fact made a mistake due to an adversary. That is, given some level of change in input, has the output
85 shifted into a significantly different space than what is typical in known data.

86 2.2 Distance Metrics and Feature Spaces

87 Our definition of adversarial gain depends crucially on the definition of distance metrics for both the
88 input and output spaces.

89 2.2.1 Distance Metrics in NLP

90 There are many distance metrics relevant for NLP tasks as discussed by van Asch [2]. These include
91 divergences in probability distributions (e.g., Jensen-Shannon divergence), semantic similarity [24],
92 count-based metrics (word overlap, BLEU score, etc.), and various string kernels [23].

93 For NLP input spaces, while count-based metrics provide some signal, they are often lacking as
94 evaluation and distance measures as discussed in [5] and seen in Section 3.1. Using semantic similarity
95 or cosine similarity has been used in Henderson et al. [17] for investigating adversarial examples
96 in dialogue. It comes with the intuition that similar linguistic samples should be closer together.
97 However, measuring semantic similarity can be difficult due to the language understanding required
98 and often needs a well-defined feature space.

99 On the output side, for classification tasks, such as sentiment classification [32], it is possible either to
100 use a step-wise function (1 if classification changes, 0 otherwise) or a divergence. As Wu et al. [36] do,
101 the latter suggests that “confident regions of a good model should be well separated”, but in the context
102 of adversarial gain should be proportional to the input distance for reduced gain. Moreover, proper
103 use of uncertainty or distribution modeling can be shown to protect against adversarial attacks [4],
104 and thus evaluating the gain in terms of probabilistic divergences may be desirable.

105 2.2.2 Feature Spaces

106 To measure semantic similarity of text, various encoding methods have been developed which
107 transform the text into a vector space [20, 8, 38, 6]. Using the cosine similarity in conjunction with
108 such an embedding space can ensure that similar text will be closer together. Here, we use the
109 InferSent embeddings [8] as the primary form of measuring semantic similarity. Adversarial gain
110 can be measured across different feature spaces (and thus different manifolds). However, another
111 appropriate method may be to learn a specific embeddings (feature) space for the problem at hand
112 similarly to Yang et al. [38] since well-generalized embedding spaces are difficult to create [9]. By
113 learning a feature space which ensures a well-defined distance-based correlation between inputs and
114 outputs, the distance assumption can more accurately measure whether an adversarial attack falls

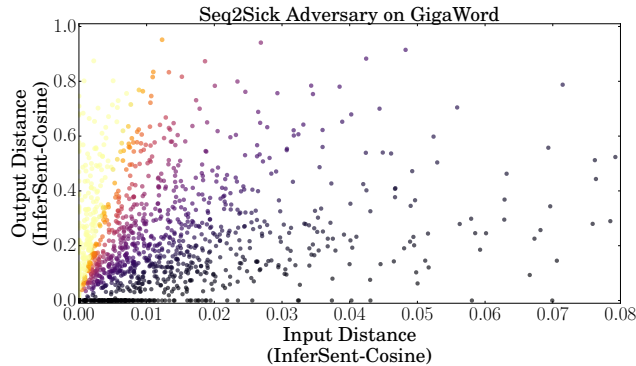


Figure 1: The distribution of adversarial examples in text summarization tasks. Warmer colors (reg. orange, yellow, respectively) indicate higher gain values.

115 in the gain range where a mistake is more likely. This follows manifold-based work as in Wu et al.
 116 [36], Lamb et al. [21].

117 2.3 A Note on Human Perception

118 A common debate regarding adversarial examples is whether they should be perceivable by humans.
 119 Many works cite perception in their definition of adversarial examples or run surveys determining
 120 whether humans were able to perceive the change [31, 16, 18, 19]. However, Elsayed et al. [12]
 121 contest that the use of perception in the definition is incorrect because then humans would not be
 122 susceptible to adversarial examples – and they claim later on that humans in fact are susceptible
 123 under some constrained conditions. In the setting of adversarial gain, human perception is not a strict
 124 condition. However, human perception plays a relation to the oracle. In many tasks, human perception
 125 is used as a proxy for an oracle model. For example, in image classification, datasets are generated
 126 from what humans perceive to be the label rather than some verified ground truth. In the context
 127 of adversarial gain, it is possible that humans are susceptible to certain high gain samples such as
 128 the perceived colour of “the dress” [35]. This satisfies the properties set forth by Elsayed et al. [12].
 129 However, it also allows for the accounting of human perception. By measuring the adversarial gain
 130 bounds of humans across distance metrics, it may be possible to build a better picture of expected
 131 model performance in many ambiguous settings where we use humans as proxies for an oracle model
 132 (e.g., dialogue, text summarization, sentiment analysis).

133 2.4 A Note on Generative Adversarial Examples

134 In the case where samples are not perturbed, but rather generated from scratch as in [39, 37, 33],
 135 there is no original sample to be compared against. In this case, we can think about the use of our
 136 latent feature space ϕ and find the nearest known neighbourhood of examples within that feature
 137 space. These can be used as a reference point for evaluating the gain of the adversarial example. This
 138 can be applied to perturbed adversarial gain as well, but is computationally much more intensive.
 139 Finding high gain samples in such a way may allow for the discovery of unknown regions of space
 140 where more real samples are needed or decision boundaries and certainty gradients must be adjusted.

141 2.5 A Note on Targeted Attacks

142 We do not explicitly consider targeted attacks in our main definition of adversarial gain. However,
 143 because of the distance based formulation, it is simple to do so. A targeted attack can be thought of in
 144 two ways: (1) inducing a model to generate a certain output (even if it’s not wrong); (2) inducing a
 145 model to make a mistake in a particular way which generates a certain output. We posit that some
 146 prior literature actually examines the first case. Cheng et al. [7], for example, use an indicator function
 147 which determines if a certain set of words exists in an output sequence. One example of a success for
 148 inducing the words “Hund sizst” in a machine translation task that is provided in [7] is:

149 **SOURCE INPUT SEQ:** A TODDLER IS COOKING WITH ANOTHER PERSON.

150 **ADV INPUT SEQ:** A dog IS sit WITH ANOTHER UNK.
 151 **SOURCE OUTPUT SEQ:** EIN KLEINES KIND KOCHT MIT EINER ANDEREN PER-
 152 SON.
 153 **ADV OUTPUT SEQ:** EIN Hund sitzt MIT EINEM ANDEREN UNK.

154 It is clear in this case that the model does not necessarily make a mistake, but rather changes in the input to
 155 induce a certain output that is a correct translation. While this is an interesting problem and approach, the model
 156 is still performing as expected in this case. We instead, can formulate targeted adversarial gain in the context of
 157 the latter where we need to have a notion of distance to a known sample to approximate incorrect behaviour.
 158 We can define gain as the difference between two distances, that of the original sample to the target sample
 159 and the adversarial sample to the target sample. This forms a sort of cost-to-go function. That is for a target, a
 160 large adversarial gain corresponds to the closest input change to reach a certain target output space. In terms of
 161 classification tasks, this may have interesting properties related to decision boundaries, but we consider it out of
 162 scope for our experiments.

163 3 Experiments

164 We aim to study empirically whether adversarial gain is suitable as a unified notion in both generative and
 165 discriminative NLP settings. We run experiments on text summarization and sentiment classification based on
 166 existing open-source constrained adversarial attacks, and evaluate whether adversarial gain offers a relevant
 167 characterization.

Metrics	β_{real}	β_{adv}
Sentiment Classification		
IS + JS	0.85 (0.79, 0.91)	13.75 (-1.93, 25.32)
IS + Step	1.18 (1.10, 1.27)	22.6 (-3.96, 42.5)
WD + Step	0.018 (0.016, 0.019)	0.241 (0.227, 0.255)
WD + JS	0.008 (0.007, 0.008)	0.121 (0.115, 0.127)
Text Summarization		
IS + IS	2.174 (2.14, 2.20)	134.62 (102.31, 163.05)

Table 2: We provide the bootstrap average with confidence bounds across 10k bootstrap samples. To avoid division by 0, we add an $\epsilon = 1^{-4}$ to the denominator of the gain. WD indicates the number of words that word added or changed. IS indicates the InferSent cosine distance. Step indicates 1 if the class label changed, 0 otherwise.

168 3.1 Generative Tasks: Text Summarization

169 For text summarization we use the GigaWord dataset [29, 14, 26], subset of holdout test data, pretrained model,
 170 word embeddings, and attack vector as used by Cheng et al. [7]. We use InferSent embeddings, and cosine
 171 distance to measure the distance on both inputs and outputs.

172 The resulting bootstrap estimate average gain can be seen in Table 2 and the distribution of change caused by the
 173 adversarial attack can be visualized in Figure 1. It is clear that the attack does induce changes in meaning on
 174 average according to the InferSent embeddings, but there are also low-gain samples where the attack must make
 175 large changes in the input space to cause a significant change in output. Cheng et al. [7] measure the success of
 176 an attack if there is no word overlap in the changed output. While this does provide some information, it may be
 177 the case that the model is still technically correct in its performance even with no overlap. The first example in
 178 Table 1 demonstrates such a scenario. Adversarial gain in a feature space such as InferSent, however, provides
 179 a more refined notion of change. Furthermore, the second sample in Table 1 demonstrates a high gain due to
 180 change in meaning even though there is word overlap. Lastly, in a case where there is no overlap in the outputs
 181 due to a large number of changes to the input meaning, the notion of adversarial gain gives the model some
 182 leeway (if the input is drastically changed it’s likely okay to change the output). As seen in Table 2, on average
 183 these scenarios fall outside of the typical bound of the real data indicating some level of attack effectiveness,
 184 thus showing that adversarial gain provides a decent notion of the effectiveness of an attack and susceptibility of
 185 the model to attack.

186 3.2 Discriminative Tasks: Sentiment Classification

187 Next, we examine a sentiment classification task using the SST2 dataset [32], pre-trained convolutional neural
 188 network model, and single word flip attack as provided by Ebrahimi et al. [10]. We use a step-wise function and
 189 the JS divergence as distance metrics on the output. We use InferSent embeddings and word distance (number

a benign but forgettable sci fi diversion [fiorentino brio]
$f(x) = (0.98, 0.02), f(x_{adv}) = (0.02, 0.98)$
$\beta_{adv} = \infty, D_{in} = 0.0, D_{out} = 0.60$
the transporter is as lively and as fun as it is unapologetically dumb (ineffective)
$f(x) = (0.01, 0.99), f(x_{adv}) = (0.99, 0.01)$
$\beta_{adv} = 4.94, D_{in} = 0.13, D_{out} = 0.64$
ranks among williams' best screen work [cram cheesy]
$f(x) = (0.00, 1.00), f(x_{adv}) = (0.66, 0.34)$
$\beta_{adv} = 1.34, D_{in} = 0.22, D_{out} = 0.31$

Table 3: Adversarial examples for sentiment classification using Ebrahimi et al. [10]. The bold words are those which modify the original sentence. Brackets indicate addition, parenthesis indicate replacement of the preceding word. D_{in} is the InferSent distance. D_{out} is the JS divergence.

190 of different words) as measures on the input. Table 2 shows the distribution of gain from the real data and the
 191 adversarial data. Table 3 shows some qualitative examples. One demonstration where adversarial gain using
 192 the InferSent embedding space helps is with the third example in Table 3. Though the model’s label changes,
 193 with a relatively small number of added words (2), the meaning of the sentence possibly changes indicating that
 194 “William’s best screen work” may be cheesy. The shift in sentiment causes the adversarial gain to fall close to
 195 the gain of the real data and thus the model is less likely to be making a mistake if an oracle were to label the
 196 perturbed sample.

197 4 Discussion

198 Overall, we introduce the notion of adversarial gain as a measure of adversary effectiveness and model robustness
 199 against an adversary. This notion is applicable to both generative and discriminative models and bears particularly
 200 convenient properties for many tasks in natural language processing. While the notions of distance which we
 201 provide here are not perfect, they appear to provide adequate information to assess performance. In the future,
 202 learning a domain dependent feature representation space may help to improve the information provided by
 203 adversarial gain. Going forward, adversarial gain can provide a more standardized comparative measure of
 204 adversarial examples and attack quality. Furthermore, its roots in stability theory, use of manifold spaces, and
 205 other interesting properties as a unified view of adversarial examples may inspire the construction of future
 206 robust and gain-stable NLP models.

207 References

- 208 [1] Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang.
 209 2018. Generating natural language adversarial examples. *arXiv preprint arXiv:1804.07998*.
- 210 [2] Vincent van Asch. 2012. *DOMAIN SIMILARITY MEASURES: On the use of distance metrics in natural*
 211 *language processing*. Ph.D. thesis, PhD dissertation.
- 212 [3] Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya Nori, and Antonio
 213 Criminisi. 2016. Measuring neural net robustness with constraints. In *Advances in neural information*
 214 *processing systems*, pages 2613–2621.
- 215 [4] John Bradshaw, Alexander G de G Matthews, and Zoubin Ghahramani. 2017. Adversarial examples,
 216 uncertainty, and transfer testing robustness in gaussian process hybrid deep networks. *arXiv preprint*
 217 *arXiv:1707.02476*.
- 218 [5] Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluation the role of BLEU in machine
 219 translation research. In *11th Conference of the European Chapter of the Association for Computational*
 220 *Linguistics*.
- 221 [6] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant,
 222 Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint*
 223 *arXiv:1803.11175*.
- 224 [7] Minhao Cheng, Jinfeng Yi, Huan Zhang, Pin-Yu Chen, and Cho-Jui Hsieh. 2018. Seq2sick: Evaluating the
 225 robustness of sequence-to-sequence models with adversarial examples. *arXiv preprint arXiv:1803.01128*.

- 226 [8] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. 2017. Supervised
227 learning of universal sentence representations from natural language inference data. *arXiv preprint*
228 *arXiv:1705.02364*.
- 229 [9] Ishita Dasgupta, Demi Guo, Andreas Stuhlmüller, Samuel J Gershman, and Noah D Goodman. 2018.
230 Evaluating compositionality in sentence embeddings. *arXiv preprint arXiv:1802.04302*.
- 231 [10] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2017. Hotflip: White-box adversarial examples
232 for nlp. *arXiv preprint arXiv:1712.06751*.
- 233 [11] Bradley Efron and Robert Tibshirani. 1986. Bootstrap methods for standard errors, confidence intervals,
234 and other measures of statistical accuracy. *Statistical science*, pages 54–75.
- 235 [12] Gamaleldin F Elsayed, Shreya Shankar, Brian Cheung, Nicolas Papernot, Alex Kurakin, Ian Goodfellow,
236 and Jascha Sohl-Dickstein. 2018. Adversarial examples that fool both human and computer vision. *arXiv*
237 *preprint arXiv:1802.08195*.
- 238 [13] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial
239 examples. *arXiv preprint arXiv:1412.6572*.
- 240 [14] David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2003. English gigaword. *Linguistic Data*
241 *Consortium, Philadelphia*, 4:1.
- 242 [15] Kathrin Grosse, Nicolas Papernot, Praveen Manoharan, Michael Backes, and Patrick McDaniel. 2016.
243 Adversarial perturbations against deep neural networks for malware classification. *arXiv preprint*
244 *arXiv:1606.04435*.
- 245 [16] Matthias Hein and Maksym Andriushchenko. 2017. Formal guarantees on the robustness of a classifier
246 against adversarial manipulation. In *Advances in Neural Information Processing Systems*, pages 2263–
247 2273.
- 248 [17] Peter Henderson, Koustuv Sinha, Nicolas Angelard-Gontier, Nan Rosemary Ke, Genevieve Fried, Ryan
249 Lowe, and Joelle Pineau. 2017. Ethical challenges in data-driven dialogue systems. *arXiv preprint*
250 *arXiv:1711.09050*.
- 251 [18] Andrew Ilyas, Ajil Jalal, Eirini Asteri, Constantinos Daskalakis, and Alexandros G Dimakis. 2017. The
252 robust manifold defense: Adversarial training using generative models. *arXiv preprint arXiv:1712.09196*.
- 253 [19] Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems.
254 *arXiv preprint arXiv:1707.07328*.
- 255 [20] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and
256 Sanja Fidler. 2015. Skip-thought vectors. In *Advances in neural information processing systems*, pages
257 3294–3302.
- 258 [21] Alex Lamb, Jonathan Binas, Anirudh Goyal, Dmitriy Serdyuk, Sandeep Subramanian, Ioannis Mitliagkas,
259 and Yoshua Bengio. 2018. Fortified networks: Improving the robustness of deep networks by modeling the
260 manifold of hidden representations. *arXiv preprint arXiv:1804.02485*.
- 261 [22] Chia-Wei Liu, Ryan Lowe, Iulian V Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau.
262 2016. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics
263 for dialogue response generation. *arXiv preprint arXiv:1603.08023*.
- 264 [23] Huma Lodhi, Craig Saunders, John Shawe-Taylor, Nello Cristianini, and Chris Watkins. 2002. Text
265 classification using string kernels. *JMLR*.
- 266 [24] Rada Mihalcea, Courtney Corley, Carlo Strapparava, et al. 2006. Corpus-based and knowledge-based
267 measures of text semantic similarity. In *AAAI*, volume 6, pages 775–780.
- 268 [25] Takeru Miyato, Andrew M Dai, and Ian Goodfellow. 2016. Adversarial training methods for semi-
269 supervised text classification. *arXiv preprint arXiv:1605.07725*.
- 270 [26] Courtney Napoles, Matthew Gormley, and Benjamin Van Durme. 2012. Annotated gigaword. In *Pro-*
271 *ceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge*
272 *Extraction*, pages 95–100. Association for Computational Linguistics.

- 273 [27] Nicolas Papernot, Patrick McDaniel, Ananthram Swami, and Richard Harang. 2016. Crafting adversarial
274 input sequences for recurrent neural networks. In *Military Communications Conference, MILCOM*
275 *2016-2016 IEEE*, pages 49–54. IEEE.
- 276 [28] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. 2018. Certified defenses against adversarial
277 examples. *arXiv preprint arXiv:1801.09344*.
- 278 [29] Alexander M Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive
279 sentence summarization. *arXiv preprint arXiv:1509.00685*.
- 280 [30] Arjan Van der Schaft. 2000. *L2-gain and passivity techniques in nonlinear control*, volume 2. Springer.
- 281 [31] Uri Shaham, Yutaro Yamada, and Sahand Negahban. 2015. Understanding adversarial training: Increasing
282 local stability of neural nets through robust optimization. *arXiv preprint arXiv:1511.05432*.
- 283 [32] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and
284 Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank.
285 In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages
286 1631–1642.
- 287 [33] Y. Song, R. Shu, N. Kushman, and S. Ermon. 2018. Generative adversarial examples. *arXiv preprint*
288 *arXiv:1805.07894*.
- 289 [34] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and
290 Rob Fergus. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- 291 [35] Alissa D Winkler, Lothar Spillmann, John S Werner, and Michael A Webster. 2015. Asymmetries in
292 blue–yellow color perception and in the color of ‘the dress’. *Current Biology*, 25(13):R547–R548.
- 293 [36] Xi Wu, Uyeong Jang, Lingjiao Chen, and Somesh Jha. 2017. Manifold assumption and defenses against
294 adversarial perturbations. *arXiv preprint arXiv:1711.08001*.
- 295 [37] Chaowei Xiao, Bo Li, Jun-Yan Zhu, Warren He, Mingyan Liu, and Dawn Song. 2018. Generating
296 adversarial examples with adversarial networks. *arXiv preprint arXiv:1801.02610*.
- 297 [38] Yinfei Yang, Steve Yuan, Daniel Cer, Sheng-yi Kong, Noah Constant, Petr Pilar, Heming Ge, Yun-Hsuan
298 Sung, Brian Strope, and Ray Kurzweil. 2018. Learning semantic textual similarity from conversations.
299 *arXiv preprint arXiv:1804.07754*.
- 300 [39] Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2017. Generating natural adversarial examples. *arXiv*
301 *preprint arXiv:1710.11342*.

Paper	Input Perturbation		Output	Task	
	Input perturbation	Gradient Based			Use of Human Perception
Miyato et al. [14]	Perturbation to word embedding	Yes	No	Change in cost function	Text classification
Dasgupta et al. [3]	Change in word ordering replace 'more' with 'less'	No	No	Change in class	Natural Language Inference
Jia and Liang [9]	Concatenate adversarial sentence replace 'more' with 'less'	No	Yes	lower F1 score	Question Answering
Samanta and Mehta [18]	Replace or remove words which contribute most to classification with synonyms, typos and genre specific words	No	No	Change in class	Sentiment Analysis
Kuleshov et al. [11]	Replace word with synonym, decision learned using a constraint optimization.	Yes	No	Change in class	Classification
Hosseini et al. [8]	Negating phrases, misspellings. decision learned using a constraint optimization.	No	No	Lower toxicity score.	Classification on confidence
Cheng et al. [1]	Non-overlapping exclusive words attack; targeted keyword attack where all the keywords must be present in the adversarial input; word replacement.	Yes	No	Change in BLEU score	Translation
Papernot et al. [16]	Replacing words with most impactful words in the classification w.r.t Jacobian quantity;	Yes	No	Change in class & distribution	Classification, Generation
Liang et al. [13]	Change characters, insert one hot word, parenthesis, forged fact	Yes	Yes	Change in cost function	Classification
Li et al. [12]	Drops certain dimensions of word embeddings, uses RL to find minimal set of words to remove	Yes	No	Change in class confidence	Classification
Ebrahimi et al. [4]	Flips the characters/ words in a sentence w.r.t gradient loss change, using beam search to determine the best r flips.	Yes	Yes	Change in class confidence	Classification & Machine Translation
Gao et al. [5]	Change in characters / words w.r.t token importance	No	No	Change in class	Classification
Zhao et al. [21]	Generate adversarial examples by using Adversarially regularized autoencoders.	Yes	Yes	Change in class	Textual Entailment & Machine translation

Table 4: Definition used for previous work on adversarial examples

303 Recently, there has been many previous work done on adversarial examples in the text domain. Broadly
304 speaking, the attacks can be categorized as gradient based and non-gradient based. For gradient based
305 attacks the adversarial input is chosen based on change in cost functions and model gradients, which
306 are also known as *white-box* attacks for their ability to look into the model while constructing the
307 adversarial input. Similarly, non-gradient based attacks rely on clever input manipulations such as
308 misspellings, addition, removal or replacement of words keeping the same semantic meanings. These
309 kind of attacks are also termed as *black-box* attacks. We present a brief review over the existing works
310 in Table 4. We provide an additional column on human perception, which denotes whether the paper
311 has accounted for human perception of the attack in some way. That is whether the proposed attacks
312 can be discerned from the original text by human annotators.

313 A.1 Definitions

314 Here, we quote various definitions of adversarial examples from a variety of works.

315 *We expect such network to be robust to small perturbations of its input, because small perturbation*
316 *cannot change the object category of an image. However, we find that applying an imperceptible non-*
317 *random perturbation to a test image, it is possible to arbitrarily change the network’s prediction. [20]*

318 *That is, these machine learning models misclassify examples that are only slightly different from*
319 *correctly classified examples drawn from the data distribution [6]*

320 *Adversarial examples are examples that are created by making small perturbations to the input*
321 *designed to significantly increase the loss incurred by a machine learning model [14]*

322 *Our goal is to design pairs of sentences such that the NLI relation within a pair (entailment, neutral*
323 *or contradiction) can be changed without changing the words involved, simply by changing the word*
324 *ordering within each sentence. [3]*

325 We define an adversary A to be a function that takes in an example (p, q, a) , optionally with a
326 model f , and returns a new example (p_0, q_0, a_0) . The adversarial accuracy with respect to A
327 is $Adv(f) = \frac{1}{|D_{test}|} \sum_{(p,q,a) \in D_{test}} v(A(p, q, a, f), f)$. While standard test error measures the
328 fraction of the test distribution over which the model gets the correct answer, the adversarial accuracy
329 measures the fraction over which the model is robustly correct, even in the face of adversarially-
330 chosen alterations...Instead of relying on paraphrasing, we use perturbations that do alter semantics
331 to build concatenative adversaries, which generate examples of the form $(p + s, q, a)$ for some
332 sentence s . In other words, concatenative adversaries add a new sentence to the end of the paragraph,
333 and leave the question and answer unchanged. [9]

334 An adversarial sample can be defined as one which appears to be drawn from a particular class by
335 humans (or advanced cognitive systems) but fall into a different class in the feature space. [18]

336 maliciously crafted inputs that are undetectable by humans but that fool the algorithm into producing
337 undesirable behavior [11]

338 Adversarial examples are inputs to a predictive machine learning model that are maliciously designed
339 to cause poor performance [4]

340 One type of the vulnerabilities of machine learning algorithms is that an adversary can change the
341 algorithm output by subtly perturbing the input, often unnoticeable by humans. [8]

342 Adversarial attack on deep neural networks (DNNs) aims to slightly modify the inputs of DNNs and
343 mislead them to make wrong predictions [1]

344 For a given sample x and a trained DNN classifier model F , the attacker aims to craft an adversarial
345 sample $x^* = x + x$ by adding a perturbation x to x , such that $F(x^*) \neq F(x)$...In order to maintain the
346 utility of a text sample, we perturb the sample not only by directly modifying its words, but also
347 inserting new items (words or sentences) or removing some original ones from it. [13]

348 A.2 Adversarial Gain Perspectives of Prior Work

349 Here we examine various works and how they can fit into the adversarial gain perspective. We
350 already demonstrate how [1] and [4] can be measured in terms of adversarial gain. Rather than
351 non-overlapping text in [1], we can examine the semantic change of the output. Similarly, we can
352 examine how well the noise preserves the meaning of the input sentence in both cases. If the semantic
353 shift is too far, this discounts the shift in output it causes.

354 Generally, most text-based adversarial attacks constrain their inputs by in some way changing words
355 while retaining meaning. This includes negation [8], misspelling [8, 18, 13], changing word order [3],
356 replacing with with synonyms [1], or simply perturbing the word embeddings [14]. In many cases,
357 such constraints will preserve the meaning of the original text, but often too strict constraints can
358 result in lower success. For example, in [4] the word-based replacement with a strict synonym
359 constraint resulted in a low success rate in adversarial examples. In other cases, preservation of word
360 meaning is not guaranteed. In fact, prior work has used samples from the generated attacks posed
361 as surveys to determine whether meaning is preserved [9], but this has not typically been done in a
362 systematic way and Jia and Liang [9] found that in some cases meaning was not preserved. In another
363 example, negation of phrases does not preserve meaning and thus a model could be totally correct
364 in changing its output. In all attacks, it is possible to evaluate preservation of meaning by using a
365 well-defined embedding space (such as [2] as a start) and the cosine distance. The use of such a
366 distance as we do as part of adversarial gain, allows attacks to change meaning and account for this
367 when evaluating the change of the model output.

368 In evaluating the results of an adversarial attack, there are many measures used. For classification
369 tasks, a change in class label is typically used as a success criterion [21, 5, 11, 18, 3]. In other cases,
370 notions such as changes in some scoring or cost function are used [1, 9, 8, 14]. However, due to the
371 change in inputs if the cost relates to the original sentence and the meaning is *not preserved*, the cost
372 may be evaluating the wrong criterion without access to an oracle. Thus adversarial gain accounts
373 for this by discounting output performance changes by the distance from the input. In a well defined
374 feature space where inputs and outputs are correlated this ensures that as an adversarial input moves
375 away from its original meaning, this is accounted for in the evaluation criteria to some extent.

376 **B Extended Perspectives on Adversarial Gain**

377 Here we discuss extended properties and perspectives on adversarial gain.

378 **B.1 Possible Feature Spaces and Distance Metrics to Measure Gain**

379 There are a number of different priors that can be used to measure gain in different ways for different
380 tasks. While in the main text we examine sentiment classification tasks and text summarization,
381 others may be relevant in domains such as dialogue systems. For example, one can use sentiment
382 classification probability and the likelihood divergence (or a step function intersection) to measure
383 difference in output of a dialogue system, text summary system, or other generative model. Similarly,
384 various sentence embeddings can be used (InferSent, Doc2Vec, etc.). Word-wise word vector distance
385 can also be used. Each of these notions of adversarial gain essentially provide a different prior on
386 the stability of the systems in different ways. For example, it is likely that unless the sentiment of an
387 input to a dialogue system doesn't change dramatically, neither should the output.

388 **C Experimental Setup**

389 In our selection of text-based attacks, we examined which attacks provided easily available open-
390 source code. Many code to replicate experiments was either unavailable or we were unable to find.
391 We settled on two text-based attacks. We used the Seq2Sick attack on text summarization by Cheng
392 et al. [1] and the word-level sentiment classification attack by Ebrahimi et al. [4]. Scripts and full
393 instructions that we used to run the code from these papers is provided at: anonymized. More samples
394 with gain and distances provided can be found in the codebase provided.

395 **C.1 Text Summarization**

396 We use the pre-trained model and code for a text summarization model based on the Open
397 Neural Machine Translation toolkit (OpenNMT) [10] as provided by Cheng et al. [1] at
398 <https://github.com/cmhcbb/Seq2Sick>. We use the GigaWord corpus the authors reference from [17]
399 based on prior versions of the dataset [7, 15].

400 When we measure cosine distance, we use the inverse of cosine similarity to follow the intuition
401 the a distance metric should keep similar words closer together. Assuming that cosine similarity is
402 bounded $[0, 1]$, cosine distance is $1 - |similarity(x, y)|$.

403 **C.2 Sentiment Classification**

404 For sentiment classification we use the binary version of the SST dataset [19] called SST2.
405 This removes all neutral labels. This is the same dataset as used by [4]. We use their pro-
406 vided code for the word-level adversarial attack and SST2 pre-processing scripts found at:
407 <https://github.com/AnyiRao/WordAdver> and <https://github.com/AnyiRao/SentDataPre>. We use the
408 pre-trained convolutional neural network classification model provided by the authors and the attack
409 as provided in our accompanying instructions. The only change we make is that we remove the cosine
410 similarity requirement on replacement words. We do this because otherwise the attack only generates
411 attacks for 95 samples. Removing this requires generates attacks for all samples (though many are
412 not successful). We note that this allows words to be added by replacing padding characters, while
413 this differs slightly from the attack mentioned by [4], the authors there do discuss that this attack has
414 a low success rate particularly due to their restrictions. Because adversarial gain as a definition does
415 not require constraints, this allows us to consider the larger set of attacks.

416 **References**

- 417 [1] Minhao Cheng, Jinfeng Yi, Huan Zhang, Pin-Yu Chen, and Cho-Jui Hsieh. 2018. Seq2sick:
418 Evaluating the robustness of sequence-to-sequence models with adversarial examples. *arXiv*
419 *preprint arXiv:1803.01128*.

- 420 [2] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. 2017.
421 Supervised learning of universal sentence representations from natural language inference data.
422 *arXiv preprint arXiv:1705.02364*.
- 423 [3] Ishita Dasgupta, Demi Guo, Andreas Stuhlmüller, Samuel J Gershman, and Noah D Goodman.
424 2018. Evaluating compositionality in sentence embeddings. *arXiv preprint arXiv:1802.04302*.
- 425 [4] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2017. Hotflip: White-box adversarial
426 examples for nlp. *arXiv preprint arXiv:1712.06751*.
- 427 [5] Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. 2018. Black-box generation of
428 adversarial text sequences to evade deep learning classifiers. *arXiv preprint arXiv:1801.04354*.
- 429 [6] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing
430 adversarial examples. *arXiv preprint arXiv:1412.6572*.
- 431 [7] David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2003. English gigaword. *Linguistic
432 Data Consortium, Philadelphia*, 4:1.
- 433 [8] Hossein Hosseini, Sreeram Kannan, Baosen Zhang, and Radha Poovendran. 2017. Deceiving
434 google’s perspective api built for detecting toxic comments. *arXiv preprint arXiv:1702.08138*.
- 435 [9] Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension
436 systems. *arXiv preprint arXiv:1707.07328*.
- 437 [10] G. Klein, Y. Kim, Y. Deng, J. Senellart, and A. M. Rush. OpenNMT: Open-Source Toolkit for
438 Neural Machine Translation. *ArXiv e-prints*.
- 439 [11] Volodymyr Kuleshov, Shantanu Thakoor, Tingfung Lau, and Stefano Ermon. 2018. Adversarial
440 examples for natural language classification problems.
- 441 [12] Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. Understanding neural networks through
442 representation erasure. *arXiv preprint arXiv:1612.08220*.
- 443 [13] Bin Liang, Hongcheng Li, Miaoqiang Su, Pan Bian, Xirong Li, and Wenchang Shi. 2017. Deep
444 text classification can be fooled. *arXiv preprint arXiv:1704.08006*.
- 445 [14] Takeru Miyato, Andrew M Dai, and Ian Goodfellow. 2016. Adversarial training methods for
446 semi-supervised text classification. *arXiv preprint arXiv:1605.07725*.
- 447 [15] Courtney Napoles, Matthew Gormley, and Benjamin Van Durme. 2012. Annotated gigaword. In
448 *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale
449 Knowledge Extraction*, pages 95–100. Association for Computational Linguistics.
- 450 [16] Nicolas Papernot, Patrick McDaniel, Ananthram Swami, and Richard Harang. 2016. Craft-
451 ing adversarial input sequences for recurrent neural networks. In *Military Communications
452 Conference, MILCOM 2016-2016 IEEE*, pages 49–54. IEEE.
- 453 [17] Alexander M Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for
454 abstractive sentence summarization. *arXiv preprint arXiv:1509.00685*.
- 455 [18] Suranjana Samanta and Sameep Mehta. 2017. Towards crafting text adversarial samples. *arXiv
456 preprint arXiv:1707.02812*.
- 457 [19] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew
458 Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a
459 sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural
460 language processing*, pages 1631–1642.
- 461 [20] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Good-
462 fellow, and Rob Fergus. 2013. Intriguing properties of neural networks. *arXiv preprint
463 arXiv:1312.6199*.
- 464 [21] Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2017. Generating natural adversarial examples.
465 *arXiv preprint arXiv:1710.11342*.