

Alma Mater Studiorum - University of Bologna

---

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

MASTER'S DEGREE IN ARTIFICIAL INTELLIGENCE

*Final Thesis in*  
NATURAL LANGUAGE PROCESSING

**SynBA: A contextualized  
Synonym-Based adversarial Attack  
for Text Classification**

*Supervisor*

Prof. Paolo Torroni

*Co-supervisors*

Dr. Federico Ruggeri

Dr. Giulia De Poli

*Candidate*

Giuseppe Murro

---

ACADEMIC YEAR 2021-2022- THIRD SESSION



Dedica

“ The algorithms that cause AI systems to work so well are imperfect, and their systematic limitations create opportunities for adversaries to attack. At least for the foreseeable future, this is just a fact of mathematical life.” — Marcus Comiter



# Abstract

*“Happiness can be found even in the darkest of times, when one only remembers to  
turn on the light.”*  
– Dumbledore

# Thanks

First, I would like to express my deepest gratitude to Professor Luigi Di Stefano and Luca De Luigi for the guidance and support during the internship

Second, I would like to thank my parents for their moral support and for their patience.

Third, I would like to thank my girlfriend and friends for being patient with me and to put up with my complainings.

*Bologna, 06 December 2022*

Giuseppe Murro

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Topic definition . . . . .	2
1.2	Problem statement . . . . .	3
1.3	Research question . . . . .	3
1.4	Solution . . . . .	3
1.5	Thesis organization . . . . .	3
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Natural Language Processing . . . . .	4
2.1.1	Lexicon . . . . .	5
2.1.2	Word Embeddings . . . . .	6
2.1.3	Masked Language Models . . . . .	8
2.1.4	Text classification . . . . .	8
2.1.5	Sentiment analysis . . . . .	9
2.1.6	Natural language inference . . . . .	10
2.1.7	Seq2Seq . . . . .	10
2.2	Adversarial Machine Learning . . . . .	11
2.2.1	Adversarial examples . . . . .	13
2.2.2	Paradigm shift . . . . .	13
2.2.3	Taxonomy of textual adversarial attacks . . . . .	13
2.2.4	Adversarial attack methods from literature . . . . .	14
2.3	Machine Learning hardening . . . . .	14
2.3.1	Vanilla adversarial training . . . . .	14
2.3.2	Attacking to Training . . . . .	14
2.4	Text Attack . . . . .	14
2.4.1	Framework structure . . . . .	14
2.4.2	Attack components . . . . .	14
2.4.3	HuggingFace integration . . . . .	14

<b>3</b>	<b>Methodology</b>	<b>14</b>
3.1	Defined goal . . . . .	17
3.1.1	Problem to solve . . . . .	17
3.1.2	Research objective . . . . .	17
3.2	Research design . . . . .	17
3.2.1	Models to attack . . . . .	17
3.3	Proposed solution . . . . .	17
3.3.1	Intuition . . . . .	17
3.3.2	SynBA components . . . . .	17
3.3.3	Hyperparameter Tuning . . . . .	17
3.3.4	Candidates ranking calibration . . . . .	17
3.4	Evaluation metrics . . . . .	17
3.4.1	Attack metrics . . . . .	17
3.4.2	Quality metrics . . . . .	17
3.4.3	Performance metrics . . . . .	17
<b>4</b>	<b>Experimental results</b>	<b>17</b>
4.1	Data collection . . . . .	18
4.1.1	Experimental setup . . . . .	18
4.1.2	Datasets perturbed . . . . .	18
4.1.3	Model attacked . . . . .	18
4.2	Qualitative evaluation . . . . .	18
4.2.1	Results on rotten-tomatoes . . . . .	18
4.2.2	Results on imdb . . . . .	18
4.3	Quantitative evaluation . . . . .	18
4.3.1	Performances on rotten-tomatoes . . . . .	18
4.3.2	Performances on imdb . . . . .	18
4.4	Human evaluation . . . . .	18
<b>5</b>	<b>Final discussions</b>	<b>18</b>

# List of Figures

2.1	Components of NLP . . . . .	5
2.2	Bag-of-words representation of a a movie review, where only the frequency of each word is considered. . . . .	7
2.3	An example of a sequence-to-sequence model for machine translation.	11
2.4	Examples of Artificial Intelligence Attacks . . . . .	12



# List of Tables

## Chapter 1

# Introduction

*In this section we will present the summarized content of the whole thesis.*

### 1.1 Topic definition

To understand why AI systems are vulnerable to the same weakness, we must briefly examine how AI algorithms, or more specifically the machine learning techniques they employ, “learn.” Just like the reconnaissance officers, the machine learning algorithms powering AI systems “learn” by extracting patterns from data.

## 1.2 Problem statement

## 1.3 Research question

## 1.4 Solution

## 1.5 Thesis organization

**First chapter** introduces the general content about thesis and gives a short presentation of the topic, the problem and the solution we propose;

**Second chapter** a deepening about the theoretical foundations used during the stage and the project;

**Third chapter** presents the datasets used during for the training and the testing of the model;

**Fourth chapter** presents the experiments did during to develop the system;

**Fifth chapter** discusses about the results and possible future developments.

During the drafting of the essay, following typography conventions are considered:

- the acronyms, abbreviations, ambiguous terms or terms not in common use are defined in the glossary, in the end of the present document;
- the first occurrences of the terms in the glossary are highlighted like this: **word**;
- the terms from the foreign language or jargon are highlighted like this: *italics*.

# Chapter 2

## Background

*In this chapter we will present the theoretical knowledge useful to understand the content from successive chapters.*

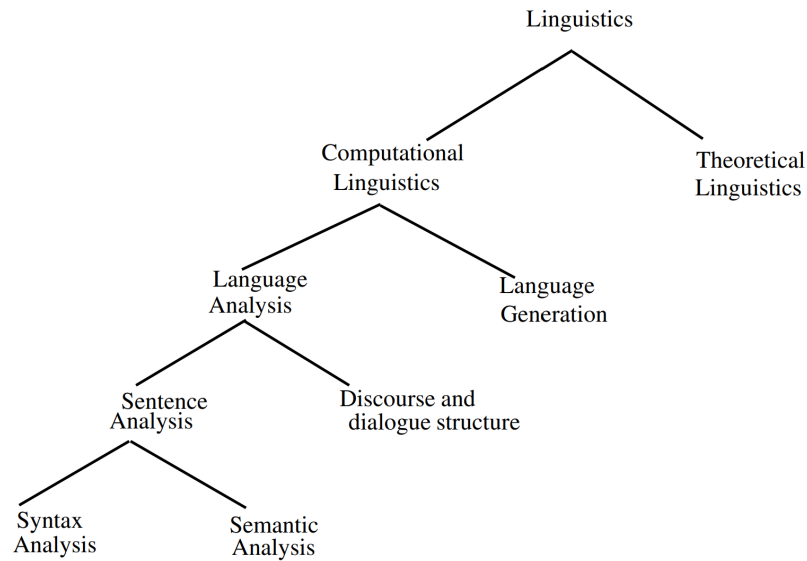
## 2.1 Natural Language Processing

The field of [Natural Language Processing \(NLP\)](#), also known as computational linguistics, is a branch of Artificial Intelligence focused on the technology of processing language. It encompasses a variety of topics, which involves the engineering of computational models and processes to solve practical problems in understanding and generating human languages. These solutions are used to build useful software.

The linguistics computational has two branches—computational linguistics and theoretical linguistics. The computational linguistics has been concerned with developing algorithms for handling a useful range of natural language as input. While the theoretical linguistics has focused primarily on one aspect of language performance, grammatical competence—how people accept some sentences as correctly following grammatical rules and others as ungrammatical. They are concerned with language universals—principals of grammar which apply to all natural languages [Cole:1996].

Computational linguistics is concerned with the study of natural language analysis and language generation. Further, the language analysis is divided into two domains, namely sentence analysis, and discourse and dialogue structure. Much more is known about the processing of individual sentences than about the determination of discourse structure. Any analysis of discourse structure requires a prerequisite as an analysis of the meaning of individual sentences. However, it is a fact that for many applications, thorough analysis of discourse is not mandatory, and the sentences can be understood without that [grishman \_ computational \_ 1986].

The sentence analysis is further divided into syntax analysis and semantic analysis. The overall objective of sentence analysis is to determine what a sentence “means” . In practice, this involves translating the natural language input into a language with simple semantics, for example, formal logic, or into a database command language. In most systems, the first stage is syntax analysis. Figure 2.1 shows the relations among different components of NLP [Chowdhary2020].



**Figure 2.1:** Components of NLP

Some of the common applications of NLP are: Classification of text into categories, Index and search large texts, Automatic translation, Information extraction, Automatic summarization, Question answering, Knowledge acquisition, and Text generations/dialogues. Some of those tasks are discussed in sections [2.1.4](#), [2.1.5](#), [2.1.6](#), [2.1.7](#).

### 2.1.1 Lexicon

Dictionaries are special texts whose subject matter is a language, or a pair of languages in the case of a bilingual dictionary. The purpose of dictionaries is to provide a wide range of information about words— etymology, pronunciation, stress, morphology, syntax, register—to give definitions of senses of words, and, in so doing, to supply knowledge not just about language, but about the world itself.

The term “dictionary” is typically related to printed wordbook for human readers. Instead “[lexicon](#)” will refer to the component of a NLP system that contains information (semantic, grammatical) about individual word strings [[Guthrie\\_ComACM96](#)].

A lexicon which provides an effective combination of traditional lexicographic information and modern computing is called WordNet [[miller1995wordnet](#)]. It is an online lexical database designed for use under program control. English nouns, verbs, adjectives, and adverbs are organized into sets of synonyms, each representing a lexicalized concept. WordNet contains more than 118,000 different word forms

and more than 90,000 different word senses. Approximately 40% of the words in WordNet have one or more synonyms.

The cognitive synonyms which are called synsets are presented in the database with lexical and semantic relations. WordNet includes the following semantic relations:

- **Hypernymy:** A hypernym is a word that is more general than the word in question. For example, the hypernym of “dog” is “canine” .
- **Hyponymy:** A hyponym is a word that is more specific than the word in question. For example, the hyponym of “dog” is “poodle” .
- **Synonymy:** A synonym is a word that has the same meaning as the word in question. For example, the synonym of “good” is “well” .
- **Antonymy:** An antonym is a word that has the opposite meaning of the word in question. For example, the antonym of “good” is “bad” .

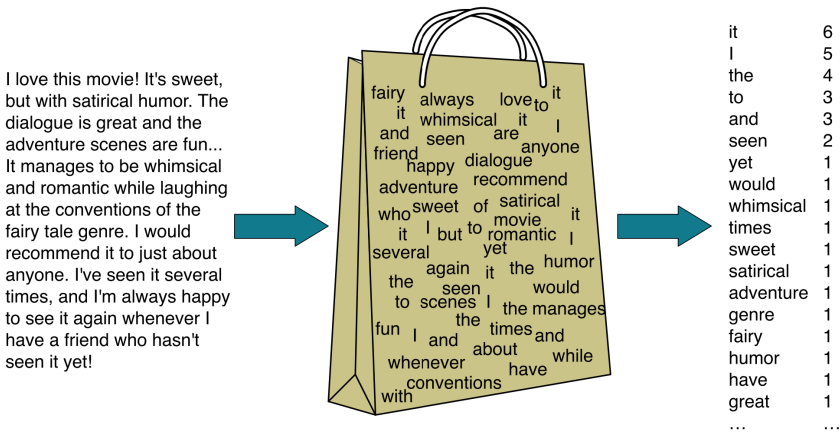
### 2.1.2 Word Embeddings

The way machine learning models process data is different from how humans do. For example, we can easily understand the text “I saw a cat” , but our models can not — they need vectors of features. Such vectors, called [word embeddings](#), are representations of words which can be fed into a model.

#### 2.1.2.1 Bag of words

A common approach to represent a text document is to use a column vector of word counts. This embedding is often called a [bag-of-words](#), because it includes only information about the count of each word, and not the order in which the words appear.

The bag-of-words representation ignores grammar, sentence boundaries, paragraphs — everything but the words. Yet the bag of words model is surprisingly effective for text classification. In the example in the Figure [2.2](#), instead of representing the word order in all the phrases like “I love this movie” and “I would recommend it” , we simply note that the word I occurred 5 times in the entire excerpt, the word it 6 times, the words love, recommend, and movie once, and so on [[Jurafsky2009](#)].



**Figure 2.2:** Bag-of-words representation of a movie review, where only the frequency of each word is considered.

2.1.2.2 Dense embeddings

Bag-of-words embeddings are sparse and long vector with dimensions corresponding to words in the vocabulary or documents in a collection. A more powerful word representation is a dense vector, where instead of mostly-zero counts, the values will be real-valued numbers that can be negative. It turns out that dense vectors work better in every NLP task than sparse vectors.

Bengio et al. [conf/nips/BengioDV00] presented a model which learned word representations using distributed representation. Authors presented a neural model which obtains word representations as to the product while training language model. The popularity of word representation methods are due to two famous models, Word2Vec [mikolov2013efficient] and GloVe [pennington2014glove].

2.1.2.3 Contextual embeddings

To address the issue of polysemous and the context-dependent nature of words, we need distinguish the semantics of words in different contexts.

Contextualised word embeddings are variable vector that are dependent on the context in which the word is used. So, representations of a given word are multiple and are directly computed from their context. The context of a word is usually composed by the words surrounding it.

These contextualized representations are set to the hidden states of a deep neural model, which is trained as a language model. By running the language model,

we obtain contextualized word representations, which can then be used as the base layer in a supervised neural network for any task. This approach yields significant gains over pretrained word embeddings on several tasks, presumably because the contextualized embeddings use unlabeled data to learn how to integrate linguistic context into the base layer of the supervised neural network.

### 2.1.3 Masked Language Models

A [Masked language model \(MLM\)](#) is a pre-training technique which first masks out some tokens from the input sentences and then trains the model to predict the masked tokens by the rest of the tokens. A special [MASK] token is used to replace some words randomly into the original text.

Masked Language Modelling is usually solved as classification problem. We feed the masked sequences to a neural encoder whose output vectors are further fed into a softmax classifier to predict the masked token.

The most popular MLM is [BERT \[devlin2018bert\]](#), which is a bidirectional encoder representation from a particular deep learning architecture called [transformer \[vaswani2017attention\]](#). It uses self-supervised training on the masked language modeling and next sentence prediction tasks to learn/produce contextual representations of words.

Concurrently, there are multiple research proposing different enhanced versions of MLM to further improve on BERT. Instead of static masking, RoBERTa [\[liu2019roberta\]](#) improves BERT by dynamic masking. While other models aim to optimize BERT's performance, DistilBERT has a different goal. Its target is to reduce the large size and enhance the speed of BERT while still keeping as much strength as possible. DistilBERT [\[sanh2019distilbert\]](#) reduces the size of  $BERT_{BASE}$  by 40%, enhances the speed by 60% while retaining 97% of its capabilities. ALBERT [\[lan2019albert\]](#) also reduces the model size of BERT, it does not have to trade-off the performance. Compared to DistilBERT, which uses BERT as the teacher for its distillation process, ALBERT is trained from scratch (just like BERT).

### 2.1.4 Text classification

Classification lies at the heart of both human and machine intelligence. Deciding what letter, word, or image has been presented to our senses, recognizing faces or voices, sorting mail, assigning grades to homeworks; these are all examples of assigning a category to an input. In this section we introduce text classification, the task of assigning a label or category to an entire text or document.

Given a text document, assign it a discrete label  $y \in Y$ , where  $Y$  is the set of possible labels. Text classification has many applications, from spam filtering to the analysis of electronic health records, or the categorization of news articles.

Classification is essential for tasks below the level of the document as well. An example of this is period disambiguation (deciding if a period is the end of a sentence or part of a word), or word tokenization (deciding if a character should be a word boundary). Even language modeling can be viewed as classification: each word can be thought of as a class, and so predicting the next word is classifying the context-so-far into a class for each next word. A part-of-speech tagger classifies each occurrence of a word in a sentence as, e.g., a noun or a verb.

The goal of classification is to take a single observation, extract some useful features, and thereby classify the observation into one of a set of discrete classes.

One method for classifying text is to use handwritten rules. There are many areas of language processing where handwritten rule-based classifiers constitute a state-of-the-art system, or at least part of it. Rules can be fragile, however, as situations or data change over time, and for some tasks humans aren't necessarily good at coming up with the rules. Most cases of classification in language processing are instead done via supervised machine learning, where an algorithm learn how to map from an observation to a correct output [Jurafsky2009].

Many kinds of machine learning algorithms are used to build classifiers. Formerly, statistical and machine learning approaches, such as naïve Bayes, k-nearest neighbors, hidden Markov models, conditional random fields (CRFs), decision trees, random forests, and support vector machines, were widely used to design classifiers. However, during the past several years, there has been a wholesale transformation, and these approaches have been entirely replaced, or at least enhanced, by neural network models [surveyNlpDeepLearning].

### 2.1.5 Sentiment analysis

A popular application of text classification is sentiment analysis, the extraction of sentiment, the positive or negative orientation that a writer expresses toward some object. A review of a movie, book, or product on the web expresses the author's sentiment toward the product, while an editorial or political text expresses sentiment toward a candidate or political action. Extracting consumer or public sentiment is thus relevant for fields from marketing to politics. [Jurafsky2009]

The simplest version of sentiment analysis is a binary classification task, and the words of the review provide excellent cues. Consider, for example, the following



phrases extracted from positive and negative reviews of movies and restaurants. Words like great, richly, awesome, and pathetic, and awful and ridiculously are very informative cues:

- + ...*zany characters and richly applied satire, and some great plot twists*
- *It was pathetic. The worst part about it was the boxing scenes...*
- + ...*awesome caramel sauce and sweet toasty almonds. I love this place!*
- ...*awful pizza and ridiculously overpriced...*

The area of sentiment analysis it is becoming increasingly popular and utilizing deep learning. Applications are varied, including product research, futures prediction, social media analysis, and classification of spam [ZhengWG18]. Good results were obtained using an ensemble, including both LSTMs and CNNs [Cliche17]. But the current trend in state-of-the-art models in all application areas is to use pretrained stacks of transformer units in some configuration, whether in encoder-decoder configurations or just as encoders.

### 2.1.6 Natural language inference

The task of [Natural Language Inference \(NLI\)](#), also known as recognizing textual entailment, asks a system to evaluate the relationships between the truth-conditional meanings of two sentences or, in other words, decide whether one sentence follows from another. The relationship can be entailment, contradiction, or neutral.

Specifically, natural language inference (NLI) is concerned with determining whether a natural language hypothesis  $h$  can be inferred from a premise  $p$ , as depicted in the following example from [Manning2009NaturalLI], where the hypothesis is regarded to be entailed from the premise.

$p$ : *Several airlines polled saw costs grow more than expected, even after adjusting for inflation.*

$h$ : *Some of the companies in the poll reported cost increases.*

### 2.1.7 Sequence-to-Sequence models

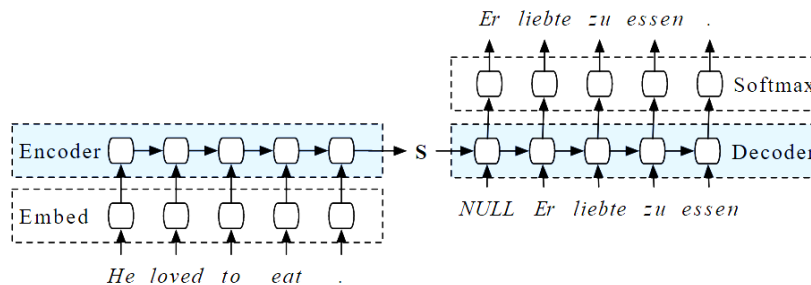
All the tasks we have discussed so far are classification-based, where the input is a text and the output is a label. However, there are many tasks where the input and output are both sequences of tokens. For example, machine translation, summarization, and question answering are all tasks where we want to generate a sequence in human-like language as output.

A [Sequence to Sequence \(seq2seq\)](#) model is a special class of [Recurrent Neural Network \(RNN\)](#) architectures that can be used to solve these tasks.

In the general case, input sequences and output sequences have different lengths (e.g. machine translation) and the entire input sequence is required in order to start predicting the target. This requires a more advanced setup:

- An *encoder* processes the input sequence and returns its own internal state. This vector is called the context vector.
- A *decoder* is a neural network that takes the context vector as input and outputs a sequence of tokens. It is trained to predict the next characters of the target sequence, given previous characters of the target sequence.

An example of this architecture is shown in Figure 2.3.



**Figure 2.3:** An example of a sequence-to-sequence model for machine translation.

## 2.2 Adversarial Machine Learning

Deep Learning algorithms have achieved the state-of-the-art performance in many tasks. However, the interpretability of deep neural networks is still unsatisfactory as they work as black boxes, which means it is difficult to get intuitions from what each neuron exactly has learned. One of the problems of the poor interpretability is evaluating the robustness of deep neural networks.

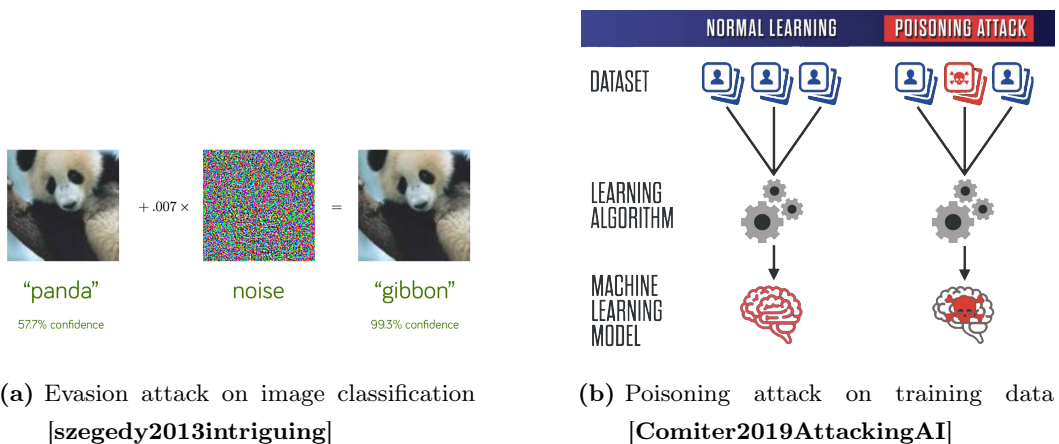
Adversarial Machine Learning is a collection of techniques to train neural networks on how to spot intentionally misleading data or behaviors. This differs from the standard classification problem in machine learning, since the goal is not just to spot “bad” inputs, but preemptively locate vulnerabilities and craft more flexible learning algorithms.

The objective of an adversary could be to attempt to manipulate either the data collection or the data processing in order to corrupt the target model, thus tampering the original output.

Unlike traditional cybersecurity attacks, these weaknesses are not due to mistakes made by programmers or users. They are just shortcomings of the current state-of-the-art methods. Put more bluntly, the algorithms that cause AI systems to work so well are imperfect, and their systematic limitations create opportunities for adversaries to attack. At least for the foreseeable future, this is just a fact of mathematical life [Comiter2019AttackingAI].

Two main types of AI attacks can be defined according to the time at which the attack happens [journals/corr/abs-1810-00069]:

- **Adversarial Attacks (Evasion):** this is the most common type of attack in the adversarial setting. The adversary tries to evade the system by adjusting malicious samples during testing phase. This setting does not assume any influence over the training data. In Figure 2.4a is depicted how adding an imperceptible and carefully constructed noise to the input originally recognized as “panda” with 57.7% confidence, we can change the classification output given by the same neural network toward another target (in the example “gibbon” with 99.3% confidence).
- **Data Poisoning Attacks:** This type of attack, known as contamination of the training data, is carried out at training phase of the machine learning model. An adversary tries to inject skilfully crafted samples to poison the system in order to compromise the entire learning process. The Figure 2.4b shows as in normal machine learning (left), the learning algorithm extracts patterns from a dataset, and the “learned” knowledge is stored in the machine learning model—the brain of the system. In a poisoning attack (right), the attacker changes the training data to poison the learned model.



**Figure 2.4:** Examples of Artificial Intelligence Attacks

In this thesis, we will focus on the first type of attack, the Adversarial Attacks. Over past few years, researchers [goodfellow2014explaining, szegedy2013intriguing] used small unperceivable perturbations to evaluate the robustness of deep neural networks and found that they are not robust to these perturbations.

### 2.2.1 Adversarial examples

Szegedy et al. [szegedy2013intriguing] first evaluated the state-of-the-art deep neural networks used for image classification with small generated perturbations on the input images. They found that the image classifiers were fooled with high probability, but human judgment is not affected. The perturbed image pixels were named *adversarial examples* and this notation is later used to denote all kinds of perturbed samples in a general manner. Formally, adversarial example  $x'$  is an example created via worst-case perturbation of the input to a deep learning model. An ideal deep neural network would still assign correct class  $y$  (in the case of classification task) to  $x'$ , while a victim deep neural network would have high confidence on wrong prediction of  $x'$ .  $x'$  can be formalized as:

$$\begin{aligned} x' &= x + \eta, & f(x) &= y, & x &\in X, \\ f(x') &\neq y, \\ \text{or } f(x') &= y', & y &\neq y' \end{aligned} \tag{2.1}$$

where  $\eta$  is the worst-case perturbation. The goal of the adversarial attack can be deviating the label to incorrect one ( $f(x') \neq y$ ) or specified one ( $f(x') = y'$ ) [journals/tist/ZhangSAL20].

### 2.2.2 Paradigm shift: from CV to NLP

Adversarial examples were first proposed for attacking DNNs for object recognition in the [Computer Vision \(CV\)](#) community.

### 2.2.3 Taxonomy of textual adversarial attacks

Adversarial attacks at the testing time do not tamper with the targeted model but rather forces it to produce incorrect outputs. The effectiveness of such attacks is determined mainly by the amount of information available to the adversary about the model. Testing phase attacks can be broadly classified into either White-Box or Black-Box attacks. Before discussing these attacks, we provide a formal definition of a training procedure for a machine learning model. Let us consider a target machine learning model  $f$  is trained over input pair  $(X, y)$  from the data distribution

$\mu$  with a randomized training procedure train having randomness  $r$  (e.g., random weight initialization, dropout, etc.). The model parameters  $\theta$  are learned after the training procedure. More formally, we can write:  $\theta \leftarrow \text{train}(f, X, y, r)$  Now, let us understand the capabilities of the white-box and black-box adversaries with respect to this definition. An overview of the different threat models have been shown in Figure. 4 White-Box Attacks. In white-box attack on a machine learning model, an adversary has total knowledge about the model ( $f$ ) used for classification (e.g., type of neural network along with

#### 2.2.4 Adversarial attack methods from literature

##### 2.2.4.1 TextFooler

##### 2.2.4.2 BERT-based attacks

### 2.3 Machine Learning hardening

#### 2.3.1 Vanilla adversarial training

#### 2.3.2 Attacking to Training

### 2.4 Text Attack

#### 2.4.1 Framework structure

#### 2.4.2 Attack components

#### 2.4.3 HuggingFace integration

## Chapter 3

# Methodology

*In this chapter we will present the datasets created and used for the visual odometry.*

- Introduction

- Research Design
- Research Questions and Hypotheses
- Setting and Sample
- Data Collection
- Data Analysis
- Conclusion

demonstration of fit between methods chosen and research question(s) rationale for choosing materials, methods and procedures details of materials, equipment and procedures that will allow others to: replicate experiments understand and implement technical solutions



### 3.1 Defined goal

#### 3.1.1 Problem to solve

#### 3.1.2 Research objective

### 3.2 Research design

#### 3.2.1 Models to attack

### 3.3 Proposed solution

#### 3.3.1 Intuition

#### 3.3.2 SynBA components

##### 3.3.2.1 Search Method

##### 3.3.2.2 Transformation

##### 3.3.2.3 Constraints

##### 3.3.2.4 Goal Function

#### 3.3.3 Hyperparameter Tuning

#### 3.3.4 Candidates ranking calibration

### 3.4 Evaluation metrics

#### 3.4.1 Attack metrics

#### 3.4.2 Quality metrics

#### 3.4.3 Performance metrics

## Chapter 4

# Experimental results

*In this chapter we will discuss about different models and different prediction strategies.*



## 4.1 Data collection

### 4.1.1 Experimental setup

### 4.1.2 Datasets perturbed

### 4.1.3 Model attacked

## 4.2 Qualitative evaluation

### 4.2.1 Results on rotten-tomatoes

### 4.2.2 Results on imdb

## 4.3 Quantitative evaluation

### 4.3.1 Performances on rotten-tomatoes

### 4.3.2 Performances on imdb

## 4.4 Human evaluation

# Chapter 5

# Final discussions

*In this chapter we will discuss the results achieved, future developments and personal comments.*

A clear answer to your research question or hypothesis  
Summary of the main findings or argument  
Connections between your findings or argument to other research  
Explanation and significance of the findings  
Implications of the findings  
Limitations of the research and methodology  
Recommendations for future research

Summary of Findings  
Limitations of the research  
Suggestions for Future Research  
Conclusion