



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

BACHELOR THESIS

Jakub Hejhal

**Exploring the vulnerabilities of
commercial AI systems against
adversarial attacks**

Department of Theoretical Computer Science and Mathematical Logic

Supervisor of the bachelor thesis: Mgr. Roman Neruda, CSc.

Study programme: Computer Science

Study branch: General Computer Science

Prague 2021

I declare that I carried out this bachelor thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date
Author's signature

I would like to thank Mgr. Roman Neruda, CSc. for his patience, helpful advice and guidance throughout my thesis work. I would also like to thank my family, my friends and my girlfriend for supporting me and for providing me with love and care that allowed me to push through occasionally difficult times.

Title: Exploring the vulnerabilities of commercial AI systems against adversarial attacks

Author: Jakub Hejhal

Department: Department of Theoretical Computer Science and Mathematical Logic

Supervisor: Mgr. Roman Neruda, CSc., Department of Theoretical Computer Science and Mathematical Logic

Abstract: Abstract. TODO

Keywords: Machine learning Deep learning Adversarial attack Black-box

Contents

1	Introduction	3
	Introduction	3
1.1	The rise of artificial neural networks	3
1.2	Concerns	3
1.3	Adversarial attacks	4
2	Preliminaries	5
2.1	Machine learning	5
2.2	Deep neural networks	5
2.3	DNN training	5
2.4	Convolutional neural networks	5
2.5	Adversarial attacks on DNNs	5
2.6	Threat models	5
3	Related Work	6
3.1	Title of the first subchapter of the third chapter	6
4	Our approach	7
4.1	Adapting off-the-shelf attack algorithms to partial-information setting	7
4.1.1	Object-organism binary classification mapping	7
4.2	PoC black-box GVision attacks	9
4.2.1	TREMBA	9
4.2.2	RayS	9
4.2.3	SquareAttack	9
4.2.4	Sparse-RS	9
4.3	The need for query-efficient attack	9
4.4	Leveraging transferability	10
4.4.1	Transfer attacks provide better seeds for blackbox optimization	10
4.4.2	Train, validate, test	10
4.4.3	Local training and validation is cheap	10
4.4.4	Multiple candidates save queries	10
4.5	The need for attack pipeline	11
4.5.1	Possibility of multiple blackbox workers	11
4.5.2	The need for unified attack and model API	11
4.6	AdvPipe	11
4.6.1	The vision	11
4.6.2	The reality	11
4.7	Implementation	12
4.7.1	Wrapping whitebox and blackbox models	12

5	Experiments	13
5.1	Blackbox PoC	13
5.1.1	TREMBA	13
5.1.2	RayS	13
5.1.3	SquareAttack	13
5.2	Local transferability experiments	13
5.2.1	Choice of dataset	13
5.2.2	Choice of local models	13
5.2.3	Baseline	14
5.2.4	Baseline	14
5.2.5	Baseline	14
5.2.6	Baseline	14
5.3	Transferability evaluation on GVision	14
5.3.1	Choice of evaluation metrics	14
5.3.2	Wordcloud	14
	Conclusion	15
	Bibliography	16
	List of Figures	18
	List of Tables	19
	List of Abbreviations	20
A	Attachments	21
A.1	First Attachment	21

1. Introduction

1.1 The rise of artificial neural networks

In recent years there has been an enormous surge in applications of artificial intelligence technologies based on neural networks to various fields. One of the most significant milestones that kickstarted today's AI revolution has undoubtedly been the year 2012. ImageNet Large Scale Visual Recognition Challenge (ILSVRC), which hand-crafted specialized image-labeling systems have previously dominated, was won by AlexNet with its CNN architecture.

Artificial neural networks inspired by their biological analog had been known and researched for a long time. Even though they enjoyed much enthusiasm initially, there has been a time period (called an "AI winter" by some) when they had been neglected as an unpromising direction towards general intelligence. More classical approaches like SVM and rule-based AI systems showed better performance and computational efficiency on AI benchmarks of the time.

What changed the game has been the available computational power that came with Moore's law and the usage of GPUs, mainly their parallel nature in accelerating matrix multiplication operations that neural networks use heavily.

Another factor that helped the rise of neural networks has been the availability of large datasets like ImageNet, which contains 1,281,167 images for training and 50,000 images for validation organized in 1,000 categories.

The availability of large amounts of labeled and unlabeled data, sometimes referred to as "Big data," is only getting better. A large part of our lives has moved to the virtual space thanks to the internet. Businesses had started to realize the value of the enormous amounts of traffic generated every day, and they are increasingly trying to figure out how to take advantage of it.

What was previously limited to academic circles had quickly become mainstream. Artificial neural networks have proven to be very versatile and have quickly been successfully applied to a wide range of problems. New neural network architectures and new training regimes allowed training deeper networks, which gave rise to a new field of machine learning called "Deep learning."

Deep neural networks have shown state-of-the-art performance in machine translation, human voice synthesis and recognition, drug composition, particle accelerator data analysis, recommender systems, algorithmic trading, reinforcement learning, and many other areas.

1.2 Concerns

Large-scale deployment of neural network systems has been criticized by many for their inherent unexplainability. It is often hard to pinpoint why neural network behaves in some way and why it makes certain decisions. One problem is the role of training data, where possible biases may be negatively and unexpectedly reflected in the behavior of the AI system. Another problem is the performance on out-of-distribution examples, where network inference occurs on different kinds of data than used in the training stage.

Those concerns lead people to study the robustness of AI systems. It turned out that image recognition CNN networks are especially susceptible to the so-called adversarial attacks, where the input is perturbed slightly, but the output of the network changes wildly. Similar kinds of attacks have since been demonstrated in other areas like speech recognition, natural language processing, or reinforcement learning.

1.3 Adversarial attacks

This vulnerability of neural networks has led to a cat-and-mouse game of various attack strategies and following defenses proposed to mitigate them.

Neural networks can be attacked at different stages:

- training
- testing
- deployment

Training attacks exploit training dataset manipulation, sometimes called dataset poisoning, to change the behavior of the resulting trained network.

Testing attacks do not modify trained neural network, but often use the knowledge of the underlying architecture to craft adversarial examples which fool the system.

Deployment attacks deal with real black-box AI systems, where the exact details of the system are usually hidden from the attacker. Nevertheless, partly because similar neural network architectures are used in the same classes of problems, some vulnerabilities in testing scenarios can still be exploited in deployment, even though the exact network parameters, architecture, and output are unknown to the attacker.

The purpose of this thesis is to explore the applicability of certain classes of testing attacks on real-world deployed AI systems. For simplicity, many kinds of SOTA adversarial attacks have only been explored in the testing regime but have not been applied to truly black-box systems.

The main aim of the thesis will be to test different types of testing attacks on AI SaaS providers like Google Vision API, Amazon recognition, or Clarify. Understandably, attacking an unknown system will be more challenging than attacking a known neural network in the testing stage. We will try to measure this attack difficulty increase. This information could prove helpful in selecting the most promising attack to a specific situation.

Many SOTA testing attacks were not designed to attack specific deployed systems, so some attacks will need to be slightly changed to be used. We will explore different ways to modify existing attacks and evaluate them.

If some of those services are proven to be vulnerable, this would have a massive impact on all downstream applications using those SaaS APIs. For instance, content moderation mechanisms, which rely mainly on automatic detection, could be circumvented.

2. Preliminaries

In this section we briefly introduce and explain the necessary theoretical background, upon which we build later on.

We first define the concepts of Machine learning (ML) and Deep neural networks (DNNs). We explain, how DNNs can learn patterns from data by using powerful gradient-based stochastic optimization algorithm called Stochastic gradient descent (SGD). Then we talk about a subset of DNNs that perform very well on image data called Convolutional neural networks (CNNs). Finally, we describe the systemic vulnerability of DNNs and we define and explain adversarial attacks and defenses, adversarial examples, adversarial robustness and different adversarial attack threat models.

2.1 Machine learning

Goodfellow et al. [2016]

2.2 Deep neural networks

2.3 DNN training

2.4 Convolutional neural networks

2.5 Adversarial attacks on DNNs

2.6 Threat models

3. Related Work

3.1 Title of the first subchapter of the third chapter

TODO: rewrite your review papernotes.md here.

4. Our approach

4.1 Adapting off-the-shelf attack algorithms to partial-information setting

Cloud-based image classifiers don't usually classify input images into a fixed number of classes. They instead output variable-length list of probable labels with scores. And what's worse, those scores aren't even probabilities, because they don't sum up to one.

Most of current score-based SoTA adversarial attacks assume that we have access to all output logits of the attacked target classifier. If we want to use them, we need to map somehow the cloud's variable-length score-output to fixed-length vector, which will simulate logits output of a standard CNN classifier.

4.1.1 Object-organism binary classification mapping

To simplify the experiments, we define a simple benchmarking attack scenario:

**Given an image containing a living thing, fool the target into
classifying it as a non-living object.**

This choice makes our simulated classifier a binary one. It should assign each input image (x) an organism score $o_{organism}(x)$ and object score $o_{object}(x)$. This 2-D score vector $(o_{object}(x), o_{organism}(x))$ is further denoted by $o(x)$ for simplicity.

We chose this split, because we perform majority of our experiments on ImageNet validation dataset (Deng et al. [2009]) and ImageNet is relatively balanced between those two semantic categories.

Why ImageNet? Despite the dataset being quite old, it is still the most heavily used dataset in the research community and the majority of freely available pretrained models are pretrained on it.

Furthermore, when $\|C\| = 2$ (where C is a set of output categories), targeted and untargeted attack scenarios don't differ anymore and are neatly unified.

Adapting the attack algorithm to a different attack objective only requires swapping the label mapping layer.

Imagenet category mapping

Classic ImageNet dataset contains real-world photos, each corresponding to one and only one classification category out of 1000 possible categories. Each ImageNet category c corresponds to a unique wordnet synset $w(c)$. These synsets are rather specific, but we can take a look at their set of hypernyms $h(w(c))$. If this hypernym set $h(w(c))$ contains the *organism* synset, it should be an organism, otherwise c is probably an object.

Written more rigorously, we map each ImageNet category c into $\{organism, object\}$ using the following mapping $m_{local}(c)$:

$$m_{local}(c) = \begin{cases} organism & organism \in h(w(c)) \\ object & organism \notin h(w(c)) \end{cases}$$

Cloud label mapping

This situation isn't so clear-cut in the case of general labels returned by cloud classifier. Labels might not even be single words, but whole sentences. We therefore resort to a more powerful label classification method and use a GPT-2 transformer (Radford et al. [2019]) for this matter. More specifically, we use the HuggingFace (Wolf et al. [2020]) zero-shot classification pipeline to encode text labels into embedding vector space and then compute their similarity $s(l_{cloud}, l_{gold})$ with carefully chosen set of organism labels $L_{organism} = \{animal, species\}$ and with a set of object labels $L_{object} = \{object, instrument\}$. The resulting binary cloud label mapping $m_{cloud}(l_{cloud})$ is defined as follows:

$$m_{cloud}(l_{cloud}) = \begin{cases} organism & \arg \max_{c \in (L_{organism} \cup L_{object})} s(l_{cloud}, c) \in L_{organism} \\ object & \arg \max_{c \in (L_{organism} \cup L_{object})} s(l_{cloud}, c) \in L_{object} \end{cases}$$

From now on, by $m(c)$ we mean either $m_{local}(c)$ or $m_{cloud}(l_{cloud})$ where the distinction wouldn't make any difference.

Computing the adversarial loss

There is one more step we have to do to transparently simulate a binary classifier with 2 output logits.

By passing the model outputs through the separation mapping $m(c)$ we obtain two score sets: $S_{organism}$ and S_{object} . In the local case:

$$\|S_{organism}\| + \|S_{object}\| = \|L_{ImageNet}\| = 1000$$

In the case of a general cloud classifier these sets have variable sizes and one or both can be even empty.

There are multiple sensible ways to produce output logits vector $o(x)$ that would roughly correspond to an organism binary classifier output and which could be attacked using standard untargeted adversarial attacks.

Just to name a few choices for $o(x)$ that could intuitively work:

1. top-1 score: $\max S$
2. sum of logits: $\sum S$
3. sum of un-normalized probabilities: $\sum \exp(S)$
4. log-sum of probabilities: $\log \sum \text{softmax}(S)$
5. ...

But in the end we want to achieve a misclassification of the original organism image x_{org} .

If we go with 1) and produce output vector $o(x) = (\max S_{object}, \max S_{organism})$, misclassification is achieved when $o(x)_0 > o(x)_1$. We can define a margin loss objective $\mathcal{L}_{margin}(x, \kappa) = \max S_{organism}(x) - \max S_{object}(x) + \kappa$ with additional

parameter κ , by which we adjust our requirement for the degree of misclassification.

Another hint that this might be a solid choice comes from the Carlini and Wagner [2017], where they discuss the choice of optimization objective. They try a number of different alternatives, but in the end they conclude that optimizing the margin loss works the best, so we stick with it.

4.2 PoC black-box GVision attacks

We first explored the viability and sample-efficiency of current SoTA black-box attacks. We ran the following against Google Vision API image classifier.

- 4.2.1 TREMBA (Huang and Zhang [2020])
- 4.2.2 RayS (Chen and Gu [2020])
- 4.2.3 SquareAttack (Andriushchenko et al. [2020])
- 4.2.4 Sparse-RS (Croce et al. [2020])

TODO: describe each attack briefly and show sample images

4.2.1 TREMBA

4.2.2 RayS

4.2.3 SquareAttack

4.2.4 Sparse-RS

We go into more details in the Experiments 5

4.3 The need for query-efficient attack

In the previous section 4.2 we empirically showed, that Google Vision API isn't 100% robust to iterative blackbox attacks. But although the previously mentioned blackbox attacks are often successful in producing adversarial image, query count to the target may be often unacceptably high. Huge query stress to the target is troublesome for several reasons:

- High cost - 1.5\$ per 1000 queries
- Raising suspicion
- Often unrealistic in practical setting

The problem is that these blackbox attacks (with the exception of TREMBA) mostly rely on random search and don't make use of the gradient similarity of various CNN models. The high dimensionality of the input data doesn't make the blackbox optimization task easy. Current SoTA blackbox attacks that don't use

any gradient priors are already at the query efficiency limit with their medium queries being often less than 100 (but that of course depends on the precise threat-model under which the attack is evaluated). Even though the median in the hundreds is amazing when compared to early attempts like Chen et al. [2017] which required queries on the order of 10^4 , it is still not satisfactory for a practical use.

4.4 Leveraging transferability

After these early experiments that proved the concept, we focused our attention on transferability, which has a huge potential to save queries.

4.4.1 Transfer attacks provide better seeds for blackbox optimization

Even if the locally-produced adversarial images don't transfer directly to the target, Suya et al. [2020] showed, that the output of transfer-attack can provide better starting seeds for blackbox optimization attacks and improve their query efficiency, which basically adds a degree of freedom to the blackbox optimization. They also discuss different prioritization strategies, as the the number of seeds produced isn't limited by target queries and we can therefore afford to produce as many candidate starting points as we like.

4.4.2 Train, validate, test

There is a weak analogy between the crafting of an adversarial example and the training of machine learning model. ML model weights are first fitted against specified loss constraint. This constraint is (among other things) a function of training data. The weights are then validated and checked against overfitting on a slightly different constraint, which now depends on a validation dataset. When all is good, model is happily deployed to production.

With a bit of imagination, ML model weights correspond to pixel values of an adversarial image. The pixels are first trained by gradient descent on training loss provided by a surrogate model. They are validated against ensemble set of diverse indepenent classifiers, and when the foolrate is good, they are sent for test evaluation to the cloud.

4.4.3 Local training and validation is cheap

We want to offload the cloud query-stress to local simulation as much as possible. An attacker can often afford to spend orders of magnitude more queries to local surrogates and validation models than to the actual target.

4.4.4 Multiple candidates save queries

Iterative black-box attacks usually have query-distrubutions which are tail-heavy. In other words, the median queries needed to create a successful adversarial image are much lower than the average queries.

Let's imagine an attack scenario, where we want to submit a photo to a platform with automatic content moderation mechanism. Querying the target hundreds of times would certainly attract unwanted attention and our heavy queries can quickly trigger human evaluation. If our primary goal is to craft only one adversarial image and as much as possible evade detection, having multiple candidate images would give us another degree of freedom and it could potentially mitigate the heavy-tail problem. This approach can be in principle transparently combined with the multiple-seed candidate suggestions mentioned in 4.4 using the same prioritization candidate scoring mechanism.

4.5 The need for attack pipeline

We argued in 4.4 that combining multiple whitebox and blackbox attack approaches could create more powerful attack as well as giving us more freedom and flexibility to tailor this combination to a specific attack scenario constraints. As of now, there isn't any general whitebox/blackbox attack pipeline which would combine different algorithms and allow us attacking cloud services in a practical way.

4.5.1 Possibility of multiple blackbox workers

We can also imagine running multiple different attacks in parallel and having some meta-controller orchestrating individual attack algorithms such that we minimize queries to the target and efficiently make use of the additional degrees of freedom.

4.5.2 The need for unified attack and model API

There are several frameworks unifying whitebox/blackbox attacks. To name a few, there is FoolBox (Rauber et al. [2020]) or AutoAttack (Croce and Hein [2020]).

Although they are excellent at testing the robustness of local models, they don't give us the flexibility we need to implement all the pipeline features mentioned previously. They cannot be used without some modification to attack cloud models and their optimization attacks cannot be cooperatively scheduled step by step, which is what would be required for effective multi-attack orchestration.

4.6 AdvPipe

4.6.1 The vision

Solves all our problems. At least in theory.

4.6.2 The reality

Solves some of our problems. Like 10%.

TODO: Make some excuses why you didn't make it in time.

TODO: Make some flowchart of what is actually working.

4.7 Implementation

4.7.1 Wrapping whitebox and blackbox models

All models (cloud, local) are wrapped as PyTorch `torch.nn.Modules`. This way they can be used in a plug-and-play manner and passed easily to existing attack algorithms.

Let $L(label)$

Wrapping cloud models

5. Experiments

5.1 Blackbox PoC

Here we go into more technical details about previously mentioned blackbox attacks we initially tried.

- TREMBA
- RayS
- SquareAttack
- Sparse-RS

5.1.1 TREMBA

5.1.2 RayS

This one is hard label attack and doesn't use the continuous loss from GVision.

5.1.3 SquareAttack

SquareAttack L2

Show an image of nice cat.

SquareAttack Linf

Comparison with local success-rate

5.2 Local transferability experiments

Here go all the different transfer-matrices

5.2.1 Choice of dataset

5.2.2 Choice of local models

We performed all our experiments on the following pretrained PyTorch ImageNet models.

- ResNet-18, ResNet-50 (He et al. [2015])
- ResNeXt-50 (32x4d) (Xie et al. [2017])
- Wide-ResNet-50-2 (Zagoruyko and Komodakis [2017])
- Squeezenet (Iandola et al. [2016])
- DenseNet-121 (Huang et al. [2018])

- EfficientNet-b0 (Tan and Le [2020])
- EfficientNet-b0 adversarially trained (Tramèr et al. [2020])

Apart from EfficientNets, all models were taken from the torchvision.models Python package. For the EfficientNets we used github.com/lukemelas/EfficientNet-PyTorch reimplementation, because the original implementation uses in TensorFlow.

5.2.3 Baseline

5.2.4 Baseline

5.2.5 Baseline

5.2.6 Baseline

5.3 Transferability evaluation on GVision

TODO: just run the images against GVision

5.3.1 Choice of evaluation metrics

5.3.2 Wordcloud

Conclusion

Bibliography

- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *ECCV*, 2020.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks, 2017.
- J. Chen and Quanquan Gu. Rays: A ray searching method for hard-label adversarial attack. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, 2017.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks, 2020.
- Francesco Croce, Maksym Andriushchenko, Naman D. Singh, Nicolas Flammarion, and Matthias Hein. Sparse-rs: a versatile framework for query-efficient sparse black-box adversarial attacks. *ArXiv*, abs/2006.12834, 2020.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks, 2018.
- Z. Huang and Tong Zhang. Black-box adversarial attack with transferable model-based embedding. *ArXiv*, abs/1911.07140, 2020.
- Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 10.5mb model size, 2016.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- Jonas Rauber, Roland S. Zimmermann, M. Bethge, and W. Brendel. Foolbox native: Fast adversarial attacks to benchmark the robustness of machine learning models in pytorch, tensorflow, and jax. *J. Open Source Softw.*, 5:2607, 2020.

- Fnu Sua, Jianfeng Chi, David Evans, and Y. Tian. Hybrid batch attacks: Finding black-box adversarial examples with limited queries. *ArXiv*, abs/1908.07000, 2020.
- Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks, 2020.
- Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses, 2020.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface’s transformers: State-of-the-art natural language processing, 2020.
- Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks, 2017.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks, 2017.

List of Figures

List of Tables

List of Abbreviations

A. Attachments

A.1 First Attachment