

DEMOTIVATE ADVERSARIAL DEFENSE IN REMOTE SENSING

Adrien Chan-Hon-Tong¹, Gaston Lenczner^{1,2}, Aurelien Plyer¹

¹ONERA/DTIS, Université Paris-Saclay, FR-91123 Palaiseau, France

²Alteia, FR-31400 Toulouse, France

ABSTRACT

Convolutional neural networks are currently the state-of-the-art algorithms for many remote sensing applications such as semantic segmentation or object detection. However, these algorithms are extremely sensitive to over-fitting, domain change and adversarial examples specifically designed to fool them. While adversarial attacks are not a threat in most remote sensing applications, one could wonder if strengthening networks to adversarial attacks could also increase their resilience to over-fitting and their ability to deal with the inherent variety of worldwide data. In this work, we study both adversarial retraining and adversarial regularization as adversarial defenses to this purpose. However, we show through several experiments on public remote sensing datasets that adversarial robustness seems uncorrelated to geographic and over-fitting robustness.

1. INTRODUCTION

World seen from remote sensing sensors can exhibit a great variability in appearance and the algorithms should be robust to these domain changes. For example, a roof segmentation module [1] should deal with different roof appearances regardless of the country, weather, density area or time of the day. Currently, a deep learning module trained on some part of the world could perform very well on neighborhood regions while having poor performance far away [2]. This is an issue for both high resolution datasets which cover limited areas due to the cost of data and for low resolution datasets which either deal with high level label only or cover limited areas due to the ground-truth annotation cost. Deep neural networks are also notoriously prone to over-fit on training data. In such scenarios, data augmentation and regularization are common efficient ways to strengthen learning algorithms.

On the other hand, deep neural networks are also extremely weak against adversarial attacks [3] which aim to find optimal examples to fool the algorithms. Figure 1 shows an example of the impact such attack. Adversarial defense has been widely studied in the literature to protect the networks from these attacks and these approaches usually rely

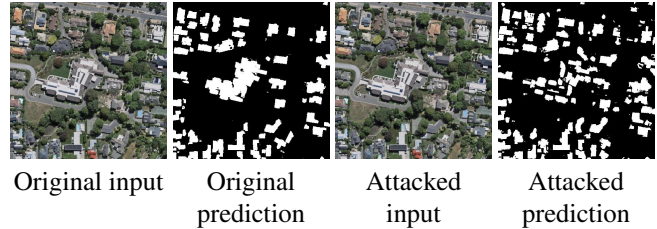


Fig. 1. Adversarial attack on a network trained on AIRS [1] for building segmentation without adversarial defense

on adversarial retraining [4] or adversarial regularization [5]. However, contrary to autonomous driving, adversarial attacks are not a real threat in most remote sensing applications which are neither in real time nor safety critical (e.g. land cover mapping) and due to the large physical changes implied. In this context, adversarial defense could be seen as being hard data augmentation and regularization and thus still be relevant by benefiting to geographic and over-fitting robustness.

1.1. Contributions

In this paper, we focus on evaluating whether adversarial defense is relevant to strengthen neural networks for remote sensing applications without adversarial examples threats. As a result, we found that adversarial robustness seems uncorrelated to geographical and over-fitting ones. Precisely, the contributions of this paper are:

- An evaluation of the impact of adversarial defense on robustness to geographic transfer, adversarial attacks and over-fitting.
- We practically study both adversarial retraining and adversarial regularization on semantic segmentation and detection tasks on public remote sensing datasets.
- We show that adversarial robustness seems uncorrelated to transfer and over-fitting robustness.

Code is available¹ and contains a toolbox for small perturbations and adversarial attacks.

¹https://github.com/achanhon/Lispchitz_penalty

Thanks to the BPI AI4GEO project <http://ai4geo.eu/> for funding.
Corresponding author: gaston.lenczner@alteia.com

1.2. Related works

Fooling learning algorithms, and specifically deep neural networks, has been extensively studied in literature. Small agnostic modification of the inputs [6] can already highly decrease accuracy of these algorithms. Adversarial perturbations [3] make it even worse for the algorithms. At test time, it is thus possible to design a specific marginal perturbation such as a targeted network eventually predicts different outputs on original and disturbed input. Adversarial examples threats have been recently studied in remote sensing [7]. However, adversarial perturbation is not a threat in many non safety-critical remote sensing applications as very large physical changes are required to produce remote sensing adversarial examples. Yet, one could still wonder if geographic generalization or over-fitting resistance could be helped by adversarial defenses.

There are currently two main adversarial defenses. First, adversarial retraining [4] consists in finetuning the network using worse possible adversarial example at each iteration. Similar to a data augmentation process, this can lead to provable defense [8] and corresponds somehow to learn the network in a support vector margin framework. Second, adversarial regularization [5] relies on Lipschitz constants to lower the neural networks adversarial sensibility.

2. METHODOLOGY

2.1. Purpose

Given the major influence of adversarial attacks, an idea could be to use adversarial defense as a way to increase general robustness in addition to the adversarial one. Our overall purpose is thus to evaluate under adversarial defense the correlation between adversarial robustness and robustness to geographic transfer and over-fitting. To better apprehend the relevance of this approach, we increase the amount of training data with additional data for comparison. Hence, we also analyze the impact of a larger training database on adversarial attacks.

2.2. Attack and defense

2.2.1. Adversarial retraining

We consider the Fast Sign Gradient Method (FSGM) adversarial attack [9] and base our adversarial retraining scheme on this algorithm. FSGM is implemented by maximizing the loss over the image. The loss is averaged over the pixels in segmentation and over the targets in detection. Formally, given an image x , its associated ground-truth y , a loss \mathcal{L} and a trained neural network f_θ parameterized by θ , FSGM considers the solution $\varepsilon = \lambda \text{sign}(\nabla_x \mathcal{L}(f_\theta(x), y))$ to the following problem :

$$\max_{\varepsilon / \|\varepsilon\|_1 \leq \delta} \mathcal{L}(f_\theta(x + \varepsilon), y) \quad (1)$$

δ is fixed such that only a dense modification of each pixel by 2 over 255 is possible. Consistently in all of our experiments, FSGM attack is more effective than agnostic attacks like blur or pepper and salt noise.

2.2.2. Adversarial regularization

A function f from \mathbb{R}^I to \mathbb{R}^J is said K-Lipschitz for norms $\|\cdot\|_I$ and $\|\cdot\|_J$ if there exists a constant $K \in \mathbb{R}$ such that, for all $x, y \in \mathbb{R}^I$, $\|f(x) - f(y)\|_J \leq K\|x - y\|_I$. As piece-wise infinitely derivable functions, ReLU-based deep network are K-Lipschitz functions and could thus be naively expected to be smooth. However, both K and I can be very large: indeed, K is estimated to be higher than 50000 for first Alexnet layers in L_2 norm [10]. As sufficiently small Lipschitz coefficient would lead to less adversarial sensibility [11], we implement an idea close to [5] to this aim. Mainly, we add the following regularization term to the classical cross entropy: $\|\nabla_x \|f(x) - \hat{f}\|_2\|_1$ where f is the last feature map of the encoder part of the network and \hat{f} the average of $f(x)$ on the batch.

Code is available for implementation details.

2.3. Experimental set-up

2.3.1. Data

Experiments have been conducted on multiple public remote sensing datasets for both object detection and segmentation tasks: AIRS [1], INRIA [12], ISPRS (Potsdam and Vaihingen) [13], IEEE GRSS DFC [14] and VEDAI [15] datasets both for segmentation and detection. All data are gray-scaled to be invariant to color channels. This is especially relevant for Vaihingen which is an IR-R-G dataset, all the others being RGB but with different distributions.

These datasets provide a wide range of use cases: AIRS and INRIA address building detection while VEDAI is only for vehicle detection. On the other hand, DFC and ISPRS are smaller but more versatile since they consist in multi class semantic segmentation. Finally, they also provide a large range of resolutions, going from 5 cm/pixel for Potsdam to 30 cm/pixel for INRIA. The impact of resolution is also considered in the experiments.

2.3.2. Backbone

For detection, detectors are based on the SSD framework [16] with either VGG or ResNet backbone. For segmentation, experiments rely on a UNet [17] architecture but consistent results have been observed with Deeplabv3+.

2.3.3. Metric

Intersection over Union (IoU) is used to evaluate the performances in all the segmentation experiments.

Train \ Test	AIRS <i>naive</i>	AIRS <i>adv</i>	AIRS <i>retr</i>
AIRS	74	85	72
AIRS <i>adv</i>	63	77	68
ISPRS	58	80	62
ISPRS <i>adv</i>	44	69	54

Table 1. Root segmentation IoU depending on train/test data. Test data can be with (*adv*) or without adversarial attacks. Training can be *naive*, with adversarial defense (*retr*) or with extra data (*add*).

For detection, the classic way to match predicted bounding boxes and ground-truth bounding boxes is an Intersection over Union above some fixed threshold. However, since scale is not an issue when resolution is known, we follow [18] and match predictions and ground-truth when centers are closer than 1.5 meter. Then, we rely on F and G-scores to evaluate detection performances: F-score is the harmonic mean of precision and recall while G-score is the geometric one. Precision is the number of true alarms divided by the number of alarms and recall is the number of true alarms divided by the number of true objects.

3. RESULTS

3.1. Segmentation results

Segmentation results are presented in Table 1. In addition to these results, we observed an IoU of 80% on ISPRS when the network is trained on the ISPRS train set. We use adversarial retraining as the adversarial defense for this experiment since we observed better results than adversarial regularization. The additional training set consists of INRIA and DFC datasets.

Consistently with adversarial works, the performances of the naive classifier decrease importantly under adversarial attack: they respectively drop by 11% and 14% on AIRS and ISPRS. Regarding these attacks, adversarial defense allows to make the network very resilient. Indeed, IoU only decreases by 6% on average under adversarial attacks for the defended network against 12.5% for the naive one.

However, this adversarial robustness does not provide any geographic robustness. Indeed, IoU highly decreases when the model trained on the train set of AIRS is applied on ISPRS Vaihingen. This phenomenon is independent of the use of adversarial examples: IoU goes from 74% to 58% for the naive network and from 72% to 62% for the defended one.

Inversely, adding additional data from two other datasets to the training set decreases overfitting and transfer capability. Indeed, performance increases both on the test set AIRS and on ISPRS, going respectively from 74% & 58% to 85% &

Model \ Data	VEDAI <i>32×32</i>	DFC <i>48×48</i>	ISPRS <i>64×64</i>	VEDAI <i>adv</i>	DFC <i>adv</i>	ISPRS <i>adv</i>
VGG	53	66	78	30	51	66
VGG <i>reg</i>	53	65	80	36	57	74
VGG <i>add</i>	62	71	85	31	49	64
RESNET	39	60	80	22	24	50
RESNET <i>reg</i>	43	49	81	20	29	69
RESNET <i>add</i>	57	65	82	21	26	52

Table 2. G-score performances of a naive detector, an adversarial defended (*reg*) one and one with extra data (*add*) on clean and adversarial (*adv*) datasets. Target size in different dataset is indicated in italic for completeness.

80%. It thus almost catches up on ISPRS with the control network trained on ISPRS. Nonetheless, it does not protect against adversarial attacks: IoU still drops under attack by about 10% both on AIRS and ISPRS.

Hence, since additional data increases geographic robustness only while adversarial training increases adversarial robustness only, we conclude that adversarial and geographic robustness are uncorrelated in this experiment.

3.2. Detection results

Detection results are reported in Table 2. Contrary to the previous segmentation experiment, pure transfer leads to extremely low performance in detection, so we consider here the impact of additional data and adversarial defense on classical test set performance: we study the impact on over-fitting. Thus, for a given dataset, we train a network classically, another one with adversarial defense and a last one with additional data from ISPRS Potsdam and DFC datasets and compare them on the associated test set. Finally, we base the adversarial defense here on adversarial regularization instead of retraining for computational reasons.

Like in the previous experiment, native models are very sensitive to FGSM. Indeed, their performance under adversarial attacks are extremely low compared to their performances on clean data: we observe G-score drops up to 30%. Adding adversarial defense again importantly moderates this drop. Regularized algorithms eventually outperform their undefended version on adversarial data by 5% of G-score. However, adversarial defense has no effect on clean data. Inversely, the additional training data clearly improves the performance on clean data but does not provide any defense against FGSM attack.

To summarize, Table 2 clearly shows that adding data outperforms the baseline and the defended baseline on clean data whilst the defended baseline outperforms the other ones under adversarial attack. This shows a lack of correlations between adversarial robustness and robustness to over-fitting.

4. CONCLUSION

This article aims to evaluate whether adversarial defense would be relevant data augmentation in a context where adversarial threats would be non-existent. Hence, it focuses on the search of a correlation between adversarial, geographic and over-fitting robustness. Specifically, we investigate this concern for both object detection and segmentation, with different backbones and adversarial defense.

Consistently in our experiments, adversarial defense does not improve both geographic and over-fitting robustness even though it strengthens the network against adversarial attacks. Therefore, using adversarial data as data augmentation does not seem relevant in an adversarial attack free context and under-performs compared to other methods such as the use of additional data. Inversely, adversarial defense greatly improves performance under adversarial attack, contrary to the use of extra datasets. To conclude, it seems that adversarial defense framework has little relevance in remote sensing applications without adversarial examples threats.

5. REFERENCES

- [1] Qi Chen, Lei Wang, Yifan Wu, Guangming Wu, Zhiling Guo, and Steven L Waslander, "Aerial imagery for roof segmentation: A large-scale dataset towards automatic mapping of buildings," *arXiv preprint*, 2018.
- [2] Javiera Castillo-Navarro, Nicolas Audebert, Alexandre Boulch, Bertrand Le Saux, and Sébastien Lefèvre, "What data are needed for semantic segmentation in earth observation?," in *JURSE*. IEEE, 2019, pp. 1–4.
- [3] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, "Explaining and harnessing adversarial examples.," *arXiv preprint*, 2014.
- [4] Uri Shaham, Yutaro Yamada, and Sahand Negahban, "Understanding adversarial training: Increasing local stability of supervised models through robust optimization," *Neurocomputing*, vol. 307, pp. 195–204, 2018.
- [5] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky, "Spectrally-normalized margin bounds for neural networks," in *NeurIPS*, 2017, pp. 6240–6249.
- [6] Dan Hendrycks and Thomas Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," *arXiv preprint*, 2019.
- [7] Li Chen, Guowei Zhu, Qi Li, and Haifeng Li, "Adversarial example in remote sensing image recognition," *arXiv preprint*, 2019.
- [8] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan, "Theoretically principled trade-off between robustness and accuracy," *arXiv preprint*, 2019.
- [9] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami, "The limitations of deep learning in adversarial settings," in *EuroS&P*. IEEE, 2016, pp. 372–387.
- [10] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus, "Intriguing properties of neural networks," *arxiv technical report*, 2013.
- [11] Todd Huster, Cho-Yu Jason Chiang, and Ritu Chadha, "Limitations of the lipschitz constant as a defense against adversarial examples," in *ECML-PKDD*. Springer, 2018, pp. 16–29.
- [12] Bohao Huang, Kangkang Lu, Nicolas Audebert, Andrew Khalel, Yuliya Tarabalka, Jordan Malof, Alexandre Boulch, Bertrand Le Saux, Leslie Collins, Kyle Bradbury, et al., "Large-scale semantic classification: outcome of the first year of inria aerial image labeling benchmark," in *IGARSS*. IEEE, 2018, pp. 6947–6950.
- [13] Franz Rottensteiner, Gunho Sohn, Jaewook Jung, Markus Gerke, Caroline Baillard, Sebastien Benitez, and Uwe Breitkopf, "The ISPRS benchmark on urban object classification and 3d building reconstruction," *ISPRS Annals*, vol. 1, no. 1, pp. 293–298, 2012.
- [14] Adrien Lagrange, Bertrand Le Saux, Anne Beaupere, Alexandre Boulch, Adrien Chan-Hon-Tong, Stéphane Herbin, Hicham Randrianarivo, and Marin Ferecatu, "Benchmarking classification of earth-observation data: From learning explicit features to convolutional networks," in *IGARSS*. IEEE, 2015, pp. 4173–4176.
- [15] Sebastien Razakarivony and Frederic Jurie, "Vehicle detection in aerial imagery: A small target detection benchmark," *Journal of Visual Communication and Image Representation*, vol. 34, pp. 187–203, 2016.
- [16] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg, "SSD: Single shot multibox detector," in *ECCV*. Springer, 2016, pp. 21–37.
- [17] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *MICCAI*. Springer, 2015, pp. 234–241.
- [18] Nicolas Audebert, Bertrand Le Saux, and Sébastien Lefèvre, "Segment-before-detect: Vehicle detection and classification through semantic segmentation of aerial images," *Remote Sensing*, vol. 9, no. 4, pp. 368, Apr 2017.