

---

# [Re] Subspace Attack: Exploiting Promising Subspaces for Query-Efficient Black-box Attacks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 As part of the NeurIPS 2019 Reproducibility Challenge, we chose to attempt re-  
2 produce the attack algorithm proposed in “Subspace Attack: Exploiting Promising  
3 Subspaces for Query-Efficient Black-box Attacks”. Our reported results are better  
4 than the original paper in terms of the median number of queries per attack, but  
5 worse in terms of failure rate. A concise assessment of our implementation is also  
6 included.

## 7 1 Introduction

8 As the use of Machine Learning in services and applications has increased, it becomes important to  
9 assess its reliability and security. In the last few years, several flaws have been found in state-of-the-art  
10 Machine Learning methods and algorithms. One of the most studied flaws is the fact that many  
11 models can be attacked using the so-called adversarial examples, inputs that are minimally perturbed  
12 in order to be misclassified by the victim model.

13 In this paper we attempt to reproduce the results obtained in the paper “Subspace Attack: Exploiting  
14 Promising Subspaces for Query-Efficient Black-box Attacks” by Yan and Guo et al. [1], published  
15 among the proceedings of NeurIPS 2019. In their paper, the authors present a new kind of black-box  
16 attack, that, for the first time, ensembles attacks’ transferability and gradient estimation for Projected  
17 Gradient Descent.

18 Our report has been written for the NeurIPS 2019 Reproducibility Challenge<sup>1</sup>. It consists of a  
19 background about adversarial attacks (Section 2), our methodology (Section 3), and a concise  
20 reproducibility section (Section 4). These are followed by our results, discussion, and conclusions  
21 (Sections 5 to 7, respectively).

## 22 2 Background

### 23 2.1 Adversarial Examples

24 Adversarial examples are inputs fed into a machine learning model and mislead the model to make  
25 an incorrect prediction. Several machine learning architectures have been proven to be vulnerable  
26 in adversarial settings, where the adversarial examples can be generated with several attacking  
27 techniques. Thus, the aim of an attack is to perturb an input  $\mathbf{x} \in \mathbb{R}^n$  and to trick a victim model  
28  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ , such that it gives a wrong prediction,  $\arg \max_i f(\mathbf{x})_i \neq y$  (where  $y$  the true label).  
29 This general definition describes the so-called *untargeted* attacks. A more specific type of attack is a  
30 *targeted* one, which aims for a wrong classification of a specific label  $y' \neq y$ ,  $\arg \max_i f(\mathbf{x})_i = y'$ .

---

<sup>1</sup><https://reproducibility-challenge.github.io/neurips2019/>

31 In general, adversarial attacks can be classified into two main types: white-box and black-box. While  
 32 only the confidence score from the victim model is accessible to the latter, the former has full access  
 33 to the network parameters of the victim. This allows for an efficient use of various gradient-based  
 34 methods [2] to generate adversarial examples. This emphasises the attacker’s preference for the  
 35 white-box models. However, it is not a realistic setting.

36 Let us now consider the more realistic setting in which the attacker has no access to the trained  
 37 victim model parameters. In this case, it is not possible to compute the gradient of the model with  
 38 respect to the perturbed input. Nonetheless, several methods have been proposed to overcome this  
 39 issue [3, 4, 5, 6, 7, 8, 9, 10, 11]. Some of these methods rely on the so-called adversarial example  
 40 transferability [12], in which a surrogate model is trained with data samples labeled by the victim  
 41 model, and is then used to generate adversarial examples that work on the victim model too. The  
 42 other kind of methods make use of zeroth-order optimization methods (e.g. finite differences [5]) by  
 43 querying the victim model to estimate the gradient of the victim model to run Projected Gradient  
 44 Descent (PGD). A drawback of these methods is the need for a large number of queries for a good  
 45 estimation of the victim model’s gradient. This motivates the attempt for reducing the required  
 46 number of queries for a successful attack. The next subsection describes the solution proposed in [1],  
 47 which harnesses the features of both zeroth-order optimization and transferability.

## 48 2.2 Subspace Attack Algorithm

49 In what follows, we present the essence of the subspace attack proposed in [1]. This will include an  
 50 introduction to the principles of the *Bandit attack* proposed in [10], on which the subspace attack  
 51 heavily relies.

52 Given a classification loss function  $\mathcal{L}(x, y)$ , where  $x$  is some input and  $y$  its corresponding label. We  
 53 can formulate the adversarial attack problem as follows:

$$x' = \arg \max_{x': \|x' - x\|_p \leq \epsilon_p} \mathcal{L}(x', y) \quad (1)$$

54 Let  $g^* = \nabla_x \mathcal{L}(x, y)$  be the gradient of  $\mathcal{L}$  at  $(x, y)$ . Then the goal of the gradient estimation problem  
 55 is to find a unit vector  $\hat{g}$  that maximize the inner product  $\mathbb{E} [\hat{g}^T g^*]^2$ .

56 In [10], the authors have proven that the gradient estimation problem can be considered as a bandit  
 57 optimization problem. Bandit optimization is a tool used in online convex optimization, in which  
 58 there is an agent playing a game that includes a sequence of rounds. During each round  $t$ , the  
 59 agent must select an action and incurs in a loss  $\mathcal{L}_t$ , whose expectation across all the rounds should  
 60 be minimized. The main novelty introduced by [10] is the fact that the estimation is improved by  
 61 accumulating prior information about gradients in a latent vector  $\mathbf{g}_t$ , that is updated each round with  
 62 an estimation of the gradient of the victim model, performed via finite differences. In [10], each new  
 63 basis vector  $\mathbf{u}_t$  used for the gradient estimation is randomly sampled from a Gaussian distribution.

64 In the subspace attack proposal, the authors present a novel way to accumulate prior information  
 65 about the gradient. Instead of using the full basis of the sample’s space to estimate the true gradient,  
 66 a subspace of directions is chosen by using the gradients of some reference models. These models  
 67 are pre-trained on the same dataset as the victim and their parameters are fully known. Each round  
 68 of the bandit optimization uses the gradient of a randomly chosen reference model, with respect to  
 69 the adversarial input, as a basis for the gradient estimation. Moreover, [1] proposes a way to apply  
 70 different ratios of drop-out to the enrich the set of reference models.

71 The method proposed by [1] is presented in Algorithm 1, which we report in our paper for the sake of  
 72 clarity and ease of read.

## 73 3 Methodology

74 In this work, we implemented the subspace attack using PyTorch [13], following the algorithm  
 75 specified in [1], without looking at the source code of the official implementation of the attack. We  
 76 chose to evaluate the reproducibility of the algorithm by attempting to replicate untargeted subspace  
 77 attacks on the GDAS [14] and WRN [15] models trained on the CIFAR-10 dataset [16]. We use

---

<sup>2</sup>The expectation here is taken over the randomness of the estimation algorithm

---

**Algorithm 1:** Subspace Attack, as described in [1].

---

**Input:** A benign example  $\mathbf{x} \in \mathbb{R}^n$ , its label  $y$ , a set of  $m$  reference models  $\{f_0, \dots, f_{m-1}\}$ , a chosen attack objective function  $\mathcal{L}(\cdot, \cdot)$ , and a victim model from which the output of  $f$  can be inferred.

**Output:** An adversarial example  $\mathbf{x}_{adv}$  that fulfills  $\|\mathbf{x}_{adv} - \mathbf{x}\|_\infty \leq \epsilon$ .

- 1: Initialize the adversarial example to be crafted  $\mathbf{x}_{adv} \leftarrow \mathbf{x}$ .
  - 2: Initialize the gradient to be estimated  $\mathbf{g} \leftarrow \mathbf{0}$ .
  - 3: Initialize the drop-out/layer ratio  $p$ .
  - 4: **while** not successful **do**
  - 5:   Choose a reference model whose index is  $i$ , uniformly at random.
  - 6:   Calculate prior gradient with drop-out/layer ratio  $p$  as  $\mathbf{u} \leftarrow \frac{\partial \mathcal{L}(f_i(\mathbf{x}_{adv}, p), y)}{\partial \mathbf{x}_{adv}}$
  - 7:    $\mathbf{g}_+ \leftarrow \mathbf{g} + \tau \mathbf{u}$ ,  $\mathbf{g}_- \leftarrow \mathbf{g} - \tau \mathbf{u}$
  - 8:    $\mathbf{g}'_+ \leftarrow \mathbf{g}_+ / \|\mathbf{g}_+\|_2$ ,  $\mathbf{g}'_- \leftarrow \mathbf{g}_- / \|\mathbf{g}_-\|_2$
  - 9:    $\Delta_t \leftarrow \frac{\mathcal{L}(f(\mathbf{x}_{adv} + \delta \mathbf{g}'_+), y) - \mathcal{L}(f(\mathbf{x}_{adv} + \delta \mathbf{g}'_-), y)}{\tau \delta} \mathbf{u}$
  - 10:    $\mathbf{g} \leftarrow \mathbf{g} + \eta \mathbf{g} \Delta_t$
  - 11:    $\mathbf{x}_{adv} \leftarrow \mathbf{x}_{adv} + \eta \cdot \text{sign}(\mathbf{g})$
  - 12:    $\mathbf{x}_{adv} \leftarrow \text{Clip}(\mathbf{x}_{adv}, \mathbf{x} - \epsilon, \mathbf{x} + \epsilon)$
  - 13:    $\mathbf{x}_{adv} \leftarrow \text{Clip}(\mathbf{x}_{adv}, 0, 1)$
  - 14:   Update the drop-out/layer ratio  $p$ , following the original paper’s policy.
  - 15: **end while**
  - 16: **return**  $\mathbf{x}_{adv}$
- 

78 as reference models VGG-11/13/16/19 [17] and AlexNet [18]. In order to verify that we use  
79 the same settings – hyper-parameters and pre-trained models – as the original subspace experiments,  
80 we checked their log files and loaded the pre-trained models from their source code repository. The  
81 latter was also required since the available models by PyTorch are pre-trained on ImageNet, while  
82 our experiment uses images from CIFAR-10.

83 Aiming to further assess the performance of our algorithm implementation, we added options to  
84 keep track of some values encountered during the attack. In addition to the required number of  
85 queries for a successful attack of each image classification, we recorded in each iteration the cosine  
86 similarity between the estimated gradient and the true gradient (as done in [10]), and the norms of  
87 these gradients. Moreover, as the subspace attack’s algorithm uses only the sign of the gradient, we  
88 also keep track of the ratio of the matching signs of the gradient elements with respect to the previous  
89 estimated gradient (denoted by *the ratio of the common signs*) and the true gradient (denoted by *the*  
90 *ratio of the correctly estimated signs*).

91 Using the evaluation metrics of [1], an attack’s performance is being evaluated with respect to its  
92 failure rate and the number of queries. Since the distribution of the latter is heavy-tailed, we use the  
93 mean and the median number of queries per a successful attack.

## 94 4 Reproducibility

95 We first try to implement Algorithm 1 and run the experiment attacking GDAS. We used VGGnets  
96 and AlexNet as reference models and the hyper-parameters listed in [1] and [10]. We then expand our  
97 implementation to be used to attack WRN and Pyramidnet<sup>3</sup>, and look for better hyper-parameters.

### 98 4.1 Machine Setup, experiment duration, and budget

99 We run our experiment both on Google Cloud Platform and on Code Ocean [19]. The Ubuntu Virtual  
100 Machine on Google Cloud Platform has an Nvidia Tesla T4 GPU, 52 GB memory and 8 vCPUs. A  
101 set of 1000 attacks against GDAS using VGGs and AlexNet as reference models took us about 7h45m  
102 to run, with a total of \$7.75 spent. On Code Ocean, the Virtual Machine we used runs Ubuntu on a  
103 4-cores CPU with 60 GB of memory and an Nvidia Tesla K80 GPU. As Code Ocean is a sponsor of  
104 the Reproducibility Challenge, we have been provided with free compute time to run the experiments.

---

<sup>3</sup>An attack on Pyramidnet is implemented but we didn’t run such experiment due to large computation time.

105 A set of 1000 attacks against GDAS using VGGs and AlexNet as reference models took us about 9  
106 hours to run.

107 Even though our reported results are produced using GPUs, our implementation can run on a CPU  
108 as well. This was tested on a laptop with Ubuntu, an Intel Core i7-7500U and 8GB memory. The  
109 computation time took about 6 times longer than a GPU (in terms of iterations per second). In  
110 addition, the results obtained could be different, due to differences in floating point precision and  
111 implementations of low-level operations<sup>4</sup>. Such results are outside of the scope of this challenge, and  
112 therefore not reported.

## 113 4.2 Algorithm implementation

114 The source code of our implementation is available online<sup>5</sup>. As mentioned in section 3, our imple-  
115 mentation is done using the PyTorch framework. We use the pre-trained reference and victim models<sup>6</sup>  
116 provided with the original paper’s code repository, along with their corresponding Python classes.

117 The CIFAR-10 dataset is loaded using PyTorch’s `torchvision.datasets`, and is iterated using  
118 a `DataLoader`. The dataset is shuffled, but the seed can be fixed, for comparable results. In all  
119 our experiments, both PyTorch and NumPy seeds are set to 0, and the following values are set:  
120 `torch.backends.cudnn.deterministic = True` and `torch.backends.cudnn.benchmark`  
121 `= False`. Regarding the implementation of the attack itself, we leveraged PyTorch autograd capabili-  
122 ties to compute the gradient of the reference models.

123 Two clipping procedures of the adversarial example are presented in lines 12 and 13 of Algorithm 1.  
124 While the former was implemented using PyTorch’s functions `torch.min` and `torch.max`, in order  
125 to keep  $\mathbf{x}_{adv}$  in the region included in  $[\mathbf{x} - \epsilon, \mathbf{x} + \epsilon]$  with  $\ell_\infty$ -norm, the latter was implemented using  
126 PyTorch’s `torch.clamp`, with 0 and 1 as arguments.

127 Yan and Guo et al. specify the use of "the hinge logit-diff adversarial loss from Carlini and Wagner"  
128 [1]. In their paper [20], Carlini and Wagner list a number of possible loss functions that can be  
129 used along with a hinge term, which regulates the  $\ell_\infty$ -distance of the adversarial example from the  
130 original one. However, we did not manage to understand how to implement it as part of the subspace  
131 attack algorithm. So, we use the Cross Entropy loss without the hinge term. We understood from  
132 Algorithm 1 that the constraint into the admitted perturbation bounds is applied by the clipping step  
133 in line 12.

## 134 4.3 Hyper-parameters and experiments settings

135 The hyper-parameters used in a subspace attack<sup>7</sup> are:

- 136 •  $\tau$  : the bandit exploration.
- 137 •  $\delta$  : the finite difference probe.
- 138 •  $\eta_g$  : the OCO learning rate.
- 139 •  $\eta$  : the image learning rate.
- 140 •  $p$  : the dropout/layer ratio of the reference models.

141 In [1], it is claimed that the used hyper-parameters were the same as those used in [10]. However, we  
142 have found out that the notation was different and a bit confusing. We reconstructed a translation in  
143 Table 1 by cross-referencing the notation used in Algorithm 1 of [1] and in Algorithms 1, 2 and 3  
144 of [10]. The reader should note that we have listed two letters as translations of the  $\delta$  used in [1] –  
145  $\epsilon$  and  $\eta$ . This is due to the fact that, in Algorithm 2 of [10], an  $\epsilon$  is used as finite difference probe,  
146 but in Table 3 of Appendix C of the same paper, where hyper-parameters values are listed, the finite  
147 difference probe is listed as “ $\eta$  (Finite difference probe)”.

148 To choose the hyper-parameters’ values, we have started from those stated in [1], which should be the  
149 same used in [10], excluding  $\epsilon$  and  $\eta$ . However, as discussed in Section 6, we obtained results which

<sup>4</sup><https://pytorch.org/docs/stable/notes/randomness.html>

<sup>5</sup><https://github.com/epfl-ml-reproducers/subspace-attack-reproduction>

<sup>6</sup><https://go.epfl.ch/subspace-pretrained>

<sup>7</sup>For hyper-parameters, we use the same notation as [1]

Table 1: Translation between hyper-parameters in Subspace attack [1] and the bandit attack [10].

	Subspace Attack [1]	Bandit Attack [10]	Value [1]
Bandit exploration	$\tau$	$\delta$	1.0
Finite difference probe	$\delta$	$\epsilon/\eta$	0.1
Image $\ell_p$ learning rate	$\eta$	$h$	1/255
OCO learning rate	$\eta_g$	$\eta$	100

150 were far worse than those obtained in [1]. After exploring the published logs of the experiment<sup>8</sup> we  
 151 found out that in the original paper, a value of  $\eta_g = 0.1$ , instead of  $\eta_g = 100$  was used. The set of  
 152 hyper-parameters that yields the best results is presented in Table 2.

Table 2: The set of Hyper-parameters which yield the best results.

Hyper-parameter	Value
$\tau$ (Bandit exploration)	1.0
$\delta$ (Finite difference probe)	0.1
$\eta$ (Image $\ell_p$ learning rate)	1/255
$\eta_g$ (OCO learning rate)	0.01

153 The dropout was applied to the reference models, in the following procedure: each experiment starts  
 154 with a 0.05 dropout ratio which is increased every iteration by 0.01, until a maximum ratio of 0.5 is  
 155 reached.

156 Finally, using the above settings, each experiment included untargeted attacks on the classification of  
 157 1,000 images using a limit of 10,000 queries per attack – after which the attack is considered a failure.  
 158 Moreover, we set a maximum perturbation of  $\epsilon = 8/255$ , defined using  $\ell_\infty$ -norm.

## 159 5 Results

160 In this section, we present the assessment of our implementation as mentioned in Section 3, followed  
 161 by a comparison of the effectiveness of our subspace attacks to those reported in the original paper.  
 162 These results are discussed in Section 6.

163 In order to evaluate the performance of our algorithm and its implementation, we saved the progression  
 164 of the loss values along with the estimated and true gradients information of a set of 49 attacks  
 165 on GDAS. We required each attack to complete 5,000 gradient estimation iterations, i.e. 10,000  
 166 queries. Figures 1 and 2 present the mean value along the gradient estimation process of the following  
 167 parameters:

- 168 • the cross entropy loss of the victim model (Figure 1a)
- 169 • the cosine similarity between the estimated and the true gradient (Figure 1b)
- 170 • the ratio of correctly estimated gradient’s sign (Figure 2a)
- 171 • the ratio of common signs between subsequent estimated gradients (Figure 2b)

172 Having only a single failed attack in this run, its corresponding curve presents the exact value, while  
 173 the rest of the curves present value averaged on the 48 successful attacks and 49 attacks overall. The  
 174 fluctuations of the curves are a feature of the stochastic nature of the algorithm.

175 In addition, we scanned the values of the  $\eta_g$  hyper-parameter. However, as presented in Table 3, it  
 176 seems that  $\eta_g = 0.1$  yields the best results in terms of a trade-off between failure rate and median  
 177 number of queries. Our best results of the different attack experiments are compared with the two  
 178 baselines [1, 10] in Table 4. There, it is demonstrated that while our implementation managed to get  
 179 similar mean and better median numbers of queries as the subspace attack paper, it leads to higher  
 180 failure rates with respect to the baselines.

<sup>8</sup><https://go.epfl.ch/subspace-original-logs>

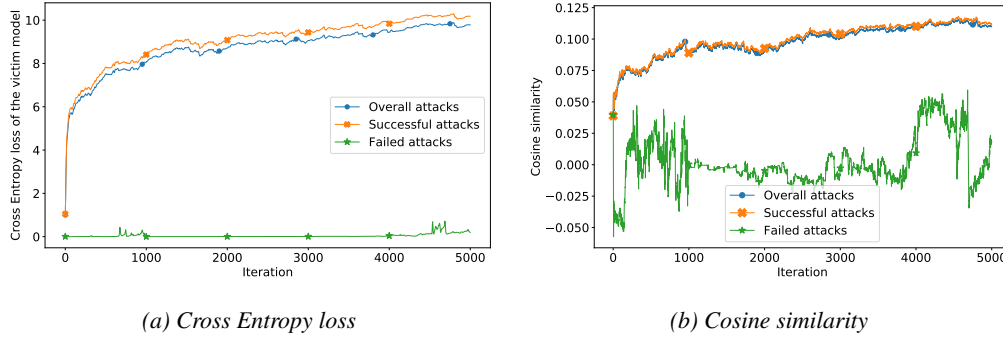


Figure 1: The mean value of the Cross Entropy loss of the true victim model (a) of cosine similarity between the estimated and the true gradient (b) as a function of the estimation iteration in our Subspace Attack implementation. The results of the successful attacks, the failed attacks, and all the attacks were averaged over 48, 1, and 49 attacks respectively.

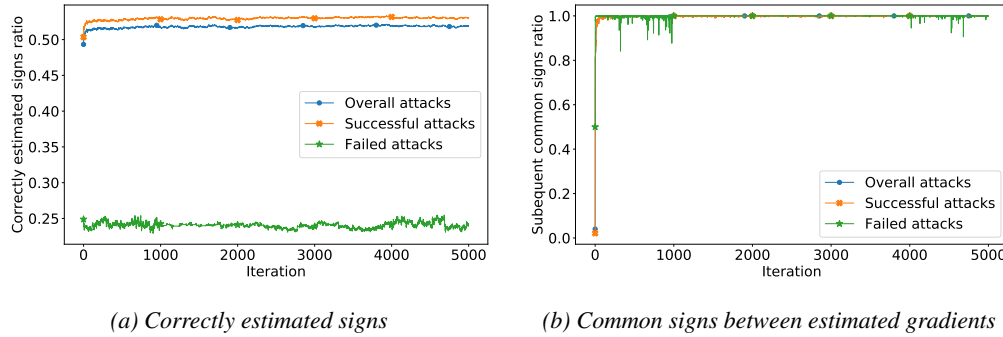


Figure 2: The mean value of the ratio of correctly estimated signs (a), and the ratio of the common signs between two subsequent estimated gradients (b) as a function of the estimation's iteration in our Subspace Attack implementation. The results of the successful attacks, the failed attacks, and the all the attacks were averaged over 48, 1, and 49 attacks respectively.

Table 3: Comparison between the results obtained with different  $\eta_g$ .

Victim	$\eta_g$	Mean Queries	Median Queries	Failure Rate
WRN	0.01	369	16	3.2%
	0.1	321	16	3.4%
	100	288	12	10.1%
GDAS	0.001	263	18	6.3%
	0.01	321	20	5.9%
	0.1	282	18	6.7%
	100	354	16	12.7%

Table 4: Comparison between the results obtained by our implementation and those of the original paper [1] and of [10]. For the results obtained with the other implementations we report the results presented in [1]. In  $\eta_g$  for the results obtained by the original paper we state 0.1\* as this is the value we found in their logs, and not the one stated in their paper.

Victim Model	Method	Ref. Models	$\eta_g$	Mean Queries	Median Queries	Failure Rate
WRN	Bandit-TD [10]	-	100	713	266	1.2%
	Subspace Attack [1]	AlexNet+VGGNets	0.1*	392	60	0.3%
	Subspace Attack (Ours)	AlexNet+VGGNets	0.1	330	14	3.7%
GDAS	Bandit-TD [10]	-	100	373	128	0.0%
	Subspace Attack [1]	AlexNet+VGGNets	0.1*	250	58	0.0%
	Subspace Attack (Ours)	AlexNet+VGGNets	0.1	282	18	6.7%

## 181 6 Discussion

182 The trends of the mean values of the loss and the cosine similarity (Figures 1a and 1b) confirms  
183 that our implementation is indeed efficiently attacking the victim model, as the loss is increased,  
184 and the directional agreement between the estimated and the true gradient improves. Furthermore,  
185 the obtained results of cosine similarity exceed the ones presented in [10], since our successful  
186 attacks demonstrate twice the cosine similarity than those reported in the reference, in the first 5,000  
187 iterations.

188 Regarding the percentage of correctly estimated signs, we see in Figure 2a that the successful  
189 attacks achieved a mean ratio of correctly estimated signs of about 55%. This is only slightly above  
190 the agreement expected by a randomly chosen vector (50%). However, the success of the attacks  
191 suggests that we are indeed exploiting the reference models to estimate only the gradient components  
192 corresponding to the most relevant subspaces used for the classification, as aimed in [1, 10]. Clearly,  
193 the reason why the failed attack did not succeed is the inability to correctly estimate the signs.  
194 Furthermore, we can observe in Figure 2b that the signs of the gradients almost stop changing after  
195 the very first few iterations. This leads us to think that the correctly estimated signs are always the  
196 same, and could correspond exactly to the most relevant subspaces.

197 Keeping the guiding rule of [1], we use the hyper-parameters of [10], excluding  $\epsilon$  and  $\eta$ . As we  
198 are performing attacks that make use of  $\ell_\infty$ -norm, we should use  $\eta_g = 100$ . This yields worse  
199 results, as seen in Table 3. Information extracted from the logs of the subspace paper mention the  
200 use of  $\eta_g = 0.1$ . The fact that our use of  $\eta_g = 0.1$  also result with better results indicates that our  
201 implementation is close to the original one.

202 Comparing our results to those of the original paper, we notice that while we managed to reduce the  
203 computational cost of the attacks, we couldn't reach the reported failure rates of either the original  
204 paper or the Bandit attack. Even though we verified that we followed every step of the algorithm,  
205 checked for differences in input models and hyper-parameters, we couldn't improve our produced  
206 failure rates. It could be that the source of the problem lies in the inconsistency of the implementation  
207 of the loss function of [20].

## 208 7 Conclusion

209 In this work, we re-implemented the algorithm of the subspace black-box adversarial attack presented  
210 by Yan et al. [1], and introduced methods for the evaluation of its performance. However, we didn't  
211 manage to match our effectiveness to that of the original paper. In terms of the reproducibility of [1],  
212 it is advisable for the authors to include a clearer statement regarding the used set of hyper-parameters.  
213 In addition, when relying on settings from another work, a corresponding mapping of the notation  
214 changes can be highly beneficial. Moreover, additional information about the used loss function  
215 could be helpful to replicate the obtained results.

## 216 Acknowledgments

217 The authors would like to thank the NeurIPS Reproducibility Challenge sponsor Code Ocean [19],  
218 for providing them with additional compute time for free.

219 **References**

- 220 [1] Y. Guo, Z. Yan, and C. Zhang, “Subspace attack: Exploiting promising subspaces for  
 221 query-efficient black-box attacks,” in *Advances in Neural Information Processing Systems 32*,  
 222 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds.  
 223 Curran Associates, Inc., 2019, pp. 3820–3829. [Online]. Available: [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8638-subspace-attack-exploiting-promising-subspaces-for-query-efficient-black-box-attacks.pdf)  
 224 [8638-subspace-attack-exploiting-promising-subspaces-for-query-efficient-black-box-attacks.](http://papers.nips.cc/paper/8638-subspace-attack-exploiting-promising-subspaces-for-query-efficient-black-box-attacks.pdf)  
 225 [pdf](http://papers.nips.cc/paper/8638-subspace-attack-exploiting-promising-subspaces-for-query-efficient-black-box-attacks.pdf)
- 226 [2] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus,  
 227 “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- 228 [3] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, “Practical black-  
 229 box attacks against machine learning,” in *Proceedings of the 2017 ACM on Asia conference on*  
 230 *computer and communications security*. ACM, 2017, pp. 506–519.
- 231 [4] Y. Liu, X. Chen, C. Liu, and D. Song, “Delving into transferable adversarial examples and  
 232 black-box attacks,” *arXiv preprint arXiv:1611.02770*, 2016.
- 233 [5] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, “Zoo: Zeroth order optimization  
 234 based black-box attacks to deep neural networks without training substitute models,”  
 235 in *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*,  
 236 ser. AISeC ’17. New York, NY, USA: ACM, 2017, pp. 15–26. [Online]. Available:  
 237 <http://doi.acm.org/10.1145/3128572.3140448>
- 238 [6] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin, “Black-box adversarial attacks with limited  
 239 queries and information,” *arXiv preprint arXiv:1804.08598*, 2018.
- 240 [7] N. Narodytska and S. Kasiviswanathan, “Simple black-box adversarial attacks on deep neural  
 241 networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops*  
 242 *(CVPRW)*. IEEE, 2017, pp. 1310–1318.
- 243 [8] A. N. Bhagoji, W. He, B. Li, and D. Song, “Practical black-box attacks on deep neural networks  
 244 using efficient query mechanisms,” in *European Conference on Computer Vision*. Springer,  
 245 2018, pp. 158–174.
- 246 [9] C.-C. Tu, P. Ting, P.-Y. Chen, S. Liu, H. Zhang, J. Yi, C.-J. Hsieh, and S.-M. Cheng, “Autozoom:  
 247 Autoencoder-based zeroth order optimization method for attacking black-box neural networks,”  
 248 in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 742–749.
- 249 [10] A. Ilyas, L. Engstrom, and A. Madry, “Prior convictions: Black-box adversarial attacks with  
 250 bandits and priors,” *ICLR 2019*, 2018. [Online]. Available: <https://arxiv.org/abs/1807.07978>
- 251 [11] C. Guo, J. R. Gardner, Y. You, A. G. Wilson, and K. Q. Weinberger, “Simple black-box  
 252 adversarial attacks,” *arXiv preprint arXiv:1905.07121*, 2019.
- 253 [12] N. Papernot, P. D. McDaniel, and I. J. Goodfellow, “Transferability in machine learning: from  
 254 phenomena to black-box attacks using adversarial samples,” *CoRR*, vol. abs/1605.07277, 2016.  
 255 [Online]. Available: <http://arxiv.org/abs/1605.07277>
- 256 [13] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,  
 257 N. Gimsheine, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani,  
 258 S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative  
 259 style, high-performance deep learning library,” in *Advances in Neural Information Processing*  
 260 *Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett,  
 261 Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: [http://papers.neurips.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
 262 [cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
- 263 [14] X. Dong and Y. Yang, “Searching for a robust neural architecture in four gpu hours,” in  
 264 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
 265 2019, pp. 1761–1770.
- 266 [15] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*,  
 267 2016.
- 268 [16] A. Krizhevsky, G. Hinton *et al.*, “Learning multiple layers of features from tiny images,”  
 269 University of Toronto, Tech. Rep., 2009.
- 270 [17] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image  
 271 recognition,” *arXiv preprint arXiv:1409.1556*, 2014.



- 272 [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional  
273 neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- 274 [19] A. Clyburne-Sherin, X. Fei, and S. A. Green, "Computational reproducibility via containers in  
275 social psychology," *Meta-Psychology*, vol. 3, 2019.
- 276 [20] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017*  
277 *IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 39–57.