

SparseFool: a few pixels make a big difference

Apostolos Modas, Seyed-Mohsen Moosavi-Dezfooli, Pascal Frossard
 École Polytechnique Fédérale de Lausanne

{apostolos.modas, seyed.moosavi, pascal.frossard}@epfl.ch

Abstract

Deep Neural Networks have achieved extraordinary results on image classification tasks, but have been shown to be vulnerable to attacks with carefully crafted perturbations of the input data. Although most attacks usually change values of many image’s pixels, it has been shown that deep networks are also vulnerable to sparse alterations of the input. However, no computationally efficient method has been proposed to compute sparse perturbations. In this paper, we exploit the low mean curvature of the decision boundary, and propose SparseFool, a geometry inspired sparse attack that controls the sparsity of the perturbations. Extensive evaluations show that our approach computes sparse perturbations very fast, and scales efficiently to high dimensional data. We further analyze the transferability and the visual effects of the perturbations, and show the existence of shared semantic information across the images and the networks. Finally, we show that adversarial training can only slightly improve the robustness against sparse additive perturbations computed with SparseFool.¹

1. Introduction

Deep neural networks (DNNs) are powerful learning models that achieve state-of-the-art performance in many different classification tasks [27, 50, 1, 22, 9], but have been shown to be vulnerable to very small, and often imperceptible, adversarial manipulations of their input data [47]. Interestingly, the existence of such adversarial perturbations is not only limited to additive perturbations [21, 12, 49] or classification tasks [8], but can be found in many other applications [48, 7, 28, 31, 39, 42, 43].

For image classification tasks, the most common type of adversarial perturbations are the ℓ_p minimization ones, since they are easier to analyze and optimize. Formally, for a given classifier and an image $\mathbf{x} \in \mathbb{R}^n$, we define as adversarial perturbation the minimal perturbation \mathbf{r} that changes

¹The code of SparseFool is available on <https://github.com/LTS4/SparseFool>, and Foolbox [41] <https://github.com/bethgelab/foolbox>.



Figure 1: Adversarial examples for ImageNet, as computed by SparseFool on a ResNet-101 architecture. Each column corresponds to different level of perturbed pixels. The fooling labels are shown below the images.

the classifier’s estimated label $k(\mathbf{x})$:

$$\min_{\mathbf{r}} \|\mathbf{r}\|_p \text{ s.t. } k(\mathbf{x} + \mathbf{r}) \neq k(\mathbf{x}), \quad (1)$$

Several methods (attacks) have been proposed for computing ℓ_2 and ℓ_∞ adversarial perturbations [47, 15, 24, 33, 6, 32]. However, understanding the vulnerabilities of deep neural networks in other perturbation regimes remains important. In particular, it has been shown [38, 45, 35, 2, 18] that DNNs can misclassify an image, when only a small fraction of the input is altered (sparse perturbations). In practice, sparse perturbations could correspond to some raindrops that reflect the sun on a “STOP” sign, but that are sufficient to fool an autonomous vehicle; or a crop-field with some sparse colorful flowers that force a UAV to spray pesticide on non affected areas. Understanding the vulnerability of deep networks to such a simple perturbation regime can further help to design methods to improve the robustness of deep classifiers.

Some prior work on sparse perturbations have been developed recently. The authors in [38] proposed the JSMA

method, which perturbs the pixels based on their saliency score. Furthermore, the authors in [45] exploit Evolutionary Algorithms (EAs) to achieve extremely sparse perturbations, while the authors in [35] finally proposed a black-box attack that computes sparse adversarial perturbations using a greedy local search algorithm.

However, solving the optimization problem in Eq. (1) in an ℓ_0 sense is NP-hard, and current algorithms are all characterized by high complexity. The resulting perturbations usually consist of high magnitude noise, concentrated over a small number of pixels. This makes them quite perceptible and in many cases, the perturbed pixels might even exceed the dynamic range of the image.

Therefore, we propose in this paper an efficient and principled way to compute sparse perturbations, while at the same time ensuring the validity of the perturbed pixels.

Our main contributions are the following:

- We propose SparseFool, a geometry inspired sparse attack that exploits the boundaries’ low mean curvature to compute adversarial perturbations efficiently.
- We show through extensive evaluations that (a) our method computes sparse perturbations much faster than the existing methods, and (b) it can scale efficiently to high dimensional data.
- We further propose a method to control the perceptibility of the resulted perturbations, while retaining the levels of sparsity and complexity.
- We analyze the visual features affected by our attack, and show the existence of some shared semantic information across different images and networks.
- We finally show that adversarial training using slightly lowers the vulnerability against sparse perturbations, but not enough to lead to more robust classifiers yet.

The rest of the paper is organized as follows: in Section 2, we describe the challenges and the problems for computing sparse adversarial perturbations. In Section 3, we provide an efficient method for computing sparse adversarial perturbations, by linearizing and solving the initial optimization problem. Finally, the evaluation and analysis of the computed sparse perturbations is provided in Section 4.

2. Problem description

2.1. Finding sparse perturbations

Most of the existing adversarial attack algorithms solve the optimization problem of Eq. (1) for $p = 2$ or ∞ , resulting in dense but imperceptible perturbations. For the case of sparse perturbations, the goal is to minimize the number of perturbed pixels required to fool the network, which

corresponds to minimizing $\|\mathbf{r}\|_0$ in Eq. (1). Unfortunately, this leads to NP-hard problems, for which reaching a global minimum cannot be guaranteed in general [3, 37, 40]. There exist different methods [34, 40] to avoid the computational burden of this problem, with the ℓ_1 relaxation being the most common; the minimization of $\|\mathbf{r}\|_0$ under linear constraints can be approximated by solving the corresponding convex ℓ_1 problem [5, 10, 36]². Thus, we are looking for an efficient way to exploit such a relaxation to solve the optimization problem in Eq. (1).

DeepFool [33] is an algorithm that exploits such a relaxation, by adopting an iterative procedure that includes a linearization of the classifier at each iteration, in order to estimate the minimal adversarial perturbation \mathbf{r} . Specifically, assuming f is a classifier, at each iteration i , f is linearized around the current point $\mathbf{x}^{(i)}$, the minimal perturbation $\mathbf{r}^{(i)}$ (in an ℓ_2 sense) is computed as the projection of $\mathbf{x}^{(i)}$ onto the linearized hyperplane, and the next iterate $\mathbf{x}^{(i+1)}$ is updated. One can use such a linearization procedure to solve Eq. (1) for $p = 1$, so as to obtain an approximation to the ℓ_0 solution. Thus, by generalizing the projection to ℓ_p norms ($p \in [1, \infty)$) and setting $p = 1$, ℓ_1 -DeepFool provides an efficient way for computing sparse adversarial perturbations using the ℓ_1 projection.

2.2. Validity of perturbations

Although the ℓ_1 -DeepFool efficiently computes sparse perturbations, it does not explicitly respect the constraint on the validity of the adversarial image values. When computing adversarial perturbations, it is very important to ensure that the pixel values of the adversarial image $\mathbf{x} + \mathbf{r}$ lie within the valid range of color images (e.g., $[0, 255]$). For ℓ_2 and ℓ_∞ perturbations, almost every pixel of the image is distorted with noise of small magnitude, so that most common algorithms usually ignore such constraints [33, 15]. In such cases, it is unlikely that many pixels will be out of their valid range; and even then, clipping the invalid values after the computation of such adversarial images have a minor impact.

This is unfortunately not the case for sparse perturbations however; solving the ℓ_1 optimization problem results in a few distorted pixels of high magnitude noise, and clipping the values after computing the adversarial image can have a significant impact on the success of the attack. In other words, as the perturbation becomes sparser, the contribution of each pixel is typically much stronger compared to ℓ_2 or ℓ_∞ perturbations.

We demonstrate the effect of such clipping operation on the quality of adversarial perturbations generated by ℓ_1 -DeepFool. For example, with perturbations computed for a VGG-16 [44] network trained on ImageNet [9], we ob-

²Under some conditions, the solution of such approximation is indeed optimal [4, 11, 16].

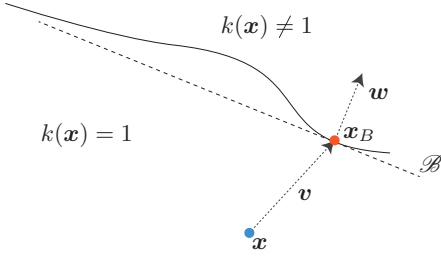


Figure 2: The approximated decision boundary \mathcal{B} in the vicinity of the datapoint \mathbf{x} that belongs to class 1. \mathcal{B} can be seen as a one-vs-all linear classifier for class 1.

served that ℓ_1 -DeepFool achieves almost 100% of fooling rate by perturbing only 0.037% of the pixels on average. However, clipping the pixel values of adversarial images to $[0, 255]$ results in a fooling rate of merely 13%. Furthermore, incorporating the clipping operator inside the iterative procedure of the algorithm does not improve the results. In other words, ℓ_1 -DeepFool fails to properly compute sparse perturbations. This underlies the need for an improved attack algorithm that natively takes into account the validity of generated adversarial images, as proposed in the next sections.

2.3. Problem formulation

Based on the above discussion, sparse adversarial perturbations are obtained by solving the following general form optimization problem:

$$\begin{aligned} & \underset{\mathbf{r}}{\text{minimize}} && \|\mathbf{r}\|_1 \\ & \text{subject to} && k(\mathbf{x} + \mathbf{r}) \neq k(\mathbf{x}) \\ & && \mathbf{l} \preceq \mathbf{x} + \mathbf{r} \preceq \mathbf{u}, \end{aligned} \quad (2)$$

where $\mathbf{l}, \mathbf{u} \in \mathbb{R}^n$ denote the lower and upper bounds of the values of $\mathbf{x} + \mathbf{r}$, such that $l_i \leq x_i + r_i \leq u_i$, $i = 1 \dots n$.

To find an efficient relaxation to problem (2), we focus on the geometric characteristics of the decision boundary, and specifically on its curvature. It has been shown [13, 14, 20] that the decision boundaries of state-of-the-art deep networks have a quite low mean curvature in the neighborhood of data samples. In other words, for a datapoint \mathbf{x} and its corresponding minimal ℓ_2 adversarial perturbation \mathbf{v} , the decision boundary at the vicinity of \mathbf{x} can be locally well approximated by a hyperplane passing through the datapoint $\mathbf{x}_B = \mathbf{x} + \mathbf{v}$, and a normal vector \mathbf{w} (see Fig. 2).

Hence, we exploit this property and linearize the optimization problem (2), so that sparse adversarial perturbations can be computed by solving the following box-

Algorithm 1: LinearSolver

Input: image \mathbf{x} , normal \mathbf{w} , boundary point \mathbf{x}_B , projection operator Q .

Output: perturbed point $\mathbf{x}^{(i)}$

- 1 Initialize: $\mathbf{x}^{(0)} \leftarrow \mathbf{x}$, $i \leftarrow 0$, $S = \{\}$
 - 2 **while** $\mathbf{w}^T(\mathbf{x}^{(i)} - \mathbf{x}_B) \neq 0$ **do**
 - 3 $\mathbf{r} \leftarrow \mathbf{0}$
 - 4 $d \leftarrow \arg \max_{j \notin S} |w_j|$
 - 5 $r_d \leftarrow \frac{|\mathbf{w}^T(\mathbf{x}^{(i)} - \mathbf{x}_B)|}{|w_d|} \cdot \text{sign}(w_d)$
 - 6 $\mathbf{x}^{(i+1)} \leftarrow Q(\mathbf{x}^{(i)} + \mathbf{r})$
 - 7 $S \leftarrow S \cup \{d\}$
 - 8 $i \leftarrow i + 1$
 - 9 **end**
 - 10 **return** $\mathbf{x}^{(i)}$
-

constrained optimization problem:

$$\begin{aligned} & \underset{\mathbf{r}}{\text{minimize}} && \|\mathbf{r}\|_1 \\ & \text{subject to} && \mathbf{w}^T(\mathbf{x} + \mathbf{r}) - \mathbf{w}^T \mathbf{x}_B = 0 \\ & && \mathbf{l} \preceq \mathbf{x} + \mathbf{r} \preceq \mathbf{u}. \end{aligned} \quad (3)$$

In the following section, we provide a method for solving the optimization problem (3), and introduce SparseFool, a fast yet efficient algorithm for computing sparse adversarial perturbations, which linearizes the constraints by approximating the decision boundary as an affine hyperplane.

3. Sparse adversarial perturbations

3.1. Linearized problem solution

In solving the optimization problem (3), the computation of the ℓ_1 projection of \mathbf{x} onto the approximated hyperplane does not guarantee a solution. For a perturbed image, consider the case where some of its values exceed the bounds defined by \mathbf{l} and \mathbf{u} . Thus, by readjusting the invalid values to match the constraints, the resulted adversarial image may eventually not lie onto the approximated hyperplane.

For this reason, we propose an iterative procedure, where at each iteration we project only towards one single coordinate of the normal vector \mathbf{w} at a time. If projecting \mathbf{x} towards a specific direction does not provide a solution, then the perturbed image at this coordinate has reached its extrema value. Therefore, at the next iteration, this direction should be ignored, since it cannot contribute any further to finding a better solution.

Formally, let S be a set containing all the directions of \mathbf{w} that cannot contribute to the minimal perturbation. Then,

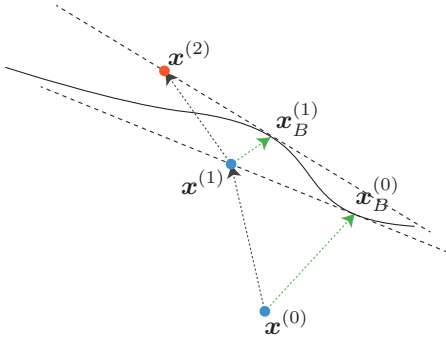


Figure 3: Illustration of SparseFool algorithm. With green we denote the ℓ_2 -DeepFool adversarial perturbations computed at each iteration. In this example, the algorithm converges after 2 iterations, and the total perturbation is $r = \mathbf{x}^{(2)} - \mathbf{x}^{(0)}$.

the minimal perturbation \mathbf{r} is updated through the ℓ_1 projection of the current datapoint $\mathbf{x}^{(i)}$ onto the estimated hyperplane as:

$$r_d \leftarrow \frac{|\mathbf{w}^T(\mathbf{x}^{(i)} - \mathbf{x}_B)|}{|w_d|} \cdot \text{sign}(w_d), \quad (4)$$

where d is the index of the maximum absolute value of \mathbf{w} that has not already been used

$$d \leftarrow \arg \max_{j \notin S} |w_j|. \quad (5)$$

Before proceeding to the next iteration, we must ensure the validity of the values of the next iterate $\mathbf{x}^{(i+1)}$. For this reason, we use a projection operator $Q(\cdot)$ that readjusts the values of the updated point that are out of bounds, by projecting $\mathbf{x}^{(i)} + \mathbf{r}$ onto the box-constraints defined by \mathbf{l} and \mathbf{u} . Hence, the new iterate $\mathbf{x}^{(i+1)}$ is updated as $\mathbf{x}^{(i+1)} \leftarrow Q(\mathbf{x}^{(i)} + \mathbf{r})$. Note here that the bounds \mathbf{l} , \mathbf{u} are not limited to only represent the dynamic range of an image, but can be generalized to satisfy any similar restriction. For example, as we will describe later in Section 4.2, they can be used to control the perceptibility of the computed adversarial images.

The next step is to check if the new iterate $\mathbf{x}^{(i+1)}$ has reached the approximated hyperplane. Otherwise, it means that the perturbed image at the coordinate d has reached its extrema value, and thus we cannot change it any further; perturbing towards the corresponding direction will have no effect. Thus, we reduce the search space, by adding to the forbidden set S the direction d , and repeat the procedure until we reach the approximated hyperplane. The algorithm for solving the linearized problem is summarized in Algorithm 1.

Algorithm 2: SparseFool

Input: image \mathbf{x} , projection operator Q , classifier f .

Output: perturbation \mathbf{r}

- 1 Initialize: $\mathbf{x}^{(0)} \leftarrow \mathbf{x}$, $i \leftarrow 0$
 - 2 **while** $k(\mathbf{x}^{(i)}) = k(\mathbf{x}^{(0)})$ **do**
 - 3 $\mathbf{r}_{\text{adv}} = \text{DeepFool}(\mathbf{x}^{(i)})$
 - 4 $\mathbf{x}_B^{(i)} = \mathbf{x}^{(i)} + \mathbf{r}_{\text{adv}}$
 - 5 $\mathbf{w}^{(i)} = \nabla f_{k(\mathbf{x}_B^{(i)})}(\mathbf{x}_B^{(i)}) - \nabla f_{k(\mathbf{x}^{(i)})}(\mathbf{x}_B^{(i)})$
 - 6 $\mathbf{x}^{(i+1)} = \text{LinearSolver}(\mathbf{x}^{(i)}, \mathbf{w}^{(i)}, \mathbf{x}_B^{(i)}, Q)$
 - 7 $i \leftarrow i + 1$
 - 8 **end**
 - 9 **return** $\mathbf{r} = \mathbf{x}^{(i)} - \mathbf{x}^{(0)}$
-

3.2. Finding the point \mathbf{x}_B and the normal \mathbf{w}

In order to complete our solution to the optimization problem (3), we now focus on the linear approximation of the decision boundary. Recall from Section 2.3 that we need to find the boundary point \mathbf{x}_B , along with the corresponding normal vector \mathbf{w} .

Finding \mathbf{x}_B is analogous to computing (in a ℓ_2 sense) an adversarial example of \mathbf{x} , so it can be approximated by applying one of the existing ℓ_2 attack algorithms. However, not all of these attacks are proper for our task; we need a fast method that finds adversarial examples that are as close to the original image \mathbf{x} as possible. Recall that DeepFool [33] iteratively moves \mathbf{x} towards the decision boundary, and stops as soon as the perturbed data point reaches the other side of the boundary. Therefore, the resulting perturbed sample usually lies very close to the decision boundary, and thus, \mathbf{x}_B can be very well approximated by $\mathbf{x} + \mathbf{r}_{\text{adv}}$, with \mathbf{r}_{adv} being the corresponding ℓ_2 -DeepFool perturbation of \mathbf{x} . Then, if we denote the network’s classification function as f , one can estimate the normal vector to the *decision boundary* at the datapoint \mathbf{x}_B as:

$$\mathbf{w} := \nabla f_{k(\mathbf{x}_B)}(\mathbf{x}_B) - \nabla f_{k(\mathbf{x})}(\mathbf{x}_B). \quad (6)$$

Hence, the decision boundary can now be approximated by the affine hyperplane $\mathcal{B} \triangleq \{\mathbf{x} : \mathbf{w}^T(\mathbf{x} - \mathbf{x}_B) = 0\}$, and sparse adversarial perturbations are computed by applying Algorithm 1.

3.3. SparseFool

However, although we expected to have a one-step solution, in many cases the algorithm does not converge. The reason behind this behavior lies on the fact that the decision boundaries of the networks are only *locally flat* [13, 14, 20]. Thus, if the ℓ_1 perturbation moves the datapoint \mathbf{x} away from its neighborhood, then the “local flatness” property

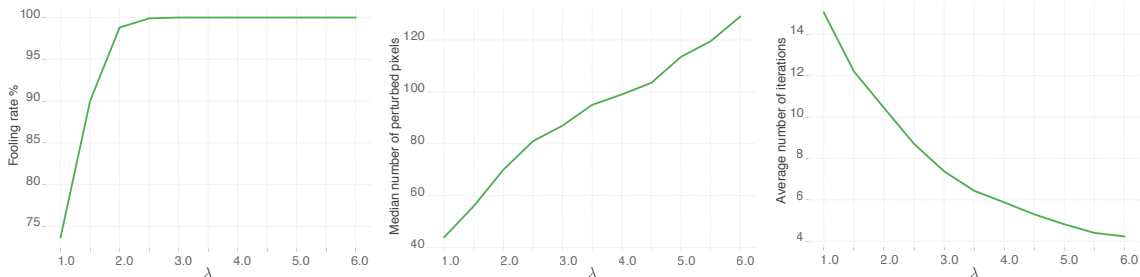


Figure 4: The fooling rate, the sparsity of the perturbations, and the average iterations of SparseFool for different values of λ , on 4000 images from ImageNet dataset using an Inception-v3 [46] model.

might be lost, and eventually the perturbed point will not reach the other side of the decision boundary (projected onto a “curved” area).

We mitigate the convergence issue with an iterative method, namely SparseFool, where each iteration includes the linear approximation of the decision boundary. Specifically, at iteration i , the boundary point $\mathbf{x}_B^{(i)}$ and the normal vector $\mathbf{w}^{(i)}$ are estimated using ℓ_2 -DeepFool based on the current iterate $\mathbf{x}^{(i)}$. Then, the next iterate $\mathbf{x}^{(i+1)}$ is updated through the solution of Algorithm 1, having though $\mathbf{x}^{(i)}$ as the initial point, and the algorithm terminates when $\mathbf{x}^{(i)}$ changes the label of the network. An illustration of SparseFool is given in Fig. 3, and the algorithm is summarized in Algorithm 2.

However, we observed that instead of using the boundary point $\mathbf{x}_B^{(i)}$ at the step 6 of SparseFool, better convergence (misclassification) can be achieved by going further into the other side of the boundary, and find a solution for the hyperplane passing through the datapoint $\mathbf{x}^{(i)} + \lambda(\mathbf{x}_B^{(i)} - \mathbf{x}^{(i)})$, where $\lambda \geq 1$. Specifically, as shown in Fig. 4, this parameter is used to control the trade-off between the fooling rate, the sparsity, and the complexity. Values close to 1, lead to sparser perturbations, but also to lower fooling rate and increased complexity. On the contrary, higher values of λ lead to fast convergence – even one step solutions –, but the resulted perturbations are less sparse. Since λ is the *only parameter* of the algorithm, it can be easily adjusted to meet the corresponding needs in terms of fooling rate, sparsity, and complexity.

Finally, note that \mathcal{B} corresponds to the boundary between the adversarial and the estimated true class, and thus it can be seen as an affine binary classifier. Since at each iteration the adversarial class is computed as the closest (in an ℓ_2 sense) to the true one, we can say that SparseFool operates as an untargeted attack. Even though, it can be easily transformed to a targeted one, by simply computing at each iteration the adversarial example – and thus approximating the decision boundary – of a specific category.

4. Experimental results

4.1. Setup

We test SparseFool on deep convolutional neural network architectures with the 10000 images of the MNIST [26] test set, 10000 images of the CIFAR-10 [23] test set, and 4000 randomly selected images from the ILSVRC2012 validation set. In order to evaluate our algorithm and compare with related works, we compute the fooling rate, the median perturbation percentage, and the average execution time. Given a dataset \mathcal{D} , the fooling rate measures the efficiency of the algorithm, using the formula: $|\{\mathbf{x} \in \mathcal{D} : k(\mathbf{x} + \mathbf{r}_\mathbf{x}) \neq k(\mathbf{x})\}|/|\mathcal{D}|$, where $\mathbf{r}_\mathbf{x}$ is the perturbation that corresponds to the image \mathbf{x} . The median perturbation percentage corresponds to the median percentage of the pixels that are perturbed per fooling sample, while the average execution time measures the average execution time of the algorithm per sample³.

We compare SparseFool with the implementation of JSMA attack [38]. Since JSMA is a targeted attack, we evaluate it on its “untargeted” version, where the target class is chosen at random. We also make one more modification at the success condition; instead of checking if the predicted class is equal to the target one, we simply check if it is different from the source one. Let us note that JSMA is not evaluated on ImageNet dataset, due to its huge computational cost for searching over all pairs of candidates [6]. We also compare SparseFool with the “one-pixel attack” proposed in [45]. Since “one-pixel attack” perturbs exactly κ pixels, for every image we start with $\kappa = 1$ and increase it till “one-pixel attack” finds an adversarial example. Again, we do not evaluate the “one-pixel attack” on the ImageNet dataset, due to its high computational cost for high dimensional images.

4.2. Results

Overall performance. We first evaluate the performance of SparseFool, JSMA, and “one-pixel attack” on differ-

³All the experiments were conducted on a GTX TITAN X GPU.

ent datasets and architectures. The control parameter λ in SparseFool was set to 1 and 3 for the MNIST and CIFAR-10 datasets respectively. We observe in Table 1 that SparseFool computes 2.9x sparser perturbations, and is 4.7x faster compared to JSMA for the MNIST dataset. This behavior remains similar for the CIFAR-10 dataset, where SparseFool computes on average, perturbations of 2.4x higher sparsity, and is 15.5x faster. Notice here the difference in the execution time: JSMA becomes much slower as the dimension of the input data increases, while SparseFool’s time complexity remains at very low levels.

Then, in comparison to “one-pixel attack”, we observe that for the MNIST dataset, our method computes 5.5x sparser perturbations, and is more than 3 orders of magnitude faster. For the CIFAR-10 dataset, SparseFool still manages to find very sparse perturbations, but less so than the “one-pixel attack” in this case. The reason is that our method does not solve the original ℓ_0 optimization problem, but it rather computes sparse perturbations through the ℓ_1 solution of the linearized one. The resulting solution is often suboptimal, and may be optimal when the datapoint is very close to the boundary, where the linear approximation is more accurate. However, solving our linearized problem is fast, and enables our method to efficiently scale to high dimensional data, which is not the case for the “one-pixel attack”. Considering the tradeoff between the sparsity of the solutions and the required complexity, we choose to sacrifice the former, rather than following a complex exhaustive approach like [45]. In fact, our method is able to compute sparse perturbations 270x faster, and by requiring 2 orders of magnitude less queries to the network, than the “one-pixel attack” algorithm.

Finally, due to the huge computational cost of both JSMA and “one-pixel attack”, we do not use it for the large ImageNet dataset. In this case, we instead compare SparseFool with an algorithm that randomly selects a subset of elements from each color channel (RGB), and replaces their intensity with a random value from the set $V = \{0, 255\}$. The cardinality of each channel subset is constrained to match SparseFool’s per-channel median number of perturbed elements; for each channel, we select as many elements as the median, across all images, of SparseFool’s perturbed elements for this channel. The performance of SparseFool for the ImageNet dataset is reported in Table 2, while the corresponding fooling rates for the random algorithm were 18.2%, 13.2%, 14.5%, and 9.6% respectively. Observe that the fooling rates obtained for the random algorithm are far from comparable with SparseFool’s, indicating that the proposed algorithm cleverly finds very sparse solutions. Indeed, our method is consistent among different architectures, perturbing on average 0.21% of the pixels, with an average execution time of 7 seconds per sample.

To the best of our knowledge, we are the first to provide

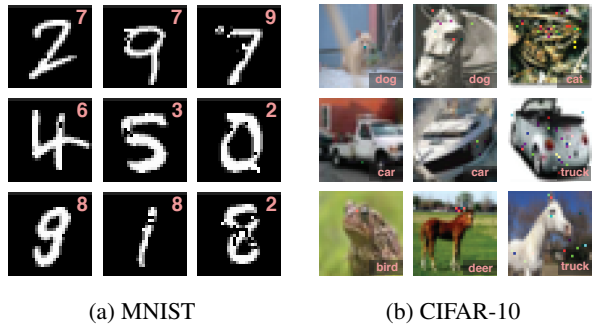


Figure 5: Adversarial examples for (a) MNIST and (b) CIFAR-10 datasets, as computed by SparseFool on LeNet and ResNet-18 architectures respectively. Each column corresponds to different level of perturbed pixels.

an adequate sparse attack that efficiently achieves such fooling rates and sparsity, and at the same time scales to high dimensional data. “One-pixel attack” does not necessarily find good solutions for all the studied datasets, however, SparseFool – as it relies on the high dimensional geometry of the classifiers – successfully computes sparse enough perturbations for all three datasets.

Perceptibility. In this section, we illustrate some adversarial examples generated by SparseFool, for three different levels of sparsity; highly sparse perturbations, sparse perturbations, and somewhere in the middle. For the MNIST and CIFAR-10 datasets (Figures 5a and 5b respectively), we observe that for highly sparse cases, the perturbation is either imperceptible or so sparse (i.e., 1 pixel) that it can be easily ignored. However, as the number of perturbed pixels increases, the distortion becomes even more perceptible, and in some cases the noise is detectable and far from imperceptible. The same behavior is observed for the ImageNet dataset (Fig. 1). For highly sparse perturbations, the noise is again either imperceptible or negligible, but as the density increases, it becomes more and more visible.

To eliminate this perceptibility effect, we focus on the lower and upper bounds of the values of an adversarial image \hat{x} . Recall from Section 2.3 that the bounds l, u are defined in way such that $l_i \leq \hat{x}_i \leq u_i, i = 1 \dots n$. If these bounds represent the dynamic range of the image, then \hat{x}_i can take every possible value from this range, and the magnitude of the noise at the element i can reach some visible levels. However, if the perturbed values of each element lie close to the original values x_i , then we might prevent the magnitude from reaching very high levels. For this reason, assuming a dynamic range of $[0, 255]$, we explicitly constrain the values of \hat{x}_i to lie in a small interval $\pm\delta$ around x_i , such that $0 \leq x_i - \delta \leq \hat{x}_i \leq x_i + \delta \leq 255$.

The resulted sparsity for different values of δ is depicted

Dataset	Network	Acc. (%)	Fooling rate (%)			Perturbation (%)			Time (sec)		
			SF	JSMA	1-PA	SF	JSMA	1-PA	SF	JSMA	1-PA
MNIST	LeNet [25]	99.14	99.93	95.73	100	1.66	4.85	9.43	0.14	0.66	310.2
CIFAR-10	VGG-19	92.71	100	98.12	100	1.07	2.25	0.15	0.34	6.28	102.7
	ResNet18 [17]	92.74	100	100	100	1.27	3.91	0.2	0.69	8.73	167.4

Table 1: The performance of SparseFool (SF), JSMA [38], and “one-pixel attack” (1-PA) [45] on the MNIST and the CIFAR-10 datasets.³ Note that due to its high complexity, “one-pixel attack” was evaluated on only 100 samples.

Network	Acc. (%)	Fooling rate (%)	Pert. (%)	Time (sec)
VGG-16	71.59	100	0.18	5.09
ResNet-101	77.37	100	0.23	8.07
DenseNet-161	77.65	100	0.29	10.07
Inception-v3	77.45	100	0.14	4.94

Table 2: The performance of SparseFool on the ImageNet dataset, using the pre-trained models provided by PyTorch.³

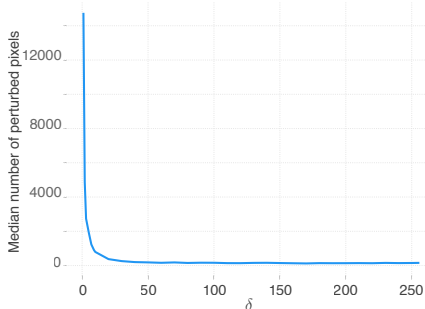


Figure 6: The resulted sparsity of SparseFool perturbations for $\pm\delta$ around the values of x , for 100 samples from ImageNet on a ResNet-101 architecture.

in Fig. 6. The higher the value, the more freedom we give to the perturbations, and for $\delta = 255$ we exploit the whole dynamic range. Observe though that after $\delta \approx 25$, the sparsity levels seem to remain almost constant, which indicates that we do not need to use the whole dynamic range. Furthermore, we observed that the average execution time per sample of SparseFool from this value onward seems to remain constant as well, while the fooling rate is always 100% regardless δ . Thus, by selecting appropriate values for δ , we can control the perceptibility of the resulted perturbations, retain the sparsity around a sufficient level. The influence of δ on the perceptibility and the sparsity of the perturbations is demonstrated in Fig. 7.

Transferability and semantic information. We now study if SparseFool adversarial perturbations can generalize across different architectures. For the VGG-16, ResNet-

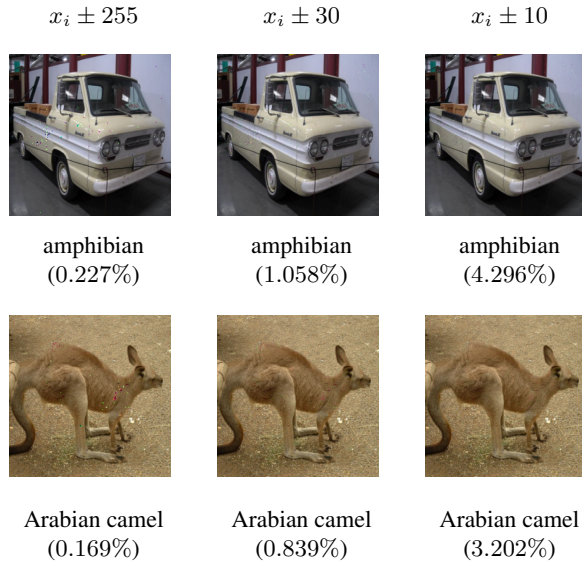


Figure 7: The effect of δ on the perceptibility and the sparsity of SparseFool perturbations. The values of δ are shown on top of each column, while the fooling label and the percentage of perturbed pixels are written below each image.

101, and DenseNet-161 [19] architectures, we report in Table 3 the fooling rate of each model when fed with adversarial examples generated by another one. We observe that sparse adversarial perturbations can generalize only to some extent, and also, as expected [29], they are more transferable from larger to smaller architectures. This indicates that there should be some shared semantic information between different architectures that SparseFool exploits, but the perturbations are mostly network dependent.

Focusing on some animal categories of the ImageNet dataset, we observe that indeed the perturbations are concentrated around “important” areas (i.e., head), but there is not a consistent pattern to indicate specific features that are the most important for the network (i.e., eyes, ears, nose etc.); in many cases the noise also spreads around different parts of the image. Examining now if semantic information is shared across different architectures (Fig. 8), we observe that in all the networks, the noise consistently lies around

	VGG16	ResNet101	DenseNet161
VGG16	100%	10.8%	8.2%
ResNet101	25.3%	100%	12.1%
DenseNet161	28.2%	17.5%	100%

Table 3: Fooling rates of SparseFool perturbations between pairs of models for 4000 samples from ImageNet. The row and column denote the source and target model respectively.

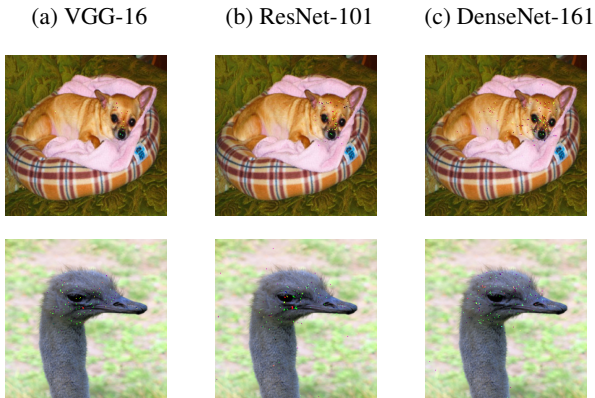


Figure 8: Shared information of the SparseFool perturbations for ImageNet, as computed on three different architectures. For all networks, the first row image was classified as “Chihuahua” and misclassified as “French Bulldog”, and second row image as “Ostrich” and “Crane” respectively.

the important areas of the image, but the way it concentrates or spreads is different for each network.

For the CIFAR-10 dataset, we observe that in many cases of animal classes, SparseFool tends to perturb some universal features around the area of the head (i.e., eyes, ears, nose, mouth etc.), as shown in Fig. 9a. Furthermore, we tried to understand if there is a correlation between the perturbed pixels and the fooling class. Interestingly, we observed that in many cases, the algorithm was perturbing those regions of the image that correspond to important features of the fooling class, i.e., when changing a “bird” label to a “plane”, where the perturbation seems to represent some parts of the plane (i.e., wings, tail, nose, turbine). This behavior becomes even more obvious when the fooling label is a “deer”, where the noise lies mostly around the area of the head in a way that resembles the antlers.

Robustness of adversarially trained network. Finally, we study the robustness to sparse perturbations of an adversarially trained ResNet-18 network on the CIFAR-10 dataset, using ℓ_∞ perturbations. The accuracy of this more robust network is 82.17%, while the training procedure and its overall performance are similar to the ones provided in [30]. Compared to the results of Table 1, the average time dropped to 0.3 sec, but the perturbation percent-



(a) Universal features (b) Fooling class features

Figure 9: Semantic information of SparseFool perturbations for the CIFAR-10 dataset, on a ResNet-18 architecture. Observe that the perturbation is concentrated on (a) some features around the area of the face, and (b) on areas that are important for the fooling class.

age increased to 2.44%. In other words, the adversarially trained network leads to just a slight change in the sparsity, and thus training it on ℓ_∞ perturbations does not significantly improve its robustness against sparse perturbations.

5. Conclusion

Computing adversarial perturbations beyond simple ℓ_p norms is a challenging problem from different aspects. For sparse perturbations, apart from the NP-hardness of the ℓ_0 minimization, one needs to ensure the validity of the adversarial example values. In this work, we provide a novel geometry inspired sparse attack algorithm that is fast and can scale to high dimensional data, where one can also have leverage over the sparsity of the resulted perturbations. Furthermore, we provide a simple method to improve the perceptibility of the perturbations, while retaining the levels of sparsity and complexity. We also note that for some datasets, the proposed sparse attack alters features that are shared among different images. Finally, we show that adversarial training does not significantly improve the robustness against sparse perturbations computed with SparseFool. We believe that the proposed approach can be used to further understand the behavior and the geometry of deep image classifiers, and provide insights for building more robust networks.

Acknowledgements

We thank Mireille El Gheche and Stefano D’Aronco for the fruitful discussions. This work has been supported by a Google Faculty Research Award, and the Hasler Foundation, Switzerland, in the framework of the ROBERT project. We also gratefully acknowledge the support of NVIDIA Corporation with the donation of the GTX Titan X GPU used for this research.

References

- [1] S. Abu-El-Haija, N. Kothari, J. Lee, P. Natsev, G. Toderici, B. Varadarajan, and S. Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016.
- [2] A. Bibi, M. Alfadly, and B. Ghanem. Analytic expressions for probabilistic moments of pl-dnn with gaussian input. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [3] T. Blumensath and M. E. Davies. Iterative thresholding for sparse approximations. *Journal of Fourier Analysis and Applications*, 14(5):629–654, 2008.
- [4] E. Candès, M. Rudelson, T. Tao, and R. Vershynin. Error correction via linear programming. *46th Annual IEEE Symposium on Foundations of Computer Science (FOCS'05)*, pages 668–681, 2005.
- [5] E. Candès and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.
- [6] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. *IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2017.
- [7] N. Carlini and D. Wagner. Audio adversarial examples: Targeted attacks on speech-to-text. *IEEE Security and Privacy Workshops (SPW)*, 2018.
- [8] M. Cisse, Y. Adi, N. Neverova, and J. Keshet. Houdini: Fooling deep structured prediction models. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 6977–6987, 2017.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F.-F. Li. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [10] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.
- [11] D. L. Donoho and M. Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization. *National Academy of Sciences*, 100(5):2197–2202, 2003.
- [12] A. Fawzi and P. Frossard. Manitest: Are classifiers really invariant? *British Machine Vision Conference (BMVC)*, 2015.
- [13] A. Fawzi, S.-M. Moosavi-Dezfooli, and P. Frossard. The robustness of deep networks: A geometrical perspective. *IEEE Signal Processing Magazine*, 34(6):50–62, 2017.
- [14] A. Fawzi, S.-M. Moosavi-Dezfooli, P. Frossard, and S. Soatto. Empirical study of the topology and geometry of deep networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [15] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations (ICLR)*, 2015.
- [16] R. Gribonval and M. Nielsen. Sparse representations in unions of bases. *IEEE Transactions on Information Theory*, 49(12):3320–3325, 2003.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [18] M. Hein and M. Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial manipulation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [19] G. Huang, Z. Liu, L. Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [20] S. Jetley, N. A. Lord, and P. H. Torr. With friends like these, who needs adversaries? *arXiv preprint arXiv:1807.04200*, 2018.
- [21] C. Kanbak, S.-M. Moosavi-Dezfooli, and P. Frossard. Geometric robustness of deep networks: analysis and improvement. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [22] I. Krasin, T. Duerig, N. Alldrin, A. Veit, S. Abu-El-Haija, S. Belongie, D. Cai, Z. Feng, V. Ferrari, V. Gomes, A. Gupta, D. Narayanan, C. Sun, G. Chechnik, and K. Murphy. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2016.
- [23] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [24] A. Kurakin, I. J. Goodfellow, and S. Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- [25] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [26] Y. LeCun and C. Cortes. Mnist handwritten digits database. *Dataset available from <http://yann.lecun.com/exdb/mnist/>*, 2010.
- [27] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *The European Conference on Computer Vision (ECCV)*, pages 740–755, 2014.
- [28] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and M. Sun. Tactics of adversarial attack on deep reinforcement learning agents. *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3756–3762, 2017.
- [29] Y. Liu, X. Chen, C. Liu, and D. Song. Delving into transferable adversarial examples and black-box attacks. *arXiv preprint arXiv:1611.02770*, 2016.
- [30] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [31] J. H. Metzen, M. C. Kumar, T. Brox, and V. Fischer. Universal adversarial perturbations against semantic image segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2774–2783, 2017.
- [32] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard. Universal adversarial perturbations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [33] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

- [34] M. Nagahara, D. E. Quevedo, and J. Ostergaard. Sparse packetized predictive control for networked control over erasure channels. *IEEE Transactions on Automatic Control*, 59(7):1899–1905, 2014.
- [35] N. Narodytska and S. P. Kasiviswanathan. Simple black-box adversarial attacks on deep neural networks. *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017.
- [36] B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM Journal on Computing*, 24(2):227–234, 1995.
- [37] M. Nikolova. Description of the minimizers of least squares regularized with ℓ_0 -norm. uniqueness of the global minimizer. *SIAM Journal on Imaging Sciences*, 6(2):904–937, 2013.
- [38] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami. The limitations of deep learning in adversarial settings. *IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 372–387, 2016.
- [39] N. Papernot, P. McDaniel, A. Swami, and R. Harang. Crafting adversarial input sequences for recurrent neural networks. *IEEE Military Communications Conference (MILCOM)*, pages 49–54, 2016.
- [40] A. Patrascu and I. Necoara. Random coordinate descent methods for ℓ_0 regularized convex optimization. *IEEE Transactions on Automatic Control*, 60(7):1811–1824, 2015.
- [41] J. Rauber, W. Brendel, and M. Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models. *arXiv preprint arXiv:1707.04131*, 2017.
- [42] A. Rozsa, M. Günther, E. M. Rudd, and T. E. Boulton. Are facial attributes adversarially robust? *International Conference on Pattern Recognition (ICPR)*, pages 3121–3127, 2016.
- [43] A. Rozsa, M. Günther, E. M. Rudd, and T. E. Boulton. Facial attributes: Accuracy and adversarial robustness. *Pattern Recognition Letters*, 2017.
- [44] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, 2015.
- [45] J. Su, D. V. Vargas, and S. Kouichi. One pixel attack for fooling deep neural networks. *arXiv preprint arXiv:1710.08864*, 2017.
- [46] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [47] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations (ICLR)*, 2014.
- [48] P. Tabacof, J. Tavares, and E. Valle. Adversarial images for variational autoencoders. *arXiv preprint arXiv:1612.00155*, 2016.
- [49] C. Xiao, J.-Y. Zhu, B. Li, W. He, M. Liu, and D. Song. Spatially transformed adversarial examples. In *International Conference on Learning Representations (ICLR)*, 2018.
- [50] Y. Zhang, K. Lee, and H. Lee. Augmenting supervised neural networks with unsupervised objectives for large-scale image classification. In *International Conference on Machine Learning (ICML)*, volume 48, pages 612–621, 2016.

A. SparseFool adversarial examples

In this section, we provide some supplementary adversarial examples generated by SparseFool for different datasets. These examples correspond to three different levels of sparsity: (a) highly sparse, (b) very sparse, and (c) sparse perturbations. The corresponding results for the ImageNet, CIFAR-10, and MNIST datasets are shown in Figures [10 – 12] respectively. Observe that for highly and very sparse perturbations, the noise is either imperceptible, or sparse enough to be negligible. However, as the number of perturbed pixels is increased, the noise may become quite perceptible.

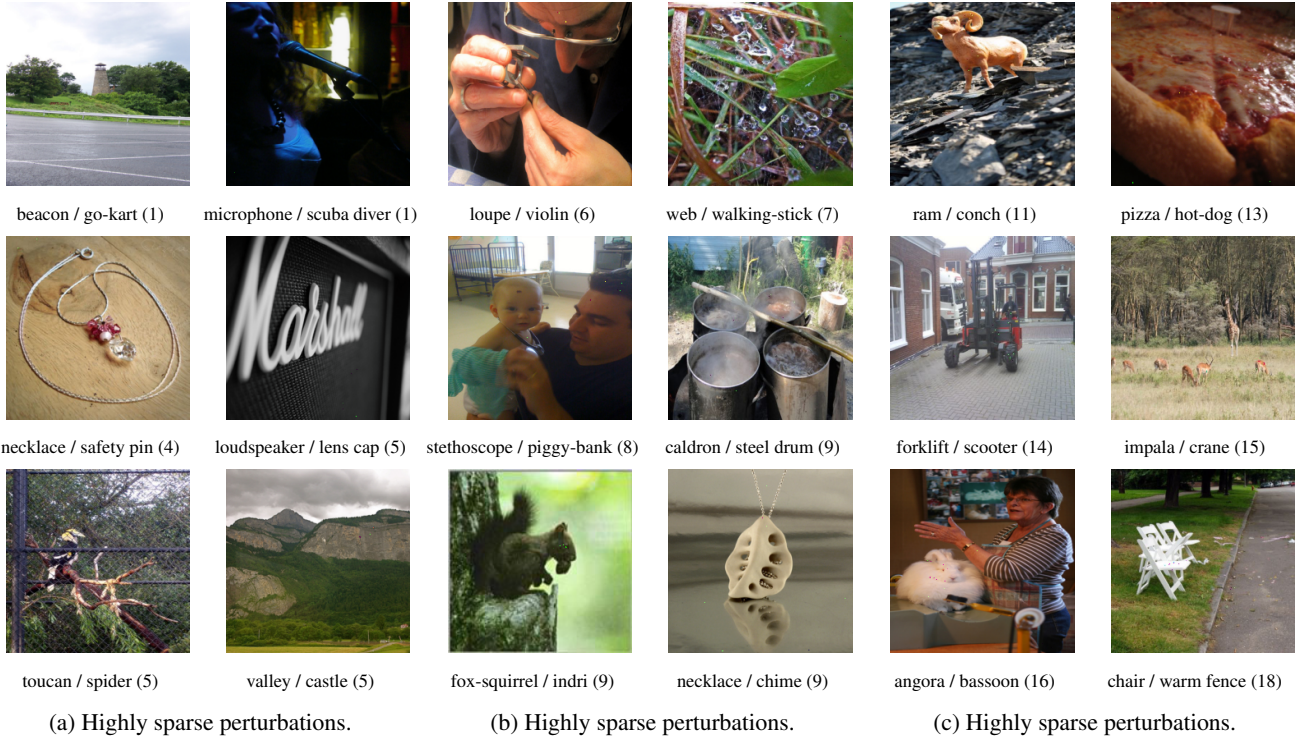


Figure 10: SparseFool adversarial examples for the ImageNet dataset for different levels of sparsity. The predicted label is shown above the image, the fooling one below, and the number of perturbed pixels is written inside the parentheses.

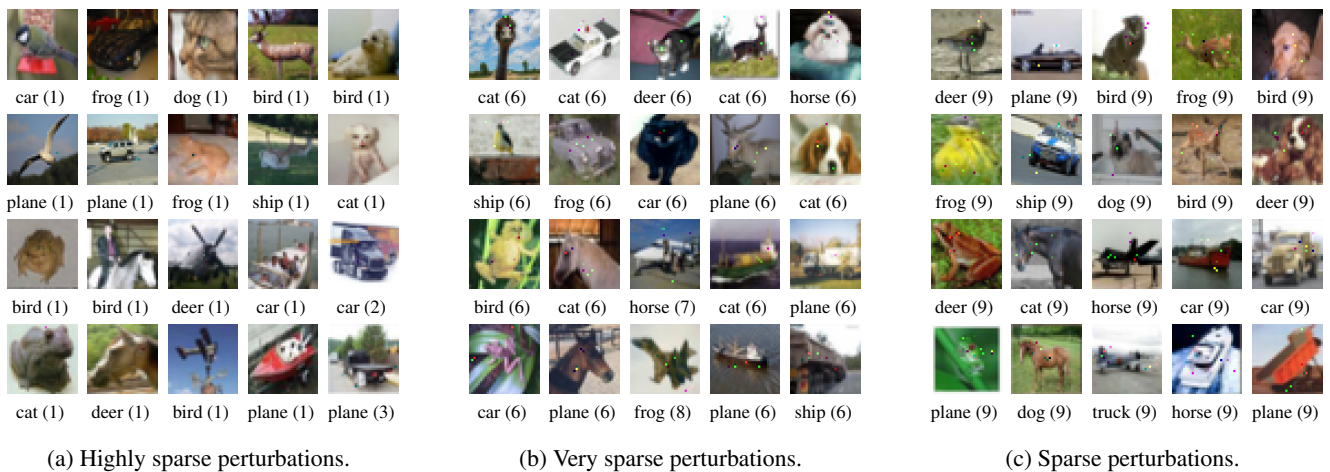


Figure 11: SparseFool adversarial examples for the CIFAR-10 dataset for different levels of sparsity. The fooling label is shown below the image, and the number of perturbed pixels is written inside the parentheses.

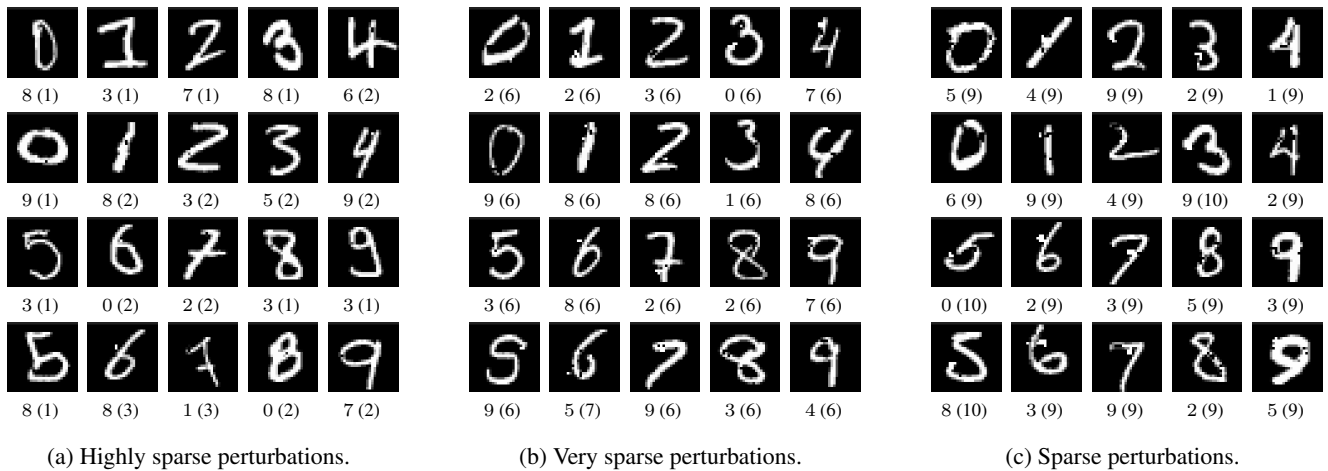


Figure 12: SparseFool adversarial examples for the MNIST dataset for different levels of sparsity. The fooling label is shown below the image, and the number of perturbed pixels is written inside the parentheses.

B. Controlling the perceptibility of the perturbations

We now present in Figure 13 some adversarial examples on the ImageNet dataset, when we control the perceptibility of the perturbations computed by SparseFool. This can be done by properly constraining the values of the perturbed image to lie $\pm\delta$ around the values of the original image.

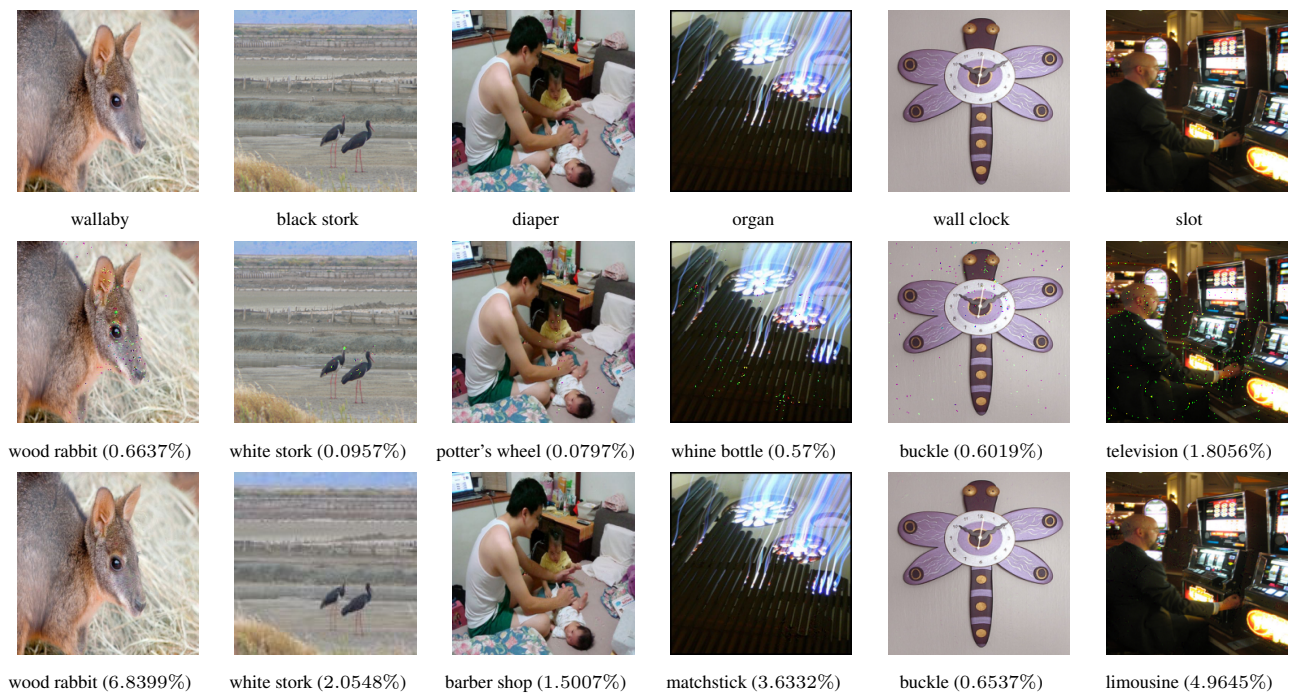


Figure 13: Controlling the perceptibility of the perturbations computed by SparseFool. First row: the original images and the predicted labels below them. Second row: adversarial examples generated by SparseFool, using the whole dynamic range. Third row: adversarial examples generated by SparseFool, constraining the noise ± 10 around the image values. For the adversarial examples, the fooling label is shown below the images, and the percentage of perturbed pixels is written inside the parentheses. Note that the fooling label might change, since SparseFool is operating as an untargeted attack.

C. The control parameter λ

In this section we provide the effect of the control parameter λ on (a) the fooling rate, (b) the sparsity of the perturbations, and (c) the convergence of SparseFool. The effect of different values of λ for different networks on the MNIST, CIFAR-10, and ImageNet datasets is shown in Figures [14 – 16] respectively.

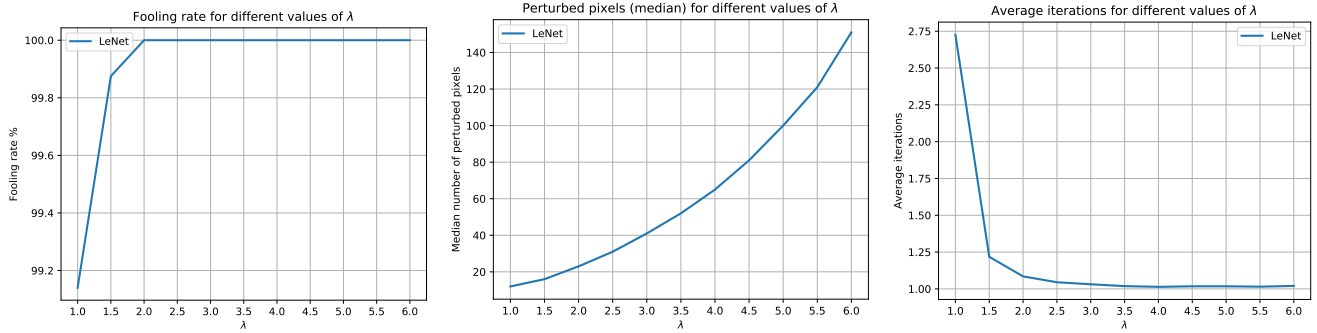


Figure 14: The effect of λ on SparseFool for a LeNet model trained on the MNIST dataset.

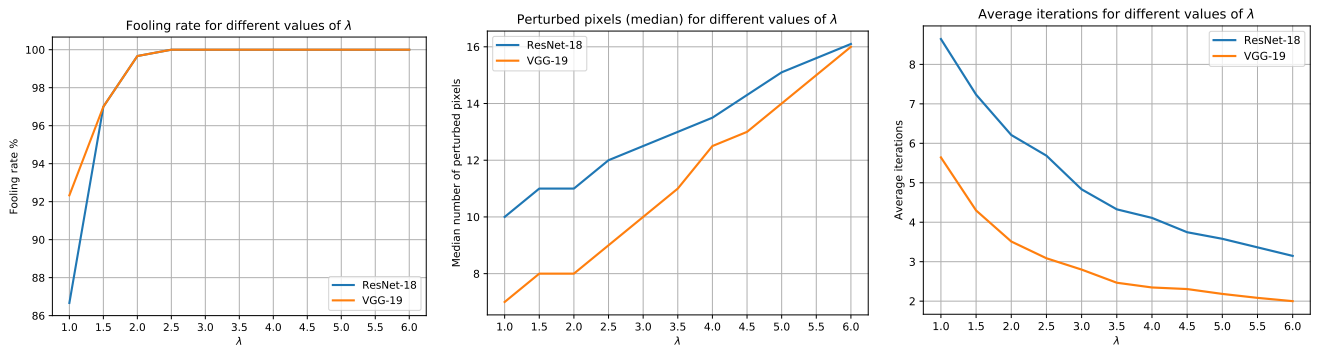


Figure 15: The effect of λ on SparseFool for different networks trained on the CIFAR-10 dataset.

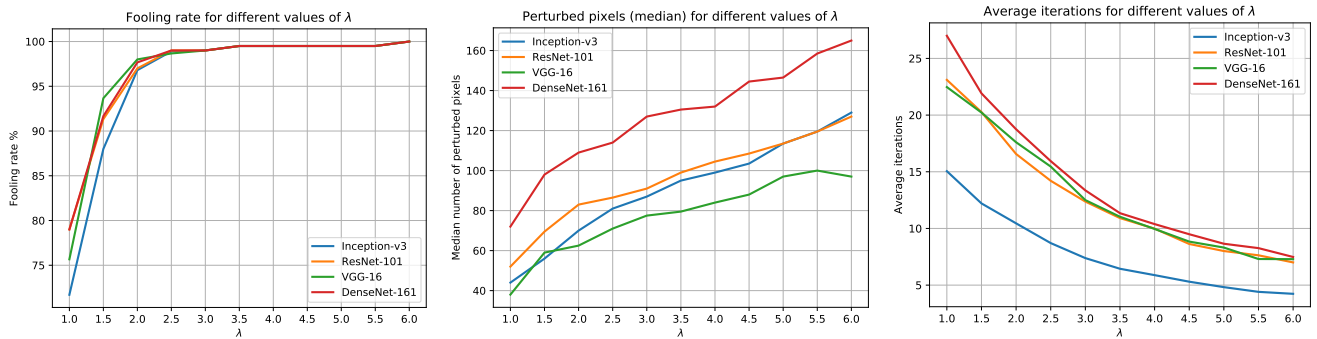


Figure 16: The effect of λ on SparseFool for different networks trained on the ImageNet dataset.

D. Adversarial examples compared to related methods

In this section we illustrate adversarial examples generated by SparseFool, compared to the corresponding ones computed by JSMA and “One-pixel attack”. The results obtained for the MNIST and CIFAR-10 datasets are depicted in Figure 17 and Figure 18 respectively.

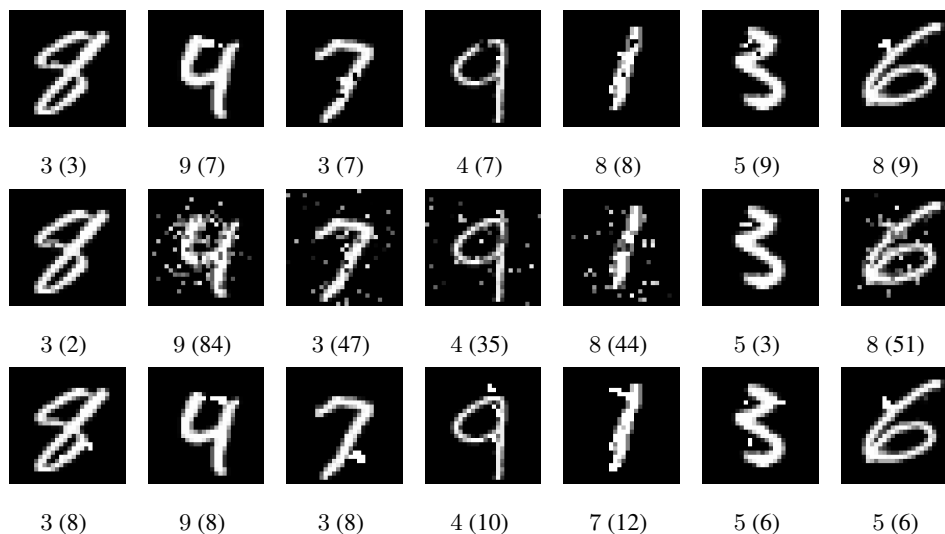


Figure 17: MNIST adversarial examples generated by (a) SparseFool (first row), (b) “One pixel attack” (second row), and (c) JSMA (third row). The fooling label is shown below each image, and the number of perturbed pixels is written inside the parentheses.

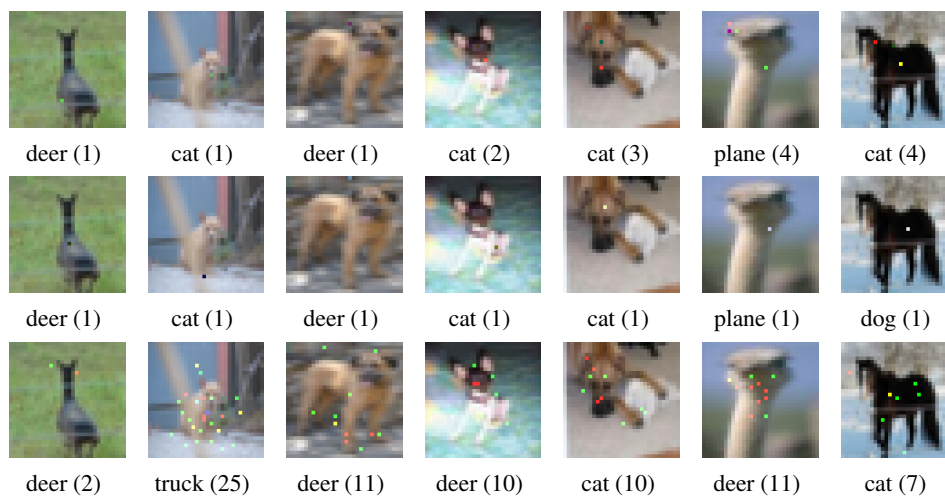


Figure 18: CIFAR-10 adversarial examples generated by (a) SparseFool (first row), (b) “One pixel attack” (second row), and (c) JSMA (third row). The fooling label is shown below each image, and the number of perturbed pixels is written inside the parentheses.