

Semantic Adversarial Attacks: Parametric Transformations That Fool Deep Classifiers

Ameya Joshi Amitangshu Mukherjee Soumik Sarkar Chinmay Hegde*

Iowa State University

{ameya, amimukh, soumiks, chinmay}@iastate.edu

Abstract

Deep neural networks have been shown to exhibit an intriguing vulnerability to adversarial input images corrupted with imperceptible perturbations. However, the majority of adversarial attacks assume global, fine-grained control over the image pixel space. In this paper, we consider a different setting: what happens if the adversary could only alter specific attributes of the input image? These would generate inputs that might be perceptibly different, but still natural-looking and enough to fool a classifier. We propose a novel approach to generate such “semantic” adversarial examples by optimizing a particular adversarial loss over the range-space of a parametric conditional generative model. We demonstrate implementations of our attacks on binary classifiers trained on face images, and show that such natural-looking semantic adversarial examples exist. We evaluate the effectiveness of our attack on synthetic and real data, and present detailed comparisons with existing attack methods. We supplement our empirical results with theoretical bounds that demonstrate the existence of such parametric adversarial examples.

1. Introduction

The existence of adversarial inputs for deep neural network-based classifiers has been well established by several recent works [5, 10, 16, 17, 58, 41]. The adversary typically confounds the classifier by adding an *imperceptible* perturbation to a given input image, where the range of the perturbation is defined in terms of bounded pixel-space ℓ_p -norm balls. Such adversarial “attacks” appear to catastrophically affect the performance of state-of-the-art classifiers [1, 22, 23, 54].

Pixel-space norm-constrained attacks reveal interesting

*This work was supported in part by NSF grants CCF-1750920, CNS-1845969, DARPA AIRA grant PA-18-02-02, AFOSR YIP Grant FA9550-17-1-0220, an ERP grant from ISU, a GPU gift grant from NVIDIA corp., and faculty fellowships from the Black and Veatch Foundation.



Figure 1. Examples of semantic adversarial attacks with a single modifiable attribute. The first and third columns are original images. Semantic adversarial examples (Columns 2 and 4) are generated by optimizing over the manifold of parametric (attribute) generative models to fool deep classifiers; specifically for adversarial facial attributes for Col. 2 and illumination for Col. 4.

insights about generalization properties of deep neural networks. However, imperceptible attacks are certainly not the only means available to an adversary. Consider an input example that comprises salient, invariant features along with modifiable attributes. An example would be an image of a face, which consists of invariant features relevant to the identity of the person, and variable attributes such as hair color and presence/absence of eyeglasses. Such adversarial examples, though perceptually distinct from the original input, appear natural and acceptable to an oracle or a human observer but would still be able to subvert the classifier. Unfortunately, the large majority of adversarial attack methods do not port over to such natural settings.

A systematic study of such attacks is paramount in safety-critical applications that deploy neural classifiers, such as face-recognition systems or vision modules of autonomous vehicles. These systems are required to be immune to a limited amount of variability in input data, particularly when these variations are achieved through natural means. Therefore, a method to generate adversarial examples using natural perturbations, such as facial attributes in the case of face images, or different weather conditions for autonomous navigation systems, would shed further in-

sights into the real-world robustness of such systems. We refer to such perceptible attacks as “semantic” attacks.

This setting fundamentally differs from existing attack approaches and has been (largely) unexplored thus far. Semantic attacks utilize nonlinear generative transformations of an input image instead of linear, additive techniques (such as image blending). Such complicated generative transformations display higher degrees of nonlinearity in corresponding attacks, the effects of which warrant further investigation. In addition, the role of the number of modifiable attributes (parameters in the generative models) in the given input is also an important point of consideration.

Contributions: We propose and rigorously analyze a framework for generating adversarial examples for a deep neural classifier by modifying *semantic* attributes.

We leverage generative models such as Fader Networks [30] that have semantically meaningful, tunable attributes corresponding to parameters into a continuous bounded space that implicitly define the space of “natural” input data. Our approach exploits this property by treating the range space of these attribute models as a manifold of semantic transformations of an image.

We pose the search for adversarial examples on this *semantic manifold* as an optimization problem over the parameters conditioning the generative model. Using face image classification as a running test case, we train a variety of parametric models (including Fader Networks and Attribute GANs), and demonstrate the ability to generate semantically meaningful adversarial examples using each of these models. In addition to our empirical evaluations, we also provide a theoretical analysis of a simplified semantic attack model to understand the capacity of parametric attacks that typically exploit a significantly lower dimensional attack space compared to the classical pixel-space attacks.

Our specific contributions are as follows:

1. We propose a novel optimization based framework to generate semantically valid adversarial examples using parametric generative transformations.
2. We explore realizations of our approach using variants of multi-attribute transformation models: Fader Networks [30] and Attribute GANs [20] to generate adversarial face images for a binary classifier trained on the CelebA dataset [37]. Some of our modified multi-attribute models are non-trivial and may be of independent interest.
3. We present an empirical analysis of our approach and show that increasing the *dimensionality* of the attack space results in more effective attacks. In addition, we investigate a sequence of increasingly nonlinear attacks, and demonstrate that a higher degree of *nonlinearity* (surprisingly) leads to weaker attacks.
4. Finally, we provide a preliminary theoretical analysis by providing upper bounds for the *classification error* for a simplified surrogate model under adversarial condition [52].

This analysis supports our empirical observations regarding the dimensionality of the attack space.

We demonstrate the effectiveness of our attacks on simple deep classifiers trained over complex image datasets; hence, our empirical comparisons are significantly more realistic than popular attack methods such as FGSM [16] and PGD [29, 39] that primarily have focused on simpler datasets such as MNIST [32] and CIFAR. Our approach also presents an interesting use-case for multi-attribute generative models which have been used solely as visualization tools thus far.

Outline: We begin with a review of relevant literature in Section 2. We describe our proposed framework, *Semantic Adversarial Generation*, in section 3. In Section 4 we describe two variants of our framework to show different methods of ensuring the semantic constraint. We provide empirical analysis of our work in Section 5. We further present empirical analysis and theoretical qualification on the dimensionality of the parametric attack space in Section 6, and conclude with possible extensions in Section 7.

2. Related Work

Due to space constraints coupled with the large amount of recent progress in the area of adversarial machine learning, our discussion of related work is necessarily incomplete. We defer a more detailed discussion to the appendix.

Our focus is on *white box, test-time* attacks on deep classification systems; other families of attacks (such as backdoor attacks, data poisoning schemes, and black-box attacks) are not directly relevant to our setting, and we do not discuss those methods here.

Adversarial Attacks: Evidence that deep classifiers are susceptible to imperceptible adversarial examples can be attributed to Szegedy *et al.* [58]. Goodfellow *et al.* [16] and Kurakin *et al.* [29] extend this line of work using the Fast Gradient Sign Method (FGSM) and its iterative variants. Carlini and Wagner [5] devise state-of-the-art attacks under various pixel-space l_p norm-ball constraints by proposing multiple adversarial loss functions. Athalye *et al.* [1] further analyze several defense approaches against pixel-space adversarial attacks, and demonstrate that most existing defenses can be surpassed by approximating gradients over defensively trained models.

Such attacks perturb the pixel-space under an imperceptibility constraint. On the contrary, we approach the problem of generating adversarial examples that have perceptible yet semantically valid modifications. Our method considers a smaller ‘parametric’ space of modifiable attributes that have physical significance.

Parametric Adversarial Attacks: Parametric attacks are a recently introduced class of attacks in which the attack space is defined by a set of parameters rather than the pixel

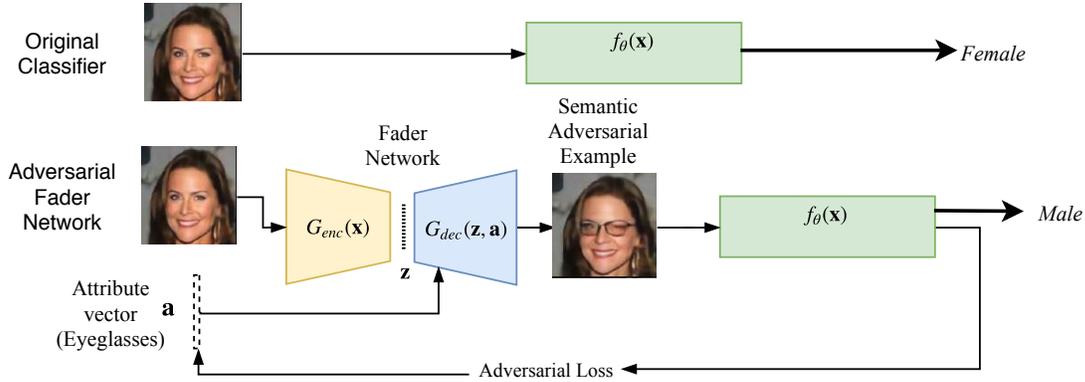


Figure 2. An single-attribute **Adversarial Fader Network**. The semantic adversarial attack framework optimizes an adversarial loss to generate an adversarial direction. Backpropagating the adversarial direction through the Fader Network with respect to the attribute vector, \mathbf{a} , ensures that the adversarial example is only generated for that specific attribute. Here, the adversarial algorithm generates eyeglasses on a face of a Female by optimizing \mathbf{a} , thus forcing the gender classifier to misclassify the image as Male.

space. Such approaches result in more “natural” adversarial examples as they target the image formation process instead of the pixel space. Recent works by Athalye *et al.* [2] and Liu *et al.* [35] use optimization over geometric surfaces in 3D space to create adversarial examples. Zhang *et al.* [71] demonstrate the existence of adversarially designed textures that can camouflage vehicles. Zhao *et al.* [72] generate adversarial examples by using the parametric input latent space of GANs[18]. Xiao *et al.* [65] employ spatial transforms to perturb image geometry for creating adversarial examples. Sharif *et al.* [55] propose a generative model to alter images of faces with eyeglasses in order to confound a face recognition classifier. Contrary to these methods, we consider the inverse approach of using a pre-trained multi-attribute generative model to transform inputs over multiple attributes for generating adversarial examples.

Song *et al.* [57] optimize over the latent space of a conditional GAN to generate unrestricted adversarial examples for a gender classifier. While our approach is thematically similar, we fundamentally differ in the context of being able to generate adversarial counterparts for given test samples while providing a finer degree of control using multi-attribute generative models. We discuss relevant literature regarding such multi-attribute generative models below.

Attribute-Based Conditional Generative Models: Generative Adversarial Networks (GAN) [18] are a popular approach for the generation of samples from a real-world data distribution. Recent advancements [49, 36, 64, 6] in GANs allow for creation of high quality realistic images. Chen *et al.* [6] introduce the concept of a attribute learning generative model where visual features are parametrized by an input vector.

Perarnau *et al.* [48] use a Conditional Generative Adversarial Network [40] and an encoder to learn the attribute

invariant latent representation for attribute editing. Fader Networks [30] improve upon this using an auto-encoder with a latent discriminator. He *et al.* [20] argue that such an attribute invariant constraint is too constrictive and replace it an attribute classification constraint and a reconstruction loss instead to alter only the desired attributes preserving attribute-excluding features. These models are primarily used for generation of a large variety of facial images. We provide a secondary (and perhaps practical) use case for such attribute models in the context of understanding generalization properties of neural networks.

3. Semantic Attacks

Conceptually, producing an adversarial semantic (“natural”) perturbation of a given input depends on two algorithmic components: (i) the ability to navigate the manifold of parametric transformations of an input image, and (ii) the ability to perform optimization over this manifold that maximizes the classification loss with respect to a given target model. We describe each component in detail below.

Notation: We assume a white-box threat model, where the adversary has access to a target model $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \{0, 1\}^c$ and the gradients associated with it. The model classifies an input image, \mathbf{x} into one of c classes, represented by a one-hot output label, \mathbf{y} . In this paper, we focus on binary classification models ($c = 2$) while noting that our framework transparently extends to multi-class models. Let $G(\mathbf{x}, \mathbf{a}) : \mathbb{R}^d \times \mathbb{R}^k \rightarrow \mathbb{R}^d$ denote parametric transformations, conditioned on a parameter vector, \mathbf{a} . Here, each element of \mathbf{a} (say, a_i) is a real number that corresponds to a specific semantic attribute. For example, a_0 may correspond to facial hair, with a value of zero (or negative) denoting absence and a positive value denoting presence of hair on a given face example. We define a semantic adversar-

Algorithm 1 Adversarial Parameter Optimization

Require: \mathbf{x}_0 : Input image, \mathbf{a}_0 : Initial attribute vector, $E(\cdot)$: Attribute encoder, $G(\cdot, \cdot)$: Pre-trained parametric transformation model, $f(\cdot)$: Target classifier, \mathbf{y} : Original label

- 1: $h_0 \leftarrow f(\mathbf{x}), l_{\text{adv}} \leftarrow \infty, i \leftarrow 0$
- 2: success = 0
- 3: **while do** $l_{\text{adv}} \neq 0$ and $i \leq \text{MaxIter}$
- 4: $\bar{\mathbf{a}} \leftarrow E(\mathbf{a})$
- 5: $h_i \leftarrow f(G(\mathbf{x}_i, \bar{\mathbf{a}}_i))$
- 6: $l_{\text{adv}} \leftarrow L_{\text{adv}}(\mathbf{y}, h_i)$
- 7: $\mathbf{a}_{i+1} \leftarrow \text{BackProp} \{ \mathbf{a}_i, \nabla l_{\text{adv}}(f(G(\mathbf{x}_i, E(\mathbf{a}_i)))) \}$
- 8: $\tilde{\mathbf{x}} \leftarrow G(\mathbf{x}, E(\mathbf{a}_{i+1}))$
- 9: **if** $f(\mathbf{x}) \neq f(\tilde{\mathbf{x}})$ **then**
- 10: return success, $\tilde{\mathbf{x}}$
- 11: **end if**
- 12: $i \leftarrow i + 1$
- 13: **end while**

ial attack as the deliberate process of transforming an input image, \mathbf{x} via a parametric model to produce a new example $\tilde{\mathbf{x}} = G(\mathbf{x}, \mathbf{a})$ such that $f(\tilde{\mathbf{x}}) \neq f(\mathbf{x})$.

3.1. Parametric Transformation Models

First, let us consider the problem of generating semantic transformations of a given input example. In order to create semantically transformed examples, the defined parametric generative model $G(\cdot)$ should satisfy two properties: $G(\cdot)$ should reconstruct the invariant data in an image, and $G(\cdot)$ should be able to independently perturb the semantic attributes while minimally changing the invariant data.

The parametric transformation model therefore, is trained to reconstruct the original example while disentangling the semantic attributes. This involves conditioning the generative model on a set of parameters corresponding to the modifiable attributes. The *semantic parameter vector* consists of these parameters and is input to the parametric model to control the expression of semantic attributes.

We argue that the range-space of such a model approximates the manifold of the semantic transformations of input images. Therefore, the transformation model can be used a projection operator to ensure that a solution to an optimization problem will lie in the set of semantic transformations of an input image. We also observe that the semantic parameter vectors will be much lower in dimension than the original image.

In this paper, we consider two variants of such conditional generative models: Fader Networks [30] and AttributeGANs (AttGAN) [20].

3.2. Adversarial Parameter Optimization

The problem of generating a semantic adversarial example essentially can be thought of as finding the right set of attributes that a classifier is adversarially susceptible to. In our approach, we model this as an optimization problem

over the semantic parameters.

The generation of adversarial examples is generally modelled as an optimization problem that can be broken down into two sub problems: (1) Optimization of an adversarial loss over the target network to find the direction of an adversarial perturbation. (2) Projection of the adversarial vector on the viable solution-space.

In the first step, we optimize over an adversarial loss, L_{adv} . We model the second step as a projection of the adversarial vector onto the range space of a parametric transformation model. This is achieved by cascading the output of the transformation function to the input of our target network. The optimization problem can then be solved by back-propagating over both the network and the transform. We also modify the Carlini-Wagner untargeted adversarial loss [5] as shown in equation 1 to include our semantic constraint:

$$\max \left(0, \max_{t \neq i} (f(\tilde{\mathbf{x}})_t) - f(\tilde{\mathbf{x}})_i \right) \quad (1)$$

s.t. $\tilde{\mathbf{x}} = G(\mathbf{x}, \mathbf{a})$

where i is the original label index and t are the class label indices for any of the other classes.

In comparison to the grid search method presented in Zhao *et al.* [72] and Engstrom *et al.* [12], our optimization algorithm scales better. In addition, we create semantic adversarial transformations with multiple attributes for a specific input allowing for a fine-grained analysis of the generalization capacities of the target model.

4. Semantic Transformations

While our semantic attack framework is applicable to any parametric transformation model that enables gradient computations, we instantiate it by constructing adversarial variants of two recently proposed generative models: Fader networks [30] and AttributeGANs (AttGAN) [20].

4.1. Adversarial Fader Network

A Fader Network [30] is an encoder-decoder architecture trained to modify images with continuously parameterized attributes. They achieve this by learning an invariance over the encoded latent representation while disentangling the semantic information of the images and attributes. The invariance of the attributes is learnt by an adversarial training step in the latent space with the help of a latent discriminator which is trained to identify correct attributes corresponding to each training sample.

Using our framework, we can adapt any pre-trained Fader Network to model the manifold of semantic perturbations of a given input. We note that minor adjustments are needed in our setting, since the parameter vector required by the approach of [30] requires each scalar attribute, a_i ,

to be represented by a tuple, $(1 - a_i, a_i)$. Since there is a one-to-one mapping between the two representations, we can project any real-valued parameter vector \mathbf{a} into this tuple form via an additional, fixed affine transformation layer. Given this extra “attribute encoding” step, all gradient computations proceed as before. We quantitatively study the effect of allowing the attacker access to single or multiple semantic attributes. In particular, we construct three approaches for generating semantic adversarial examples: (i) A single attribute Fader Network; (ii) A multi-attribute Fader Network; and (iii) A cascaded sequence of single attribute Fader Networks.

Single Attribute Attack: For the single attribute attack, we use the range-space of a pre-trained, single attribute Fader Network to constrain our adversarial attack. The single attribute attack constrains an attacker to only modify a specified attribute for all the images. In the case of face images, such attributes might include presence/absence of eyeglasses, hair color, and nose shape.

In our experiments, we present examples of attacks on a gender classifier using three separate single attributes: eyeglasses, age, and skin complexion. Fig. 2 describes the mechanism of a single-attribute adversarial Fader Network that generates an adversarial example by adding eyeglasses.

Multi-Attribute Attack: Similar to the single-attribute case, we may also use pre-trained multi-attribute Fader Networks to model cases where the adversary has access to multiple modifiable traits.

A limitation of multi-attribute Fader Networks lies in the difficulty of their training. This is because a Fader Network is required to learn disentangled representations of the attributes while in practice, semantic attributes cannot be perfectly decoupled. We resolve this using a novel conditional generative model described as follows.

Cascaded Attribute Attack: We propose a novel method to simulate multi-attribute attacks by stage-wise concatenation pre-trained single attribute Fader networks. The benefit is that the computational burden of learning disentangled representations is now removed.

Each single-attribute model exposes a attribute latent vector. During execution of Alg. 1 we jointly optimize over all the attribute vectors. The optimal adversarial vector is then segmented into corresponding attributes for each Fader Network to generate an adversarial example.

4.2. Adversarial AttGAN

A second encoder-decoder architecture [20], known as AttGAN, achieves a similar goal as Fader Networks of editing attributes by manipulating the encoded latent representation; however, AttGAN disentangles the semantic attributes from the underlying invariances of the data by considering both the original and the flipped labels while training. This is achieved by training a latent discriminator and

Attack Type	Attributes	Accuracy of target (%)	of model Random Sampling (%)
Single Attribute Attack	A1	52.0	89.00
	A2	35.0	96.00
	A3	14.0	90.00
Multi Attribute Attack	A1,A5,A6	3.00	89.00
	A2,A5,A6	1.00	81.00
	A1,A2,A7	3.00	87.00
Cascaded Multi Attribute Attack	A1-A2-A3	18.00	55.6
	A2-A3-A4	20.00	93.00
Multi Attribute AttGAN Attack	A1,A2,A6,A8,A10	70.40	32.80
	A1,A2,A6,A8,A9,A10	39.40	40.40

Table 1. Performance of the Semantic Adversarial Example under multiple implementations. Legend for attributes: A1-Eyeglasses, A2-Age, A3-Nose shape, A4-Eye shape, A5-Chubbiness, A6-Pale Skin, A7-Smiling, A8-Mustache, A9-Eyebrows, A10-Hair color. As the number of attributes increase, semantic attacks are more effective. Our optimization-based attack fares better as compared to worst-of-10 random sampling [12], showing the former’s efficacy at finding semantic adversarial examples.

classifier pair to classify both the original and the transformed image to ensure invariance.

In order to generate semantic adversarial examples using AttGAN, we use a pretrained generator conditioned on 13 attributes. The attribute vector in this case, is encoded to be a perturbation of the original sequence of attributes for the image. We consider the two sets of attributes listed in Table 1 to generate adversarial examples. In our experience, the AttGAN architecture provides a more stable reconstruction, thus allowing for more modifiable parameters.

5. Experimental Results

We showcase our semantic adversarial attack framework using a binary (gender) classifier as the target model trained on the CelebA dataset [37]. While we restrict ourselves to results on binary classifiers on faces in this paper, additional results with multi-class classifiers on the Berkeley Deep Drive dataset [69] can be found in the appendix (refer Fig. 8). All experiments were performed on a single workstation equipped with an NVidia Titan Xp GPU in PyTorch [47] v1.0.0.¹ We train the classifier using the ADAM optimizer [26] over the categorical cross-entropy loss. The training data is augmented with random horizontal flipping to ensure that the classifier does not overfit. The target model achieves a (standard) accuracy of 99.7% on the test set (10% of the dataset).

Our goal is to break this classifier model using semantic attacks. To do so, we use a subset of 500 randomly selected images from the test set. Each image is transformed by our algorithm using the various parametric transformation families described in Section 4. Our metric of comparison for all adversarial attacks is the target model accuracy on the

¹Code and models: https://github.com/ameya005/Semantic_Adversarial_Attacks

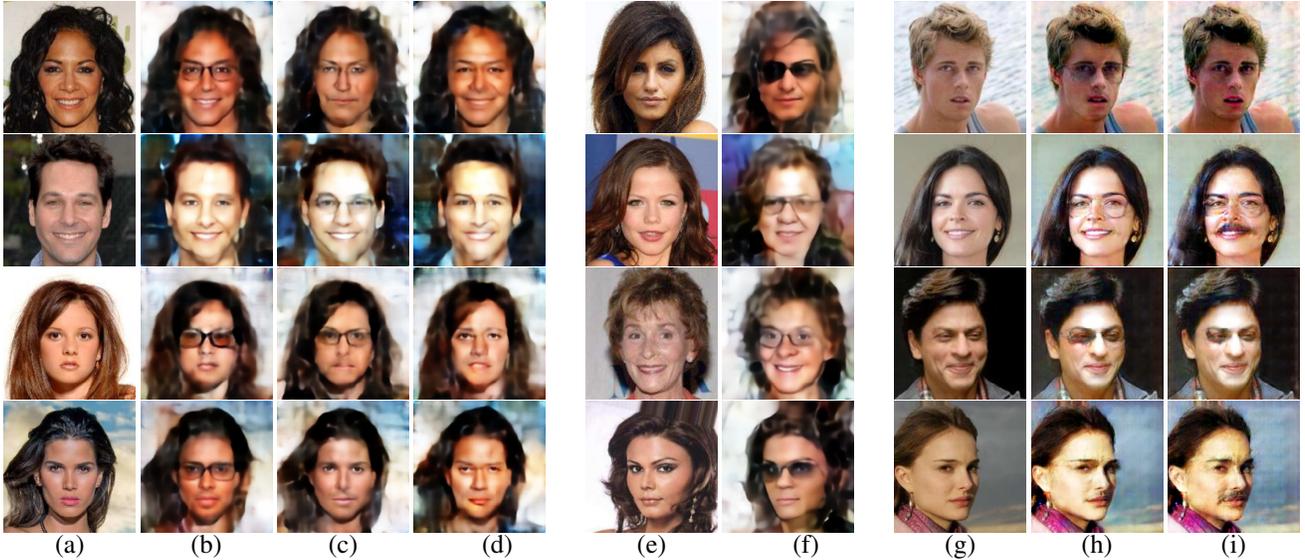


Figure 3. Semantic adversarial examples generated with multiple attribute semantic models as in table 1. Columns (a), (e) and (g) are original images. Columns (b){Attribute category: A1,A5,A6} (c){Attribute category: A1,A2,A7} and (d){Attribute category: A2,A5,A6} show examples generated using multi-attribute Fader Networks as semantic transforms. Examples in (f){Attribute category: A1-A2-A3} were generated using cascaded single attribute Fader Network. Columns (h){Attribute category: A1,A2,A6,A8,A10} and (i){Attribute category: A1,A2,A6,A8,A9,A10} are images transformed using an AttGAN with 5 and 6 attributes respectively. Additional results of semantic attacks on multi-class classifiers for traffic scenes are provided in the appendix.

generated adversarial test set.

Adversarial Fader Networks: We consider the three approaches documented in section 4.1. For every image in our original test set, we generate adversarial examples by optimizing the adversarial loss in equation 1 with respect to the corresponding attribute parameters.

In the cases of single-attribute and cascaded sequential attacks, we use the pre-trained single-attribute models provided by Lample *et al.* [30] to represent the manifold of semantic transformations. For the multi-attribute attack, we train 3 multi-attribute Fader Networks with the attributes presented in Table 1. We create an adversarial test set for each our approaches as described in Section 4.1 using our algorithm as defined in Algorithm 1.

Our experiments show that Adversarial Fader Networks successfully generate examples that confound the binary classifier in all cases; see Table 1. Visual adversarial examples are displayed in Fig. 1 and Fig. 3. We also observe that multi-attribute attacks outperform single-attribute attacks, which conforms with intuition; a more systematic analysis of the effect of the number of semantic attributes on attack performance is provided below in Section 6.

Adversarial AttGAN: We perform a similar set of experiments using the multi-attribute AttGAN implementation of He *et al.* [20]. We record the performance over two experiments: one using 5 attributes, and the second using 6 attributes, as seen in Table 1. We observe a significant improvement in performance as the number of semantic at-

tributes increases (in particular, adding the eyebrows attribute results in nearly a 30% drop in model accuracy).

Comparison with parameter-space sampling: We compare our method with a previously-proposed approach that investigates parametric attacks *et al.* [12]. They propose picking s random samples from the parameter space and choose the adversarial example generated by the sample giving the worst cross entropy loss (we use $s = 10$).

We showcase the results in Table 1, and observe that in all cases (but one), our semantic adversarial attack algorithm outperforms random sampling. In addition, the table also reveals that random examples in the range of Fader Networks or AttGANs are mostly classified correctly. This suggests that the target model is generally invariant to the low reconstruction error incurred by the parametric transformation models².

Comparison with pixel-space attacks: In addition to our analyses described above, we also compare our attacks with the state-of-the-art Carlini-Wagner(CW) l_∞ -attack [5] as well as several other attack techniques [16, 29, 12] in Table 2. To ensure fair comparison, we consider the maximum l_∞ distance over our multi-attribute attacks as the bound parameter ϵ for all pixel-norm based attacks. From the table, we observe that the CW attack is extremely effective; on

²We do not compare our work with other approaches such as the Differentiable Renderer [35] and 3D adversarial attacks [70], since these papers expect oracle access to a 3D rendering environment. We also do not compare with Song *et al.* [57] since they generate adversarial examples from scratch, whereas our attack targets specific inputs.

Attack ($\epsilon = 1.74$)	Accuracy(%)
Single Att. Semantic Attack	14.01
Multi Att. Semantic Attack	1.00
FGSM [16]	91.6
PGD [39, 29]	26.2
CW- l_∞ [5]	0.00
Spatial [12]	41.00

Table 2. Comparison of adversarial attacks with other attack strategies. A lower target accuracy corresponds to a better attack. The pixel space attacks are allowed to generate adversarial examples under the l_∞ distance corresponding to our best performing multi-attribute attack model. Observe that semantic attacks are comparable to the state of the art pixel-space attack.

the other hand, our semantic attacks are able to outperform other methods such as FGSM [17] and PGD [39].

We also compare our approach to *Spatial Attacks* of [12], which uses a grid search over affine transformations of an input to generate adversarial examples; l_∞ constraints do not apply here, and instead we use default parameters provided in [12]. Our proposed attack methods are considerably more successful. We provide additional detailed experiments on binary and multi-class classifiers for other attributes as well as other datasets in the appendix.

6. Analysis: Impact of Dimensionality

From our experiments, we observe that limiting the adversary to a low-dimensional, semantic parametric transformation of the input leads to less-effective attacks than pixel-space attacks (at least when the same loss is optimized). Moreover, single-attribute semantic attacks are more powerful than multi-attribute attacks. This observation makes intuitive sense: the dimension of the manifold of perturbed inputs effectively represents the *capacity* of the adversary, and hence a greater number of degrees of freedom in the perturbation should result in more effective attacks. In pixel-space attacks, the adversary is free to search over a high-dimensional l_p -ball centered around the input example, which is perhaps why l_p -norm attacks are so hard to defend against [1].

In this section, we provide experimental and theoretical analysis that precisely exposes the impact of the dimensionality of the attribute parameters. While our analysis is stylized and not directly applicable to deep neural classifiers, it constitutes a systematic first attempt towards upper bounds on what a semantically constrained adversary can possibly hope to achieve.

6.1. Synthetic Experiments

We propose and analyze the following synthetic setup which enables explicit control over the dimension of the semantic perturbations.

Data: We construct a dataset of $n = 500$ samples from a mixture of Gaussians (MoG) with 10 components (denoted

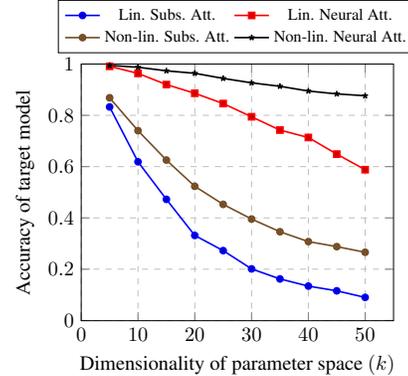


Figure 4. Effect of dimensionality of the parametric attack space. Considering subspace and rank constrained transforms to generate adversarial examples, note that the target model accuracy decreases as the dimensionality of the attack space increases. The additive attack (surrogate to PGD) is more effective than multiplicative attack (similar to our approach) over all values of k .

by \mathbb{P}_d) defined over $(\mathbf{x}, y) \in \mathbb{R}^d \times \{\pm 1\}$. Each data sample is obtained by uniformly sampling one of the mixture component means, and then adding random Gaussian noise with standard deviation $\sigma \leq \sqrt{d}$. The component means are chosen as 10 randomly selected images (1 for each digit) from the MNIST dataset [32] rescaled to 10×10 (i.e., the ambient dimension is $d = 100$).

Target Model: We artificially define two classes: the first class containing images generated from digits 0-4 and the second class containing images from samples 5-9. We train a simple two-layer fully connected network, $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \{\pm 1\}$ as the target model. The classifier is trained by optimizing cross-entropy using ADAM [26] for 50 epochs, resulting in training accuracy of 100%, validation accuracy of 99.8%, and test accuracy of 99.6%.

Parametric Transformations: We consider a stylized transformation function, $G(\mathbf{x}, \delta) : \mathbb{R}^d \times \mathbb{R}^k \rightarrow \mathbb{R}^d$. We study the effect of varying k for two specific parametric transformation models.

Subspace attacks: We first consider an additive (linear) attack model. Here, the manifold of semantic perturbations is constrained to lie a k -dimensional subspace spanned by an arbitrary matrix $\mathbf{U} \in \mathbb{R}^{d \times k}$, whose columns are assumed to be orthonormal, and $\delta \in \mathbb{R}^k$

$$G(\mathbf{x}, \delta) := \tilde{\mathbf{x}} = \mathbf{x} + \mathbf{U}\mathbf{U}^T \delta \quad (2)$$

Neural attacks: Next, we consider a *multiplicative* attack model. Here the manifold of perturbations corresponds to a rank- k transformation of the input.

$$G(\mathbf{x}, \delta) := \tilde{\mathbf{x}} = \mathbf{U} \cdot \text{diag}(\delta) \cdot \mathbf{U}^T \mathbf{x} \quad (3)$$

Here, \mathbf{U} and δ follow the definition presented earlier. This transformation can be interpreted as the action of a shallow (two-layer) auto-encoder network with k hidden neurons with scalar activations parameterized by δ .



Figure 5. Semantically transformed single-attribute examples which are classified correctly by the target model but show severe artifacts. This shows that neural networks are immune to significant changes in the semantic domain unlike the pixel domain.

Nonlinear ReLU variants: We also consider each of the above two attacks in the *rectified* setting where the transformation is passed through a rectified linear unit: $\tilde{\mathbf{x}} = \text{ReLU}(G(\mathbf{x}, \delta))$.

Results: We analyse the effect of the dimensionality of the attack space (k) by considering the performance of the subspace and neural attacks on the target binary classifier. Fig. 4 shows the comparison of the constrained attacks for the linear and non-linear cases.

We infer the following: (i) as expected, increasing dimensionality of the semantic attack space leads to less accurate target models; (ii) adding a non-linearity to the transformation function *reduces* the viability of both subspace- and rank-constrained attacks; (iii) subspace-constrained attacks are more powerful than neural attacks. In general, the degree of “nonlinearity” in the transformation model appears to be inversely proportional to the power of the corresponding semantic attack. We believe this phenomenon is somewhat surprising, and defer further analysis to future work.

6.2. Theory

In the case of subspace attacks, we explicitly derive upper bounds on the generalization behavior of target models. Our derivation follows the recent approach of Schmidt *et al.* [52], who consider a simplified version of the data model defined in Section 6.1 and bound the performance of a linear classifier in terms of its *robust classification error*.

Def. 6.1 (Robust Classification Error). Let $\mathbb{P}_d : \mathbb{R}^d \times \{\pm 1\} \rightarrow \mathbb{R}$ be a distribution and let \mathcal{S} be any set containing \mathbf{x} . Then the \mathcal{S} -robust classification error of any classifier $f : \mathbb{R}^d \rightarrow \{\pm 1\}$ is defined as $\beta = P_{(\mathbf{x}, y) \sim \mathbb{P}_d} [\exists \tilde{\mathbf{x}} \in \mathcal{S} : f(\tilde{\mathbf{x}}) \neq y]$.

Using this definition, we analyze the efficacy of subspace attacks on a simplified linear classifier trained using a mixture of two spherical Gaussians. Consider a dataset with samples $(\mathbf{x}, y) \in \mathbb{R}^d \times \{\pm 1\}$ sampled from a mixture of two Gaussians with component means $\pm \theta^*$ and standard deviation $\sigma \leq \sqrt{d}$. We assume a linear classifier $f_{\hat{\mathbf{w}}}$, defined by the unit vector $\hat{\mathbf{w}}$, as $f_{\hat{\mathbf{w}}}(\mathbf{x}) = \text{sign}(\langle \hat{\mathbf{w}}, \mathbf{x} \rangle)$.

Let $\mathcal{S}_\epsilon = \{\tilde{\mathbf{x}} \mid \tilde{\mathbf{x}} = \mathbf{x} + \mathbf{U}\mathbf{U}^T\delta, \|\mathbf{x} - \tilde{\mathbf{x}}\|_\infty \leq \epsilon\}$. Assuming that the target classifier is well-trained (i.e., $\hat{\mathbf{w}}$ is sufficiently well-correlated with the true component mean θ^*), we can upper bound the probability of error incurred by the classifier when subjected to any subspace attack.

Theorem 1 (Robust classification error for subspace attacks). *Let $\hat{\mathbf{w}}$ be such that $\langle \hat{\mathbf{w}}, \theta^* \rangle \geq k\|\mathbf{U}\|_{\infty, 1}\|\hat{\mathbf{w}}^T\mathbf{U}\|_\infty\epsilon$. Then, the linear classifier $f_{\hat{\mathbf{w}}}$ has a \mathcal{S}_ϵ -robust classification error upper bounded as:*

$$\beta \leq \exp\left(-\frac{(\langle \hat{\mathbf{w}}, \theta^* \rangle - k\|\mathbf{U}\|_{\infty, 1}\|\hat{\mathbf{w}}^T\mathbf{U}\|_\infty\epsilon)^2}{2\sigma^2}\right) \quad (4)$$

The proof is deferred to the appendix, but we provide some intuition. Lemma 20 of [52] recovers a similar result, albeit with the \sqrt{k} term in the exponent being replaced by \sqrt{d} . This is because they only consider bounded ℓ_∞ -perturbations in pixel-space, and hence their bound on the robust classification error scales exponentially according to the *ambient dimension* d , while our bound is expressed in terms of the *number of semantic attributes* $k \ll d$. A natural next step would be to derive *sample complexity* bounds analogous to [52] but we do not pursue that direction here.

7. Discussion and Conclusions

We conclude with possible obstacles facing our approach and directions for future work. We have provided evidence that there exist adversarial examples for a deep neural classifier that may be perceptible, yet are semantically meaningful and hence difficult to detect. A key obstacle is that parameters associated with semantic attributes are often difficult to decouple. This poses a practical challenge, as it is difficult to train a conditional generative model with independent latent semantic dimensions. However, the success of recent efforts in this direction, including Fader Networks [30], AttGans [20], and StarGANs [8] demonstrate promise of our approach: any newly developed conditional generative models can be used to mount a semantic attack using our framework.

Despite the existence of semantic adversarial examples, we have found that enforcing semantic validity confounds the adversary’s task, and that target models are generally able to classify a significant subset of the examples generated under our semantic constraint. Fig. 5 are examples of images generated with severe artifacts, yet that are successfully classified. This presents the question: is “naturalness” a strong defense?

This intuition is the premise of a recent defense strategy called DefenseGAN [51]. Indeed, our approach can be viewed as converse of this strategy: DefenseGAN uses the range-space of a generative model (specifically, a GAN) to *defend* against pixel-space attacks, while conversely, we use the same principle to *attack* trained target models. A closer look into the interplay between the two approaches is worthy of future study.

References

- [1] Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, 2018. 1, 2, 7
- [2] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *ICML*, 2018. 3
- [3] Stephen Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, 2004.
- [4] Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017.
- [5] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *2017 IEEE Symposium on Security and Privacy (SP)*, 2017. 1, 2, 4, 6, 7
- [6] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NeurIPS*, 2016. 3
- [7] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning. *arXiv preprint, abs/1712.05526*, 2017.
- [8] Yunje Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 8
- [9] Ali Dabouei, Sobhan Soleymani, Jeremy M. Dawson, and Nasser M. Nasrabadi. Fast geometrically-perturbed adversarial faces. *WACV*, 2019.
- [10] Sumanth Dathathri, Stephan Zheng, Sicun Gao, and RM Murray. Measuring the Robustness of Neural Networks via Minimal Adversarial Examples. In *NeurIPS-W*, volume 35, 2017. 1
- [11] Roberto Rey de Castro and Herschel A Rabitz. Targeted non-linear adversarial perturbations in images and videos. *arXiv preprint, abs/1809.00958*, 2018.
- [12] Logan Engstrom, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. A rotation and a translation suffice: Fooling cnns with simple transformations. *arXiv preprint, abs/1712.02779*, 2017. 4, 5, 6, 7
- [13] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Xiaodong Song. Robust physical-world attacks on deep learning visual classification. *CVPR*, 2018.
- [14] Alhussein Fawzi, Hamza Fawzi, and Omar Fawzi. Adversarial vulnerability for any classifier. In *NeurIPS*, 2018.
- [15] Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Analysis of classifiers’ robustness to adversarial perturbations. *Machine Learning*, 107, 2018.
- [16] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015. 1, 2, 6, 7
- [17] Ian J. Goodfellow. Defense against the dark arts: An overview of adversarial example security research and future research directions. *arXiv preprint, abs/1806.04169*, 2018. 1, 7
- [18] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 3
- [19] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Bad-nets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint, abs/1708.06733*, 2017.
- [20] Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. Attgan: Facial attribute editing by only changing what you want. *arXiv preprint*, 2017. 2, 3, 4, 5, 6, 8
- [21] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Lecture 6a, overview of mini-batch gradient descent.
- [22] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In *PMLR*, volume 80, 2018. 1
- [23] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. Prior convictions: Black-box adversarial attacks with bandits and priors. *arXiv preprint, abs/1807.07978*, 2018. 1
- [24] Takuhiro Kaneko, Kaoru Hiramatsu, and Kunio Kashino. Generative attribute controller with conditional filtered generative adversarial networks. *CVPR*, 2017.
- [25] Taeksoo Kim, Byoungjip Kim, Moonsoo Cha, and Jiwon Kim. Unsupervised visual attribute transfer with reconfigurable generative adversarial networks. *arXiv preprint, abs/1707.09798*, 2017.
- [26] Diederik Kingma and Jimmy Ba. Adam: a method for stochastic optimization (2014). In *ICLR*, 2015. 5, 7
- [27] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint, abs/1312.6114*, 2014.
- [28] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *JMLR*, volume 70, 2017.
- [29] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint, abs/1607.02533*, 2017. 2, 6, 7
- [30] Guillaume Lample, Neil Zeghidour, Nicolas Usunier, Antoine Bordes, Ludovic Denoyer, et al. Fader networks: Manipulating images by sliding attributes. In *NeurIPS*, 2017. 2, 3, 4, 6, 8
- [31] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *ICML*, 2016.
- [32] Yann LeCun and Corinna Cortes. MNIST handwritten digit database, 2010. 2, 7
- [33] Mu Li, Wangmeng Zuo, and David Zhang. Convolutional network for attribute-driven and identity-preserving human face generation. *arXiv preprint, abs/1608.06434*, 2016.
- [34] Mu Li, Wangmeng Zuo, and David Zhang. Deep identity-aware transfer of facial attributes. *arXiv preprint, abs/1610.05586*, 2016.
- [35] Hsueh-Ti Derek Liu, Michael Tao, Chun-Liang Li, Derek Nowrouzezahrai, and Alec Jacobson. Beyond pixel norm-balls: Parametric adversaries using an analytically differentiable renderer. In *ICLR*, 2019. 3, 6

- [36] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *NeurIPS*, 2017. 3
- [37] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, 2015. 2, 5
- [38] Yongyi Lu, Yu-Wing Tai, and Chi-Keung Tang. Attribute-guided face generation using conditional cycleGAN. In *ECCV*, 2018.
- [39] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018. 2, 7
- [40] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint*, abs/1411.1784, 2014. 3
- [41] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. *CVPR*, 2017. 1
- [42] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. DeepFool: A simple and accurate method to fool deep neural networks. *CVPR*, 2016.
- [43] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Jonathan Uesato, and Pascal Frossard. Robustness via curvature regularization, and vice versa. In *CVPR*, 2019.
- [44] Konda Reddy Mopuri, Utkarsh Ojha, Utsav Garg, and R. Venkatesh Babu. Nag: Network for adversary generation. *CVPR*, 2018.
- [45] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier GANs. In *ICML*, 2017.
- [46] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. *EuroS&P*, 2016.
- [47] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NeurIPS-W*, 2017. 5
- [48] Guim Perarnau, Joost van de Weijer, Bogdan Raducanu, and Jose M. Álvarez. Invertible conditional GANs for image editing. *arXiv preprint*, abs/1611.06355, 2016. 3
- [49] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint*, abs/1511.06434, 2016. 3
- [50] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models. *arXiv preprint arXiv:1707.04131*, 2017.
- [51] Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *ICLR*, 2018. 8
- [52] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In *NeurIPS*, 2018. 2, 8
- [53] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *NeurIPS*, 2018.
- [54] Ali Shafahi, W. Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. Are adversarial examples inevitable? In *ICLR*, 2019. 1
- [55] Mahmood Sharif, Sruti Bhagavatula, Lujio Bauer, and Michael K. Reiter. Adversarial generative nets: Neural network attacks on state-of-the-art face recognition. *arXiv preprint*, abs/1801.00349, 2018. 3
- [56] Wei Shen and Rujie Liu. Learning residual images for face attribute manipulation. *CVPR*, 2017.
- [57] Yang Song, Rui Shu, Nate Kushman, and Stefano Ermon. Constructing unrestricted adversarial examples with generative models. In *NeurIPS*, 2018. 3, 6
- [58] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 1, 2
- [59] O. Tange. Gnu parallel - the command-line power tool. *login: The USENIX Magazine*, 36(1):42–47, Feb. 2011.
- [60] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. *arXiv preprint*, abs/1705.07204, 2017.
- [61] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In *NeurIPS*, 2018.
- [62] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Clean-label backdoor attacks, 2019.
- [63] Paul Upchurch, Jacob R. Gardner, Geoff Pleiss, Robert Pless, Noah Snaveley, Kavita Bala, and Kilian Q. Weinberger. Deep feature interpolation for image content changes. *CVPR*, 2017.
- [64] Xiaolong Wang and Abhinav Gupta. Generative image modeling using style and structure adversarial networks. In *ECCV*, 2016. 3
- [65] Chaowei Xiao, Jun-Yan Zhu, Bo Li, Warren He, Mingyan Liu, and Dawn Xiaodong Song. Spatially transformed adversarial examples. *arXiv preprint*, abs/1801.02612, 2018. 3
- [66] Huang Xiao, Battista Biggio, Blaine Nelson, Han Xiao, Claudia M. Eckert, and Fabio Roli. Support vector machines under adversarial label contamination. *Neurocomputing*, 160, 2015.
- [67] Han Xiao, Huang Xiao, and Claudia M. Eckert. Adversarial label flips attack on support vector machines. In *ECAI*, 2012.
- [68] Taihong Xiao, Jiapeng Hong, and Jinwen Ma. Dna-gan: Learning disentangled representations from multi-attribute images. *arXiv preprint*, abs/1711.05415, 2018.
- [69] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving video database with scalable annotation tooling. *arXiv preprint arXiv:1805.04687*, 2018. 5
- [70] Xiaohui Zeng, Chenxi Liu, Yu-Siang Wang, Weichao Qiu, Lingxi Xie, Yu-Wing Tai, Chi-Keung Tang, and Alan Loddon Yuille. Adversarial attacks beyond the image space. *arXiv preprint*, abs/1711.07183, 2017. 6
- [71] Yang Zhang, Hassan Foroosh, Philip David, and Boqing Gong. Camou: Learning physical vehicle camouflages to adversarially attack detectors in the wild. In *ICLR*, 2019. 3

- [72] Zhengli Zhao, Dheeru Dua, and Sameer Singh. Generating natural adversarial examples. In *ICLR*, 2018. 3, 4
- [73] Shuchang Zhou, Taihong Xiao, Yi Yang, Dieqiao Feng, Qinyao He, and Weiran He. Genegan: Learning object transfiguration and attribute subspace from unpaired data. *arXiv preprint*, abs/1705.04932, 2017.
- [74] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *ICCV*, 2017.