

## Machine learning notes

# Contents

Notation . . . . .	3
<b>I Notes</b>	<b>9</b>
<b>1 Interplay of adversarial robustness and generalization</b>	<b>10</b>
1.1 Introduction . . . . .	10
1.2 Definitions and notation . . . . .	11
1.3 Adversarial example definitions . . . . .	12
1.4 The manifold hypothesis . . . . .	13
1.5 Properties of adversarial examples . . . . .	14
1.6 Finding adversarial examples . . . . .	18
1.6.1 Attack objectives . . . . .	18
1.6.2 Common attacks . . . . .	19
1.7 Improving adversarial robustness . . . . .	20
1.7.1 Adversarial training and empirical adversarial risk . . . . .	21
1.8 Adversarial robustness and generalization . . . . .	21
1.8.1 A trade-off between robustness and generalization . . . . .	22
1.8.2 Non-robust features . . . . .	22
1.8.3 Training with on-manifold adversarial examples . . . . .	23
1.9 Conclusion . . . . .	24
<b>2 Kompromis između generalizacije i robusnosti kod dubokih modela za razumijevanje slika</b>	<b>25</b>
2.1 Uvod . . . . .	25
2.2 Osnovni pojmovi . . . . .	25

<i>CONTENTS</i>	2
2.2.1 Osnovne oznake . . . . .	25
2.2.2 Nadzirano učenje . . . . .	26
2.2.3 Duboko učenje i konvolucijski modeli . . . . .	26
2.3 Neprijateljska robusnost i neprijateljski primjeri . . . . .	28
2.4 Osnovne definicije i oznake u vezi neprijateljskih primjera . . . . .	28
2.5 Svojstva neprijateljskih primjera kod dubokih modela . . . . .	28
2.6 Pronalaženje neprijateljskih primjera . . . . .	28
2.7 Robusnost i generalizacija . . . . .	28
2.7.1 Postizanje robusnosti . . . . .	28
2.7.2 Kompromis između robusnosti i generalizacije . . . . .	28
2.8 Eksperimenti . . . . .	28
2.8.1 Postizanje robusnosti . . . . .	28
2.8.2 Smanjivanje kompromisa između robusnosti i generalizacije kod klasifikacije slika . . . . .	28
2.8.3 Robusno učenje modela za semantičku segmentaciju . . . . .	28
2.9 Zaključak . . . . .	28
<b>3 Pieces of text</b>	<b>29</b>
3.0.1 Robustness evaluation . . . . .	30

## Notation

### Objects

Variables are generally denoted by italic serif letters. Most constants are denoted by upright serif letters. Random variables are underlined. Vectors and sequences are denoted by lowercase bold letters. Matrices and multidimensional-arrays are denoted by uppercase bold letter. Sets are denoted by uppercase blackboard-bold letters. Latin or Greek letters can be used for any type of object.

$a, A, \theta$	Variable (commonly scalar or function)
$\mathbf{a}, \boldsymbol{\theta}$	Vector or sequence (commonly column vector)
$\mathbf{A}, \boldsymbol{\Theta}$	Matrix or multidimensional array
$\mathbb{A}$	Set or multiset
$a, A, \theta$	Constant
$\mathbf{a}, \boldsymbol{\theta}$	Vector or sequence constant
$\mathbf{A}, \boldsymbol{\Theta}$	Matrix or multidimensional array constant
$\mathbb{A}$	Set constant
$\underline{a}, \underline{A}, \underline{\theta}$	Random variable
$\underline{\mathbf{a}}, \underline{\boldsymbol{\theta}}$	Random vector or sequence
$\underline{\mathbf{A}}, \underline{\boldsymbol{\Theta}}$	Random matrix or multidimensional array
$\underline{\mathbb{A}}$	Random set or multiset
$a, \text{riječ}$	Textual label not representing an object

### Constants

$\{\}$	Empty set
$e$	The number that satisfies $\frac{d}{dx}e^x = e^x$
$\mathbf{0}$	Null-vector
$\mathbf{e}_i$	$i$ -th canonical basis vector
$\mathbf{1}$	The sum of all canonical basis vectors
$\mathbf{I}, \mathbf{I}_n$	Identity matrix (with $n$ rows/columns)
$\mathbb{N}, \mathbb{Z}, \mathbb{R}, \mathbb{C}$	A standard set

$\mathbb{R}_{\geq 0}, \mathbb{R}_{> 0}$  The set of all non-negative/positive real numbers

### Defining sets and arrays

$a..b$	Shorthand notation for $a, \dots, b$
$\{a..b\}$	A subset of integers from $a$ to $b$
$\{f(a): P(a)\}, \{f(a)\}_{P(a)}$	A set with elements defined by a function $f$ and a predicate $P$
$\{f(a)\}_a$	A set with elements defined by a function $f$ and variables $a$ from an implicitly defined set
$\{a_1..a_n\}, \{a_i\}_{i=1..n}$	A set with $n$ elements
$[x_1, \dots, x_n]$	A row vector
$[a_i]_i, [a_{i,j}]_{i,j}, [a_{i,j,k}]_{i,j,k}$	A multidimensional array with an implicit or undefined number of elements
$[a, b)$	A semi-closed interval

### Subscript and superscript

U donjem i gornjem indeksu oznake mogu biti oznake drugih matematičkih objekata ili slova ili riječi koje ne predstavljaju matematičke objekte. Redni brojevi elemenata vektora ili višedimenzionalnih nizova se, ako nije određeno drugačije, pišu u donjem indeksu oznake vektora u uglatim zagradama. Npr.  $i$ -ti element vektora  $\mathbf{a} = [a_1, \dots, a_n]^T$  je  $\mathbf{a}_{[i]} = a_i$ . Indeksi kod  $n$ -dimenzionalnih nizova mogu biti i vektori iz  $\mathbb{N}^n$ , ili kombinacije vektora manje dimenzije sa skalarima.

$a_d^g$	Varijabla s oznakama u donjem i gornjem indeksu
$\mathbf{a}_{[i]}$	$i$ -ti element vektora $\mathbf{a}$
$\mathbf{a}_{[i_1:i_2]}$	Vektor kojeg čine elementi $\mathbf{a}_{[i_1]}, \mathbf{a}_{[i_1+1]}, \dots, \mathbf{a}_{[i_2]}$
$\mathbf{a}_{[(i_1..i_n)]}$	Vektor kojeg čine elementi $\mathbf{a}_{[i_1]}, \mathbf{a}_{[i_2]}, \dots, \mathbf{a}_{[i_n]}$
$\mathbf{A}_{[i,j]}$	Element $i, j$ matrice $\mathbf{A}$
$\mathbf{A}_{[i,:]}$	$i$ -ti redak matrice $\mathbf{A}$
$\mathbf{A}_{[:,i_1:i_2,j]}$	2-D odsječak 3-D niza $\mathbf{A}$
$\mathbf{A}_{[i]}$	Element $\mathbf{A}_{[i_{[1]}, \dots, i_{[n]}]}$ $n$ -D niza

$\mathbf{A}_{[i_1:i_2]}$	Podniz $\mathbf{A}_{[i_{1[1]}:i_{2[1]}, \dots, i_{1[n]}:i_{2[n]}]}$ $n$ -D niza
$\mathbf{A}_{[i_1:i_2;:]}$	Podniz $\mathbf{A}_{[i_{1[1]}:i_{2[1]}, \dots, i_{1[n-1]}:i_{2[n-1]}, :]}$ $n$ -D niza

### Operacije linearne algebre i operacije s nizovima

$\langle \mathbf{a}   \mathbf{b} \rangle, \mathbf{a}^\top \mathbf{b}$	Skalarni produkt
$\mathbf{a} \mathbf{b}^\top$	Vanjski produkt
$\mathbf{a} \odot \mathbf{b}$	Umnožak po elementima; Hadamardov produkt
$\mathbf{a} \oslash \mathbf{b}$	Dijeljenje po elementima
$\mathbf{a}^{\odot b}$	Potenciranje po elementima
$\mathbf{A} \mathbf{B}$	Matrično množenje
$\mathbf{A}^{-1}$	Inverz matrice
$\mathbf{A}^\top$	Transponiranje
$\text{diag}(\mathbf{a})$	Dijagonalna matrica kojoj dijagonalu čini vektor $\mathbf{a}$
$\det(\mathbf{A})$	Determinanta matrice $\mathbf{A}$
$\ \mathbf{a}\ _2$	$L^2$ -norma vektora $\mathbf{a}$
$\ \mathbf{a}\ _p$	$L^p$ -norma vektora $\mathbf{a}$
$\ \mathbf{A}\ _p$	Matrična $L^p$ -norma matrice $\mathbf{A}$
$\ \mathbf{A}\ _F$	Frobeniusova norma matrice $\mathbf{A}$
$\mathbf{a} \# \mathbf{b}$	Konkatenacija vektora (stupaca) $\mathbf{a} \in \mathbb{R}^n$ i $\mathbf{b} \in \mathbb{R}^m$ u vektor iz $\mathbb{R}^{n+m}$
$\mathbf{A} \# \mathbf{B}$	Konkatenacija nizova po prvoj dimenziji
$\text{vec}(\mathbf{A})$	Funkcija koja preslikava niz iz $\mathbb{R}^{d_1 \times \dots \times d_n}$ u $\mathbb{R}^{d_1 \dots d_n}$
$\dim(\mathbf{a})$	Dimenzija vektora
$\dim(\mathbf{A})$	Vektor dimenzija niza; $[d_1, \dots, d_n]$ za $\mathbf{A} \in \mathbb{R}^{d_1 \times \dots \times d_n}$

### Diferencijalni račun

$\frac{dy}{dx}, \frac{d}{dx} f(x)$	Derivacija $y = f(x)$ po $x$
------------------------------------	------------------------------

$\frac{\partial y}{\partial x}, \frac{\partial}{\partial x} f(x)$	Parcijalna derivacija $y = f(x)$ po $x$
$\nabla_{\mathbf{x}} y, \nabla_{\mathbf{x}} f(x), \left(\frac{\partial y}{\partial \mathbf{x}}\right)^T$	Gradijent $y = f(\mathbf{x})$ po $\mathbf{x}$
$\nabla_{\mathbf{X}} y, \nabla_{\mathbf{X}} f(x)$	Gradijent $y = f(\mathbf{x})$ po $\mathbf{X}$
$\frac{\partial^2 y}{\partial \mathbf{x} \partial \mathbf{x}^T}, \mathbf{H}_f(\mathbf{x}), \mathbf{H}$	Hessijan iz $\mathbb{R}^{n \times n}$ za $f: \mathbb{R}^n \rightarrow \mathbb{R}$ i $y = f(\mathbf{x})$
$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}, \mathbf{J}_f(\mathbf{x}), \mathbf{J}$	Jakobijeva matrica iz $\mathbb{R}^{m \times n}$ za $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ i $\mathbf{y} = f(\mathbf{x})$
$\int_A f(x) dx, \int_{x \in A} f(x)$	Određeni integral funkcije $f(x)$ po $x \in A$
$\int f(x) dx, \int_x f(x)$	Određeni integral funkcije $f(x)$ po $x \in A$ , gdje je $A$ implicitan

## Teorija vjerojatnosti

Svakoj slučajnoj varijabli  $\underline{a}$  jednoznačno je dodijeljena jedna razdioba  $p(\underline{a})$  (ili  $P(\underline{a})$ ) i funkcija gustoće vjerojatnosti (koja može biti poopćena funkcija)  $p_{\underline{a}}(a) = p(\underline{a} = a)$ .  $P(A)$  označava vjerojatnost događaja  $A$ , a  $P_{\underline{a}}$  funkciju vjerojatnosti slučajne varijable  $\underline{a}$ . Mogući su i kraći zapisi  $p(a)$  i  $P(a)$ , gdje se po slovu koje označava vrijednost pretpostavlja slučajna varijabla označena istim slovom bez serifa. Mogu se koristiti i druge oznake za funkciju vjerojatnosti ili funkciju gustoće vjerojatnosti.

$(\underline{a} \mid \underline{b} = b), (\underline{a} \mid b)$	Uvjetna slučajna varijabla
$(\underline{a}, \underline{b})$	Združena slučajna varijabla
$\underline{a} \perp \underline{b}$	Slučajne varijable $\underline{a}$ i $\underline{b}$ su nezavisne
$\underline{a} \not\perp \underline{b}$	Slučajne varijable $\underline{a}$ i $\underline{b}$ su zavisne
$\underline{a} \perp \underline{b} \mid \underline{c}$	Slučajne varijable $\underline{a}$ i $\underline{b}$ su uvjetno nezavisne uz poznat ishod slučajne varijable $\underline{c}$
$\underline{a} \not\perp \underline{b} \mid \underline{c}$	Slučajne varijable $\underline{a}$ i $\underline{b}$ su uvjetno zavisne uz poznat ishod slučajne varijable $\underline{c}$
$p, q$	Razdioba ili funkcija gustoće vjerojatnosti
$A$	Događaj
$\{R(\underline{a})\}$	Događaj definiran predikatom slučajne varijable $\underline{a}$

$P(\{R(\underline{a})\}), P(R(\underline{a}))$	Vjerojatnost događaja $\{R(\underline{a})\}$
$P(\underline{a}), p(\underline{a}), \mathcal{D}$	Razdioba slučajne varijable $\underline{a}$ ; $P$ ako je $\underline{a}$ diskretna slučajna varijabla, a $p$ ako nije ili ako se ne zna
$P(\underline{a} = a), P_{\underline{a}}(a), P(a)$	Vjerojatnost događaja $\{\underline{a} = a\}$
$p(\underline{a} = a), p_{\underline{a}}(a), p(a)$	Gustoća vjerojatnosti događaja $\{\underline{a} = a\}$
$p_{\underline{a} b}(a), p(a   b)$	Gustoća vjerojatnosti događaja $\{\underline{a} = a   \underline{b} = b\}$
$p_{\underline{a}, \underline{b}}(a, b), p(a, b)$	Gustoća vjerojatnosti događaja $\{\underline{a} = a, \underline{b} = b\}$
$\underline{a} \sim q, p(\underline{a}) = q$	Slučajna varijabla $\underline{a}$ ima razdiobu $q$
$\underline{a} \sim \mathbb{A}$	Slučajna varijabla $\underline{a}$ ima takvu razdiobu da svi elementi (multi)skupa $\mathbb{A}$ imaju vjerojatnost proporcionalnu višestrukosti ( $\frac{1}{ \mathbb{A} }$ za običan skup)
$a \sim q$	$a$ se izvlači iz razdiobe $q$
$a \sim \underline{a}, a \sim p(\underline{a})$	$a$ se izvlači iz razdiobe $p(\underline{a})$
$\mathbf{E}_{\substack{a \sim \underline{a}}} f(a), \mathbf{E}_{\underline{a}} f(a)$	Očekivanje funkcije slučajne varijable $\underline{a}$
$\mathbf{D}_{\substack{a \sim \underline{a}}} f(a), \mathbf{D}_{\underline{a}} f(a)$	Disperzija (varijanca) funkcije slučajne varijable $\underline{a}$
$\text{Cov}(\underline{a}, \underline{b})$	Kovarijanca
$\mathcal{N}(\mu, \sigma^2)$	Normalna razdioba s očekivanjem $\mu$ i varijancom $\sigma^2$
$\mathcal{U}(\mathbb{A})$	Uniformna razdioba nad skupom $\mathbb{A}$

## Information theory

$I(\mathbb{A})$	Information content of event $\mathbb{A}$
$H(\underline{a})$	Entropy
$h(\underline{a})$	Differential entropy
$I(\underline{a}; \underline{b})$	Mutual information
$H(\underline{a}   \underline{b})$	Conditional entropy
$H_{\underline{b}}(\underline{a})$	Cross entropy
$D_{\underline{b}}(\underline{a})$	Relative entropy (Kullback-Leibler divergence)



**Grafovi**

$\text{pa}_G(a)$	Skup čvorova roditelja čvora $a$ u grafu $G$
$\text{ch}_G(a)$	Skup čvorova djece čvora $a$ u grafu $G$
$\text{pred}_G(a)$	Skup čvorova prethodnika čvora $a$ u grafu $G$
$\text{succ}_G(a)$	Skup čvorova nasljednika čvora $a$ u grafu $G$

**Other**

$\mathcal{A} \rightarrow \mathcal{B}$	A set of functions with domain $\mathcal{A}$ and codomain $\mathcal{B}$
$f \in (\mathcal{A} \rightarrow \mathcal{B})$	Function definition; a function mapping from its domain $\mathcal{A}$ to its codomain $\mathcal{B}$
$x \mapsto f(x)$	Function definition; a function mapping $x$ from its domain to $f(x)$ in its codomain
$f + g$	Sum of functions
$fg$	Product of functions
$f * g$	Convolution of functions
$\langle f g \rangle$	Scalar product of functions
$ \mathcal{A} $	Set cardinality
$\delta(\cdot)$	Dirac delta
$\llbracket \cdot \rrbracket$	Iverson bracket; $\llbracket P \rrbracket = \begin{cases} 1, & P \equiv \top \\ 0, & P \equiv \perp \end{cases}$
$f[\mathcal{A}]$	Image of $\mathcal{A}$ (subset of the domain) under $f$ ; $f[\mathcal{A}] := \{f(a) \mid a \in \mathcal{A}\}$
$f^{-1}[\mathcal{B}]$	Inverse image (preimage) of $\mathcal{B}$ (subset of the codomain) under $f$ ; $f^{-1}[\mathcal{B}] := \{a \mid f(a) \in \mathcal{B}\}$

# **Part I**

## **Notes**

# Chapter 1

## Interplay of adversarial robustness and generalization

### 1.1 Introduction

Although common machine learning algorithms achieve human-level performance in many tasks, when given out-of-distribution, domain-shifted, corrupted or slightly modified examples from the training data distribution, they can often make overconfident and incorrect predictions [Hendrycks and Gimpel, 2016, Ganin and Lempitsky, 2015, Nguyen et al., 2015, Hendrycks and Dietterich, 2019, Engstrom et al., 2017, Szegedy et al., 2013]. Perhaps most surprisingly, by perturbing input examples (e.g. images) even imperceptibly for humans, common state-of-the-art algorithms can be made to significantly change their predictions both for examples from and outside of the training data distribution [Szegedy et al., 2013, Goodfellow et al., 2014b]. Such perturbed examples are called *adversarial examples*. The existence of adversarial examples indicates that common algorithms are probably performing well for somewhat wrong reasons, without actually *understanding data*.

Although adversarial examples exist on other tasks, image classification algorithms with deep, mostly convolutional, models is the problem considered in most research on adversarial examples. In the following text, the problem of image classification will be considered mostly as well.

Some evidence suggests that there is a trade-off between robustness and generalization [Tsipras et al., 2018, Madry et al., 2017, Su et al., 2018] with current algorithms, which is counter-intuitive because a hypothesis which optimally generalizes would have no adversarial examples. The question remains whether it is feasible or tractable with respect to computation or amount of data to implement such algorithms. A small but interesting step in this direction has

been made by [Stutz et al., 2018].

A brief overview over a subset of the research on adversarial examples will be given with the goals of improving intuition and understanding of adversarial examples, understanding how robustness and generalization of deep models relate, and getting closer to knowing whether it is possible to improve both robustness and generalization on real-world datasets.

## 1.2 Definitions and notation

To prevent some ambiguities, we define some basic machine learning concepts as they will be used here:

- Hypothesis – a mapping from inputs to predictions. A hypothesis will be denoted  $h(\mathbf{x})$  or  $h(\mathbf{x}; \boldsymbol{\theta})$ , where  $\mathbf{x} \in \mathbb{X}$  is an input vector and  $\boldsymbol{\theta}$  a fixed vector of model parameters. For discriminative models  $h(\mathbf{x}) = \mathbb{P}(\underline{y} \mid \mathbf{x})$ , where  $\underline{y} \mid \mathbf{x}$  is a conditioned random variable, short for  $\underline{y} \mid \mathbf{x} = \mathbf{x}$ . Random variables are underlined.
- Model – a set of hypotheses  $\mathcal{H}$  or a function  $h$  defined up to some free parameters  $\boldsymbol{\theta}$ :  $\mathcal{H} = \{\mathbf{x} \mapsto h(\mathbf{x}; \boldsymbol{\theta})\}_{\boldsymbol{\theta}}$ .
- Machine learning algorithm (short: algorithm) – usually a model together with inductive bias and an objective function.
- Classifier – a hypothesis that outputs a categorical distribution.
- The true hypothesis – the best hypothesis given knowledge of the true underlying distribution. Usually, it can be known only for synthetic data.
- Risk – the expected error on the data distribution. For discriminative models it can often be expressed as

$$R(h, \mathcal{D}) := \mathbf{E}_{(\mathbf{x}, y) \sim \mathcal{D}} L(y, h(\mathbf{x})), \quad (1.1)$$

where  $L$  is the loss function, which is usually proportional to negative log-likelihood  $\mathbb{P}(y \mid \mathbf{x}, \boldsymbol{\theta})$ .

- Empirical risk – the error on the empirical distribution  $p_{\mathcal{D}}$  (a uniform distribution over the set of observed examples  $\mathcal{D}$ ):  $R_E(h, \mathcal{D}) := R(h, p_{\mathcal{D}})$ . The objective of a machine learning algorithm is to minimize a combination of empirical risk and structural risk, which usually represents the prior hypothesis probability of the hypothesis (parameters)  $\mathbb{P}(\boldsymbol{\theta})$ .

Here are some common terms related to adversarial examples defined:

- Attack – an algorithm that generates adversarial examples.
- Defense – an idea used to make an algorithm resistant to adversarial examples.

- Natural example – an example from the distribution that the training data was sampled from.
- Adversarial perturbation – the difference between an adversarial example and the natural example it is based on. An adversarial example is based on some natural example iff it is found in its neighbourhood.
- Targeted attack – an attack with the goal that the hypothesis outputs some desired prediction.
- Non-targeted attack – an attack with the goal that the hypothesis outputs any misprediction.
- Threat model – a set of constraints on attacks. A defense should assume a threat model under which it is meant to work.
- Adversarial robustness (short: robustness) – either the performance of a hypothesis or algorithm under some threat model or how large a perturbations must be in order to turn natural examples into adversarial examples.
- Adversarial training – training that uses adversarial examples generated with an attack besides natural examples.
- Standard training – training that uses original examples from the training set without an attack.

Because an adversarial example can be defined in more ways, depending on the usefulness of the definition, and the definition is important for making precise conclusions, the definitions of an adversarial example will be discussed in section 1.3.

### 1.3 Adversarial example definitions

A broadly accepted definition of an adversarial example is that it is *an input designed to fool a hypothesis into producing a misprediction*. To make it more precise and useful, the "designed" part can be replaced with the constraint that adversarial examples are close to examples with correct predictions, which serves as a distinction between adversarial examples and examples far from decision boundaries with mispredicted labels causing standard generalization error. We get the following definition.

**Definition 1 (adversarial example)** *An adversarial example is an input for which the following holds:*

1. *It is close to an input with a correct prediction.*
2. *The hypothesis produces a misprediction.*

By this definition, the set of adversarial examples is a function of the hypothesis, the true hypothesis, and a neighbourhood function. A misprediction can be defined

via some distance between distributions and a threshold. In classification, a misprediction is usually defined as misclassification, i.e. the class with the highest predicted probability being different from true class.

As will be explained, a more practical, but less consistent definition is the following.

**Definition 2 (practical adversarial example)** *A (practical) adversarial example is an input for which the following holds:*

1. *It is close to an input with a correct prediction.*
2. *The hypothesis produces a different prediction than for the input it is based on (the input from point 1).*

Again, the difference between predictions can be defined via a distance between distributions. In classification, predictions of a hypothesis being different usually means that the class with the highest probability differs.

The difference from the first definition is that in the second definition knowledge of the true label is substituted with the prediction for the input with the correct prediction. This is useful for finding adversarial examples because all constraints can be known, but it is possible for a pair of natural examples differing in the label to have overlapping neighbourhoods and thus to be inside each other's neighbourhoods and thus be adversarial examples of each other. By the second definition, Even the true hypothesis can have adversarial examples close to decision boundaries. The first definition is more consistent in this sense, but then, knowing whether something is an adversarial example requires knowing the true hypothesis (i.e. it requires having solved the problem we are trying to solve with machine learning). In the following text, we will refer to the first definition as *the consistent definition* and the second one as *the practical definition*.

There are also some different or broader definitions of adversarial examples. Some authors consider out-of-distribution examples producing high-confidence predictions [Gal and Smith, 2018] or, most broadly, any inputs that fool the hypothesis [Brown et al., 2018].

## 1.4 The manifold hypothesis

A useful assumption in many machine learning problems is that the distribution of high-dimensional data is approximately concentrated in low-dimensional manifolds [Goodfellow et al., 2016]. A manifold can be described as a set of points related through neighbourhoods that locally resembles a euclidean space. In machine learning, the term tends to be used loosely, usually denoting a connected set of points that can be approximated well with a small number of degrees of freedom, each corresponding to a local direction of variation, embedded in a higher-dimensional space. [Goodfellow et al., 2016] provide the following evidence.

Firstly, the fact that the chance a natural example will be sampled from a simple distribution in a high-dimensional space is negligible indicates that the data distribution is highly concentrated. Secondly, we can imagine transformations that represent connections between similar examples and can be applied to traverse the manifold (e.g. for images).

One of the goals of generative models is often to model the data manifold with a lower-dimensional representation. This assumption provides useful intuition for generative models, dimensionality reduction, the understanding of adversarial examples, and some adversarial defenses.

## 1.5 Properties of adversarial examples

Here some properties, phenomena and hypotheses regarding adversarial examples are described, some of which are more or less supported by evidence. Besides the hypotheses described here, there is also a quite predictive and well supported by evidence – the *non-robust features* hypothesis [Ilyas et al., 2019, Tsipras et al., 2018] described in section 1.8, which might not be in contradiction with most of the hypotheses presented here. A broader overview can be found in e.g. Serban and Poll [2018].

**Rareness and closeness to natural examples** The existence of adversarial examples shows that near almost every input, there are misclassified inputs nearby. Evidence also shows that they are rare, i.e. they can not be easily found with random search in the neighbourhood of a natural example. This initially lead to the hypothesis that existence of adversarial examples is caused by high nonlinearity of deep models and that they are located in *low-probability pockets* in the data manifold [Szegedy et al., 2013].

**Rareness of non-adversarial examples in the proximity of adversarial examples** An adversarial example is unlikely to be turned into a non-adversarial example by applying a random perturbation to it. ??? Is this really true given randomized smoothing?

**Local linearity** However, Goodfellow et al. [2014b] have found that it is usually enough to know the locally linear behaviour (gradient with respect to the input) to reliably generate adversarial examples along directions of increasing the linear approximation. This is visualized in figure 1.1. They show that the direction of the perturbation matters more than the position in input space, evidencing against the low-probability pockets hypothesis. They suggest the *linearity hypothesis* whereby high local linearity of models accounts for the existence of adversarial examples. Summing small perturbations in many elements of high-dimensional inputs can produce a large change in the prediction. Su et al. [2017] have shown that even

modifying only 1 pixel can often be enough to cause a misclassification, indicating that models can be highly sensitive to some input elements. [Tabacof and Valle \[2016\]](#) have found more evidence that adversarial examples span many-dimensional subsets in the input space.

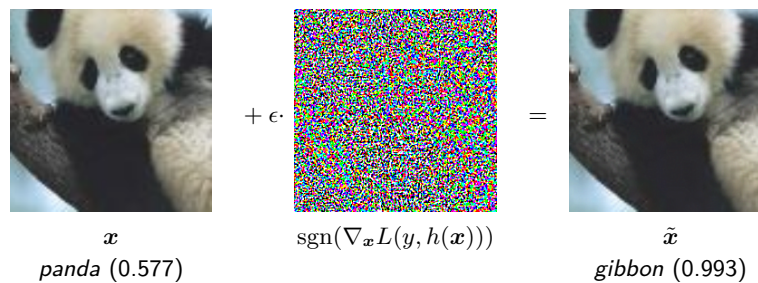


Figure 1.1: Generation of an adversarial example with FGSM, a single step attack. Italic words and numbers represent classes and confidences. The images are from [Goodfellow et al. \[2014b\]](#).

**Transferability** An interesting property is that adversarial examples generated for one model are often misclassified by other models and models trained on different datasets, i.e. they generalize across models and datasets [[Szegedy et al., 2013](#)]. This property of adversarial examples is called transferability. This suggests that most state-of-the-art deep models are similarly biased. More analysis of transferability can be found in [Papernot et al. \[2016\]](#), [Liu et al. \[2017\]](#), [Tramèr et al. \[2017\]](#).

**Universal perturbations** [Moosavi-Dezfooli et al. \[2016a\]](#) have observed that there exist perturbations that can reliably produce adversarial examples when added to almost any input and hypothesize that they exploit geometric correlations between different parts of decision boundaries.

**Boundary tilting** [Tanay and Griffin \[2016\]](#) hypothesize that adversarial examples exist when the decision boundary lies close to the submanifold of sampled data. They suggest that adversarial examples might be occurring along low-variance directions of the data where it is close to the manifold and that robustness could be improved with regularization. This is illustrated in figure 1.2.

**High-dimensional space manifold geometry** [Gilmer et al. \[2018\]](#) hypothesize that the existence of adversarial examples could be a naturally occurring result of the geometry of high-dimensional data manifolds. For a simple dataset with two classes consisting of examples from a pair of high-dimensional concentric spheres, they have observed that most random points in the data distribution are both



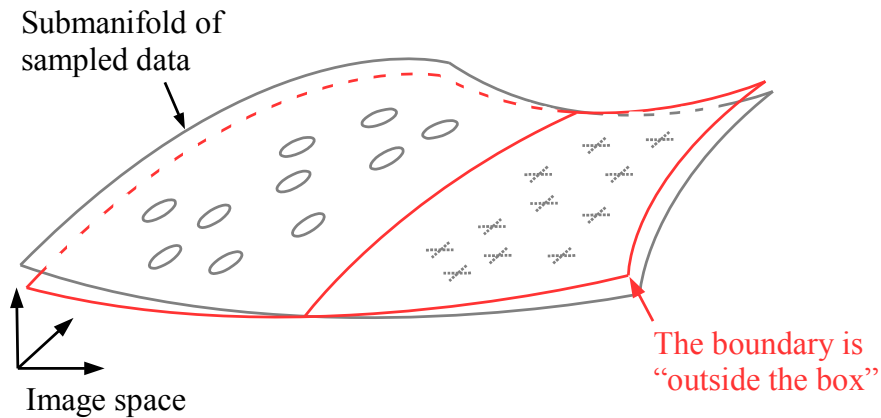


Figure 1.2: An illustration of boundary tilting from Tanay and Griffin [2016].

correctly classified and close to a misclassified point. The authors have also given negative evidence on the hypothesis that adversarial examples are off the data manifold by showing that *on-manifold* adversarial examples can exist as well.

**Adversarial examples of generative models** Adversarial examples are not just a phenomenon related to discriminative models. They have been found to exist for some generative models as well [Goodfellow et al., 2014b, Kos et al., 2018]. An example is shown in figure 1.3.

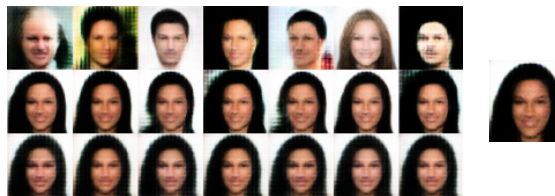


Figure 1.3: Reconstruction outputs for targeted attacks on a VAE-GAN model from Kos et al. [2018]. The rows represent of original image reconstructions (top), reconstructions of adversarial examples generated using an attack in latent space (middle) and a VAE-loss attack (bottom). The target reconstruction is on the right.

**True ambiguity of adversarial examples of robust classifiers** Qualitative assesment of adversarial examples of robust classifiers suggest that they *understand data* much better. Their adversarial examples really are ambiguous, i.e.

generated perturbations don't look like noise but are semantically meaningful to humans [Tsipras et al., 2018, Li, 2018]. This is illustrated in figures 1.4 and 1.5. An interesting observation is that adversarially trained discriminative classifiers together with an iterative attack can interpolate between and generate quite realistic examples. Tsipras et al. [2018] suggest a connection between the saddle point problem of adversarial training and GAN training [Goodfellow et al., 2014a].



Figure 1.4: Original images and adversarial examples generated with a large perturbation using an iterative non-targeted attack on an adversarially trained Restricted ImageNet [Tsipras et al., 2018] classifier from Tsipras et al. [2018].

Further interesting phenomena and hypotheses related to the nature of adversarial examples, adversarial robustness and generalization will be discussed in section 1.8.

### Adversarial reprogramming



Figure 1.5: Clean images (left) and adversarial examples generated using an iterative non-targeted attack on a generative MNIST [LeCun et al., 2015] classifier with the factorization  $p(z)p(y|z)p(x|z,y)$  (right) from Li [2018]. The adversarial examples marked in green are successful.

## 1.6 Finding adversarial examples

Let  $\mathbb{X}$  be the input space, and  $d \in (\mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}^+)$  a *distance function* for defining similarity between inputs. For each example  $x$ , we can also define its *neighbourhood* as  $B_\epsilon(x) = \{x' : d(x', x) \leq \epsilon\}$ , where  $\epsilon$  is the maximum distance from the example.

Ideally, the neighbourhood of an example  $x$  should be the set of *perceptually similar* examples that all belong to the same class as  $x$  (their true class may be at most ambiguous), but it is hard to define such a neighbourhood (as it requires knowing the true model). A practical and common way of defining the neighbourhood function for images is to have  $d$  be a  $L^p$  distance where  $p$  is usually  $\infty$  or 2. Note that, if an example is very near the true class boundary, such a neighbourhood may contain examples belonging to another class.

### 1.6.1 Attack objectives

Finding an adversarial example can be defined as a constrained optimization problem of maximizing some loss with respect to the input with the constraint that the input is in the neighbourhood  $B_\epsilon(x)$ :

$$\tilde{x} = \arg \max_{x' \in B_\epsilon(x)} L(y, h(x')), \quad (1.2)$$

where  $y$  is the true label. Let  $\hat{h}(x) := \arg \max_y h(x)_{[y]}$  denote the function that assigns the label with the highest probability to an input. An objective can also be to find the  $\tilde{x}$  closest to  $x$  such that the classifier misclassifies it [Moosavi-Dezfooli

et al., 2016b]:

$$\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}' : \mathbf{x}' \in B_\epsilon(\mathbf{x}) \wedge \hat{h}(\mathbf{x}') \neq y} d(\mathbf{x}', \mathbf{x}). \quad (1.3)$$

The described objectives, where it only matters that the adversarial example is misclassified, are objectives for *non-targeted adversarial attacks*. There are also *targeted adversarial attacks*, where the objective is to create an adversarial example such that the model classifies it as some desired target. Targeted attack objectives corresponding to equations (1.2) and (1.3) are:

$$\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}' \in B_\epsilon(\mathbf{x})} L(y_a, h(\mathbf{x}')), \quad (1.4)$$

$$\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}' : \mathbf{x}' \in B_\epsilon(\mathbf{x}) \wedge \hat{h}(\mathbf{x}') = y_a} d(\mathbf{x}', \mathbf{x}), \quad (1.5)$$

where  $y_a$  denotes the adversarial target label. The difference in equation (1.4) is loss minimization and the adversarial target label instead of the true label. The difference in equation (1.5) is the condition that the predicted labels equals the adversarial target instead of differing from the true label.

non-targeted adversarial examples can also be generated without knowledge of the true label. Instead of the true label  $y$ , the predicted label  $\hat{h}(\mathbf{x})$  can be used in equations (1.2) and (1.3). Such adversarial examples are called *virtual adversarial examples*. Miyato et al. [2017] propose the following attack objective for use in semi-supervised learning:

$$\tilde{\mathbf{x}} = \arg \min_{\mathbf{x}' \in B_\epsilon(\mathbf{x})} D((y \mid \mathbf{x}, \theta), (y \mid \mathbf{x} = \mathbf{x}', \theta)), \quad (1.6)$$

where  $D$  is some non-negative function that represents distance between distributions.

For adversarial training, non-targeted attacks should be preferred due to the *label-leaking* phenomenon [Kurakin et al., 2016b] where the learned classifier can overfit to adversarial examples and perform better on them than on natural examples, especially with attacks with a small number of iterations. For robustness evaluation with datasets that have many similar classes, non-targeted attacks can too easily fool the classifier and targeted attacks give more meaningful evaluation results [Athalye et al., 2018].

### 1.6.2 Common attacks

Being that finding adversarial examples is a constrained optimization problem, general gradient-based and black-box optimization algorithms can be used for attacks. Additionally, sometimes techniques specific to some potential defense and machine learning algorithm have to be used.

Some commonly known gradient-based attacks for finding adversarial examples are the following (using the notation from section 1.6):

- Box-constrained L-BFGS – Szegedy et al. [2013] propose to minimize  $c\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 + L(y, h(\tilde{\mathbf{x}}))$  with the constraint  $\tilde{\mathbf{x}} \in [0, 1]$  with L-BFGS, a quasi-Newton optimization method.  $c$  is a number obtained via line-search that yields adversarial examples of minimum distance.
- Fast gradient sign method (FGSM) – an attack proposed by Goodfellow et al. [2014b] that requires a single gradient computation:

$$\tilde{\mathbf{x}} = \mathbf{x} + \epsilon \nabla_{\mathbf{x}} L(y, h(\mathbf{x})), \quad (1.7)$$

where  $\epsilon$  is the  $L^\infty$ -norm of the perturbation. For  $L^2$ -constrained perturbations, Miyato et al. [2017] propose  $L^2$  normalization instead of the sign function.

- DeepFool – an iterative non-targeted attack proposed by Moosavi-Dezfooli et al. [2016b] that in each step finds the optimal solution to a linear approximation of a loss in the  $L^2$  ball  $B_\epsilon(\mathbf{x})$  using the gradient in the current adversarial input. It is faster and finds smaller perturbations than L-BFGS and stronger than FGSM.
- Projected gradient descent (PGD) [Madry et al., 2017] basic iterative method (BIM) [Kurakin et al., 2016a] – an iterative gradient-based algorithm with random initialization [Madry et al., 2017] of the perturbation from within the  $B_\epsilon(\mathbf{x})$  at the start and steps in the direction of the gradient sign:

$$\tilde{\mathbf{x}}_i = \Pi_{B_\epsilon(\mathbf{x})}(\tilde{\mathbf{x}}_{i-1} + \alpha \operatorname{sgn}(\nabla_{\tilde{\mathbf{x}}_{i-1}} L(y, h(\tilde{\mathbf{x}}_{i-1}))). \quad (1.8)$$

$\alpha$  is the step size, and  $\Pi_{B_\epsilon(\mathbf{x})}$  is the projection into the  $L^p$   $\epsilon$ -ball around  $\mathbf{x}$ .

- Carlini-Wagner (CW) attacks – Carlini and Wagner [2017b] propose attacks with similar minimal perturbation objectives as Szegedy et al. [2013] and Moosavi-Dezfooli et al. [2016b]. They modify the loss function and, to enable unconstrained optimization, they introduce change of variables  $\delta = \frac{1}{2}(\tanh(\mathbf{w}) + 1) - \mathbf{x}$ , which limits the perturbation  $\delta$  to the interval  $[0, 1]$ . With the PGD (BIM) attack, it is currently probably one of the 2 strongest attacks.

## 1.7 Improving adversarial robustness

There are different approaches (*defenses*) trying to improve adversarial robustness, most of which have been shown to actually be non-robust, but had appeared robust because they intentionally or unintentionally caused attacks that they were evaluated on not to be able to find adversarial examples [Carlini and Wagner, 2017a, Athalye et al., 2018, Uesato et al., 2018, Carlini and Wagner, 2017b]. Thus, it is important to put as much effort as needed to correctly evaluate

robustness, i.e. get an as low as possible upper bound on robustness. [Carlini et al., 2019] is a recent helpful overview on evaluation of robustness.

A broad overview of many defenses can be found in Serban and Poll [2018]. To give a few examples, some approaches use generative models to approximately project inputs to a learned data manifold (e.g. Samangouei et al. [2018]), some approaches are based on limiting the Lipschitz constant of the model to limit sensitivity to small input perturbations by regularization and model modification (e.g. Qian and Wegman [2018]), some research is looking into ways of guaranteeing robustness (e.g. Cohen et al. [2019]).

### 1.7.1 Adversarial training and empirical adversarial risk

The only defense currently believed to be effective according to Athalye et al. [2018] is adversarial training [Goodfellow et al., 2014b] with a strong attack [Madry et al., 2017], where the model is trained on adversarial examples as well as natural examples.

Madry et al. [2017] define what can be called *empirical adversarial risk* by allowing the worst-case attack to modify each the input in the empirical risk expression:

$$R_{\text{EA}}(h, \mathcal{D}) := \mathbf{E}_{(\mathbf{x}, y) \sim p_{\mathcal{D}}} \left( \max_{\tilde{\mathbf{x}} \in B_{\epsilon}(\mathbf{x})} L(y, h(\tilde{\mathbf{x}})) \right). \quad (1.9)$$

They propose PGD for the attack during training and PGD with as large a number of iterations as necessary to approximate the worst-case adversary, i.e. get a better upper bound on robustness.

Still, adversarially trained models are not robust to attacks with weaker constraints than those used for training [Schott et al., 2018]. Furthermore, because adversarial examples are generated and robustness is evaluated according to the practical definition of an adversarial example and using  $L^p$  distance as a non-ideal approximation of perceptual similarity, performance is affected [Madry et al., 2017, Tsipras et al., 2018] and there can exist misclassified examples among which are *invariance-based* adversarial examples [Jacobsen et al., 2019].

## 1.8 Adversarial robustness and generalization

Based on the practical definition of an adversarial example or similar definitions, experimental [Madry et al., 2017, Su et al., 2018] and theoretical evidence [Tsipras et al., 2018] suggests that there is a trade-off between robustness and standard generalization for current models.

Here are some recent discoveries related to adversarial robustness and generalization presented. Robustness to distributional shift and corruptions also

seems to be quite relevant [Gilmer et al., 2019, Hendrycks and Dietterich, 2019], but it will not be discussed here.

### 1.8.1 A trade-off between robustness and generalization

Madry et al. [2017], Su et al. [2017], Tsipras et al. [2018] and others have empirically observed that adversarial robustness with current algorithms requires more capacity and negatively affects generalization. Su et al. [2017] have observed that older convolutional architectures with no shortcut connections, like AlexNet [Krizhevsky et al., 2012] and VGG [Simonyan and Zisserman, 2014] seem to be inherently more robust than architectures like ResNet [He et al., 2015], DenseNet [Huang et al., 2016] and MobileNet [Howard et al., 2017] and NASNets [Zoph et al., 2018] with standard training. Madry et al. [2017] has observed that with adversarial training, more model capacity is required and natural test set performance is reduced. Furthermore, Tsipras et al. [2018] have, based on the practical definition of an adversarial example, theoretically demonstrated an aspect of the trade-off. Another reason that affects performance suggested by them is that salient features might be harder to learn and that algorithms rely on highly predictive but *non-robust* features.

### 1.8.2 Non-robust features

Based on some ideas from Tsipras et al. [2018], Ilyas et al. [2019] propose an interesting and experimentally well supported hypothesis on the nature of features that well-generalizing non-robust classifiers learn. They show that existence of adversarial examples can be directly attributed to existence of non-robust features, "features derived from patterns in the data that are highly predictive but brittle and incomprehensible to humans".

They demonstrate the predictiveness of non-robust features by:

1. constructing a dataset  $\mathcal{D}_{NR}$  where approximately the only useful features are non-robust features by turning inputs of the original dataset  $\mathcal{D}$  into adversarial examples for a classifier that was trained standardly and relabeling them with the adversarial label,
2. training a new classifier on the non-robust dataset  $\mathcal{D}_{NR}$ ,
3. testing the new classifier on the original test set where it achieves performance close to the original classifier and lower robustness.

In another experiment, they try to remove non-robust-features from inputs with the help of an adversarially trained classifier. A new classifier trained on the dataset with removed non-robust features achieves a bit lower performance and robustness not as high as the adversarially trained classifier but quite significant compared to the standardly trained classifier. Results of these experiments with the CIFAR-10 dataset [Krizhevsky, 2009] are shown in figure 1.6.



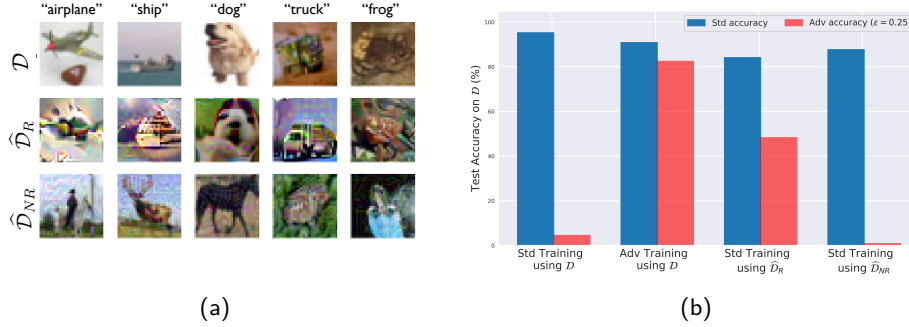


Figure 1.6: (a) Random samples from variants of the CIFAR-10 training set: the original training set; the *robust training set*  $\mathcal{D}_R$ , restricted to features used by a robust model; and the *non-robust training set*  $\mathcal{D}_{NR}$ , restricted to features relevant to a standard model (labels appear incorrect to humans). (b) Standard and robust accuracy on the CIFAR-10 test set ( $\mathcal{D}$ ) for models trained with standard training, adversarial training, and standard training on datasets  $\mathcal{D}_R$  (robust) and  $\mathcal{D}_{NR}$  (non-robust). Adapted from Ilyas et al. [2019].

### 1.8.3 Training with on-manifold adversarial examples

Stutz et al. [2018] challenge the hypothesis that there is a trade-off between robustness and generalization. They hypothesize that most adversarial examples come from directions orthogonal to the learned class manifolds and that training adversarial examples (as per a definition similar to definition 1) limited to the known or learned class manifolds (on-manifold adversarial examples) can improve generalization. They conduct experiments with a synthetic dataset with known class-invariant transformations and datasets with small images that support the hypothesis that generalization can be improved with adversarial training with on-manifold adversarial examples.

Class manifolds to which adversarial examples are restricted and on-manifold and off-manifold adversarial examples are illustrated in figure 1.7.

In one of the experiments Stutz et al. [2018] construct a synthetic dataset with a known manifold (geometric transformations of letters) in order to be able to generate exactly on-manifold adversarial examples by modifying parameters of the geometric transformations. With this dataset, they succeed in improving generalization and on-manifold robustness<sup>1</sup> with adversarial training.

In other experiments they use EMNIST [Cohen et al., 2017], FashionMNIST [Xiao et al., 2017] and CelebA [Liu et al., 2015]. In order to better approximate class manifolds and disable leaving the manifold of a class when an adversarial

<sup>1</sup>By the authors' definition of an adversarial example, which is similar to the consistent definition (definition 1), except for that there is no closeness constraint, making it equivalent to the definition of a misclassified example, on-manifold robustness essentially boils down to generalization.



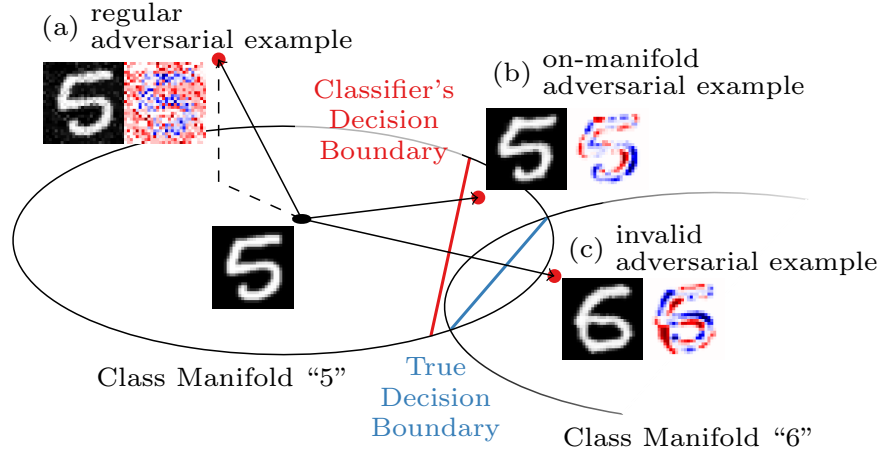


Figure 1.7: An illustration by [Stutz et al. \[2018\]](#) of class manifolds (classes "5" and "6") with a regular (off-manifold) adversarial example and an on-manifold adversarial example.

example is generated for adversarial training, they first train one variational autoencoder (VAE-GAN) per class. They perform training and evaluation analogously to the experiment with synthetic data by allowing the attack to perturb the latent representation of the autoencoder corresponding to the correct class. They measure positive correlation between robustness to on-manifold adversarial examples and generalization. They observe worse quality of on-manifold adversarial examples for the more complex dataset CelebA due to worse approximation quality of their VAE-GAN-s.

## 1.9 Conclusion

This paper presents an overview of ideas related to the existence adversarial examples, hypotheses on their existence and some of their properties. It describes general principles regarding adversarial attacks, defenses and robustness evaluation and presents examples of such algorithms. Finally, some recent discoveries and hypotheses regarding the relation between robustness and generalization are explored. Some recent results [[Stutz et al., 2018](#)] suggest that finding ways of improving both robustness generalization might be an interesting research direction to explore.

## Chapter 2

# Kompromis između generalizacije i robusnosti kod dubokih modela za razumijevanje slika

### 2.1 Uvod

### 2.2 Osnovni pojmovi

#### 2.2.1 Osnovne oznake

**Varijable.** Varijable su nakošene, a konstante uspravne. Vektori se označavaju malim podebljanim slovima ( $\mathbf{a}$ ), matrice i višedimenzionalni nizovi velikim podebljanim slovima ( $\mathbf{A}$ ), a skupovi slovima s udvostručenim linijama ( $\mathbb{A}$ ). Slučajne varijable su podvučene ( $\underline{a}$ ).

**Nizovi.**  $[x_1, \dots, x_n]$  je vektor redak,  $a..b$  i  $a : b$  su kraći zapisi za  $a, \dots, b$ . Redni brojevi elemenata vektora ili višedimenzionalnih nizova se pišu uglatoj zagradi u donjem indeksu. Npr.  $i$ -ti element vektora  $\mathbf{a} = [a_1, \dots, a_n]^\top$  je  $\mathbf{a}_{[i]} = a_i$ ,  $\mathbf{a}_{[i_1:i_2]}$  je vektor koji čine elementi  $\mathbf{a}_{[i_1]}, \dots, \mathbf{a}_{[i_2]}$ ,  $\mathbf{A}_{[i,j]}$  je element, a  $\mathbf{A}_{[i,:]}$  redak matrice.

**Linearna algebra.**  $\langle \mathbf{a} | \mathbf{b} \rangle$  ili  $\mathbf{a}^\top \mathbf{b}$  je skalarni produkt,  $\mathbf{a} \mathbf{b}^\top$  vanjski produkt,  $\mathbf{a} \odot \mathbf{b}$  umnožak po elementima,  $\mathbf{a} \oslash \mathbf{b}$  dijeljenje po elementima,  $\text{diag}(\mathbf{a})$  dijagonalna matrica kojoj dijagonalu čini  $\mathbf{a}$ ,  $\text{vec}(\mathbf{A})$  predstavljanje niza vektorom.

**Diferencijalni račun.**  $\frac{\partial y}{\partial x} = \frac{\partial}{\partial x} f(x)$  je parcijalna derivacija  $y = f(x)$  po  $x$   
 $\nabla_x y = \nabla_x f(x) = \left( \frac{\partial y}{\partial x} \right)^T$  gradijent,  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{J}_f(\mathbf{x})$  Jakobijeva matrica iz  $\mathbb{R}^{m \times n}$  za  
 $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , a  $\int_A f(x) dx = \int_{x \in A} f(x)$  određeni integral.

**Vjerojatnost.** Svako slučajnoj varijabli  $\underline{a}$  jednoznačno je dodijeljena jedna razdioba  $p(\underline{a})$  (ili  $P(\underline{a})$ ) i funkcija gustoće vjerojatnosti (koja može biti poopćena funkcija)  $p_{\underline{a}}(a) = p(\underline{a} = a)$ .  $P(A)$  označava vjerojatnost događaja  $A$ , a  $P_{\underline{a}}$  funkciju vjerojatnosti slučajne varijable  $\underline{a}$ . Mogući su i kraći zapisi  $p(a)$  i  $P(a)$ , gdje se po slovu koje označava vrijednost pretpostavlja slučajna varijabla označena istim slovom bez serifa. Mogu se koristiti i druge oznake za funkciju vjerojatnosti ili funkciju gustoće vjerojatnosti. Još neki primjeri su:  $(\underline{a} \mid \underline{b} = b)$  ili  $(\underline{a} \mid \underline{b})$  je uvjetna slučajna varijabla,  $(\underline{a}, \underline{b})$  združena slučajna varijabla,  $P(\underline{a})$ ,  $p(\underline{b})$ ,  $\mathcal{D}$  razdiobe,  $\underline{a} \sim q$  oznaka da slučajna varijabla  $\underline{a}$  ima razdiobu  $q$ ,  $a \sim q$  oznaka da je  $a$  uzorak iz razdiobe  $q$ ,  $\mathbf{E}_{a \sim \underline{a}} f(a)$  ili  $\mathbf{E}_{\underline{a}} f(a)$  očekivanje funkcije slučajne varijable  $\underline{a}$ ,  $\mathbf{D}_{a \sim \underline{a}} f(a)$  ili  $\mathbf{D}_{\underline{a}} f(a)$  disperzija (varijanca),  $\mathcal{N}(\mu, \sigma^2)$  normalna razdioba s očekivanjem  $\mu$  i varijancom  $\sigma^2$ ,  $\mathcal{U}(A)$  uniformna razdioba nad skupom  $A$ .

**Teorija informacije.**  $H(\underline{a})$  je entropija,  $h(\underline{a})$  diferencijalna entropija,  $I(\underline{a}, \underline{b})$  uzajamna informacija,  $H(\underline{a} \mid \underline{b})$  uvjetna entropija,  $H_{\underline{b}}(\underline{a})$  unakrsna entropija,  $D_{\text{KL}}(\underline{a} \parallel \underline{b})$  Kullback-Leiblerova divergencija (relativna entropija).

**Ostalo.**  $\mathbf{0}$  je nul-vektor,  $\mathbf{e}_i$   $i$ -ti vektor kanonske baze,  $\mathbf{1}$  zbroj svih vektora kanonske baze, a  $\mathbf{I}$  matrica identiteta.  $\mathcal{A} \rightarrow \mathcal{B}$  je skup funkcija s domenom  $\mathcal{A}$  i kodomenom  $\mathcal{B}$ ,  $x \mapsto g(x)$  funkcija koja preslikava  $x$  u  $g(x)$ ,  $|\mathcal{A}|$  kardinalitet

skupa,  $\llbracket \cdot \rrbracket$  Iversonovu zagrada za koju vrijedi  $\llbracket P \rrbracket := \begin{cases} 1, & P \equiv \top, \\ 0, & P \equiv \perp. \end{cases}$

## 2.2.2 Nadzirano učenje

Kapacitet, podnaučenost i prenaučenosť

Rizik i funkcija pogreške

## 2.2.3 Duboko učenje i konvolucijski modeli

Duboki unaprijedni modeli

Učenje

Algoritam propagacije pogreške unatrag

## *CHAPTER 2. KOMPROMIS IZMEĐU GENERALIZACIJE I ROBUSTNOSTI KOD DUBOKIH MODELA ZA RAZUM*

### **Stohastička optimizacija**

### **Problemi i poboljšavanje učenja i generalizacije**

### **Konvolucijski modeli**

### **Konvolucija**

### **Konvolucijski sloj**

### **Sažimanje**

## **2.3 Neprijateljska robusnost i neprijateljski primjeri**

## **2.4 Osnovne definicije i oznake u vezi neprijateljskih primjera**

## **2.5 Svojstva neprijateljskih primjera kod dubokih modela**

## **2.6 Pronalaženje neprijateljskih primjera**

## **2.7 Robusnost i generalizacija**

### **2.7.1 Postizanje robusnosti**

### **2.7.2 Kompromis između robusnosti i generalizacije**

## **2.8 Eksperimenti**

### **2.8.1 Postizanje robusnosti**

### **2.8.2 Smanjivanje kompromisa između robusnosti i generalizacije kod klasifikacije slika**

### **2.8.3 Robusno učenje modela za semantičku segmentaciju**

## **2.9 Zaključak**

## **Chapter 3**

### **Pieces of text**

### 3.0.1 Robustness evaluation

Athalye et al. [2018] state that, if there are classes that are very similar, using targeted attacks for robustness evaluation can give more meaningful results.

stronger attack targeted, stronger evaluation untargeted

#### Adversarial training

Kurakin et al. [2016b] note that by using the true label in the loss in untargeted attacks ( $y$  in equation (??)) can cause





### Distance metrics for images

Usually, an  $L^p$  distance ( $d(\mathbf{x}', \mathbf{x}) = \|\tilde{\mathbf{x}} - \mathbf{x}\|_p$ ) is used as a distance metric for adversarial examples.

**Scale-invariant norms.** Let  $\mathbf{x}$  denote some image (or perturbation) and  $\mathbf{x}_\lambda$  the same image with dimensions scaled by  $\lambda$ .  $\mathbf{x}_\lambda$  has a greater norm because it contains  $\lambda^2$  the number of pixels of the original image and every pixel is approximately effectively repeated  $\lambda^2$  times. In terms of  $\|\mathbf{x}\|_p$ , its norm can be approximated like this:

$$\begin{aligned}
 \|\mathbf{x}_\lambda\|_p &= \left( \sum_{u \in \{0 \dots \lambda H\}} \sum_{v \in \{0 \dots \lambda W\}} |\mathbf{x}_{\lambda[u,v]}|^p \right)^{\frac{1}{p}} && \text{(norm definition)} \\
 &\approx \left( \sum_{u \in \{0 \dots H\}} \sum_{v \in \{0 \dots W\}} \lambda^2 |\mathbf{x}_{[u,v]}|^p \right)^{\frac{1}{p}} && (\lambda^2 \text{ element copies}) \\
 &= \lambda^{\frac{2}{p}} \left( \sum_{u \in \{0 \dots H\}} \sum_{v \in \{0 \dots W\}} |\mathbf{x}_{[u,v]}|^p \right)^{\frac{1}{p}} \\
 &= \lambda^{\frac{2}{p}} \|\mathbf{x}\|_p. && \text{(norm definition)}
 \end{aligned}$$

If there is a perturbation in 2 different resolutions, the higher-resolution one will have a greater norm by approximately a factor of  $(\lambda^2)^{\frac{1}{p}}$ , especially if the perturbations don't have too high spatial frequencies. Hence, scale invariant equivalents of  $L^p$  norms can be defined by dividing the norm by the scale factor. If we choose the scale factor to be relative to the scale of an image with area 1,  $\lambda^2$  is  $n$ , the total number of pixels.

We can define **scale-invariant norms** like this:

$$\mathfrak{m}_p(\mathbf{x}) := n^{-\frac{1}{p}} \|\mathbf{x}\|_p, \quad (3.1)$$

This can also be expressed as the **generalized mean**  $M_p$  (also known as **power mean**) of the elementwise absolute value:

$$\mathfrak{m}_p(\mathbf{x}) = M_p(|\mathbf{x}|) := \left( \frac{1}{n} \sum_i |\mathbf{x}_{[i]}|^p \right)^{\frac{1}{p}}. \quad (3.2)$$

Such norms could probably be useful for comparison of norms between different-resolution images and different datasets, and hyperparameter choice for adversarial attacks.

Maybe something similar could be done about objects of different scale?

...

**Expectation of scale-invariant norms uniformly distributed random vectors.**

Let  $\mathbf{x}_n$  be a random vector with  $n$  independent elements  $\mathbf{x}_{n[i]} \sim \mathcal{U}([- \epsilon, \epsilon])$ .

An easy special case is when  $p = 1$ :

$$\mathbf{E}[\mathbf{m}_1(\mathbf{x}_n)] = \mathbf{E}\left[\frac{1}{n} \sum_i |\mathbf{x}_{n[i]}|\right] = \mathbf{E}[|\mathbf{x}_{n[i]}|] = \frac{\epsilon}{2}.$$

For very large  $n$ , we can approximate it like this:

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbf{E}[\mathbf{m}_p(\mathbf{x}_n)] &= \mathbf{E}\left[\lim_{n \rightarrow \infty} \left(\frac{1}{n} \sum_i |\mathbf{x}_{n[i]}|^p\right)^{\frac{1}{p}}\right] && \text{(dominated convergence)} \\ &= \mathbf{E}\left[\mathbf{E}[|\mathbf{x}_{n[0]}|^p]^{\frac{1}{p}}\right] && \text{(IID elements)} \\ &= \mathbf{E}[|\mathbf{x}_{n[0]}|^p]^{\frac{1}{p}} && \text{(redundant expectation)} \\ &= \frac{1}{2\epsilon} \left(\int_{-\epsilon}^{\epsilon} |x|^p dx\right)^{\frac{1}{p}} && \text{(definition of expectation)} \\ &= \frac{1}{\epsilon} \left(\int_0^{\epsilon} x^p dx\right)^{\frac{1}{p}} && \text{(absolute value symmetry)} \\ &= \frac{1}{\epsilon} \left(\frac{\epsilon^{p+1}}{p+1}\right)^{\frac{1}{p}}. && (3.3) \end{aligned}$$

Specially, for  $\epsilon = 1$ , it evaluates to

$$\lim_{n \rightarrow \infty} \mathbf{E}[\mathbf{m}_p(\mathbf{x}_n)] = (p+1)^{-\frac{1}{p}}. \quad (3.4)$$

TODO: Show that (3.3) is the upper bound if  $p > 1$  and the lower bound if  $p < 1$ .

By the power mean inequality, which can be proved by applying the Jensen inequality,  $\mathbf{m}_p(\mathbf{x})$  is monotonically increasing in  $p$  (for  $p > 0$ ).

For easier analysis, let's first consider a special case:  $\mathbf{x}_{n[i]} \sim \mathcal{U}([0, 1])$ .

$$\begin{aligned} \mathbf{E}[m_p(\mathbf{x}_n)] &= \mathbf{E} \left[ \left( \frac{1}{n} \sum_i |\mathbf{x}_{n[i]}|^p \right)^{\frac{1}{p}} \right] \\ &= \int \cdots \int_{[0,1]^n} \left( \frac{1}{n} \sum_i |\mathbf{x}_{[i]}|^p \right)^{\frac{1}{p}} d\mathbf{x} \\ &= n^{-\frac{1}{p}} \int \cdots \int_{[0,1]^n} \left( \sum_i |\mathbf{x}_{[i]}|^p \right)^{\frac{1}{p}} d\mathbf{x} \end{aligned}$$

**Local  $p$ -norm and hierarchical norm.** Let  $\mathbf{k}$  denote a non-negative 2-D kernel with  $\|\mathbf{k}\|_1 = 1$ , e.g. Gaussian centered at  $(0, 0)$ . Let  $\mathbf{x}$  denote an image perturbation and assume that it has a single channel for simplicity. We can define the local  $p$ -norm around a pixel  $(i, j)$  as

$$\text{LocalNorm}_{p, \mathbf{k}}(\mathbf{x})_{[i,j]} = (|\mathbf{x}|^p * \mathbf{k})_{[i,j]}^{\frac{1}{p}}, \quad (3.5)$$

where the absolute value and powering operations are elementwise.

We can then define a bi-level hierarchical  $(p_1, p_2)$ -norm as  $\|\text{LocalNorm}_{p_1, \mathbf{k}}(\mathbf{x})\|_{p_2}$ . This can be generalized to a multi-level norm  $(p_1, \dots, p_n)$ -norm by chaining multiple local norms with potentially different kernels until the last, global,  $p_n$ -norm.

Why?

# Bibliography

Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 274–283, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/athalye18a.html>.

Peter L. Bartlett, Dylan J. Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pages 6241–6250, USA, 2017. Curran Associates Inc. ISBN 978-1-5108-6096-4. URL <http://dl.acm.org/citation.cfm?id=3295222.3295372>.

Wieland Brendel and Matthias Bethge. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019. URL <https://openreview.net/forum?id=SkfMWhAqYQ>.

Tom B. Brown, Nicholas Carlini, Chiyuan Zhang, Catherine Olsson, Paul Francis Christiano, and Ian J. Goodfellow. Unrestricted adversarial examples. *CoRR*, abs/1809.08352, 2018. URL <http://arxiv.org/abs/1809.08352>.

Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, AISec '17, pages 3–14, New York, NY, USA, 2017a. ACM. ISBN 978-1-4503-5202-4. doi: 10.1145/3128572.3140444. URL <http://doi.acm.org/10.1145/3128572.3140444>.

Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 39–57, 2017b. doi: 10.1109/SP.2017.49. URL <https://doi.org/10.1109/SP.2017.49>.

Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian J. Goodfellow, Aleksander Madry, and Alexey

- Kurakin. On evaluating adversarial robustness. *CoRR*, abs/1902.06705, 2019. URL <http://arxiv.org/abs/1902.06705>.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of mnist to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/cohen19c.html>.
- Neven Elezović. *Vjerojatnost i statistika: Slučajne varijable*. Element, 2007.
- Logan Engstrom, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. A rotation and a translation suffice: Fooling cnns with simple transformations. *CoRR*, abs/1712.02779, 2017. URL <http://arxiv.org/abs/1712.02779>.
- Y. Gal and L. Smith. Sufficient Conditions for Idealised Models to Have No Adversarial Examples: a Theoretical and Empirical Study with Bayesian Neural Networks. *ArXiv e-prints*, June 2018.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37*, ICML'15, pages 1180–1189. JMLR.org, 2015. URL <http://dl.acm.org/citation.cfm?id=3045118.3045244>.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *ICLR*. OpenReview.net, 2019. URL <http://dblp.uni-trier.de/db/conf/iclr/iclr2019.html#GeirhosRMBWB19>.
- Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S. Schoenholz, Maithra Raghu, Martin Wattenberg, and Ian J. Goodfellow. Adversarial spheres. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings*, 2018. URL <https://openreview.net/forum?id=SkthILkPf>.
- Justin Gilmer, Nicolas Ford, Nicholas Carlini, and Ekin Cubuk. Adversarial examples are a natural consequence of test error in noise. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2280–2289, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/gilmer19a.html>.

- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'14, pages 2672–2680, Cambridge, MA, USA, 2014a. MIT Press. URL <http://dl.acm.org/citation.cfm?id=2969033.2969125>.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2014b. URL <http://arxiv.org/abs/1412.6572>.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. *CoRR*, abs/1706.04599, 2017. URL <http://arxiv.org/abs/1706.04599>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *CoRR*, abs/1903.12261, 2019. URL <http://arxiv.org/abs/1903.12261>.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *CoRR*, abs/1610.02136, 2016. URL <http://arxiv.org/abs/1610.02136>.
- Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications, 2017. URL <http://arxiv.org/abs/1704.04861>.
- Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks. *CoRR*, abs/1608.06993, 2016. URL <http://arxiv.org/abs/1608.06993>.
- Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features, 2019. URL <http://arxiv.org/abs/1905.02175>. cite arxiv:1905.02175.
- Jörn-Henrik Jacobsen, Jens Behrmann, Nicholas Carlini, Florian Tramèr, and Nicolas Papernot. Exploiting excessive invariance caused by norm-bounded adversarial robustness. *CoRR*, abs/1903.10484, 2019. URL <http://arxiv.org/abs/1903.10484>.

- Yiding Jiang, Dilip Krishnan, Hossein Mobahi, and Samy Bengio. Predicting the generalization gap in deep networks with margin distributions. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019. URL <https://openreview.net/forum?id=HJIQfnCqKX>.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *CoRR*, abs/1703.04977, 2017. URL <http://arxiv.org/abs/1703.04977>.
- J. Kos, I. Fischer, and D. Song. Adversarial examples for generative models. In *2018 IEEE Security and Privacy Workshops (SPW)*, pages 36–42, May 2018. doi: 10.1109/SPW.2018.00014.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. URL <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.
- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *CoRR*, abs/1607.02533, 2016a. URL <http://arxiv.org/abs/1607.02533>.
- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *CoRR*, abs/1611.01236, 2016b. URL <http://arxiv.org/abs/1611.01236>.
- Yann LeCun, Yoshua Bengio, and Geoffrey E. Hinton. Deep learning. *Nature*, 521 (7553):436–444, 2015. doi: 10.1038/nature14539. URL <https://doi.org/10.1038/nature14539>.
- Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *CoRR*, abs/1711.09325, 2017.
- Yingzhen Li. Are generative classifiers more robust to adversarial attacks? *CoRR*, abs/1802.06552, 2018. URL <http://arxiv.org/abs/1802.06552>.
- Shiyu Liang, Yixuan Li, and R. Srikant. Principled detection of out-of-distribution examples in neural networks. *CoRR*, abs/1706.02690, 2017. URL <http://arxiv.org/abs/1706.02690>.
- Qianli Liao, Brando Miranda, Andrzej Banburski, Jack Hidary, and Tomaso A. Poggio. A surprising linear relationship predicts test performance in deep networks. *CoRR*, abs/1807.09659, 2018. URL <http://arxiv.org/abs/1807.09659>.

- Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017. URL <https://openreview.net/forum?id=Sys6GJqxl>.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *CoRR*, abs/1706.06083, 2017. URL <http://arxiv.org/abs/1706.06083>.
- Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *CoRR*, abs/1704.03976, 2017. URL <http://arxiv.org/abs/1704.03976>.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. *CoRR*, abs/1610.08401, 2016a. URL <http://arxiv.org/abs/1610.08401>.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *CVPR*, pages 2574–2582. IEEE Computer Society, 2016b.
- Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012. ISBN 0262018020, 9780262018029.
- Anh Mai Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *CVPR*, pages 427–436. IEEE Computer Society, 2015. ISBN 978-1-4673-6964-0. URL <http://dblp.uni-trier.de/db/conf/cvpr/cvpr2015.html#NguyenYC15>.
- Nicolas Papernot, Patrick D. McDaniel, and Ian J. Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *CoRR*, abs/1605.07277, 2016. URL <http://dblp.uni-trier.de/db/journals/corr/corr1605.html#PapernotMG16>.
- Stanislav Pidhorskyi, Ranya Almohsen, and Gianfranco Doretto. Generative probabilistic novelty detection with adversarial autoencoders. In *NeurIPS*, pages 6823–6834, 2018.
- Haifeng Qian and Mark N. Wegman. L2-nonexpansive neural networks. *CoRR*, abs/1802.07896, 2018. URL <http://arxiv.org/abs/1802.07896>.
- Mohammad Sabokrou, Mohammad Khalooei, Mahmood Fathy, and Ehsan Adeli. Adversarially learned one-class classifier for novelty detection. In *CVPR*, pages 3379–3388. IEEE Computer Society, 2018.



- Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against adversarial attacks using generative models. *CoRR*, abs/1805.06605, 2018. URL <http://arxiv.org/abs/1805.06605>.
- Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. Towards the first adversarially robust neural network model on mnist. *CoRR*, abs/1805.09190, 2018. URL <http://arxiv.org/abs/1805.09190>.
- Alexandru Constantin Serban and Erik Poll. Adversarial examples - A complete characterisation of the phenomenon. *CoRR*, abs/1810.01185, 2018. URL <http://arxiv.org/abs/1810.01185>.
- K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- David Stutz, Matthias Hein, and Bernt Schiele. Disentangling adversarial robustness and generalization. *CoRR*, abs/1812.00740, 2018.
- Dong Su, Huan Zhang, Hongge Chen, Jinfeng Yi, Pin-Yu Chen, and Yupeng Gao. Is robustness the cost of accuracy? - A comprehensive study on the robustness of 18 deep image classification models. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XII*, pages 644–661, 2018. doi: 10.1007/978-3-030-01258-8\_39. URL [https://doi.org/10.1007/978-3-030-01258-8\\_39](https://doi.org/10.1007/978-3-030-01258-8_39).
- Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. One pixel attack for fooling deep neural networks. *CoRR*, abs/1710.08864, 2017. URL <http://arxiv.org/abs/1710.08864>.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *CoRR*, abs/1312.6199, 2013. URL <http://arxiv.org/abs/1312.6199>.
- Pedro Tabacof and Eduardo Valle. Exploring the space of adversarial images. In *2016 International Joint Conference on Neural Networks, IJCNN 2016, Vancouver, BC, Canada, July 24-29, 2016*, pages 426–433, 2016. doi: 10.1109/IJCNN.2016.7727230. URL <https://doi.org/10.1109/IJCNN.2016.7727230>.
- Thomas Tanay and Lewis D. Griffin. A boundary tilting persepective on the phenomenon of adversarial examples. *CoRR*, abs/1608.07690, 2016. URL <http://arxiv.org/abs/1608.07690>.
- Florian Tramèr, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. The space of transferable adversarial examples. *CoRR*, abs/1704.03453, 2017. URL <http://arxiv.org/abs/1704.03453>.
- Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy, 2018. URL <http://arxiv.org/abs/1805.12152>. cite arxiv:1805.12152.

- Jonathan Uesato, Brendan O'Donoghue, Pushmeet Kohli, and Aäron van den Oord. Adversarial risk and the dangers of evaluating against weak attacks. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pages 5032–5041, 2018. URL <http://proceedings.mlr.press/v80/uesato18a.html>.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. 2017.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable architectures for scalable image recognition. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 8697–8710, 2018. doi: 10.1109/CVPR.2018.00907. URL [http://openaccess.thecvf.com/content\\_cvpr\\_2018/html/Zoph\\_Learning\\_Transferable\\_Architectures\\_CVPR\\_2018\\_paper.html](http://openaccess.thecvf.com/content_cvpr_2018/html/Zoph_Learning_Transferable_Architectures_CVPR_2018_paper.html).