

Adversarial Examples

Generative Adversarial Networks



Presenter: Efstathios Chatzikyriakidis

M.Sc. in Informatics and Communications
Specialization in Computational Intelligence and Digital Media
School of Informatics, Faculty of Sciences
Aristotle University of Thessaloniki, Hellas

B.Sc. in Informatics and Communications
Specialization in Software Engineering

Department of Informatics and Communications, Faculty of Applied Technology
Technological Educational Institute of Central Macedonia, Hellas

Email : contact@efxa.org

Website : <http://www.efxa.org/>

Contributor: Christos Papaioannidis



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731667 (MULTIDRONE)



Introduction

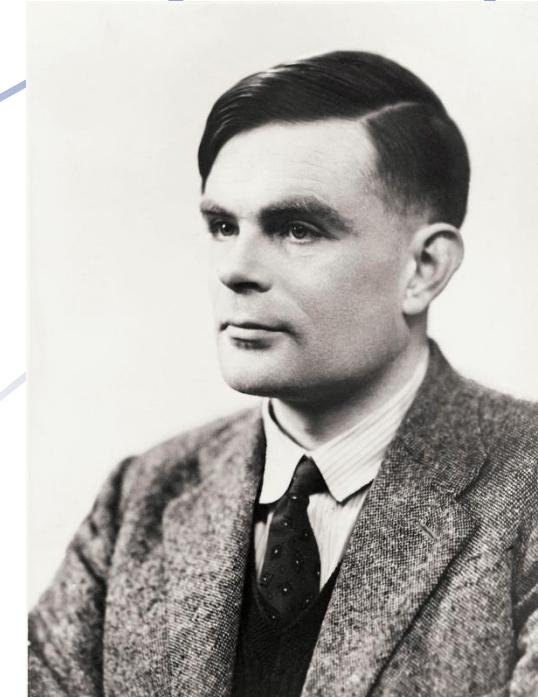
Thoughts on Machine Intelligence

MultDrone



Alan Turing, "Computing Machinery and Intelligence", 1950

- “Can machines think?” (this was a really difficult question)
- “Can machines do what we can do?” (this question was more promising)
- “Can a machine act indistinguishably from the way a thinker acts?” (Turing Test)
- So, when the machines become intelligent we will not be able to prove that they are not



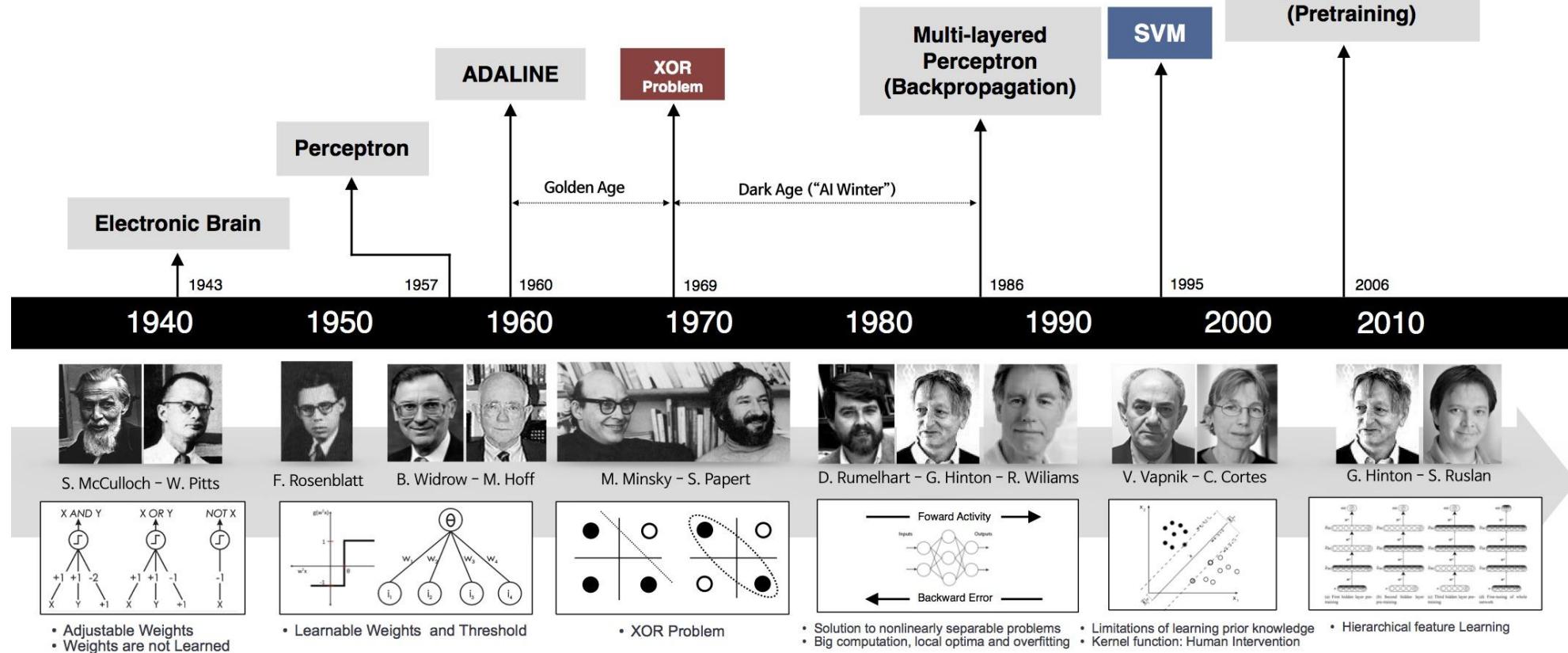
Introduction

Important history of Artificial Intelligence

MultDrone



We had a long journey... and we are still at the birth of it...



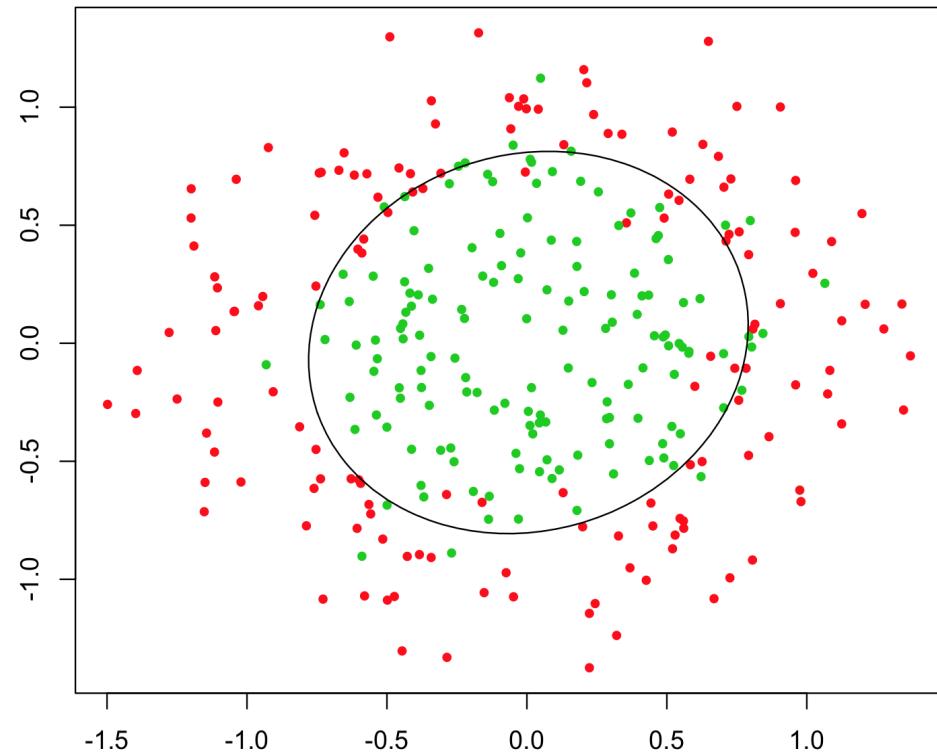
Introduction

Advantages of Artificial Neural Networks

MultDrone



Satisfactory separation of non-linear separable input data



Introduction

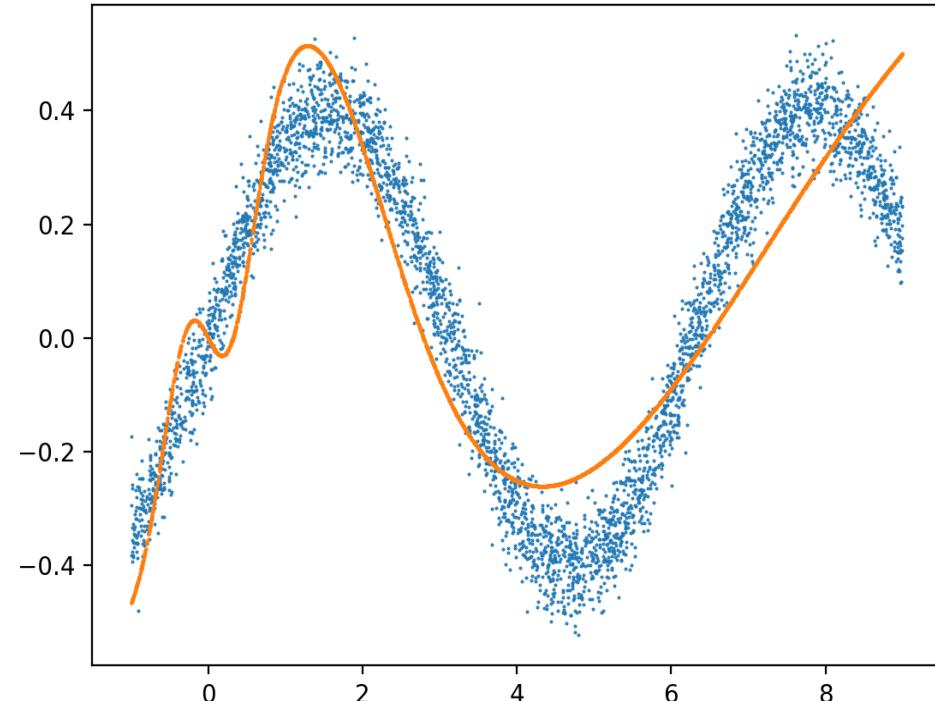
Advantages of Artificial Neural Networks

MultDrone



Universal Approximation Theorem (George Cybenko)

Satisfactory function approximation using only input data



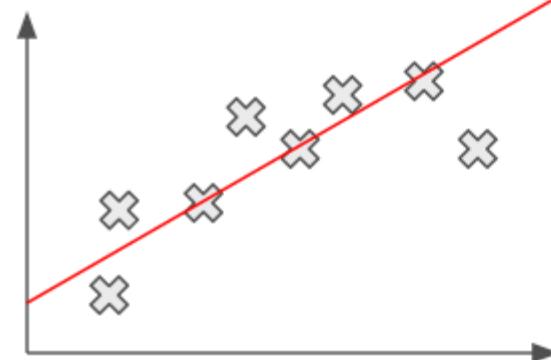
Introduction

Advantages of Artificial Neural Networks

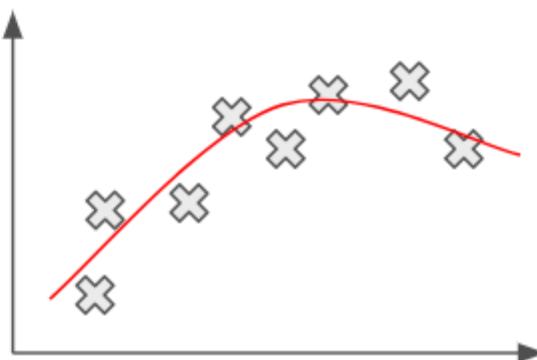
MultDrone



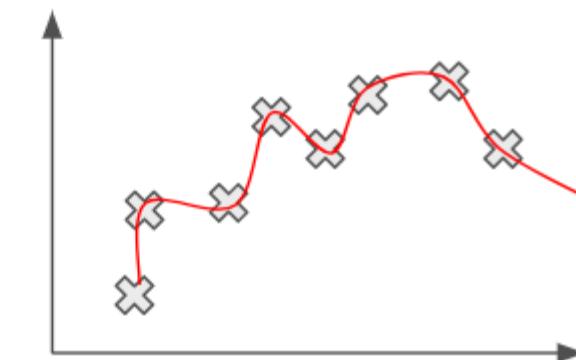
Good generalization that captures the general manifold of data



Underfitting



Optimal



Overfitting

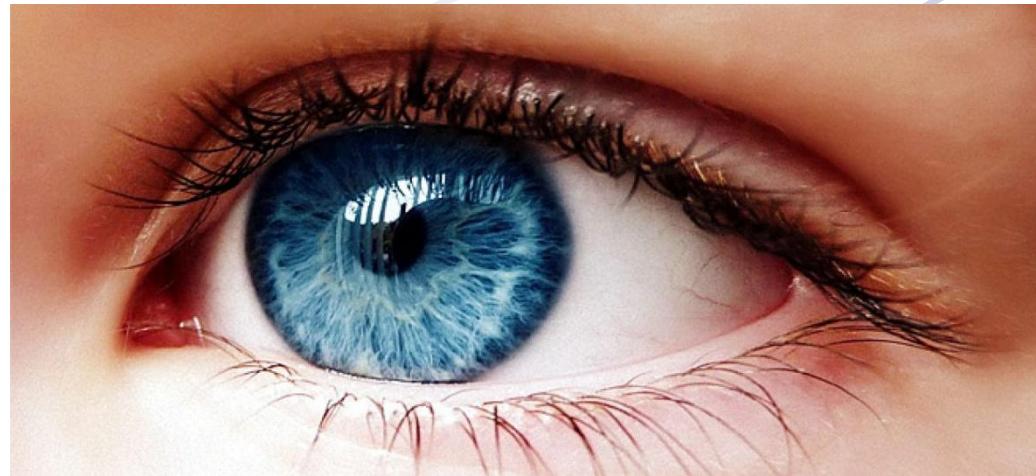
Introduction

Computer Vision and its future goal

MultDrone



- The strongest sense of most animal species
- Simulate how brains see and understand the world through vision sense



Introduction

Computer Vision and its future goal

| | |
|--|---------|
| 1. Q: Is there a person in the blue region? | A: yes |
| 2. Q: Is there a unique person in the blue region? (Label this person 1) | A: yes |
| 3. Q: Is person 1 carrying something? | A: yes |
| 4. Q: Is person 1 female? | A: yes |
| 5. Q: Is person 1 walking on a sidewalk? | A: yes |
| 6. Q: Is person 1 interacting with any other object? | A: no |
| : | |
| 9. Q: Is there a unique vehicle in the yellow region? (Label this vehicle 1) | A: yes |
| 10. Q: Is vehicle 1 light-colored? | A: yes |
| 11. Q: Is vehicle 1 moving? | A: no |
| 12. Q: Is vehicle 1 parked and a car? | A: yes |
| : | |
| 14. Q: Does vehicle 1 have exactly one visible tire? | A: no |
| 15. Q: Is vehicle 1 interacting with any other object? | A: no |
| 17. Q: Is there a unique person in the red region? | A: no |
| 18. Q: Is there a unique person that is female in the red region? | A: no |
| 19. Q: Is there a person that is standing still in the red region? | A: yes |
| 20. Q: Is there a unique person standing still in the red region? (Label this person 2) | A: yes |
| : | |
| 23. Q: Is person 2 interacting with any other object? | A: yes |
| 24. Q: Is person 1 taller than person 2? | A: amb. |
| 25. Q: Is person 1 closer (to the camera) than person 2? | A: no |
| 26. Q: Is there a person in the red region? | A: yes |
| 27. Q: Is there a unique person in the red region? (Label this person 3) | A: yes |
| : | |
| 36. Q: Is there an interaction between person 2 and person 3? | A: yes |
| 37. Q: Are person 2 and person 3 talking? | A: yes |



MultiDrone

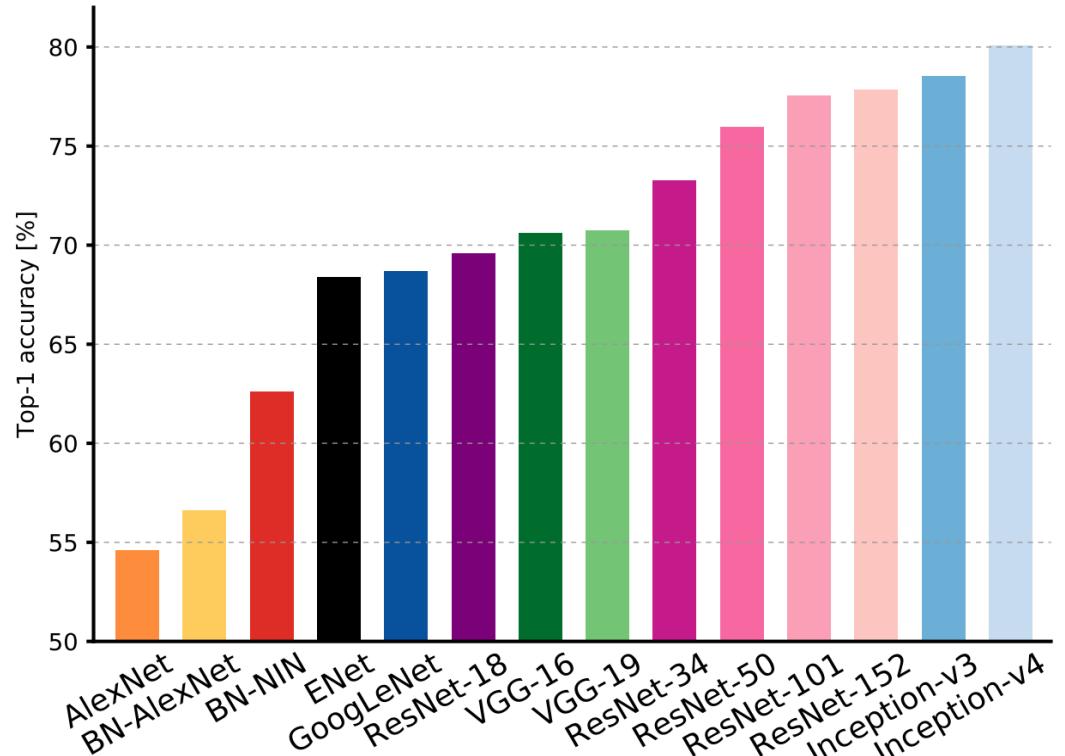


Visual Turing Test

Introduction

State-of-the-art Computer Vision

ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

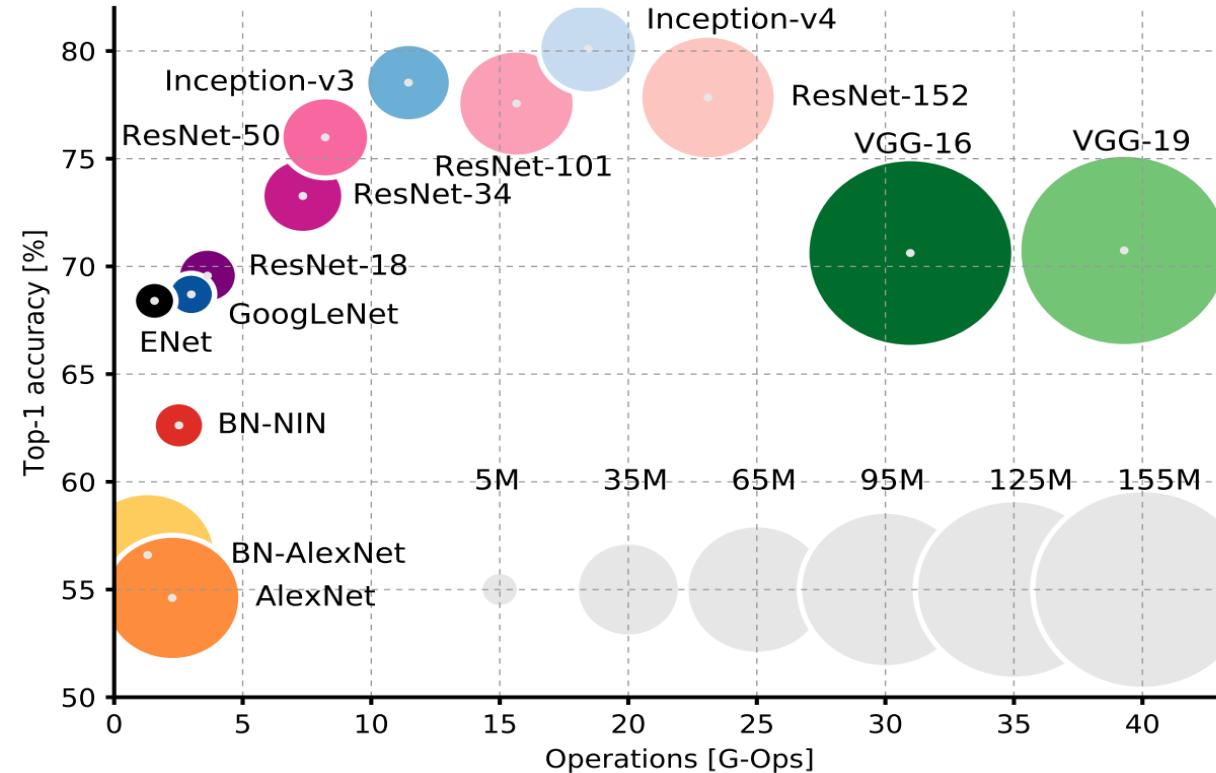


Introduction

State-of-the-art Computer Vision

MultDrone

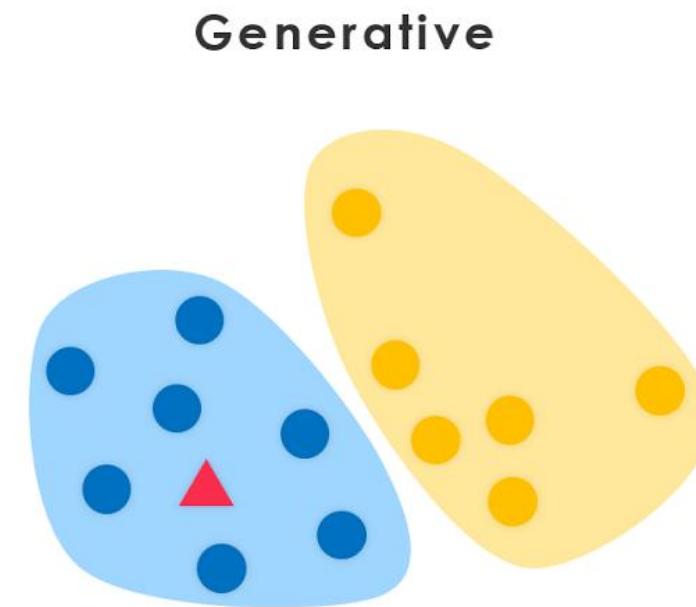
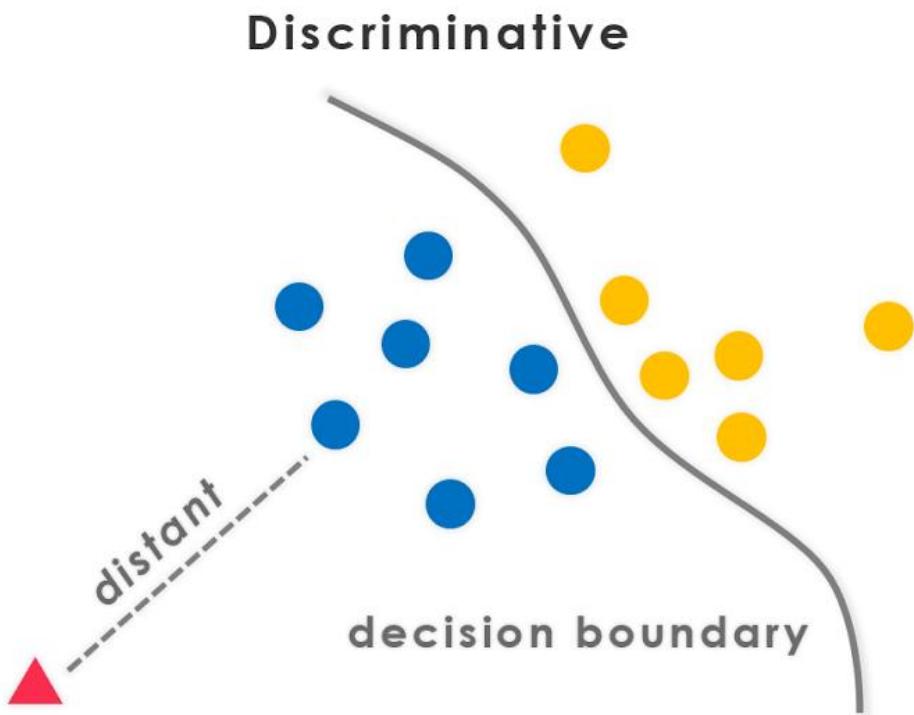
ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



Introduction

Discriminative and Generative models

MultDrone



Introduction

Discriminative and Generative models

In the past years we have focused mostly on building discriminative models for classification and regression

MultDrone



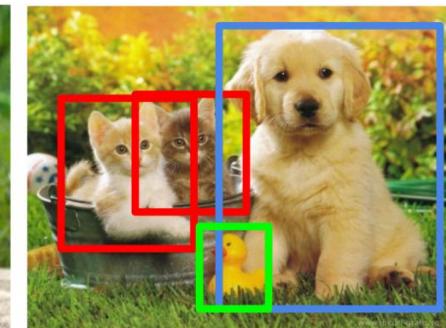
Classification



Classification + Localization



Object Detection



Instance Segmentation



CAT

CAT

CAT, DOG, DUCK

CAT, DOG, DUCK

Single object

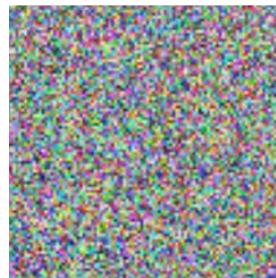
Multiple objects

Introduction

Discriminative and Generative models

Nowadays, we are focusing also on generative models for generating artificially realistic data

Noise $\sim N(0,1)$



Generative Model



Adversarial Examples

Local Generalization in Computer Vision

MultDrone



Slight pixel changes should not affect the decision of a model

There is an image and its ground truth class: x, y

A classification model f predicts the class of x : $\hat{y} = f(x)$

The prediction \hat{y} is the same as the ground truth: $y = \hat{y}$

There is an image x_p which is x perturbated by p : $x_p = x + p$

The distance of the two images is restricted by threshold e : $d(x, x_p) \leq e$

The threshold e is positive and small for imperceptibility changes

The classification model classifies the image x_p same as x : $\hat{y}_p = \hat{y}$

Adversarial Examples

Local Generalization in Computer Vision



Slight pixel changes should not affect the decision of a model



Adversarial Examples

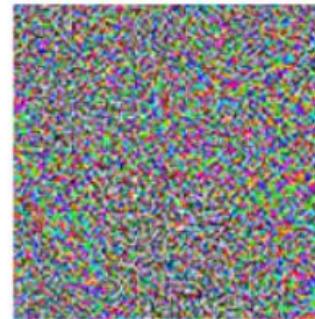
What exactly are these?



“panda”

57.7% confidence

+



=



“gibbon”

99.3% confidence

MultDrone



- Examples where the Local Generalization does not apply
- Perturbation: Optimal direction to move all the pixels so that the model will do a mistake
- Most models fail to work (LR, Softmax Regression, SVM, k-NN, Decision Trees, Neural Nets, Ensembles)

Adversarial Examples

Important related papers

Fooling a spam filter

“Adversarial Classification”
Nilesh Dalvi et al, 2004

Fooling neural networks

“Evasion Attacks Against Machine Learning at Test Time”
Battista Biggio et al, 2013

Fooling ImageNet classifiers

“Intriguing properties of neural networks”
Christian Szegedy et al, 2013

Cheap, closed form attack

“Explaining and Harnessing Adversarial Examples”
Ian J. Goodfellow et al, 2014

MultDrone



Adversarial Examples

The big question

Why adversarial examples exist and how is it feasible a model that is not overfitted and has high test/validation accuracy not to be functional in adversarial examples which are very similar to the original data?



MultDrone



Morpheus to Neo: "I imagine that right now you're feeling a bit like Alice, tumbling down the rabbit hole.", The Matrix (1999)

Adversarial Examples

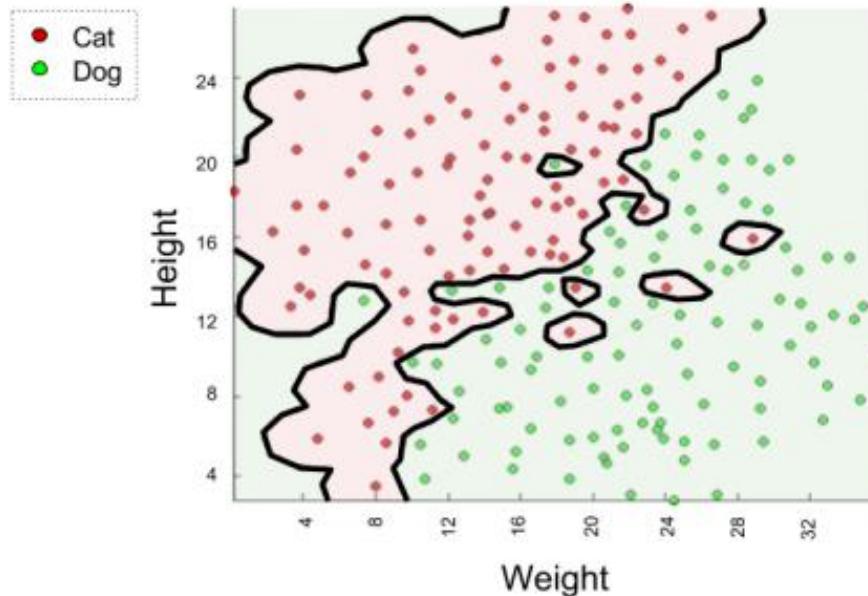
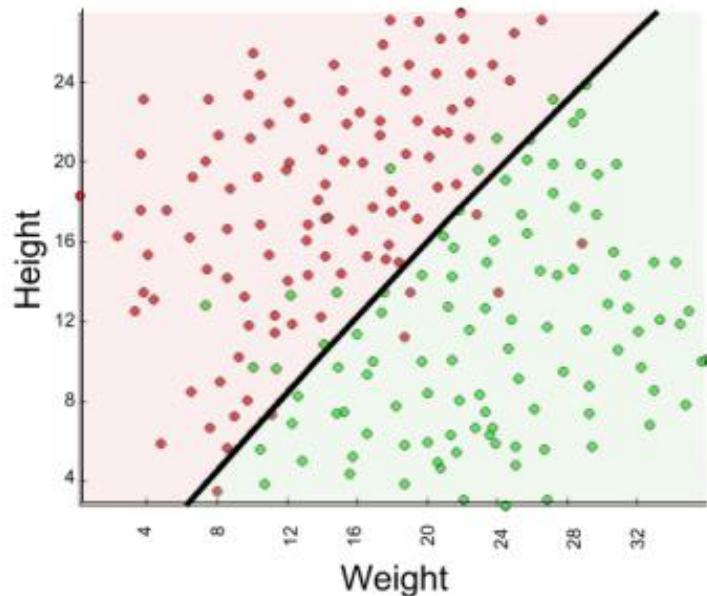
Why do they exist?

MultDrone



Is it due to overfitting?

- a) an overfitted model is sensitive to small input changes since it learns the idiosyncrasies of input

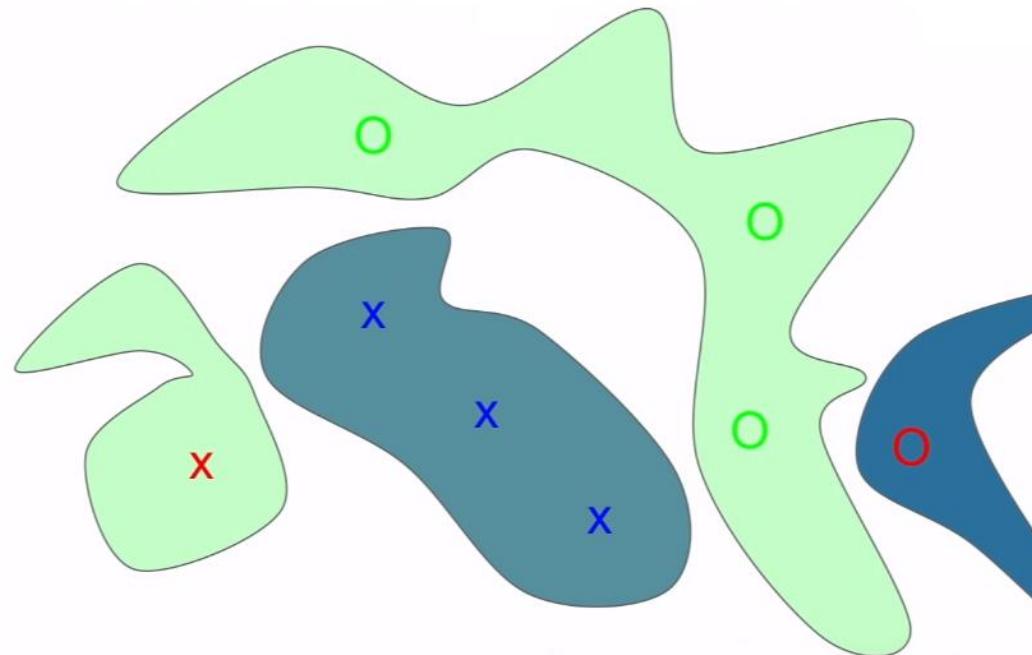


Adversarial Examples

Why do they exist?

Is it due to overfitting?

b) an overfitted model assigns some blobs of probabilities mass in unseen places of input



MultDrone



Adversarial Examples

Why do they exist?

MultDrone



If adversarial examples exist due to (a) and (b) then these are just the result of bad luck and should not be transferable between different model settings since different models give different probabilities in unseen places of input.

However, the adversarial examples are transferable between different model settings and that seems that is something more systematic than random.

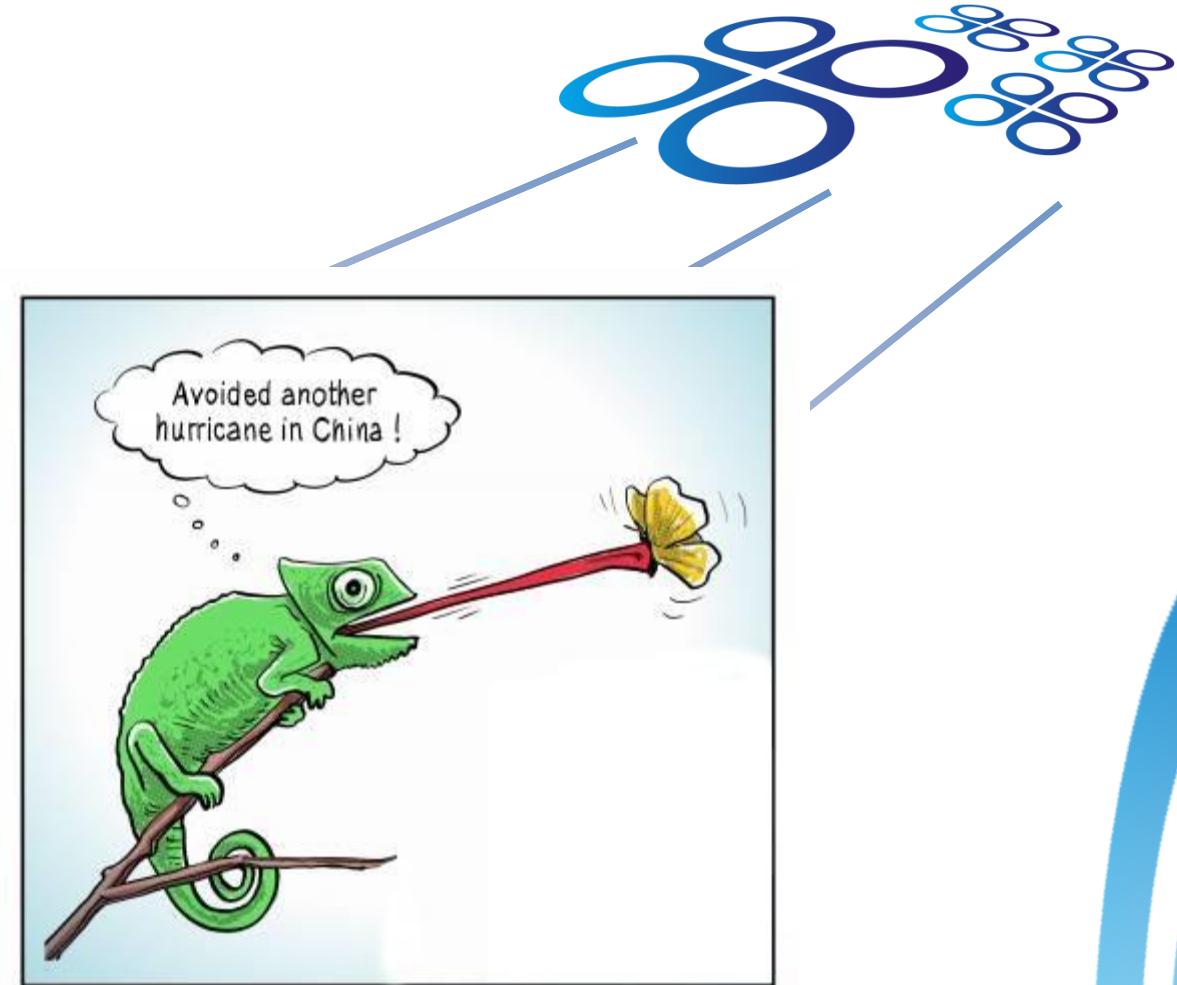
Adversarial Examples

Why do they exist?

Latest research supports that it is due to high-dimensional input and linearity of models

Many small pixel changes in high-dimensional input lead all together to a huge negative side effect

MultDrone

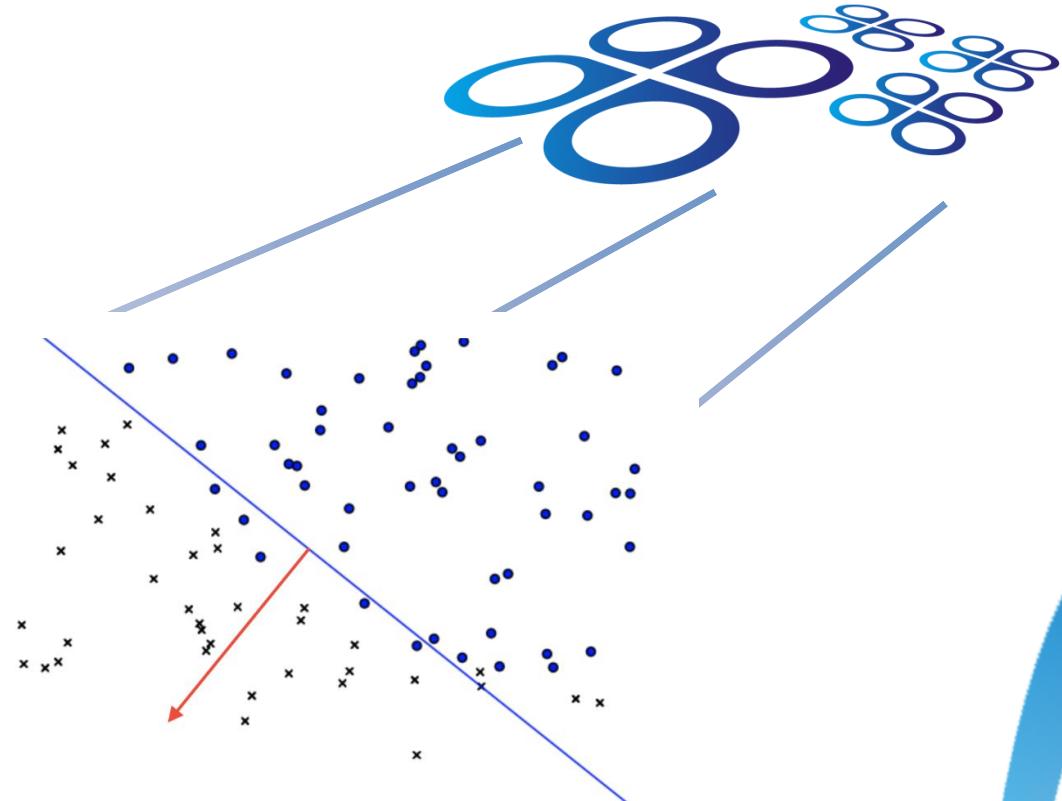


Adversarial Examples

Why do they exist?

Linear models are quite pathological outside of the region where training data is concentrated (initial experiments with shallow linear models shown that this affects greatly the $w^T x$ calculation and leads to wrong misclassification)

MultDrone



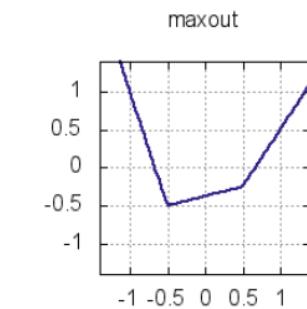
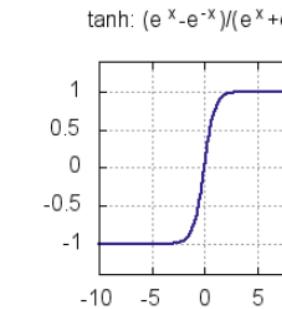
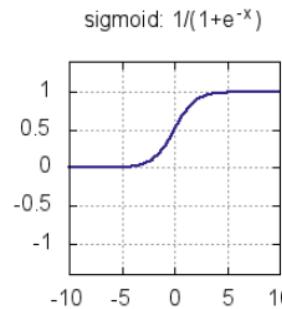
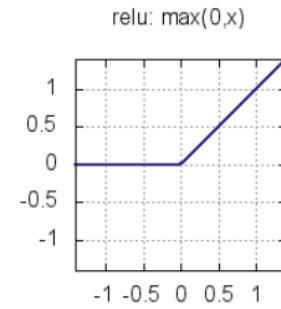
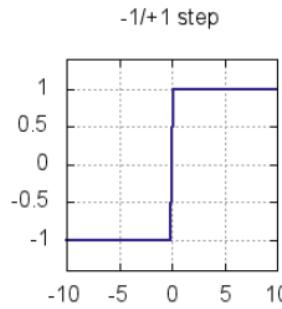
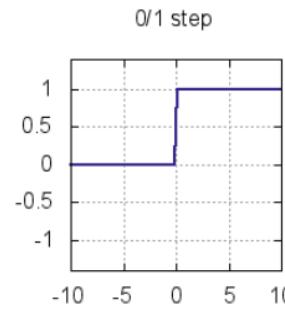
Adversarial Examples

Why do they exist?

MultDrone



Why adversarial examples exist also in non-linear models (e.g. Neural Networks)? Although by definition these are non-linear, are designed knowingly to be piecewise linear. Nowadays, state-of-the-art deep models have various piecewise linear elements.



“The fault is not in our stars, but in our activation functions, for they are piecewise linear”

“Explaining and Harnessing Adversarial Examples”, 2014
Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy



Adversarial Examples

Where does this lead us?

MultDrone



Achilles Heel



Potemkin Village

Adversarial Examples

How useful they are?

MultDrone



Adversarial examples can be used:

- to evaluate the robustness of a model
- to attack a model in order to break it
- in model training as a defensive regularization technique
- concealment / misclassification (e.g: de-detection / de-identification)

Adversarial Examples

Multiple ways to create them!



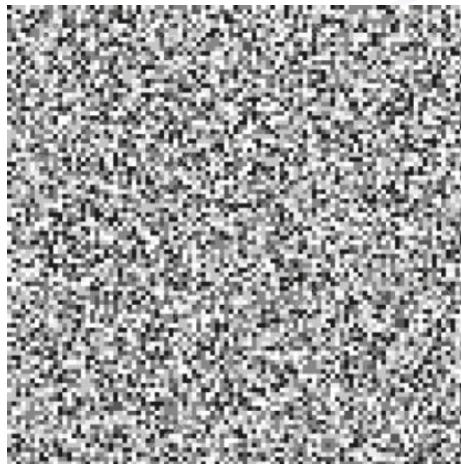
- Random / Non-random input
- Attack Method Frequency
- Adversarial Specificity
- Perturbation Scope
- Perturbation Measurement
- Realistic / Non-realistic Output
- Adversarial Falsification
- Adversary Knowledge
- Attack Method

Adversarial Examples

Random / Non-random input



- **Random:** From uniform, Gaussian (normal) or other probability distribution
- **Non-random:** Any example from the training / validation / test data set



Adversarial Examples

Attack Method Frequency

- **One-time:** Take only one time to optimize adversarial examples
- **Iterative:** Multiple times to update the adversarial examples

Adversarial Specificity

- **Targeted:** Find an input that is misclassified as a specific label
- **Non-targeted:** Find an input that is misclassified in a label just different than the ground truth

MultDrone



Adversarial Examples

Targeted / Non-targeted Definition

MultDrone



There is a sample and its ground truth class: x, y

A classification model f predicts the class of x : $\hat{y} = f(x)$

The prediction \hat{y} is the same as the ground truth: $y = \hat{y}$

There is a sample x_p which is x perturbated by p : $x_p = x + p$

The distance of the two samples is restricted by threshold e : $d(x, x_p) \leq e$

The threshold e is positive and small for imperceptibility changes

The classification model classifies the perturbated sample x_p : $\hat{y}_p = f(x_p)$

Non-targeted adversarial example constraint: $\hat{y}_p \neq \hat{y}$

Targeted adversarial example constraint: $\hat{y}_p = y_{target}$

Adversarial Examples

Perturbation Scope

- **Individual:** Each image has its own individual perturbation



x
“panda”
57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

MultDrone



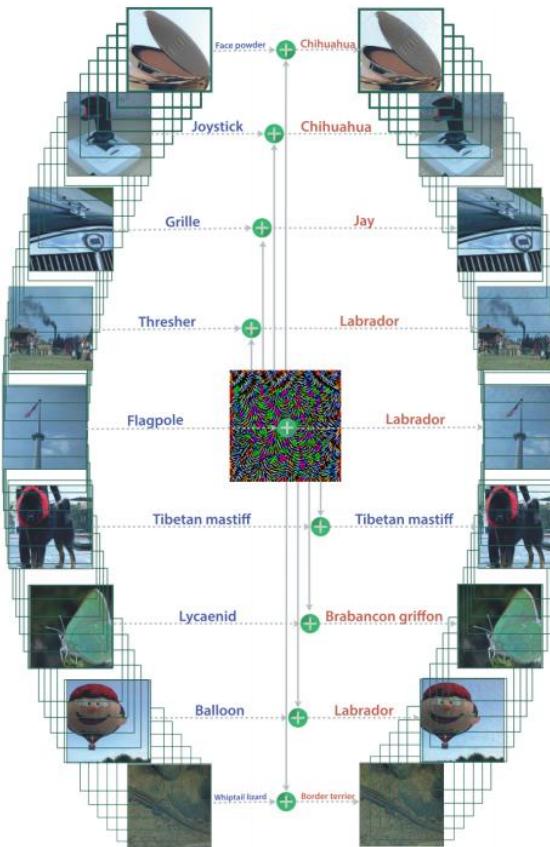
Adversarial Examples

Perturbation Scope

MultDrone



- **Universal:** A universal perturbation for the whole dataset



Adversarial Examples

Perturbation Measurement



$$\| x \|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

Use an ℓ_p -norm to calculate the magnitude of the adversarial perturbation:

- ℓ_0 -norm (The total number of non-zero values in the vector)
- ℓ_1 -norm (The sum of the absolute values in the vector)
- ℓ_2 -norm (The square root of the sum of the squares of the values in the vector)
- ℓ_∞ -norm (The maximum of the absolute values in the vector)

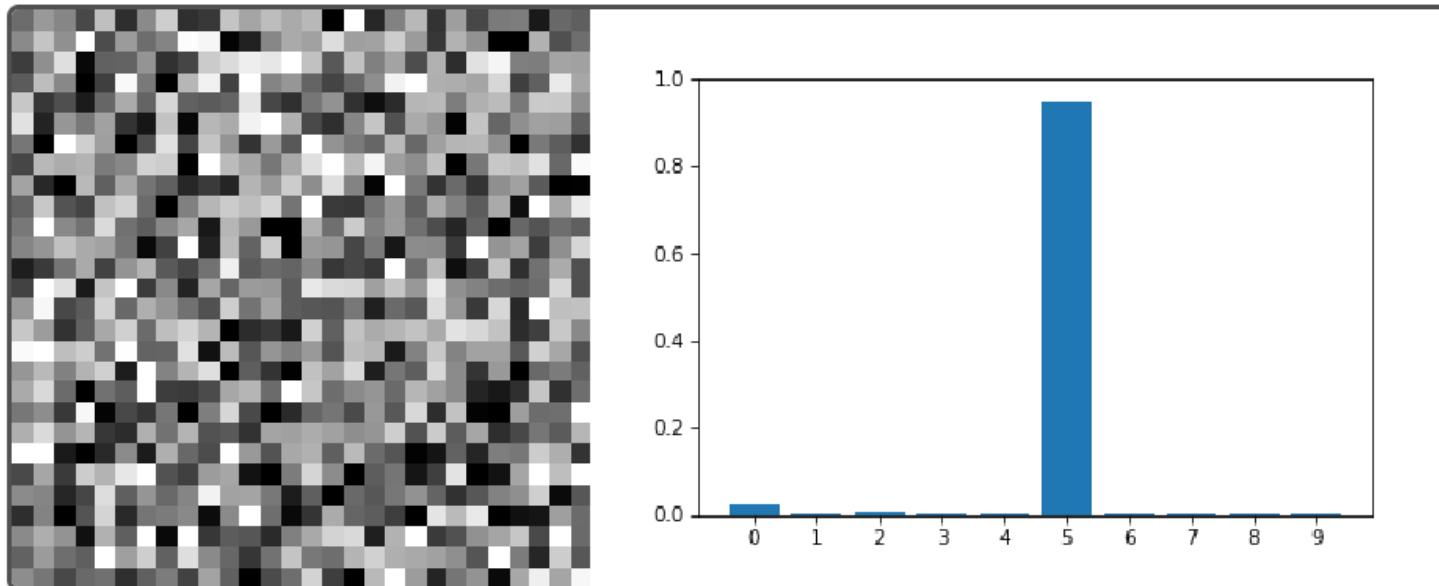
Adversarial Examples

Realistic / Non-realistic Output

MultDrone



- **Non-realistic:** Any adversarial input that fools the model



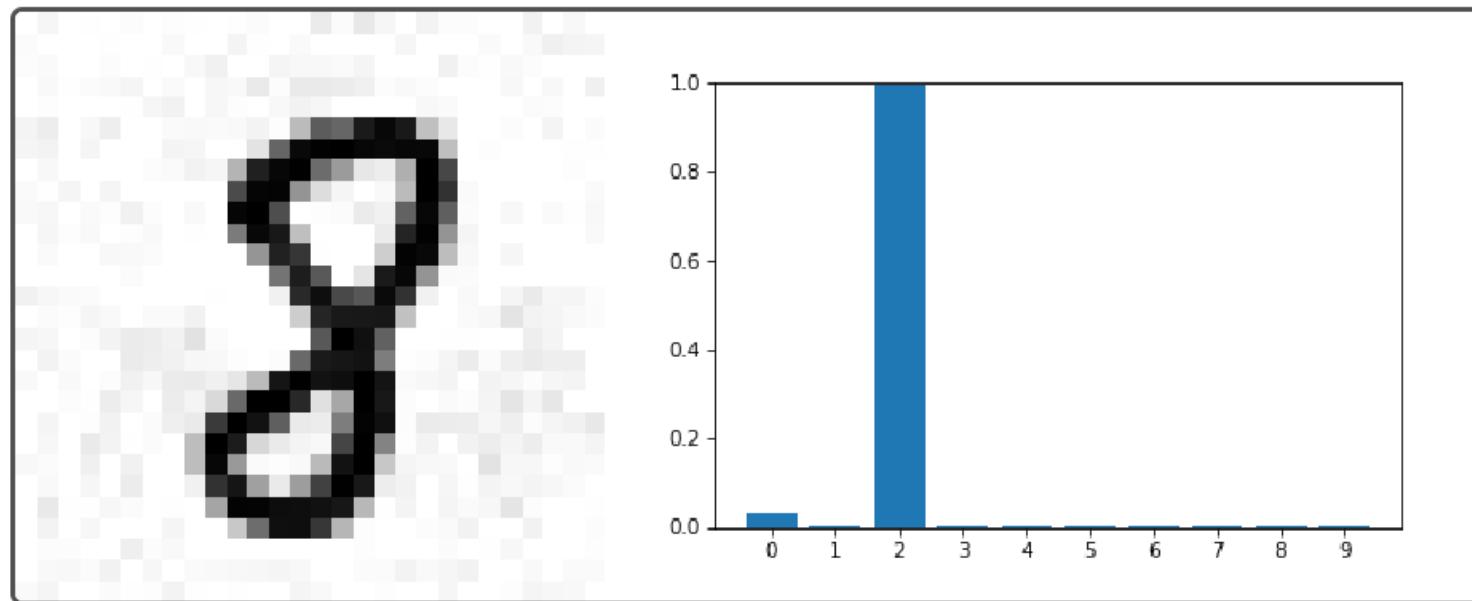
Adversarial Examples

Realistic / Non-realistic Output

MultDrone



- **Realistic:** An adversarial input that is close to a target example

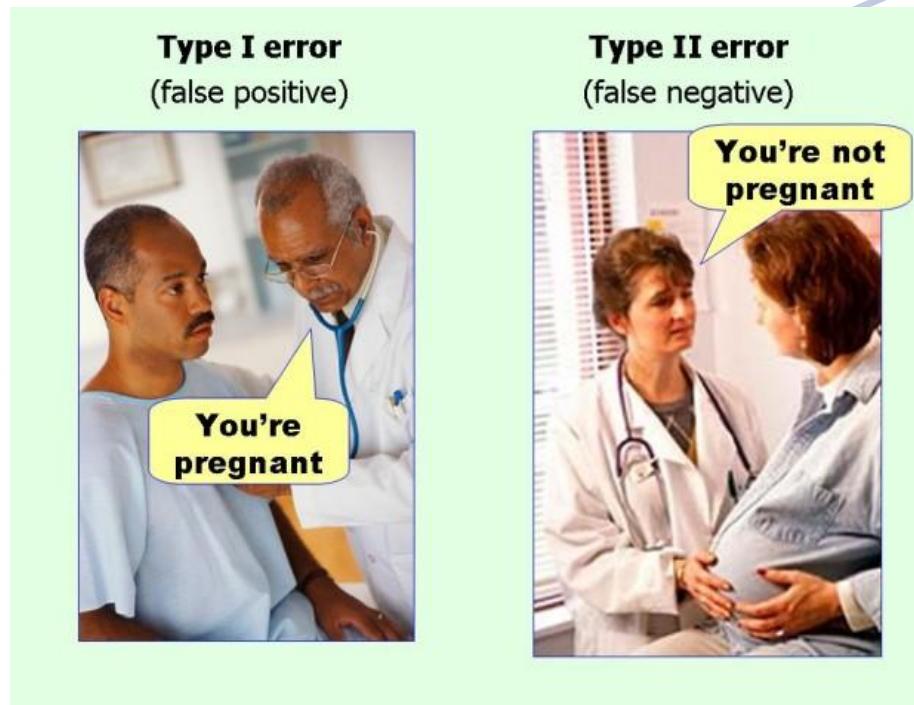


Adversarial Examples

Adversarial Falsification

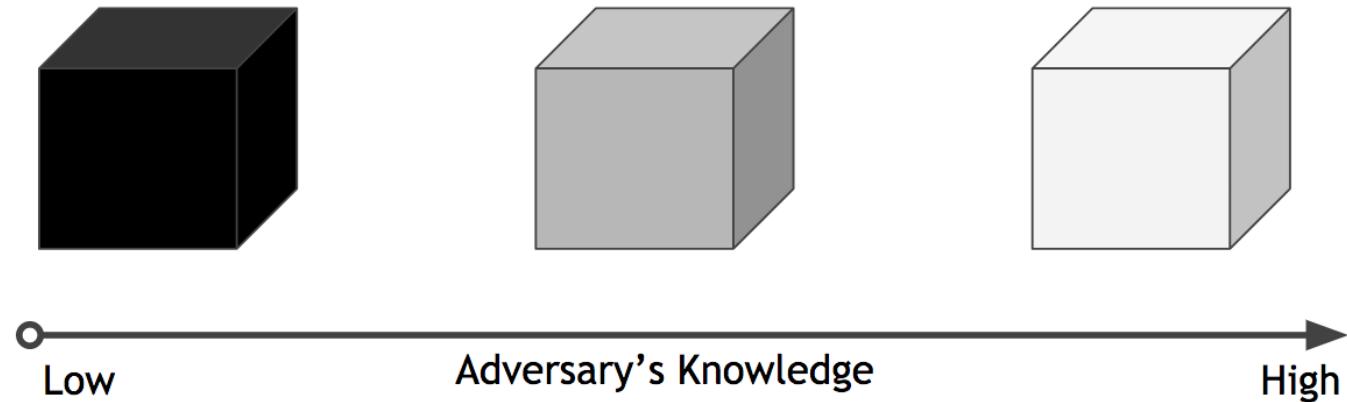


- **Type 1 error:** Negative sample classified as positive
- **Type 2 error:** Positive sample classified as negative



Adversarial Examples

Adversary Knowledge



MultDrone



- **Black-box:** Zero knowledge about the model to attack (knowing only the final classification)
- **Grey-box:** Limited knowledge about the model to attack (something between Black-box and White-box)
- **White-box:** Full knowledge about the model to attack (architecture, parameters, dataset, etc)

Adversarial Examples

Attack Methods



L-BFGS: Limited-memory BFGS

Intriguing properties of neural networks

Christian Szegedy et al, 2013

$$\text{minimize} \|x - x^{adv}\|_2^2$$

such that

$$f(x^{adv}) = y_{target}$$

$$x^{adv} \in [0, 1]^n$$

Given an input image x , the method finds a targeted adversarial valid image x^{adv} that is similar to x under squared l_2 -norm and is labeled as y_{target} by the classifier f . However, this problem can be very difficult to solve. So, instead solve the following problem:

$$\text{minimize } c \|x - x^{adv}\|_2^2 + loss_f(x^{adv}, y_{target})$$

such that $x^{adv} \in [0, 1]^n$

Adversarial Examples

Attack Methods



Let's get an idea with “Fast Gradient” Methods

Explaining and Harnessing Adversarial Examples

Ian J. Goodfellow et al, 2014



x
“panda”
57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

=

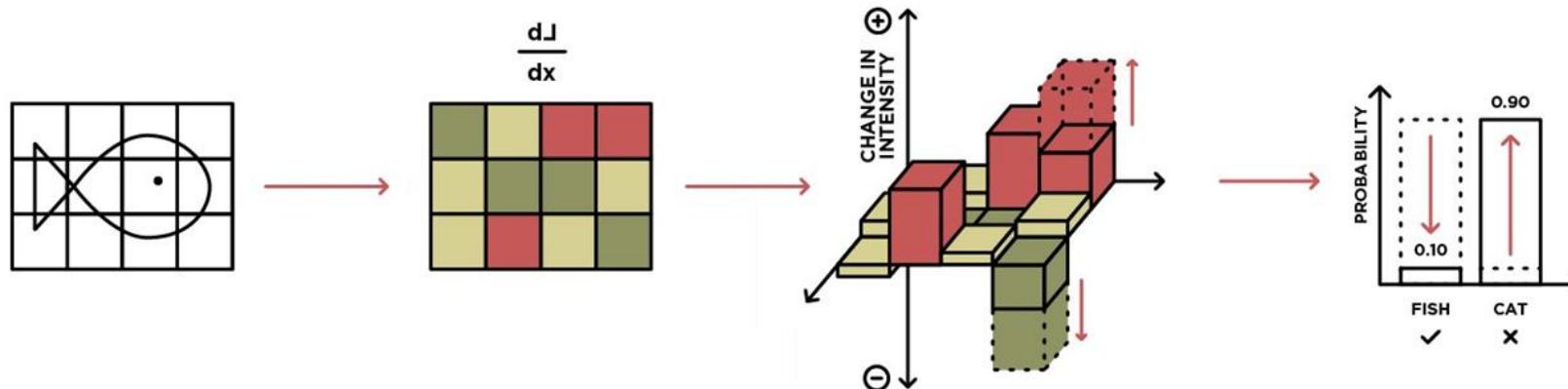


$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

Adversarial Examples

Attack Methods

- Use gradients of loss function w.r.t. input
- Gradient descent for targeted or ascent for non-targeted
- Very effective for the domain of image
- Fast and easy to compute
- ‘ ϵ ’ controls the size of the change (should be a small value)
- Can be used for run-time adversarial training
- For NNs the $\nabla_x J(\theta, x, y)$ can be calculated with backpropagation



Adversarial Examples

Attack Methods

FGSM : Fast Gradient Sign Method

- Non-targeted: One-shot method with gradient ascent
- Targeted: One-shot method with gradient descent

$$x^{adv} = x + \varepsilon \cdot \text{sign}(\nabla_x J(x, y_{true}))$$

where

x – Clean Input Image

x^{adv} – Adversarial Image

J – Loss Function

y_{true} – Model Output for x

ε – Tunable Parameter

$$x^{adv} = x - \varepsilon \cdot \text{sign}(\nabla_x J(x, y_{target}))$$

where

x – Clean Input Image

x^{adv} – Adversarial Image

J – Loss Function

y_{target} – Target Label

ε – Tunable Parameter

MultDrone



Adversarial Examples

Attack Methods

MultDrone



I-FGSM : Iterative Fast Gradient Sign Method

- Non-targeted: Iterative method with gradient ascent
- Targeted: Iterative method with gradient descent

$$x_0^{adv} = x, \quad x_{N+1}^{adv} = Clip_{x,\varepsilon}\{ x_N^{adv} + \alpha \cdot sign(\nabla_x J(x_N^{adv}, y_{true})) \}$$

where

x – Clean Input Image

x_i^{adv} – Adversarial Image at i^{th} step

J – Loss Function

y_{true} – Model Output for x

ε – Tunable Parameter

α – Step Size

$$x_0^{adv} = x, \quad x_{N+1}^{adv} = Clip_{x,\varepsilon}\{ x_N^{adv} - \alpha \cdot sign(\nabla_x J(x_N^{adv}, y_{target})) \}$$

where

x – Clean Input Image

x_i^{adv} – Adversarial Image at i^{th} step

J – Loss Function

y_{target} – Target Label

ε – Tunable Parameter

α – Step Size

Adversarial Examples

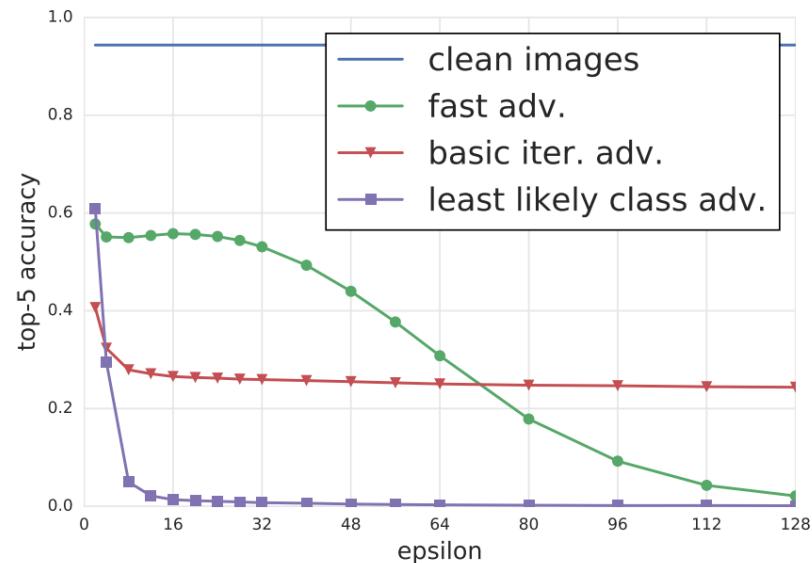
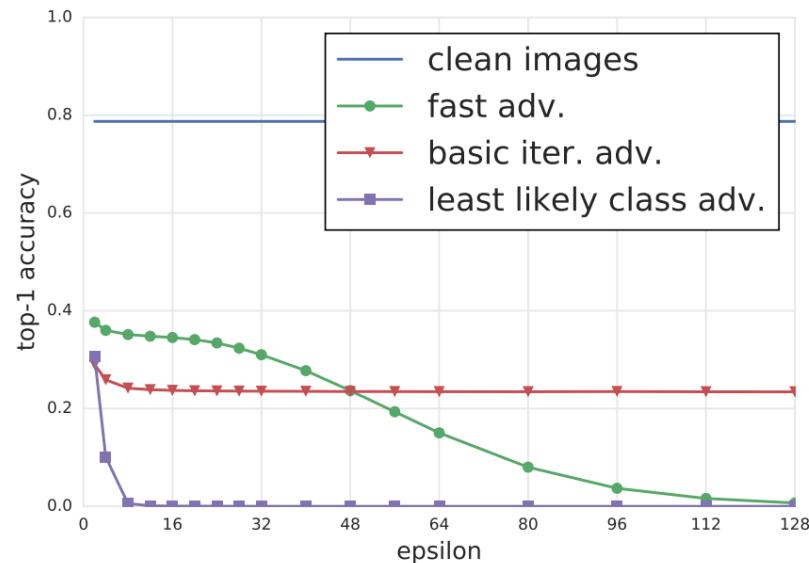
Attack Methods

MultDrone



Iterative Least-Likely Class Method

Using I-FGSM with $y_{\text{target}} = \operatorname{argmin}_y \{ p(y|x) \}$



Top-1 and Top-5 accuracy of Inception V3 under attack by different adversarial methods and different ‘e’ compared to unmodified images from the dataset. The accuracy was computed on all 50,000 validation images from the ImageNet dataset.

Adversarial Examples

Attack Methods

What would be more crazy? Of course, “Single Pixel” attack!

One pixel attack for fooling deep neural networks

Jiawei Su et al, 2017

MultDrone



ONE PIXEL ATTACK

Adversarial Examples

Attack Methods

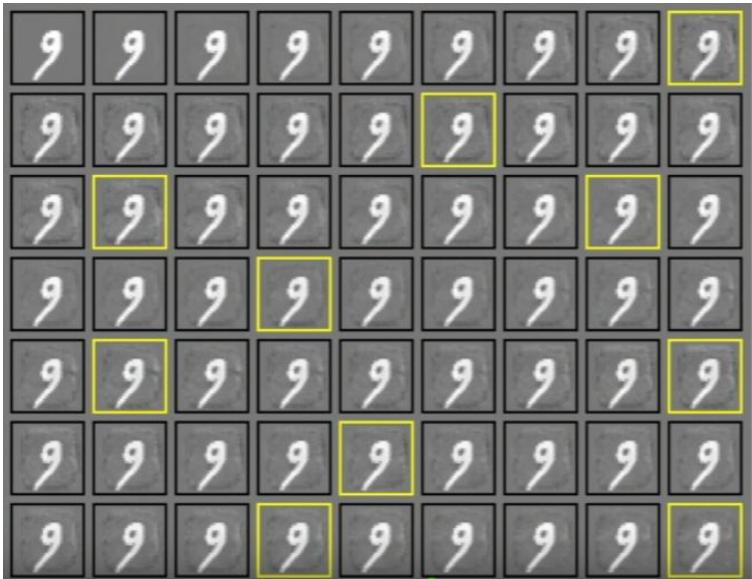
- L-BFGS (Szegedy et al, 2013)
- Fast Gradient Sign Method (Goodfellow et al, 2015)
- Basic Iterative Method (Kurakin et al, 2016)
- One Pixel (Jiawei Su et al, 2018)
- DeepFool (Moosavi-Dezfooli et al, 2015)
- Jacobian-based Saliency Map Method (Papernot et al, 2016)
- SPSA (Uesato et al, 2018)
- Carlini Wagner L2 (Carlini and Wagner, 2016)
- Feature Adversaries (Sabour et al, 2016)
- Elastic Net Method (Chen et al, 2017)
- Virtual Adversarial Method (Miyato et al, 2016)
- The Momentum Iterative Method (Dong et al, 2017)
- Projected Gradient Descent (Madry et al, 2017)

MultDrone



Adversarial Examples

Softmax regression (MNIST)



- The first image in the table is a random sample from the training data set
- The rest images are generated from top-to-bottom and left-to-right sequentially using L-BFGS
- Yellow Boxes: cases with $\geq 99\%$ probability where the model thinks is 0, 1, ..., 9
- The animal species recognize all above images as 9 but the model seems different

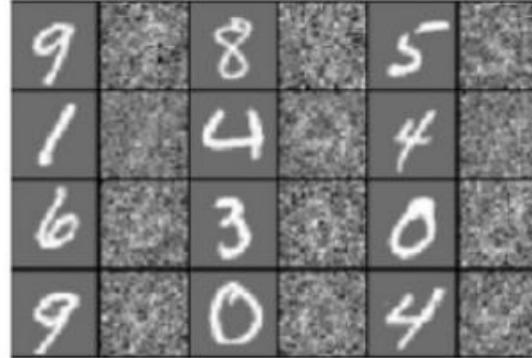
Adversarial Examples

Softmax regression (MNIST)

MultDrone



(a)



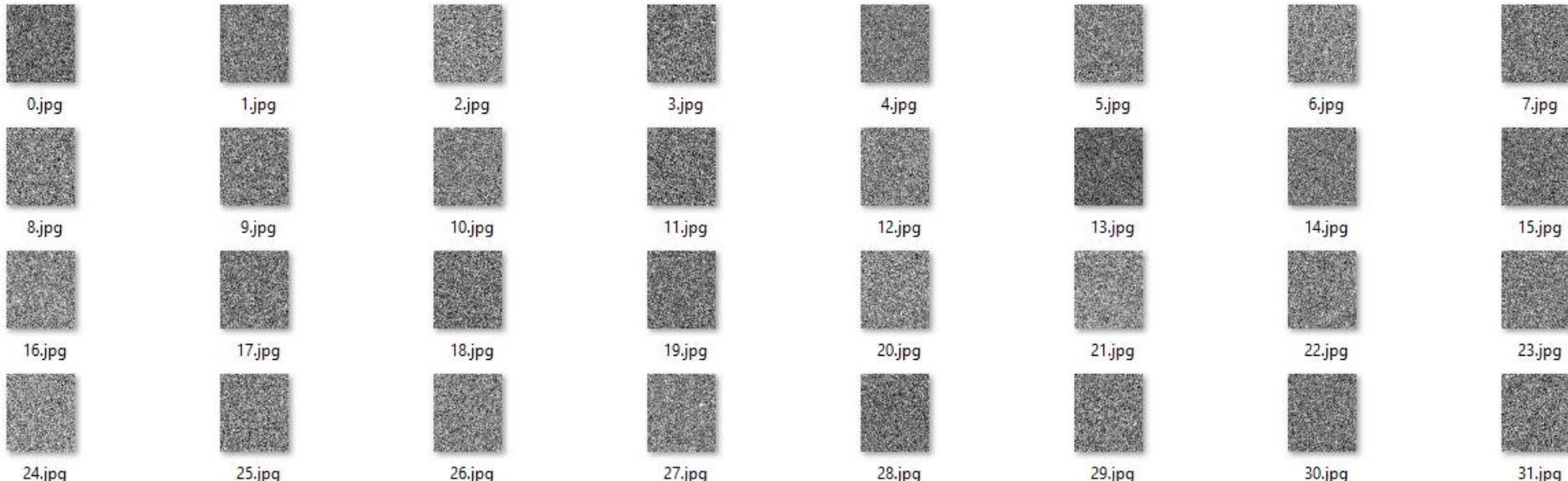
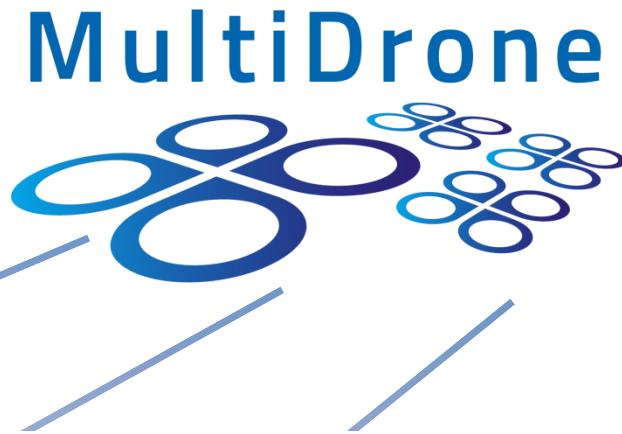
(b)

- Images in table's odd columns are random samples from the training data set
- Images in table's even columns correspond to distorted counterparts
- (a) We successfully classify images with adversarial perturbation but models fail to do it
- (b) We fail to classify images with Gaussian noise added but models do it successfully

Adversarial Examples

MLP w/ Ext. Yale Face Database B

An example of face de-identification with an MLP model
Successfully target to any class with non-realistic images

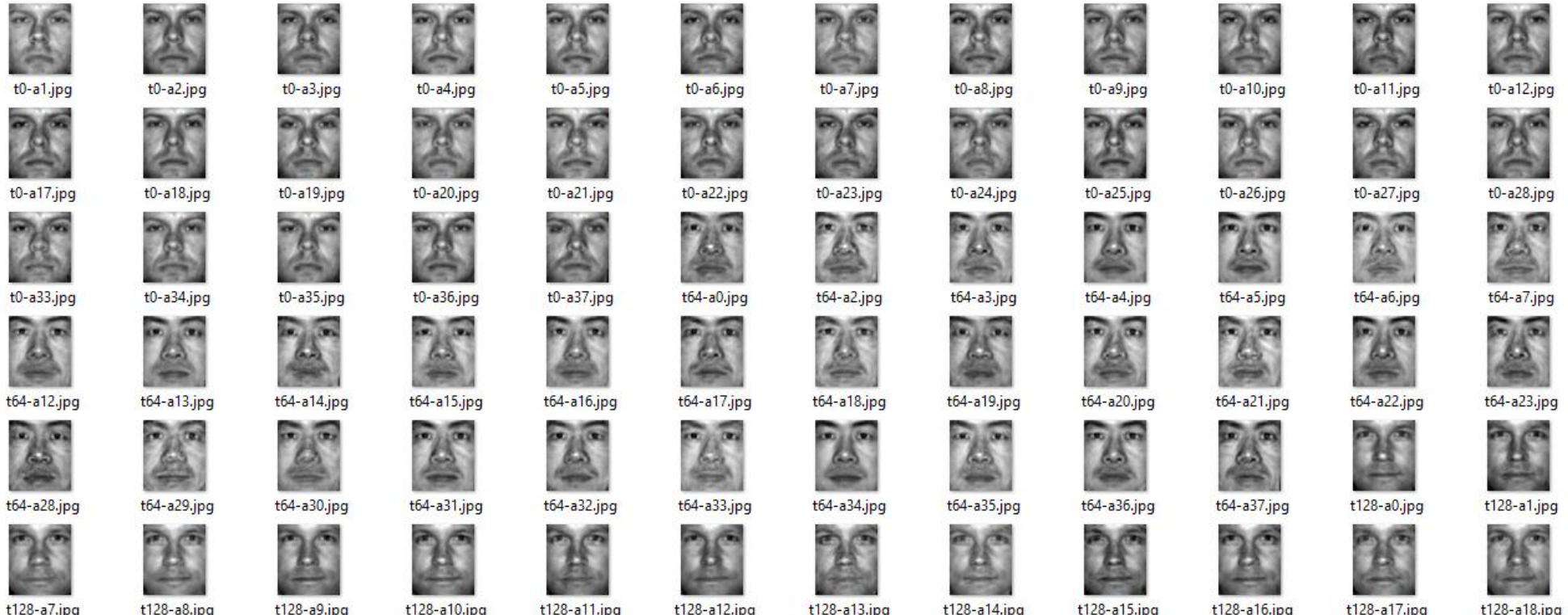


Adversarial Examples

MLP w/ Ext. Yale Face Database B

An example of face de-identification with an MLP model
Successfully target to any class with realistic images

MultDrone



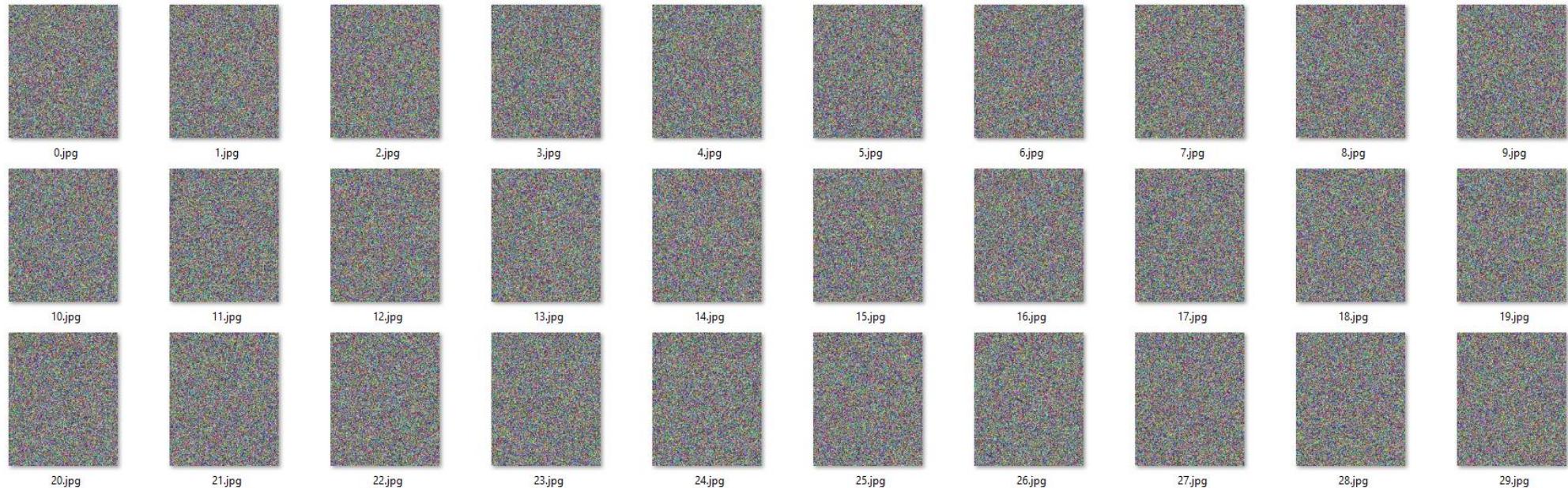
Adversarial Examples

CNN w/ CelebA

MultDrone



An example of face de-identification with a CNN model
Successfully target to any class with non-realistic images



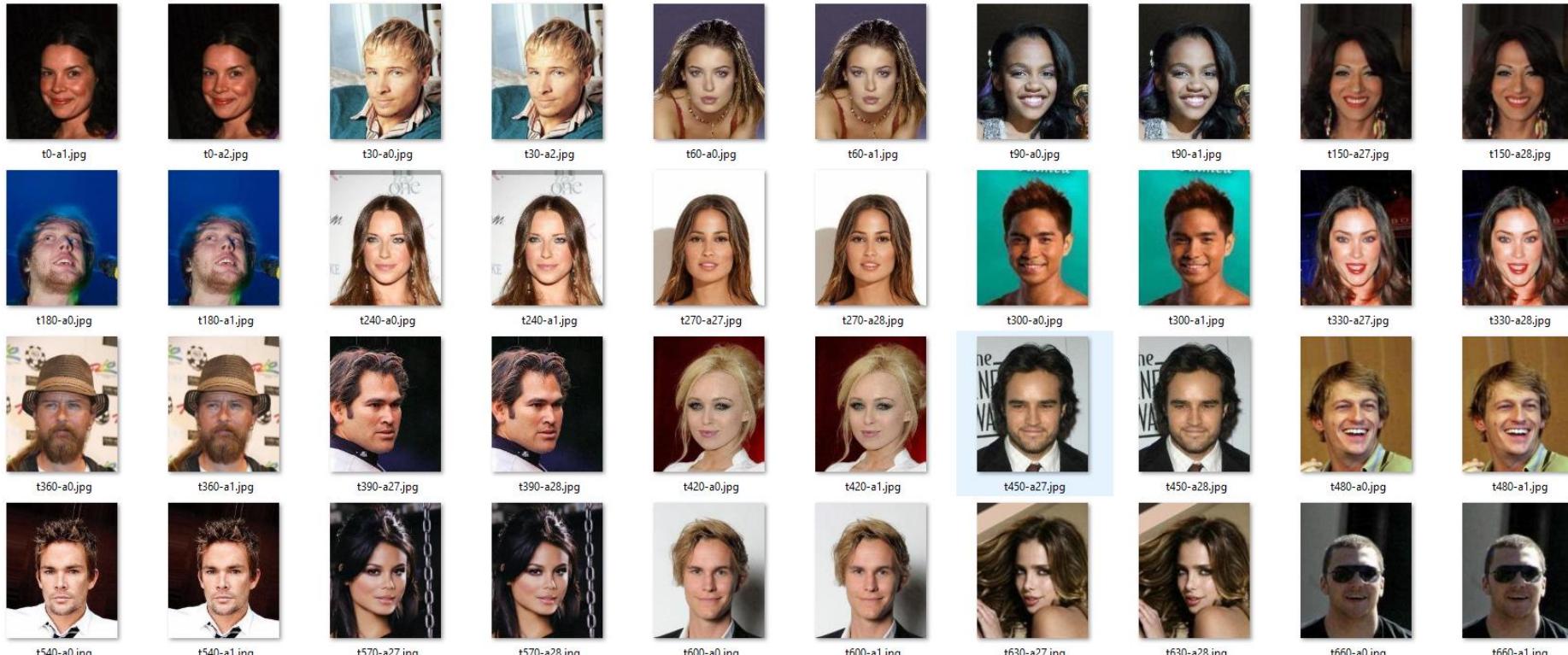
Adversarial Examples

CNN w/ CelebA

MultDrone



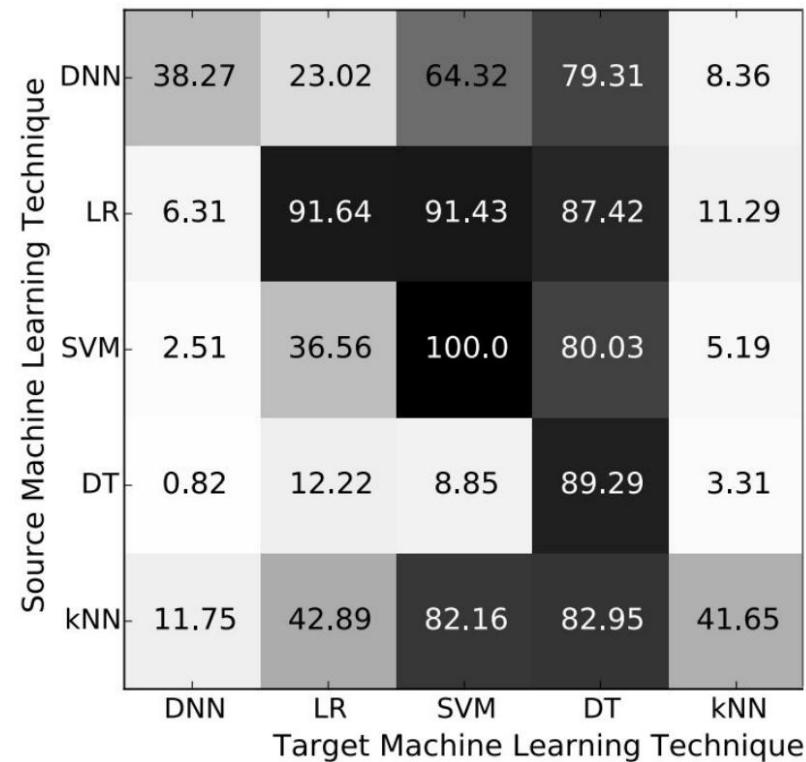
An example of face de-identification with a CNN model
Successfully target to any class with realistic images



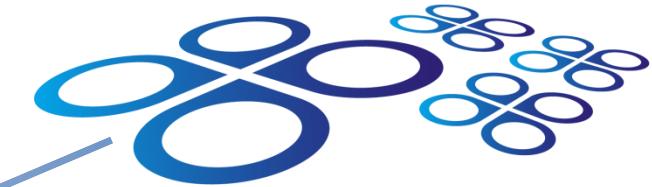
Adversarial Examples

They are transferable!

- Cross model generalization



MultDrone



Adversarial Examples

They are transferable!

- Cross training-set generalization

| Source LR | A | B | C | D | E |
|-----------|----|----|----|----|----|
| A | 98 | 95 | 95 | 95 | 95 |
| B | 95 | 98 | 95 | 95 | 94 |
| C | 94 | 94 | 98 | 95 | 95 |
| D | 94 | 95 | 95 | 98 | 95 |
| E | 95 | 95 | 95 | 95 | 98 |

| Source SVM | A | B | C | D | E |
|------------|----|----|----|----|----|
| A | 99 | 41 | 38 | 40 | 41 |
| B | 34 | 99 | 32 | 46 | 34 |
| C | 36 | 41 | 99 | 38 | 45 |
| D | 37 | 43 | 37 | 99 | 38 |
| E | 39 | 37 | 47 | 37 | 99 |

| Source DNN | A | B | C | D | E |
|------------|----|----|----|----|----|
| A | 81 | 67 | 66 | 49 | 54 |
| B | 71 | 86 | 75 | 53 | 58 |
| C | 67 | 70 | 84 | 52 | 57 |
| D | 64 | 64 | 65 | 68 | 57 |
| E | 75 | 73 | 74 | 57 | 80 |

Strong

Weak

Intermediate

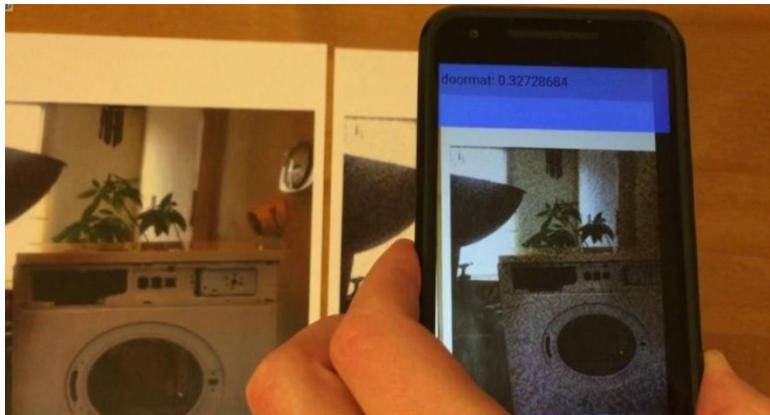
MultIDrone



Adversarial Examples

They are transferable!

- Can exist in the real world!



Adversarial examples in the physical world
Alexey Kurakin et al, 2017

“We used images taken from a cell-phone camera as an input to an Inception V3 image classification neural network. We showed that in such a set-up, a significant fraction of adversarial images crafted using the original network are misclassified even when fed to the classifier through the camera.”, Kurakin et al.

MultDrone



Adversarial Examples

Adversarial Training

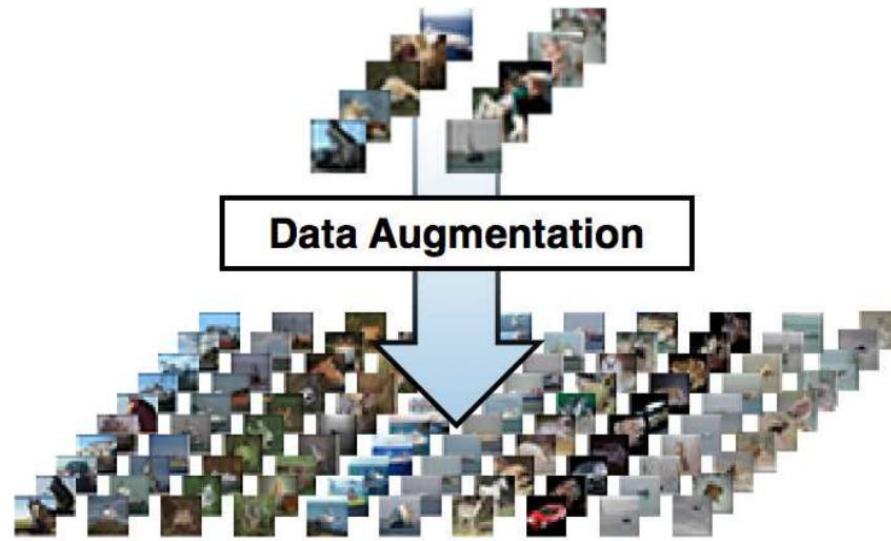
MultDrone



- The Neural Networks can approximate any function (Universal Approximation Theorem)
- We can train our Neural Networks with some adversarial examples to be more robust
- Adversarial training works as a regularization technique
- There is a need for sophisticated adversarial training algorithms with robustness by priori

Adversarial Examples

Adversarial Training



MultDrone



- Mixing training/validation/testing data sets with adversarial examples
- Same idea with data augmentation as a regularization technique

Adversarial Examples

Adversarial Training

MultDrone



- Using an adversarial loss function for single examples
- Generate fast adversarial examples at run-time with FGSM
- ‘c’ : hyperparameter coefficient to balance the weight of adversarial learning

$$J'(\theta, x, y) = c \times J(\theta, x, y) + (1 - c) \times J(\theta, x + p, y)$$

$$p = e \times \text{sign}(\nabla_x J(\theta, x, y))$$

Adversarial Examples

Adversarial Training

2016, “Adversarial Machine Learning at Scale”
Alexey Kurakin, Ian Goodfellow, Samy Bengio

Using an adversarial loss function for large scale models and datasets

MultDrone



Algorithm 1 Adversarial training of network N .

Size of the training minibatch is m . Number of adversarial images in the minibatch is k .

- 1: Randomly initialize network N
 - 2: **repeat**
 - 3: Read minibatch $B = \{X^1, \dots, X^m\}$ from training set
 - 4: Generate k adversarial examples $\{X_{adv}^1, \dots, X_{adv}^k\}$ from corresponding clean examples $\{X^1, \dots, X^k\}$ using current state of the network N
 - 5: Make new minibatch $B' = \{X_{adv}^1, \dots, X_{adv}^k, X^{k+1}, \dots, X^m\}$
 - 6: Do one training step of network N using minibatch B'
 - 7: **until** training converged
-

$$Loss = \frac{1}{(m - k) + \lambda k} \left(\sum_{i \in CLEAN} L(X_i | y_i) + \lambda \sum_{i \in ADV} L(X_i^{adv} | y_i) \right)$$

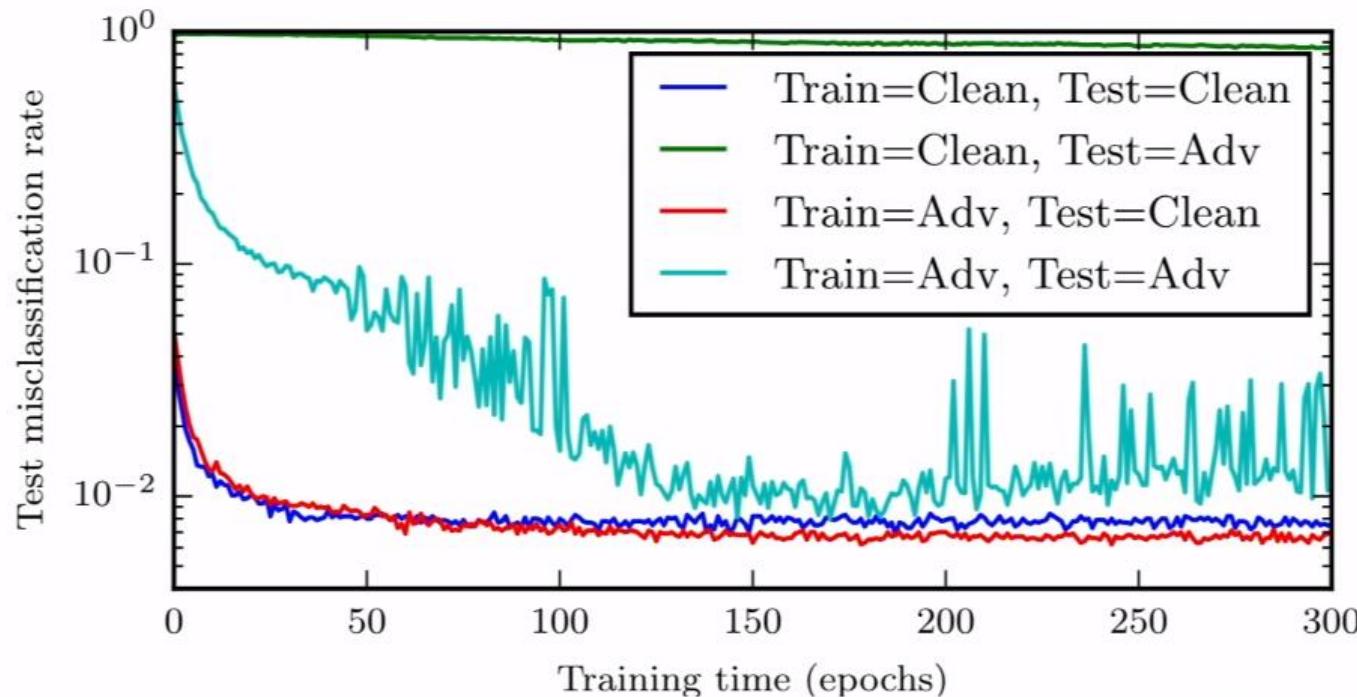
Adversarial Examples

Adversarial Training

MultDrone



Training on Adversarial Examples



Generative Adversarial Networks

What exactly are these?

MultDrone



- Discriminative models map rich, high dimensional data input to low dimensions output
- Nowadays, interest has turned quite a bit on generative models which do the opposite
- GAN: A framework for evolving adversarially a generative and a discriminative model
- Generative model as ‘Generator’ or ‘G’ / Discriminative model as ‘Discriminator’ or ‘D’

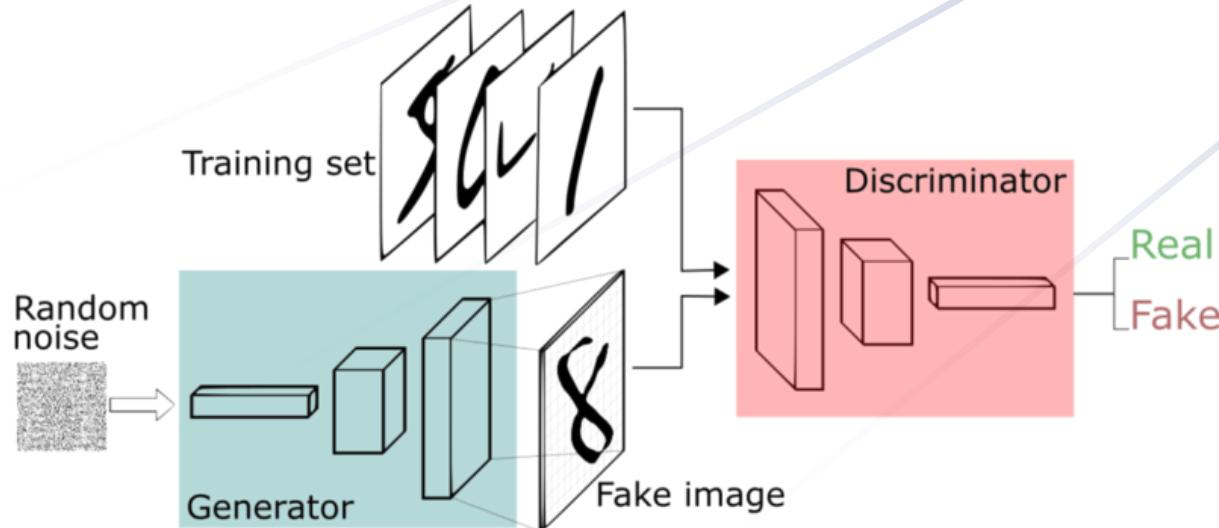
Generative Adversarial Networks

What exactly are these?

MultDrone



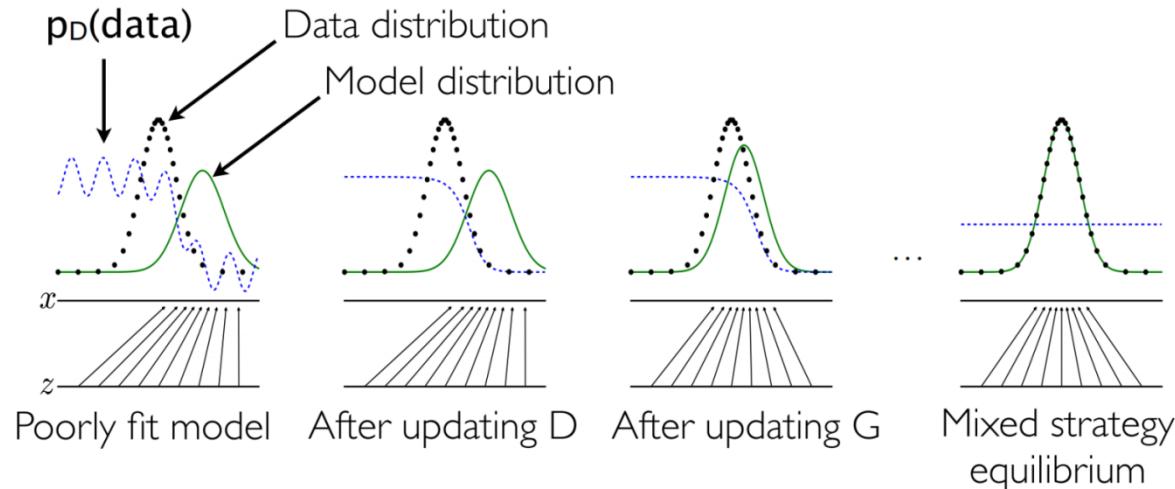
- 'G' and 'D' evolve through an adversarial process
- 'G' and 'D' are players in a two-player minimax game
- 'G' is trained to learn to generate fake data and approximate the real data distribution
- 'D' is trained to learn to separate the fake from the real data distribution



Generative Adversarial Networks

What exactly are these?

MultDrone

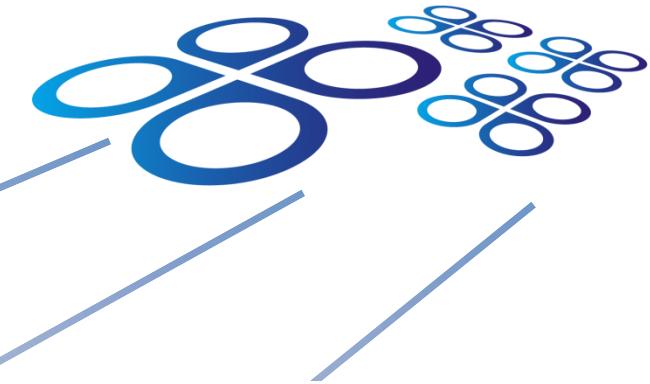


- 'G' samples random noise 'z' from a uniform / normal (Gaussian) distribution
- 'G' generates a fake sample with a single forward pass: $x = G(\Theta_G, z)$
- 'D' separates (as a binary classifier) fake from real samples: $y = D(\Theta_D, x)$
- Optimal solution: Nash Equilibrium where $y = 0.5$ probability (fake = real)

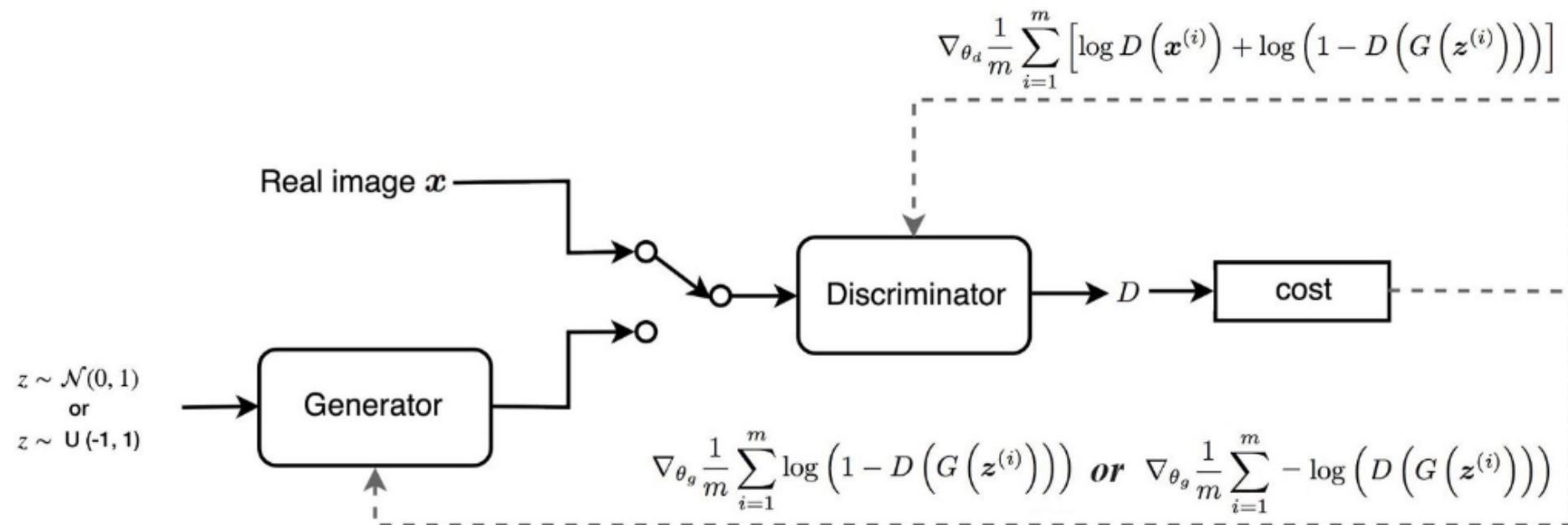
Generative Adversarial Networks

Adversarial Training

MultDrone



- ‘G’ and ‘D’ can be trained using backpropagation (in case of NNs)
- ‘G’ tries to **decrease** the correct classification of ‘D’ on the fake data
- ‘D’ tries to **increase** the correct classification for fake and real data



Generative Adversarial Networks

Adversarial Training

2014, “Generative Adversarial Networks”, Ian J. Goodfellow et al.



Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule.

Generative Adversarial Networks

Why GANs?



- First decent generative model that produces promising results
- Ability to approximate the probability distribution of input data through training
- Ability to create huge amounts of realistic data (any digital media)
- Instead of using data augmentation we can use 'G' to generate realistic samples
- Using it with unsupervised algorithms that need unlabeled data samples
- Ability to use convolution layers, autoencoders, regularization techniques
- 'G' is not trained on the real dataset as it gets only random input noise

Generative Adversarial Networks

Famous GANs variations

- Bidirectional GAN
- Conditional GAN
- Context-Conditional GAN
- CycleGAN
- Deep Convolutional GAN
- DualGAN
- InfoGAN
- Pix2Pix
- Semi-Supervised GAN
- Super-Resolution GAN



400+ named GANs variations exist today (GAN Zoo: <https://bit.ly/2sxGUHl>)

Generative Adversarial Networks

Applications (Image Synthesis)



2017, “Progressive Growing of GANs for Improved Quality, Stability, and Variation”, Tero Karras et al

Generative Adversarial Networks

Applications (Image Synthesis)

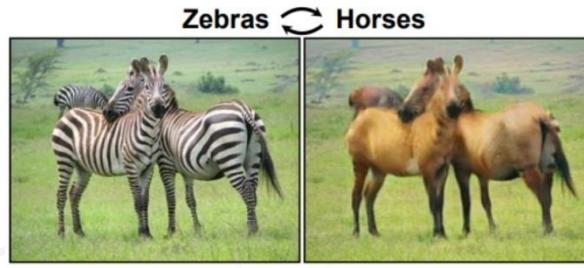
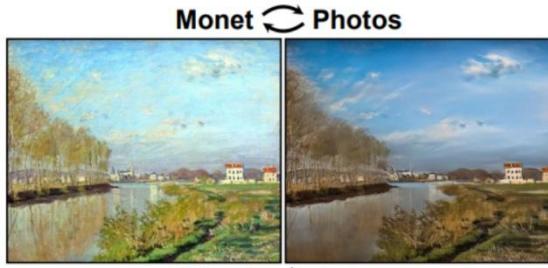


2017, "Towards the Automatic Anime Characters Creation with Generative Adversarial Networks", Yanghua Jin et al

Generative Adversarial Networks

Applications (Style Transfer)

MultDrone



2017, "Image to Image Translation using Cycle-Consistent Adversarial Neural Networks", Jun-Yan Zhu et al

Generative Adversarial Networks

Applications (Text-to-Image Synthesis)

This flower has a lot of small purple petals in a dome-like configuration



This flower is pink, white, and yellow in color, and has petals that are striped



This flower has petals that are dark pink with white edges and pink stamen



This bird is red and brown in color, with a stubby beak



The bird is short and stubby with yellow on its body



A bird with a medium orange bill white body gray wings and webbed feet



MultDrone



2017, "StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks", Han Zhang et al

Generative Adversarial Networks

Applications (Sketch-to-Image Synthesis)

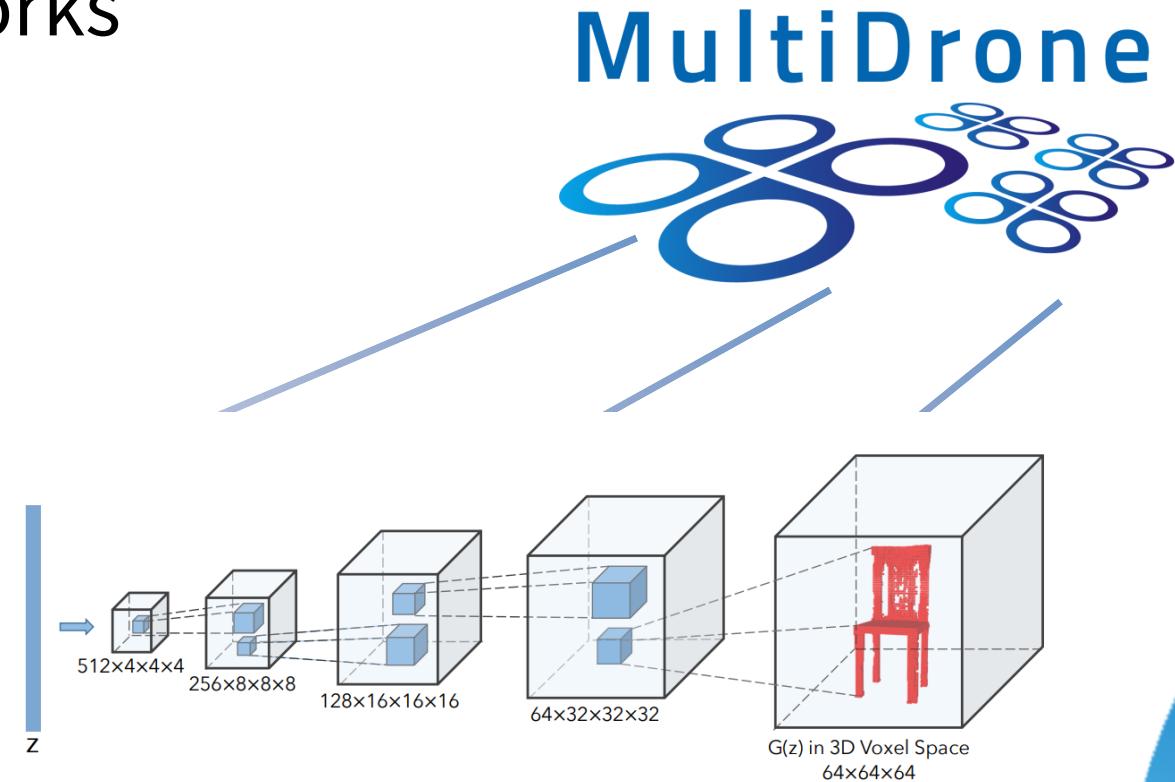
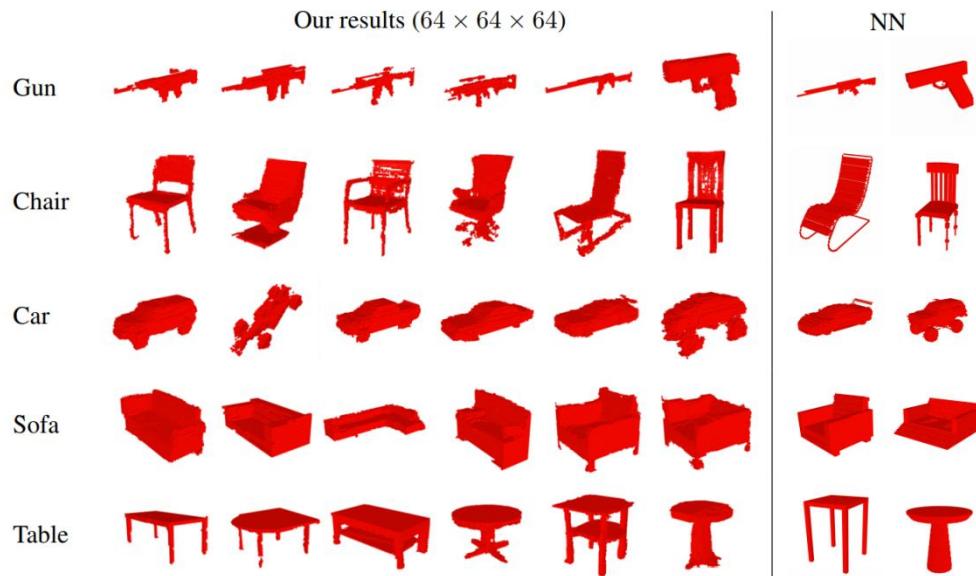
MultDrone



2016, “Image-to-Image Translation with Conditional Adversarial Networks”, Phillip Isola et al

Generative Adversarial Networks

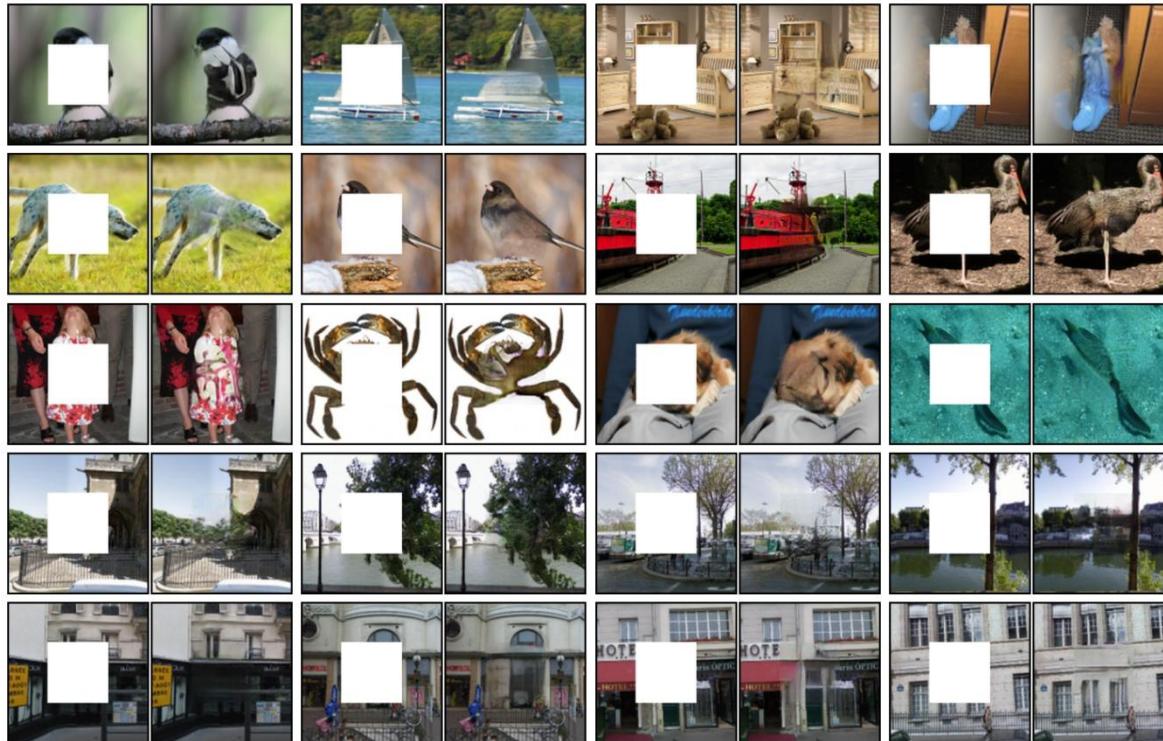
Applications (3D Object Synthesis)



2016, "Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling", Jiajun Wu et al

Generative Adversarial Networks

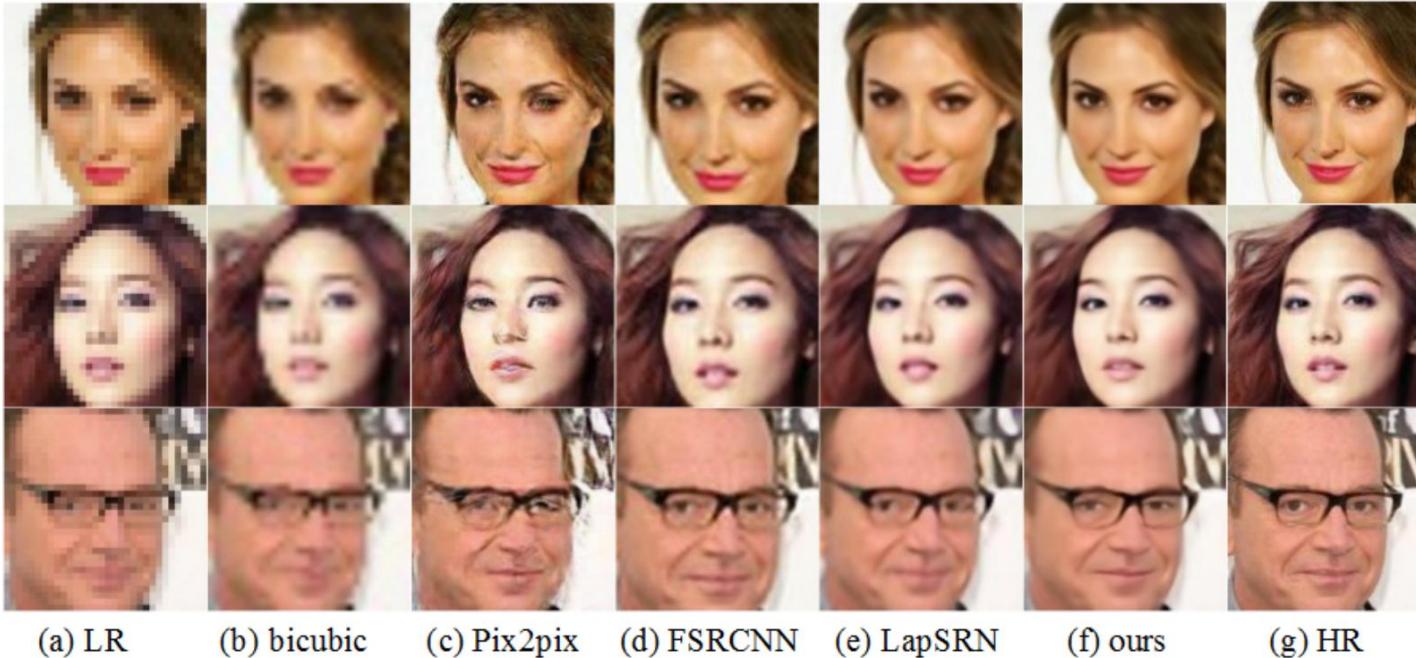
Applications (Image Gap Filling)



2016, “Context Encoders: Feature Learning by Inpainting”, Deepak Pathak et al

Generative Adversarial Networks

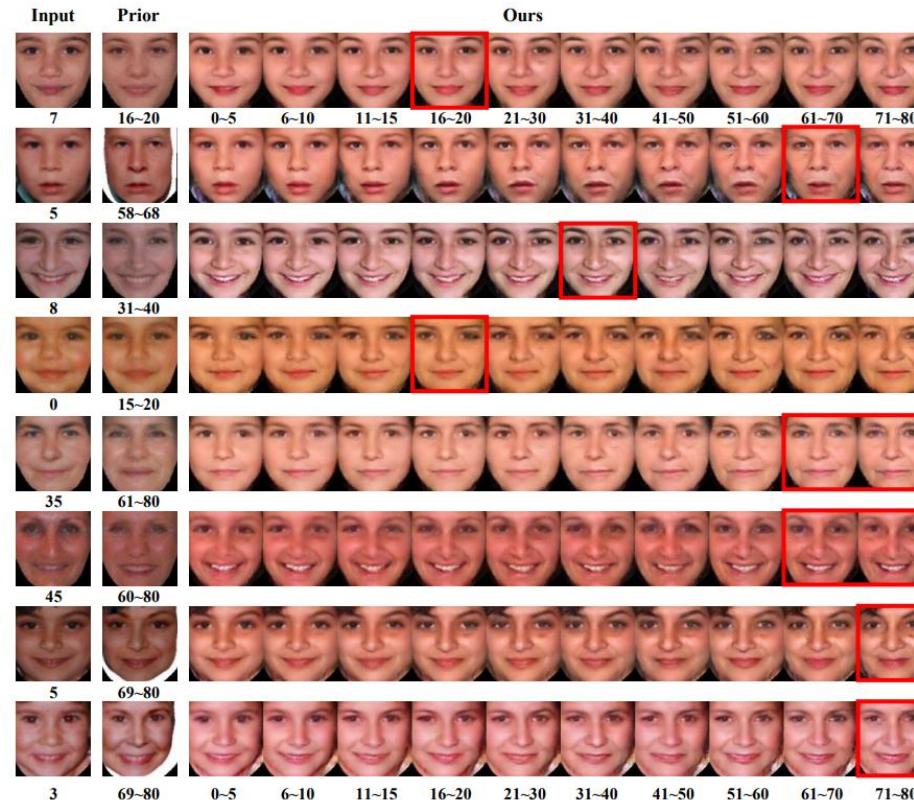
Applications (Increase Image Resolution)



2017, "High-Quality Face Image Super-Resolution Using Conditional Generative Adversarial Networks", Huang Bin et al

Generative Adversarial Networks

Applications (Progressive Face Aging)



MultiDrone



2017, "Age Progression/Regression by Conditional Adversarial Autoencoder", Zhifei Zhang et al

Generative Adversarial Networks

Applications (Music Synthesis)

MultDrone



2016, “C-RNN-GAN: Continuous recurrent neural networks with adversarial training”, Olof Mogren

Conclusions

MultDrone



- We have a long journey before we reach human-like high-level Computer Vision
- Many knowledge is created in recent years for Adversarial Examples and GANs
- 400+ named GANs variations exist today (GAN Zoo: <https://bit.ly/2sxGUHl>)
- It is the era where the Generative Models rise
- Machine Learning needs to be secure for a safe digital intelligent world
- The future of fake digital media is coming and is going to be really dangerous
- A lot of fundamental questions need to be well answered
- Exploit the vulnerabilities of Neural Networks
- Exploit the robustness of Adversarial Examples

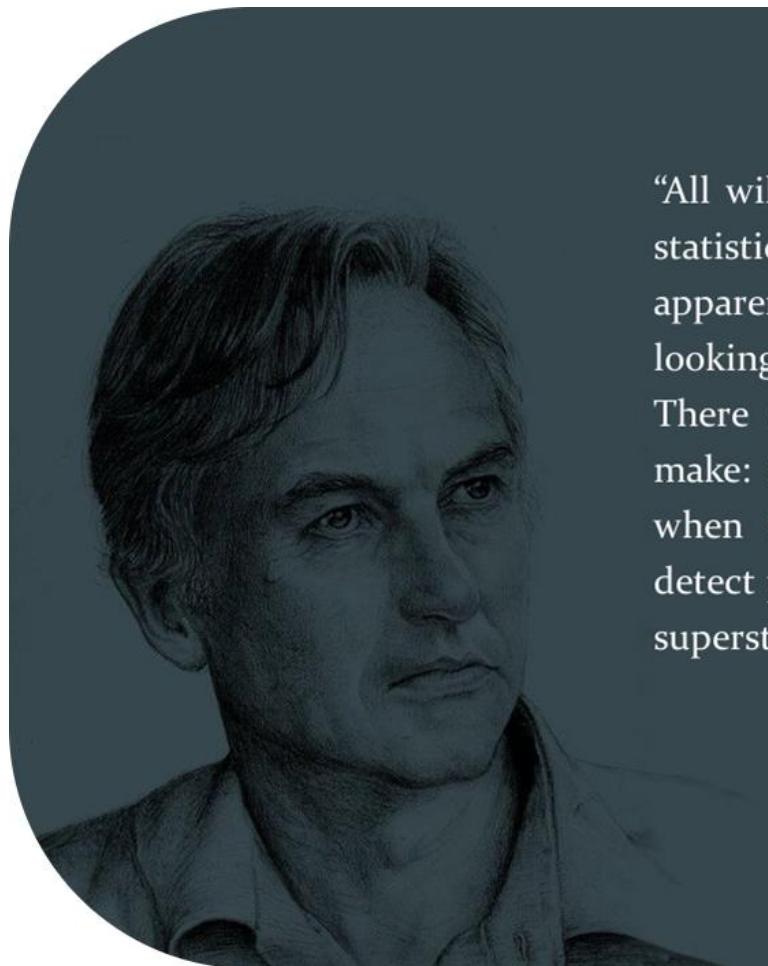
Recent Research

Adversarial Examples - GANs



- Universal adversarial perturbations (Moosavi-Dezfooli et al, 2016)
- Adversarial Machine Learning at Scale (Alexey Kurakin et al, 2016)
- Synthesizing Robust Adversarial Examples (Anish Athalye et al, 2017)
- Delving into Transferable Adversarial Examples and Black-box Attacks (Yanpei Liu et al, 2017)
- Adversarial Examples Are Not Easily Detected: Bypassing Detection Methods (Carlini et al, 2017)
- Adversarial vulnerability for any classifier (Alhussein Fawzi et al, 2018)
- Adversarial Attacks and Defences Competition (Alexey Kurakin et al, 2018)
- Adversarial Logit Pairing (Harini Kannan et al, 2018)
- Evolutionary Generative Adversarial Networks (Chaoyue Wang et al, 2018)
- Generating Adversarial Examples with Adversarial Networks (Chaowei Xiao et al, 2018)
- Robust Conditional Generative Adversarial Networks (Grigorios G. Chrysos et al, 2018)

Thank you a lot and have a nice day!



“All wild animals have to be kind of natural statisticians, looking for patterns in the apparent randomness of nature when they are looking for food or trying to avoid predators. There are two kinds of mistake they can make: they can either fail to detect pattern when there is some, or they can seem to detect pattern when there isn't any, and that's superstition.”

Richard Dawkins

MultDrone

