# Machine Learning for Deepfake Detection

Shannon K. Gallagher, PhD

skgallagher@sei.cmu.edu

Softw are Engineering Institute Carnegie Mellon University Pittsburgh, PA 15213



Carnegie Mellon University Software Engineering Institute

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

Copyright 2022 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

References herein to any specific commercial product, process, or service by trade name, trade mark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by Carnegie Mellon University or its Software Engineering Institute.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

Carnegie Mellon<sup>®</sup> and CERT<sup>®</sup> are registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM22-0747

#### We learned why we need detectors, but why do we need machine learning?



Image from DW.com. The righthand side image is an example of a deepfake used to impersonate the mayor of Kyiv. Brackets are ours.

#### **Potential Dangers:**

- >700k hours of video uploaded to YouTube daily
- Deepfake apps can be run with push of a button
- Deepfakes are generated with ML, logical then to think that we can detect them with ML
- Castle defense

#### We need scalable detectors!

**Carnegie Mellon University** Software Engineering Institute

#### First, a crash course in modeling



### Problem set-up: Is this a deepfake?

Model/Alg./Blackbox



Input



**Carnegie Mellon University** Software Engineering Institute

## Problem set-up: Now with some math

Model/Alg./Blackbox



**Carnegie Mellon University** Software Engineering Institute

# We need to first specify a *model* for f(X)

Example: Logistic Regression

$$f(X) = logit^{-1}(X\beta)$$
$$= \frac{1}{1 + e^{-X\beta}}$$

X = vector of features

 $\beta$  = vector of parameters to learn

Logistic (sigmoid or inverse logit) function



This Photo by Unknown Author is licensed under <u>CCBY</u>

**Carnegie Mellon University** Software Engineering Institute

# Often f(X) is a neural net we need to learn

 $f: \mathcal{X} \to [0,1]$  $X \mapsto f(X; \theta)$ 

X = image

 $\theta$  = parameters we need to *learn* 



This Photo by Unknown Author is licensed under <u>CCBY</u>

#### Statistical and ML models let us estimate $\theta$ from data

Data = { $(X_i, Y_i)_{i=1...n}$ } Data , REAL , FAKE , REAL , FAKE

**Carnegie Mellon University** Software Engineering Institute

# We learn by minimizing a loss function

 $\hat{Y}_i = f(\boldsymbol{X}_i, \boldsymbol{\theta})$ 

 $L(\widehat{Y}_i, Y_i)$  = distance between *predicted* and actual value

$$\widehat{\boldsymbol{\theta}} = argmin_{\boldsymbol{\theta}} \sum_{i=1...n} L\left(\widehat{Y}_{i}, Y_{i}\right)$$



# Our favorite loss is binary cross entropy

$$L(y_i, \hat{y}_i) = -[y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)]$$

If  $y_i = 1$ 

if  $\hat{y}_i$  is close to 1 then term is small if  $\hat{y}_i$  is close to 0 then term is large If  $y_i = 0$ if  $\hat{y}_i$  is close to 1 then term is large

if  $\hat{y}_i$  is close to 0 then term is small

### Goal: We want loss to be *small*

**Carnegie Mellon University** Software Engineering Institute Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University

11

#### In summary, we need only three things

- 1. We have a set of labeled data
- 2. We specify a function class that has parameters we need to learn
- 3. Using data, we minimize a loss function to estimate parameters in neural net

#### The devil is in the details

#### These two are hard problems, but we have lots of help

2. We specify a function class that has parameters we need to learn (e.g. neural net)3. Using data, we minimize a loss function to estimate parameters in neural net





**Carnegie Mellon University** Software Engineering Institute

#### So the devil is in the data

Problem 1: Overfitting



#### Solution: Train and Test Data

**Carnegie Mellon University** Software Engineering Institute

#### Train and Test data

Idea: Don't 'train' your model on all the data. Leave some for testing.



## Problem 2. Affine transformations

*Idea:* a deepfake is still a deepfake if the face is big/large, rotated, upside down, off-center, etc.



#### Solution: Augmented data

**Carnegie Mellon University** Software Engineering Institute

#### **Problem 3: Spurious features**

*Idea:* Data are biased. Sometimes the machine finds coincidental features, not real ones.



#### Solution: standardization and masking

**Carnegie Mellon University** Software Engineering Institute

#### Standardization

- Align and center faces
- Subtract the 'average'

#### AVG(REAL) - AVG(FAKE)



#### Takeaway: we want real differences to stand out

**Carnegie Mellon University** Software Engineering Institute

# Masking



#### Takeaway: we want *real* differences to stand out

**Carnegie Mellon University** Software Engineering Institute

## There's a very clear pattern in deepfake detection

- 1. Gather labeled data
- 2. Transform data to emphasize useful features and mitigate biases
- 3. Train a model on some data
- 4. Test the model on separate data

#### We call this process the deepfake detection pipeline

**Carnegie Mellon University** Software Engineering Institute Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University

20

#### Gather labeled data

**Carnegie Mellon University** Software Engineering Institute

## Gathering deepfake data is harder than it may seem

- Ethical issues
- Proprietary issues
- Accessibility issues

# **Consequence**: There are about a dozen datasets the public effectively uses for deepfake detection

**Carnegie Mellon University** Software Engineering Institute

#### Popular datasets for deepfake detection

Name	Type	Format	Labels	Size (08)	Size (#)	Resolution	GANF	Gen.	Faces?	Year	Access	Ex.
The fire Galegier												
Flickt-Faces-HQ	. Iniges	sng and json	Real	2 TB total, 90 GB for a condensed awt	210s files and 70k in condensed yst	. Variubie	Ne		Yes	2020	Google Drive	
Desplaie Detection Challenge (DFDC)	Video	Agm	Real/Fake	47008 compressed	1004+ 10± clips 8428 unique actors	1080p (mostly)	yes	1-3	yes	2020	Register on Kaggle	atductorizonal
MetFaces	images	and and (son	Real (paintings of faces)	150B	2621 files	9024x1024	No		Yes (but paintings)	2020	Seo hero	<b>E</b>
DeeperForensics	Video	<i>m</i> 04	Resylfanipulated	30008	60k videos, 500 Individuale		Yes		Yes	2020	Google furm/license	
DeepFake Detection Dataset (DFD) Note: The data is Face Porensica++	Video	hqn	Resphanipulated	-8008 compressed -2 TB new	303 original videos and 3068 manipulated	Vartuble	yes		yes	2018	Google form	

**Carnegie Mellon University** Software Engineering Institute

## Flickr-Face HQ (Real portrait photos)









![](_page_23_Picture_5.jpeg)

**Carnegie Mellon University** Software Engineering Institute

Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University [DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

## StyleGAN2 (Synthetic individuals, portrait style)

![](_page_24_Picture_1.jpeg)

![](_page_24_Picture_2.jpeg)

![](_page_24_Picture_3.jpeg)

![](_page_24_Picture_4.jpeg)

![](_page_24_Picture_5.jpeg)

**Carnegie Mellon University** Software Engineering Institute

Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University [DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

#### Deepfake Detection Challenge (Real)

![](_page_25_Picture_1.jpeg)

**Carnegie Mellon University** Software Engineering Institute

## DeepFake Detection Challenge (Fake)

![](_page_26_Picture_1.jpeg)

**Carnegie Mellon University** Software Engineering Institute

Celeb DF v2

![](_page_27_Picture_1.jpeg)

![](_page_27_Picture_2.jpeg)

![](_page_27_Picture_3.jpeg)

![](_page_27_Picture_4.jpeg)

**Carnegie Mellon University** Software Engineering Institute

Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University [DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

Abstracting a video/image into computer representation

- Inputs of dimension (W, H, C, F)
  - W = Pixel width
  - H = Pixel height
  - C = Channel (RGB)
  - F = Frame #

#### **Data Transformations**

**Carnegie Mellon University** Software Engineering Institute

Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University [DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

#### How do we extract useful features?

Thought: Only small sections of images/videos are 'deepfaked' Problem: Extract the 'area' where we think deepfake will take place

#### For us this usually translates to extracting faces

**Carnegie Mellon University** Software Engineering Institute

#### Haar cascades

![](_page_31_Picture_1.jpeg)

#### From <u>Ngo et al. (2009)</u>

**Carnegie Mellon University** Software Engineering Institute

#### **Edge Detectors**

![](_page_32_Picture_1.jpeg)

#### From this canny edge detection article

**Carnegie Mellon University** Software Engineering Institute

#### **Entirely separate Neural Nets**

![](_page_33_Picture_1.jpeg)

#### Facial boundary detection from MTCNN

**Carnegie Mellon University** Software Engineering Institute

#### Then we can augment the data

![](_page_34_Picture_1.jpeg)

#### From Zeno, Kalinovskiy, and Matveev (2021)

**Carnegie Mellon University** Software Engineering Institute

# Modeling

**Carnegie Mellon University** Software Engineering Institute Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

# Training Models – largely pre-trained Neural Nets

<u>AlexNet</u>	<u>RegNet</u>	
<u>ConvNeXt</u>	<u>ResNet</u>	
<u>DenseNet</u>	<u>ResNeXt</u>	P7/12
<b>EfficientNet</b>	ShuffleNet V2	P5/32
EfficientNetV2	<u>SqueezeNet</u>	P3/8
<u>GoogLeNet</u>	SwinTransformer	P1/2
Inception V3	VGG	
<b>MNASNet</b>	VisionTransformer	
MobileNet V2	Wide ResNet	
MobileNet V3		

![](_page_36_Figure_2.jpeg)

37

## Testing/Evaluation

**Carnegie Mellon University** Software Engineering Institute

Prototype results: data bias makes generalization hard

#### Accuracy (%) of fine-tuned ResNet

	Data Set	Celeb DF v1	Stylegan2	Stylegan3-t	Stylegan3-r	DFDC Pt. 0
_	Celeb DF v1	99.1	44.2	44.2	44.0	51.2
n	Stylegan2	24.1	98.7	52.9	48.4	57.4
	Stylegan3-t	16.7	69.7	96.7	84.0	7.0
ine	Stylegan3-r	16.9	68.0	89.0	97.2	7.0
Tra	DFDC Pt. 0	68.1	57.4	57.5	57.5	88.7

## Tested on

## Testing: How do the best models do??

#### Great\*\*

![](_page_39_Picture_2.jpeg)

0

0.42842

2

2

ZY

2Y

3 Ntechiah	(a) // 1088	
	G 0.43	1452

WM

#### \*\*In controlled scenarios

**Carnegie Mellon University** Software Engineering Institute

Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University

68 -0

## In closing

- Deepfake detection can be completely adapted to the ML modeling framework
- In theory, deepfake detection is a simple four step process
  - Data collection
  - Data transformation
  - Modeling
  - Evaluation
- But the devil is in the details

And Dr. Bernaciak will show you exactly how!

## The GAN problem

Red makes a generator to create deepfakes

Blue makes a detector

Red uses results from blue's detector to make generator better

Blue uses new red images to improve detector

Who wins?

. . .

#### GAN set-up

X' = G(Z) = generator fake image X = real image D(X) = discriminator in [0,1] Y=0, Y'=1 (0 is real, 1 is fake)

Data ={
$$(X_i, 0)_{i=1...m}, (X'_j, 1)_{j=1...n}$$
}  
 $L(x, y) = loss function$ 

![](_page_42_Figure_3.jpeg)

Fig. 18.1.1 Generative Adversarial Networks

#### Fig. from d2l.ai

**Carnegie Mellon University** Software Engineering Institute Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University [DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution

## GAN Set-up

Round 0: Generator introduces fakes

Round 1:

Discriminator turn: Use generated data to get best discriminator

$$\widehat{D}^{1}|\widehat{G}^{0} = argmin_{D}\sum_{\{i=1\}}^{m} L(D(X_{i}), 0) + \sum_{\{j=1\}}^{n} L(D(X_{j}'), 1)$$

Generator turn: Directly try to deceive discriminator

$$\widehat{G}^{1} | \widehat{D^{1}} = argmin_{G} \sum_{\{j=1\}}^{n} L(\widehat{D}^{1}(X_{j}'), 0)$$
$$= argmin_{G} \sum_{\{j=1\}}^{n} L(\widehat{D}^{1}(G(Z_{j})), 0)$$

#### Repeat

**Carnegie Mellon University** Software Engineering Institute Machine Learning for Deepfake Detection Contact: skgallagher@sei.cmu.edu © 2022 Carnegie Mellon University

44