Detecting Facial Manipulation Deepfakes

Evan Kravitz, Huazhe Xu April 21, 2020

ZUCKERBERG: WE'RE INCREASING TRANSPARENCY ON ADS ANNOUNCES NEW MEASURES TO "PROTECT ELECTIONS"

FACEBOOK



What is a deepfake?

- Synthetic image/video of a person that looks realistic to human viewers, which can be used to perpetrate fraud or spread misinformation
- Deepfakes are a form of social engineering attack
- We have focused our research on detecting facial deepfakes

Social engineering attacks

Are you aware of bill and Melinda gates foundation grant?



click on the link it will take you to their page, like the page and message them that you want to apply for a grant, okay

Face synthesis

StyleGan (2019)





FaceSwap (2016)



Deepfake FaceSwap (2020)



Jennifer Lawrence/Steve Buscemi FaceSwap using the Villain model

Face attribute

StarGAN (2018)



Facial expression

Face2Face (2016)



Protecting against deepfake

• We need a system for authenticating media





Convolutional Neural Network (CNN)



Convolutional Neural Network (CNN) cont.





Optical flow



Amerini et al., 2019



Figure 2. Optical flow for original (left) and deepfake (right) videos.

CNN's with optical flow



CNN's with self-labeled data (Li et al., 2019)

1. Generate "negative" examples that contain deepfake generation artifacts



2. Use "negative" examples to train a CNN



Forensic deepfake detection

- Forensic approach
 - Generate correlations between facial features in a video to determine "signature motion" (Agarwal et al., 2019)





Figure 3. Shown is a 2-D visualization of the 190-D features for Hillary Clinton (brown), Barack Obama (light gray with a black border), Bernie Sanders (green), Donald Trump (orange), Elizabeth Warren (blue), random people [23] (pink), and lip-sync deep fake of Barack Obama (dark gray with a black border).

Forensic deepfake detection cont.



Our contribution

- We aim to improve upon existing neural network and forensic feature models.
 - ✓ Feature augmentation and enhancement
 - ✓ Better classification model

Original labeled data

Altered labeled data



Entire YouTube 8M Dataset

Cropped faces from video frames

Face2Face manipulated video frames

Dataset cont.

- 704 videos for training (368,135 images)
- 150 videos for validation (75,526 images)
- 50 videos for testing (77,745 images)



Forensic analysis of facial landmarks



Principal Component Analysis (PCA)

- Popular technique for dimensionality reduction
- Transform feature space into orthogonal basis features, only capture most prominent features
- Fewer features \rightarrow less variance, less overfitting



Method: Random forest classifier

• Pros:

- Works with few features
- Lower variance compared to regular decision tree
- Explainable model
- Low cost



- Cons:
 - Hard to tune

Method: Support vector machine

• Pros:

- Supports non-linear decision boundaries
- Cons:
 - Hard to tune kernel and hyperparameters



Method: Neural Network with facial landmarks



Need extensive tuning

Metrics

Accuracy: (True Positive + True Negative) / total samples

Precision: True Positives / All the predicted positives

Recall: True Positives / All the actual positives

Results: in-distribution samples (small scale)

- Near perfect performance for random forest
- What does this imply? We can perfectly detect fake/real across the web if we have label for part of a clip.
- 10K training images

	SVI	M	Random Forest		NN	
Accuracy	80.0	00%	98.10%		85.12%	
Table 1: A	Table 1: Accuracy for different models					
		Random Fo	orest	NN		
Precision		98.52%		92.81%		
Recall		98.72%		85.01	%	

 Table 2: Precision and Recall for top 2 models

Results: out-of-distribution training and testing

- Both methods drops significantly
- Neural Net performs slightly better (the training accuracy for NN is 90% and for random forest 99.9%)
- Training data is too little!
- 14K training images

	SVM	Random Forest	NN
Accuracy	N/A	70.50%	73.78%

Table 1: Accuracy of Random and NN model

	Random Forest	NN
Precision	77.15%	79.23%
Recall	58.82%	63.44%

Table 2: Precision and Recall for top 2 models

Public Benchmark Results w/ ~5 times our current training data

- Larger net
- More data
- Utilize video property

Method	Info	Deepfakes	Face2Face
Sentinel		0.964	0.905
Xception	P	0.964	0.869
Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess	, Justus Thies, Matthias Nießner: FaceForensics++: Learning	to Detect Manipulate	d Facial Images. ICC
Inception Resnet V1		0.936	0.839
Nika Dogonadze, Jana Obernosterer: Deep Face Forgery Detection.	Advanced Deep Learning for Computer Vision Course at TU	м	
SCAN	C	0.909	0.825
EfficientNet-b4		0.955	0.796
Face Defence		0.945	0.766
swap_classify		0.909	0.759
KceptionNet Full Image	P	0.745	0.759
ndreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess	, Justus Thies, Matthias Nießner: FaceForensics++: Learning	to Detect Manipulate	d Facial Images. ICC
Bayar and Stamm		0.845	0.737
Belhassen Bayar and Matthew C. Stamm: A deep learning approact	h to universal image manipulation detection using a new conv	olutional layer. ACM	Workshop on Informa
Steganalysis Features		0.736	0.737
lessica Fridrich and Jan Kodovsky: Rich Models for Steganalysis of	f Digital Images. IEEE Transactions on Information Forensics a	and Security	
GAEL-Net		0.718	0.686
Recasting		0.855	0.679
Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva: Recasting n Hiding and Multimedia Security	esidual-based local descriptors as convolutional neural netwo	rks: an application to	image forgery detec
Forged face detection : fCNN		0.791	0.642
Rahmouni		0.855	0.642
Nicolas Rahmouni, Vincent Nozick, Junichi Yamagishi, and Isao Ech Security,	izen: Distinguishing computer graphics from natural images u	using convolution neu	iral networks. IEEE W

Visualized Examples



Original Image



Altered Image



Scale up & Analysis



Compare with public Benchmark



Temporal Features



CNN + Forensic Features

Thank you!