Defeating the Fakes

Vocal Fake Detection using Discrete Fourier Transforms and Spectrograms in Neural Networks

Tina Tami Lars Tijsmans

SZ

Supervisor Zeno Geradts



Tech

Deepfake queen to deliver Channel 4 Christmas message [1]

() 23 December 2020

Tech

Deepfake queen to deliver Channel 4 Christmas message [1]

() 23 December 2020



22 april 2021 19:53 Aangepast: 22 april 2021 20:00

Tech

'Deepfake is the future of content Deepfake queen to delive Christmas message [1]

(23 December 2020

Europese politici benaderd door deepfake van Russische oppositie [2]

22 april 2021 19:53 Aangepast: 22 april 2021 20:00

Tech

'Deepfake is the future of content Deepfake queen to deliver Christmas message [1]

() 23 December 2020

Europese politici benaderd door deepf van Russische oppositie [2]

22 april 2021 19:53 Aangepast: 22 april 2021 20:00

[4] https://www.independent.co.uk/celebrity-news/tom-cruise-deep-

Scarily authentic new deep fake of Tom Cruise attracts millions of views [4]

'Reality is becoming mutable'

Peony Hirwani | @peony_hirwani | Thursday 27 May 2021 17:02 | comments



navalny-europese-politic

[5] https://www.wsj.com/articles/fraudst in-unusual-cybercrime-case-11567157402

Research questions

Can Fourier transformations be used to distinguish vocal fakes from real voices in a classification model?

- How do discrete Fourier transform and spectrograms compare in accuracy?
- How do discrete Fourier transform and spectrograms compare in computational performance?
- How does the model perform on different datasets?

Related work

- 'DeepSonar: Towards Effective and Robust Detection of Al-Synthesized Fake Voices'^[1]
- 'A Machine Learning Model to Detect Fake Voice'^[2]
- 'A Fourier transform based audio watermarking algorithm'^[3]

Datasets

ASVSpoof2019



RTVCSpeech using LibriSpeech ASR corpus



Preprocessing: RTVCSpeech



Preprocessing: ASVSpoof



Real Time Voice Cloning (RTVC)



Discrete Fourier transform

$$egin{aligned} X_k &= \sum_{n=0}^{N-1} x_n \cdot e^{-rac{i2\pi}{N}kn} \ &= \sum_{n=0}^{N-1} x_n \cdot \left[\cos\!\left(rac{2\pi}{N}kn
ight) - i \cdot \sin\!\left(rac{2\pi}{N}kn
ight)
ight], \end{aligned}$$

Discrete Fourier transform



Spectrograms



https://towardsdatascience.com/understanding-audio-data-fourier-transform-fft-spectrogram-and-speech-recognition-a4072d228520



1. RTVCSpeech as training and test set

2. RTVCSpeech as training set, ASVSpoof as test set

3. ASVSpoof as training and test set

4. ASVSpoof as training set, RTVCSpeech as test set

RTVCSpeech	Real	Fake
Training set	2116	2116
Test set	528	528

ASVSpoof	Real	Fake
Training set	2064	2116
Test set	516	528

Experiments



Models

2D Convolutional Neural Network

- Image classification
- Network employs a mathematical operation called <u>convolution</u>

Fully Connected Neural Network

- All the neurons in a layer are connected to those in the next layer
- Learns features from all the combinations of the features of the previous layer

Models - 2D CNN: Convolution layer



image and computing the dot product to detect patterns

https://www.slideshare.net/Simplilearn/convolutional-neural-network-tutorial-cnn-how-cnn-works-deep-learning-tutorialsimplilearn/Simplilearn/convolutional-neural-network-tutorial-cnn-how-cnn-works-deep-learning-tutorial-simplilearn

Models – 2D CNN: MaxPooling layer





Decrease the size of the feature map by taking the largest element

https://peltarion.com/knowledge-center/documentation/modeling-view/build-an-ai-model/blocks/2d-max-pooling-block

Models - 2D CNN and FCNN : Fully connected (dense) layer



Models - 2D CNN and FCNN: ReLU (Rectified Linear Unit) activation function



Models - 2D CNN and FCNN: Sigmoid activation function



Models: 2D Convolutional Neural Network



Models: Fully Connected Neural Network



Results – spectrograms: RTVCSpeech as training set

RTVCSpeech as test set

Predicted

		Real	Fake
aı	Real	527	1
ערנמ	Fake	0	528

	Real	Fake
Accuracy	0.999	0.999
F1-score	0.999	0.999

Time on GPU: 20s for 1 epoch Time on CPU: 655s for 1 epoch

ASVSpoof as test set

		Real	Fake
lμ	Real	516	0
Actua	Fake	386	142

	Real	Fake
Accuracy	0.630	0.630
F1-score	0.728	0.424

Results – spectrograms: ASVSpoof as training set

ASVSpoof as test set

Predicted

		Real	Fake
aı	Real	514	2
АСГИ	Fake	0	528

	Real	Fake
Accuracy	0.998	0.998
F1-score	0.998	0.998

Time on GPU: 20s for 1 epoch Time on CPU: 655s for 1 epoch

RTVCSpeech as test set

		Real	Fake
al	Real	198	330
Actu	Fake	26	502

	Real	Fake
Accuracy	0.663	0.663
F1-score	0.527	0.738

Results – DFT : RTVCSpeech training set

RTVCSpeech as test set

Predicted

		Real	Fake
aı	Real	449	61
HLLU	Fake	81	467

	Real	Fake
Accuracy	0.886	0.886
F1-score	0.863	0.868

Time on GPU: 1s for 1 epoch Time on CPU: 1s for 1 epoch

ASVSpoof as test set

		Real	Fake
al	Real	175	334
Actu	Fake	123	413

	Real	Fake
Accuracy	0.563	0.563
F1-score	0.434	0.664

Results - DFT : ASVSpoof as training set

ASVSpoof as test set

Predicted

		Real	Fake
al	Real	457	70
ACTU	Fake	83	435

	Real	Fake
Accuracy	0.854	0.854
F1-score	0.857	0.850

Time on GPU: 1s for 1 epoch Time on CPU: 1s for 1 epoch

RTVCSpeech as test set

		Real	Fake
lı	Real	271	235
Actuc	Fake	152	400

	Real	Fake
Accuracy	0.634	0.634
F1-score	0.583	0.675

Spectrograms: accuracy

Training set

		RTVCSpeech	ASVSpoof
set	RTVCSpeech	0.999	0.663
Test	ASVSpoof	0.630	0.998

How do discrete Fourier transform and spectrograms compare in accuracy?

Using spectrograms results in a higher accuracy.

Discrete Fourier Transform: accuracy

Training set

		RTVCSpeech	ASVSpoof
Test set	RTVCSpeech	0.886	0.634
	ASVSpoof	0.563	0.854

Spectrograms: accuracy

Training set

		RTVCSpeech	ASVSpoof
set	RTVCSpeech	0.999	0.663
Fest	ASVSpoof	0.630	0.998

How do discrete Fourier transform and spectrograms compare in accuracy?

Using spectrograms results in a higher accuracy.

Discrete Fourier Transform: accuracy

Training set

		RTVCSpeech	ASVSpoof
Test set	RTVCSpeech	0.886	0.634
	ASVSpoof	0.563	0.854

How does the model perform on different datasets?

ASVSpoof as training set gives better results. Training and testing on the same dataset gives better results.

How do discrete Fourier transform and spectrograms compare in computational performance?

	GPU	CPU
2D CNN	20s for 1 epoch	655s for 1 epoch
FCNN	1s for 1 epoch	1s for 1 epoch

Can Fourier transformations be used to distinguish vocal fakes from real voices in a classification model? Yes, but...

Discussion and future work

- Improving the neural networks
 - Transfer learning
 - Adding explainability
 - Better vocal fakes