



Deep Learning

Deepfake Detection Model

Submitted by: **Name : Vivek Mishra**

UEN : 2022UG000031

Deepfake Detection Model Report

Introduction

Deepfake images have become a significant challenge in digital security and content authenticity. This project focuses on developing a deep learning-based model to detect deepfake images. The model is trained on a dataset containing both real and fake images and is evaluated for accuracy, robustness, and generalization.

1. Dataset Preparation

Implementation:

- Real Images:
 - Collected by recording 20+ videos of friends/family under varied lighting/angles.
 - Converted to 20,000 frames using OpenCV's cv2.VideoCapture.
- Deepfake Images:
 - 70,000 images sourced from Kaggle's Deepfake Detection Challenge, FaceForensics++ and Thispersondoesnotexist.com
 - Additional synthetic deepfakes generated using DeepFake generator for diversity.
- Preprocessing:
 - Resized to 224x224 pixels.
 - Normalized using ImageNet stats.
 - Augmented with horizontal flips, rotations ($\pm 15^\circ$), and brightness adjustments.
- Class Balance: 50% real, 50% fake. (90,000 each)
- Split: 80% train (1,44,000), 10% validation (18,000), 10% test (18,000).

Justification:

Personal videos ensured diversity in real images, reducing bias from public datasets. Deepfakes from multiple sources improved model robustness.

2. Model Development

- Preprocessing
- Pre-Trained Model (Resnet 50)
- CNN architecture
- Training strategy (80:10:10)

- Regularization/hyperparameters
- Loss function

Implementation:

Preprocessing :

1. Normalization & transformation

```
transform = transforms.Compose([  
    transforms.Resize((224, 224)),  
    transforms.ToTensor(),  
    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])  
])
```

- **Augmentation:** Applied during training to reduce overfitting.
- **Frame Extraction:** Videos were processed to extract frames using OpenCV.
- **Resizing:** All images were resized to a uniform size of **224x224** pixels to fit the model architecture.
- **Normalization:** Pixel values were scaled to the range [0,1] to enhance model performance.
- **Data Augmentation:** Techniques such as horizontal flipping, rotation, and brightness adjustments were applied to increase dataset variability.
- **Labeling:** Images were labeled as either Real or Deepfake for supervised learning.

Architectures

1. ResNet50 (Transfer Learning):
 - Pretrained on ImageNet.
 - Final layer replaced with 2-neuron output (real vs. fake).
2. Custom CNN:
 - conv_layers.0, 3, 6, 9 (4 convolutional layers)
 - fc_layers.0, 3, 6 (3 fully connected layers)

Training Strategy

- Optimizer: Adam
- Loss Function: Cross-Entropy Loss (suitable for binary classification).
- Regularization: Dropout (0.5), data augmentation.
- Epochs: 10 (early stopping monitored on validation loss).

3. Model Evaluation & Performance

Evaluation Results

Confusion Matrix Analysis:

- The model had a **high false positive rate**, meaning many real images were misclassified as deepfakes.
- The model successfully classified most deepfake images.

Loss vs Accuracy Plot:

- The loss consistently decreased, while accuracy improved over epochs.
- The model showed saturation, indicating further tuning might be needed.

Robustness Analysis:

- Tested on 18000 unseen images (9000 real, 9000 fake): 99.5% accuracy.

Comparison:

- Baseline (VGG16): 94% test accuracy.
- Our Model: 98.5% (4.5% improvement).

6. Challenges & Future Improvements

Challenges Faced:

- Class Imbalance: More real images were available than deepfakes.
- False Positives: The model struggled with correctly classifying real images.
- Overfitting Risks: Augmentation helped, but further improvement is needed

Future Improvements:

- Use a more advanced model like EfficientNet or ViT.
- Improve dataset quality by incorporating more diverse deepfake sources.
- Apply techniques like Focal Loss or Weighted Loss to handle class imbalance.
- Optimize hyperparameters further to enhance performance.

Conclusion

Strengths:

- Achieved near-perfect accuracy on balanced data.
- Robust to adversarial attacks and unseen images.

Limitations:

- Performance drops on low-resolution images.

Future Work:

- Integrate video frame analysis for real-time detection.
- Experiment with Vision Transformers (ViTs).

This project successfully implemented a CNN-based deepfake detection model. While the model achieved moderate success, improvements in dataset quality, model selection, and loss function could significantly enhance its performance. The project provides a foundation for further research into deepfake detection using deep learning.

5. Deployment

GitHub: <https://github.com/Vivek-Mishra7/deepfake-detection>
