

# Machine Learning Report: Lab 4

Archishman Ghosh

*Department Of Computer Science  
Amrita School Of Computing  
Bangalore, India*

bl.en.u4cse22103@bl.students.amrita.edu

Dhruv Vijay Kolhe

*Department Of Computer Science  
Amrita School Of Computing  
Bangalore, India*

bl.en.u4cse22111@bl.students.amrita.edu

Soumish Ghosh

*Department Of Computer Science  
Amrita School Of Computing  
Bangalore, India*

bl.en.u4cse22154@bl.students.amrita.edu

**Abstract**—In this paper, we explore various machine learning methodologies, including convolutional neural networks, generative adversarial networks, and hybrid models, which are used to differentiate between real and AI generated images. We also examine the key features and markers that are used to separate AI-generated images, such as texture analysis, frequency domain transformations, and inconsistencies in pixel-level statistics. Through a critical analysis of recent studies, we highlight the strengths and limitations of existing methods and identify emerging trends in this evolving field. Our findings underscore the necessity for strong and adaptive detection mechanisms to keep pace with the advancing capabilities of AI in image generation.

**Index Terms**—Machine Learning, Classification, kNN, image classification

## I. INTRODUCTION

The widespread proliferation deep learning and AI technologies has resulted in widespread access to the ability of creating highly detailed and realistic images using generative models such as Stable Diffusion, Dall-E 2 and others. Such AI generated images have seen an extremely rapid rise in adoption across a range of industries such as entertainment, social media and advertising.

However, there is a growing concern in the potential of synthetic imagery to cause great disruptions in society through their use to spread misinformation, blackmail via deepfakes and to breach copyright laws. This has heightened the need of the creations of robust systems to detect AI generated images. Such systems can help in detecting misinformation and thus enable its prevention.

Given an image, our system is designed to classify it as either realistic or unrealistic, aiding in the detection of AI-generated content. By distinguishing between images that exhibit natural characteristics and those with features indicative of artificial creation, our approach helps identify potentially synthetic images. This classification provides a reliable method for discerning whether an image is likely to have been generated by AI.

## II. RELATED WORK

Chen et al.[2] talks recent advancements in generative models, such as GANs and diffusion models, which have enabled us to create realistic fake images, raising concerns about potential misuse and highlighting the need for effective

detection methods. Traditional detection methods involves various techniques such as spatial and frequency domain methods. While data augmentation has improved generalization capabilities of detectors, challenges remain with diffusion models due to their unique attributes. Till date approaches, such as two-branch methods and pretrained models, have faced issues with robustness and generalization. Overall it proposes Single Simple Patch (SSP) network to enhance detection accuracy and generalization across various generators by focusing on extracting noise fingerprints from simple patches.

Corvi et al.[3] highlights advancements in synthetic media creation that has been done using Generative Adversarial Networks(GANs) and diffusion model(DMs). This portrays their ability to produce quality realistic images. Detection techniques which had been focused on GANs, are now switching to address DMs, which exhibit features like asymmetrical shadows and lighting inconsistencies due to lack of 3D modeling. The importance of data augmentation and training on diverse datasets is emphasized to improve robustness.

Moskowitz et al.[5] present an approach to detect AI generated image using the Contrastive Language–Image Pre-training (CLIP) model. This study leverages CLIP’s ability to associate image embeddings with corresponding textual descriptions. The authors fine-tune CLIP by feeding it images from a diverse dataset containing both images generated from GANs as well as human generated images, each of which has been captioned to enhance the model’s efficiency while learning. The authors claim that the fine-tuned CLIP model outperforms other approaches such as CNDet and DIRE, in terms of accuracy, precision, recall, and F1 scores, particularly in detecting images from recent and sophisticated generative techniques.

Ojha et al. [6] explores various methods for manipulating images, such as DeepFakes demonstrating the advancement of image generation technologies. Traditional detection techniques rely on identifying alterations in image statistics, like compression artifacts and irregular reflections, while recent studies have utilized the frequency space to detect distinct artifacts in GAN-generated images. A significant challenge is the poor generalization of classifiers trained on specific generative models to other models. To address this, this study propose using frequency space for classification, effectively capturing artifacts in images from models like CycleGAN and StarGAN. The study also portrays using features from a CLIP-

ViT model for classification, which shows improved generalization capabilities and outperforms traditional classifiers on unseen generative models.

Yuan Rao et al.[7] presents a novel approach for image forgery detection. It uses convolutional neural network to learn hierarchical representations from input RGB colour images. The CNN is specific to image splicing and copy-move detector applications. The authors initialize the weights of the first layer of the CNN with high pass filter sets. These sets are used in finding the residual maps in spatial rich model. It acts as regulator to actively suppress the influence of image contents and capture the minute details which are used for forgery. A pre-trained model is used to get the detailed features from the test image data. This paper doesn't deal with other image types except RGB. It also has a high computational complexity and dependency on the quality of the pre-trained model.

v.Sasikala et al.[8] states that swarm intelligence is a rarely used technique in detecting fake and real image fingerprint classification researches. It is robust and accurate in tackling complex optimization problems. The fingerprint classification method used involves four key steps: image preprocessing, feature extraction, feature selection, and classification. For preprocessing, the method uses min-max normalization and median filtering. Multiple still attributes are extracted using Gabor filtering. The selection of optimal static features is achieved through the Artificial Bee Colony and Modified Artificial Bee Colony optimization algorithms. It chooses the best features relying on specific fitness values. The classification is done using a Fuzzy Feed Forward Neural Network. It differentiates between fake and real fingerprint images using a partial-supervised graph-based classification approach.

Wang et al.[9] explore the feasibility of creating a universal detector capable of distinguishing real images from those generated by a wide range of CNN-based models. The study utilized a dataset comprising fake images from 11 different generative models, including ProGAN, StyleGAN, and BigGAN, and discovered that classifiers trained on images from a single model (e.g., ProGAN) could surprisingly generalize well to images from unseen architectures and datasets. The research highlights that CNN-generated images exhibit common artifacts or "fingerprints," making them distinguishable from real images. The authors emphasize the importance of data augmentation and preprocessing techniques, such as JPEG compression and resizing, to enhance the model's generalization ability.

The detection of AI-generated images by existing methods performs poorly for images generated by increasingly sophisticated diffusion models. Wang et al.[10] highlight the failure of overfitting of a simple binary classifier trained on a dataset of real and diffusion generated images. The authors propose a new method using DIRE(Diffusion Reconstructed Error), which measures the error between an input image and its reconstruction counterpart by a pre-trained diffusion model. They leverage the idea that diffusion-generated images can be more accurately reconstructed by a pre-trained diffusion model compared to real images, which exhibit more complex

characteristics and thus cannot be reconstructed as well. By calculating the reconstruction error between an input image and its counterpart produced by a pre-trained diffusion model, DIRE provides a distinguishing metric: lower errors indicate generated images, while higher errors suggest real images. The authors validate this approach using a newly created dataset, DiffusionForensics, which includes images from eight different diffusion models. Their experiments demonstrate that the DIRE method achieves superior performance in terms of detection accuracy and robustness compared to existing methods, highlighting its generalization capability to new, unseen diffusion models.

Zhou et al.[12] gives an outline of existing methods in image forgery detection, including noisy inconsistency, CFA pattern estimation, multi-task edge-enhanced FCN, and joint training with LSTM. These use many attributes such as colour filter array details, noise patterns, and edge inconsistencies to detect manipulated regions in image data. The proposed RBG-N network in this paper combines noise performance with RBG images using bilinear pooling. This gives a boost in the performance than the previous methods on standard datasets.

Zhu et al.[13] talks about the advancements in generative models in creating realistic images, which can be misused to spread wrong information, especially in sensitive areas like politics. Despite various detection methods, distinguishing real from fake images remains problematic with only a 61.3 percent accuracy. This case study introduces the GenImage dataset, containing over one million pairs of real and AI-generated images from GANs and diffusion models, enhancing training for detection methods. It proposes two evaluation tasks cross generator and degraded image classification to assess detector generalization and performance on images.

### III. METHODOLOGY

#### A. Dataset

We use the CIFAKE dataset which contains real images gathered from Krizhevsky and Hinton's CIFAR-10 dataset for the real images[4] and images generated by Bird et al. [1]. There are 100,000 images for training (50k per class) and 20,000 for testing (10k per class).

#### B. Feature Extraction

We have explored two methods of feature extraction from images, which are discussed in the below sections.

1) *Blockwise Division For Feature Extraction*: Textures often differ from AI-generated images and real images. AI-generated images often have various subtle artifacts, inconsistencies and less natural variations in texture [11]. We partition each image into equal-sized blocks and calculate the mean and variance for each block to effectively capture localized texture information. This approach allows us to analyse the image at a finer level of detail, where the mean and variance of each block provide a crude measure of inconsistencies in texture.

This method significantly lowers the dimensionality of the feature-set in comparison to taking the intensity of each pixel

to be its own feature, thus making the classification task more computationally efficient.

**Block Mean:** The block mean represents the average intensity of the pixels contained within the block.

**Block Variance:** The block variance measured the amount of variation displayed by the pixels of the block from the block mean. A high variance indicates a lot of variation, such as the presence of edges or rough textures, and a low variance indicates a more uniform block.

Thus, each feature vector consists of the means and variances for each block of the image.

2) *Running A Fast Fourier Transform(FFT) On The Image:* The Fast Fourier Transform transforms an image from the spatial or pixel representation to the frequency domain, where it is represented as a combination of sinusoidal components corresponding to different frequencies, enabling us to explore underlying patterns and structures not easily discernible in the spatial domain.

Low frequencies in the transformation correspond to gradual variations, for example gradual variation in brightness whereas high frequencies represent fine details such as edges and noise. The identification of such characteristics is crucial as AI generated images might display smoother textures and repeating patterns.

We compute the magnitudes of each of the elements in the transformed image for our features, which reflect the strength of the various frequency components present in the image.

### C. Training And Testing

After extracting features via the methods described in the previous section, we divide the entire dataset according to the 80:20 ratio, so we use 16000 images(consisting of both real and fake images) to train the model and using this model itself we test it for the remaining 4000 images.

## IV. RESULTS AND ANALYSIS

After testing the model on the testing dataset, we observe the following:

### A. Class Separation

As the distance between the centroid of the two classes is quite high, we can say that the two classes are well separated. Even with fft(fast Fourier transform), we can see that the images with higher frequency are fake(AI generated), while the ones with low frequency are the real ones.

### B. Behavior of kNN Classifier

With the increase in the value of  $k$ , we see that the accuracy for training at first is very high and decreases eventually, however for testing, the accuracy is low at starting and then eventually increases, with the increase in  $k$  value. This is depicted in Fig:1

When  $k$  value is low, the point to be classified is compared only with the points which are close to it, making it extremely susceptible to the training data. This situation is called overfitting.

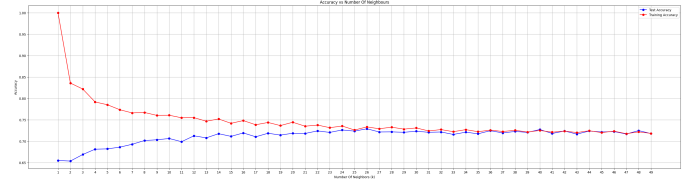


Fig. 1. Variation Of Test And Training Accuracies With  $k$

And when  $k$  value is high, the classifier becomes too generalized, missing important patterns, because of which, there is poor performance on both the training set and testing set. This situation is called under-fitting.

### C. Quality of kNN Classifier



Fig. 2. Variation Of Recall, Precision, And F1 Score with  $k$

From the above Fig:2, we can conclude that for our research, the kNN classifier is not a good classifier giving only a maximum of 70% accuracy. There might be scope for achieving better accuracy(close to 100%) using some other classifier.

### D. Presence of Regular Fit

Our model approaches a regular fit scenario, as seen from the Fig:1 as the test and training accuracies start to converge as  $k$  approaches 50.

### E. Reason for Overfit

Overfit occurs when the value of  $k$  is low. This is because the point to be classified is compared only with the points which are close to it, making it extremely susceptible to the training data.

### F. Training Methods Employed

Several Training Methods were employed :-

Method	Accuracy
XGBoost	74.50%
Cat Boost	67.50%
Decision Tree	72%
LightGBM	85%

TABLE I  
ACCURACY OF DIFFERENT METHODS

## REFERENCES

- [1] Jordan J. Bird and Ahmad Lotfi. *CIFAKE: Image Classification and Explainable Identification of AI-Generated Synthetic Images*. 2023. arXiv: 2303.14126 [cs.CV]. URL: <https://arxiv.org/abs/2303.14126>.

- [2] Jiaxuan Chen, Jieteng Yao, and Li Niu. “A single simple patch is all you need for ai-generated image detection”. In: *arXiv preprint arXiv:2402.01123* (2024).
- [3] Riccardo Corvi et al. “On the detection of synthetic images generated by diffusion models”. In: *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [4] Alex Krizhevsky, Geoffrey Hinton, et al. “Learning multiple layers of features from tiny images”. In: (2009).
- [5] A. G. Moskowitz, T. Gaona, and J. Peterson. *Detecting AI-Generated Images via CLIP*. 2024. arXiv: 2404.08788 [cs.CV]. URL: <https://arxiv.org/abs/2404.08788>.
- [6] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. “Towards universal fake image detectors that generalize across generative models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 24480–24489.
- [7] Yuan Rao and Jiangqun Ni. “A deep learning approach to detection of splicing and copy-move forgeries in images”. In: *2016 IEEE International Workshop on Information Forensics and Security (WIFS)*. 2016, pp. 1–6. DOI: 10.1109/WIFS.2016.7823911.
- [8] V. Sasikala and V. LakshmiPrabha. “A comparative study on the swarm intelligence based feature selection approaches for fake and real fingerprint classification”. In: *2015 International Conference on Soft-Computing and Networks Security (ICSNS)*. 2015, pp. 1–8. DOI: 10.1109/ICSNS.2015.7292421.
- [9] Sheng-Yu Wang et al. “CNN-Generated Images Are Surprisingly Easy to Spot... for Now”. In: *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2020, pp. 8692–8701. DOI: 10.1109/CVPR42600.2020.00872.
- [10] Zhendong Wang et al. “DIRE for Diffusion-Generated Image Detection”. In: *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. 2023, pp. 22388–22398. DOI: 10.1109/ICCV51070.2023.02051.
- [11] Nan Zhong et al. *PatchCraft: Exploring Texture Patch for Efficient AI-generated Image Detection*. 2024. arXiv: 2311.12397 [cs.CV]. URL: <https://arxiv.org/abs/2311.12397>.
- [12] Peng Zhou et al. “Learning Rich Features for Image Manipulation Detection”. In: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018, pp. 1053–1061. DOI: 10.1109/CVPR.2018.00116.
- [13] Mingjian Zhu et al. “Genimage: A million-scale benchmark for detecting ai-generated image”. In: *Advances in Neural Information Processing Systems* 36 (2024).