

# Image Classification: AI-Generated and Real Images

Brinda Muralie  
2021140  
ECE'25

brinda21140@iiitd.ac.in

Parisha Agrawal  
2021270  
CSAI'25

parisha21270@iiitd.ac.in

Sanya Madan  
2021561  
CSAI'25

sanya21561@iiitd.ac.in

## Abstract

*Differentiating between real and AI-generated images has become a complex task for humans due to the remarkable progress in AI image generation using General Adversarial Networks (GANs) and Stable Diffusion Models. Misleading images that convincingly imitate reality can cause significant harm like spreading fake news, damaging reputations, and manipulating public opinion through media. The judicial system also often relies on visual evidence. Identification and mitigation of fake images is crucial to ensure the authenticity of visual content.*

## 1. Introduction

This project aims to develop robust classification methods for distinguishing between real and AI-generated images. Using datasets such as CIFAKE<sup>2</sup> which includes both types of images, extracting discriminative features, and evaluating various models, the project seeks to address the challenge of identifying increasingly realistic AI-generated images, with potential applications in image regulation both online and in real world.

## 2. Literature Survey

- [1] **CIFAKE: Advancing AI-Generated Image Recognition**<sup>2</sup>: The paper *CIFAKE: Advancing AI-Generated Image Recognition*<sup>2</sup> by Jordan J. Bird and Ahmad Lotfi, focuses on the critical necessity of data reliability and authentication in the era of AI-generated synthetic images. The authors propose a method to enhance our ability to recognize these images through computer vision.

The study involves the creation of a synthetic dataset, CIFAKE<sup>2</sup>, which mirrors the ten classes of the CIFAR-10<sup>1</sup> dataset. The dataset is generated using latent diffusion, providing a contrasting set of images for comparison to real photographs. They used prompt modifiers to generate 6,000 images in each class.

The authors then propose the use of a Convolutional Neural Network (CNN) to classify the images into two categories: Real or Fake. After hyperparameter tuning and training of 36 individual network topologies, the optimal approach could correctly classify the images with an accuracy of 92.98%.

To further understand the classification process, the study implements explainable AI via Gradient Class Activation Mapping. This reveals that the model focuses on small visual imperfections in the background of the images for classification, rather than the actual image itself.

The CIFAKE<sup>2</sup> dataset is made publicly available to the research community for future work. This study provides a significant contribution to advancing AI-generated image recognition.

- [2] **GenImage: Advancing AI-Generated Fake Image Detection**<sup>3</sup>:

The paper *GenImage: A Million-Scale Benchmark for Detecting AI-Generated Image*<sup>3</sup> addresses the growing concern about the spread of misinformation due to the extraordinary ability of generative models to generate photographic images. The authors introduce the GenImage dataset.

The GenImage dataset has three main advantages:

1. **Plenty of Images:** It includes over one million pairs of AI-generated fake images and collected real images.
2. **Rich Image Content:** It encompasses a broad range of image classes.
3. **State-of-the-art Generators:** The dataset synthesizes images with advanced diffusion models and GANs.

The authors conducted a comprehensive analysis of the dataset and proposed two tasks for evaluating the detection method in resembling real-world scenarios.

The first task was the cross-generator image classification task which measured the performance of a detector trained on one generator when tested on others. The second task was the degraded image classification task which assessed the capability of detectors in handling degraded images such as low-resolution, blurred, and compressed images.

With the GenImage dataset, researchers can effectively expedite the development and evaluation of superior AI-generated image detectors in comparison to prevailing methodologies.

### 3. Dataset

In this section, we provide a comprehensive overview of the dataset used for our project and detail the preprocessing techniques applied to prepare the data for model training.

#### 3.1. Dataset Description

The dataset utilized in this study is known as CIFAKE<sup>2</sup>, which consists of two primary classes: real and fake images. CIFAKE<sup>2</sup> contains a total of 120,000 images, evenly distributed between synthetically-generated and real images. Key characteristics of the dataset are as follows:

- **Classes:** Both the classes in CIFAKE<sup>2</sup>, real and fake, are divided into ten classes, each representing distinct objects or entities. These classes include: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck, with 6,000 images per class.
- **Data Sources:** The real images used in the dataset were extracted from the CIFAR-10 dataset introduced by Krizhevsky and Hinton<sup>1</sup>. The synthetic images, on the other hand, were generated using Hugging Face’s Stable Diffusion Model Version 1.4. To enhance diversity within the synthetic dataset, different prompt modifiers were applied for each class label.
- **Data Split:** To facilitate training and evaluation of machine learning models, the dataset was partitioned into 100,000 training images (50,000 per class) and 20,000 testing images (10,000 per class). All images were represented in RGB format and resized to a uniform resolution of  $32 \times 32$  pixels.
- **Scaling:** It was observed that no additional scaling was required during preprocessing since all images were already resized to the optimised  $32 \times 32$  pixel dimensions during the data collection process.

#### 3.2. Data Visualisation

The t-SNE plot (Figure 1) is used for dimensionality reduction that is helpful for the visualization of high-dimensional datasets in a lower-dimensional space. The

real and fake images in our dataset are scattered throughout the plot with no clear pattern or clustering. This may be due to loss of information due to reduction in dimensions of images or the fact that it focuses on preserving local similarities, which can distort the global structure of the data. Hence t-SNE is not apt to be used in this context.

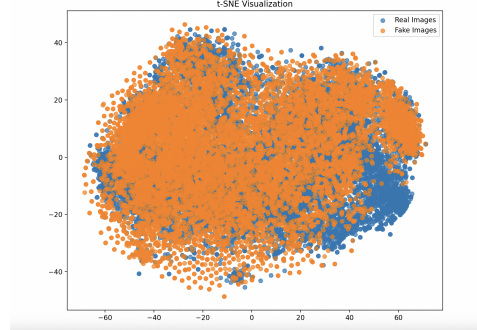


Figure 1. t-SNE Visualisation

The pixel intensity histogram (Figure 2) shows the raw pixel intensity values for all channels red, green, and blue. The x-axis represents the pixel intensity, which ranges from 0 (black) to 255 (white) for each color channel in an 8-bit image. The y-axis represents the frequency of these intensities in the images. There are some differences in the pixel intensity distributions between real and fake images which could be a useful feature for a machine learning models trying to distinguish between the two.

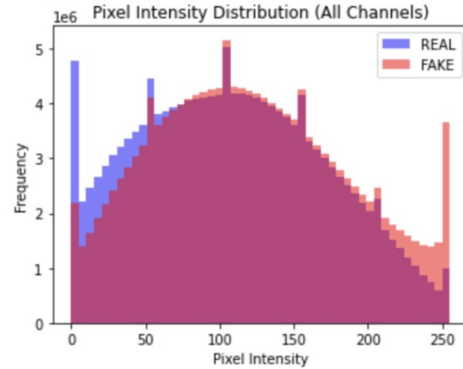


Figure 2. Pixel Intensity Histogram

The blur level histogram (Figure 3) displays the distribution of image blur across the dataset. On the x-axis, lower values signify sharper images, while higher values indicate increased blur. The y-axis represents the frequency of images corresponding to each blur level. Notably, real images tend to have lower blur levels, capturing finer details, while fake images show a concentration of higher blur levels, indicating potential differences in image quality or generation methods.

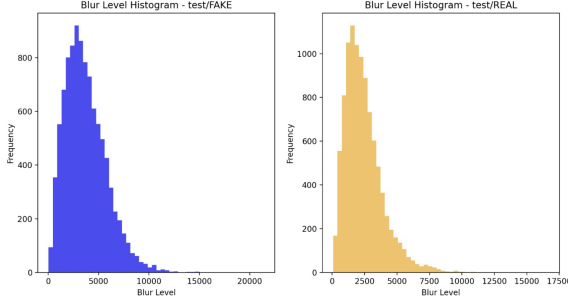


Figure 3. Blur Level Histogram

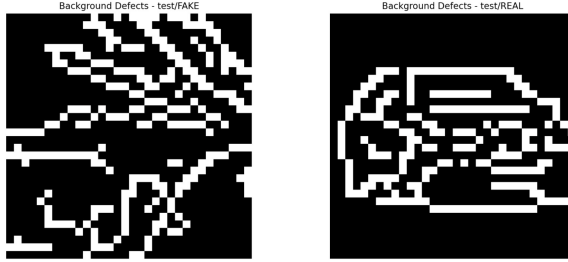


Figure 4. Background Defects

The background defects histogram (Figure 4) illustrates variations in background quality. On the x-axis, it shows the presence or absence of defects, and the y-axis represents the frequency of images with these defects. Real images exhibit fewer background defects, suggesting a more natural background, while fake images show a higher frequency of defects, indicating potential artifacts introduced during generation. These histograms serve as informative features for distinguishing real and fake images based on blur and background characteristics.

### 3.3. Data Preprocessing Techniques

After loading the data from the provided directory, and segregating images into their respective classes ('REAL' or 'FAKE'), The following steps were undertaken to preprocess the dataset in preparation for model training:

1. **Image load and array conversion:** We used the 'imread' function of the CV2 library to load images from the specified directory. Subsequently, we converted the images to numeric data using the 'numpy' library.
2. **Image Resizing:** To ensure uniformity in image dimensions, each image was resized to  $32 \times 32$  pixels using the OpenCV library (For precaution as in data provided each image was  $32 \times 32$  pixels only).
3. **Label Encoding:** Class labels were encoded to numeric values using the LabelEncoder from scikit-learn. This transformation is essential for machine learning model compatibility.

4. **No Outliers:** Using data visualisation and boxplots, we analysed that there are no outliers.

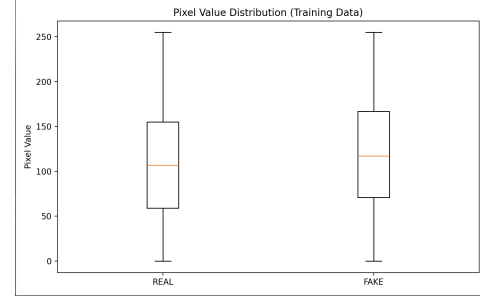


Figure 5. Box Plot

5. **Dimension Reduction using PCA:** Using PCA for dimension reduction is not useful here as the dataset has already provided the optimal reduced dimension. We had analysed this using data visualisation, TSNE as well as training model on different PCA values. Here is a table showing the number of features and the accuracy.

Number of components	Accuracy
10	0.71375
50	0.795
100	0.81175
500	0.81125
1000	0.81325
1024 ( $32 \times 32$ )	0.8135

By following these preprocessing techniques, we ensured that the dataset was suitably prepared for training and evaluation of machine learning models.

## 4. Methodology

- **Logistic Regression :** Logistic regression a fundamental machine learning algorithm used for binary classification. Performed data preprocessing, model creation and training, and evaluation on image data for binary classification (Real vs. AI generated) using both scikit-learn and pytorch. Achieved 0.6793 Validation Accuracy and Test accuracy of 0.67725.

Value	Validation	Test
Accuracy	0.6793	0.67725
Precision	0.66386	0.6926
Recall	0.7218	0.6374
F1-Score	0.6916	0.6638
Specificity	0.6370	0.6374
Confusion Matrix:	[[6393 3642] [2772 7193]]	[[6374 3626] [2829 7171]]
False Positive Rate	0.3629	0.3626

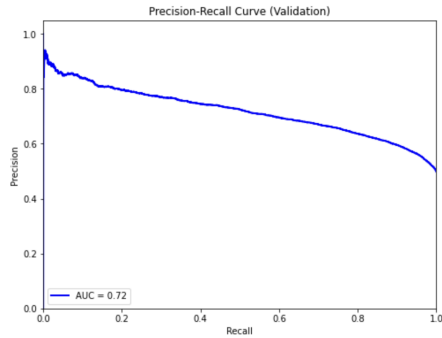


Figure 6. Precision Recall Curve Logistic Regression

- **Naïve Bayes Classifier** : Naive bayes assume features to be independent of each other. Performed Multinomial Naïve Bayes classifier on image data for binary classification (Real vs. AI generated) using scikit-learn. Achieved 0.5893 Validation Accuracy and Test accuracy of 0.59275.

Value	Validation	Test
Accuracy	0.5893	0.59275
Precision	0.59412	0.59570
Recall	0.57269	0.5773
F1-Score	0.58321	0.5863
Specificity	0.60602	0.6082
Confusion Matrix:	[[6039 3926] [4288 5747]]	[[6082 3918] [4227 5773]]
False Positive Rate	0.3939	0.3918

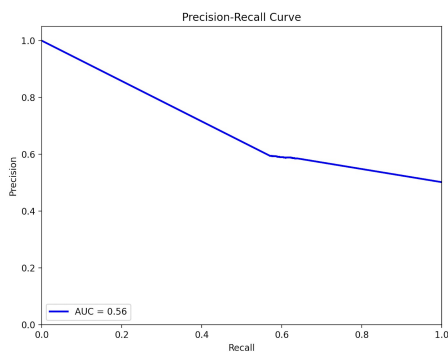


Figure 7. Precision Recall Curve Naïve Bayes

- **Decision Tree Classifier** : Decision Tree is a tree-like model representing decisions and consequences in hierarchical manner. Achieved 0.69345 Validation and 0.69765 Test accuracy on image data for binary classification (Real vs. AI generated) using scikit-learn.

Value	Validation	Test
Accuracy	0.6869	0.69745
Precision	0.68457	0.6983
Recall	0.6972	0.6953
F1-Score	0.6908	0.6967
Specificity	0.6764	0.6996
Confusion Matrix:	[[6741 3224] [3038 6997]]	[[6996 3004] [3047 6953]]
False Positive Rate	0.3235	0.3004

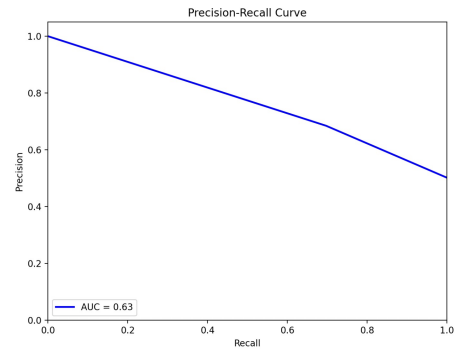


Figure 8. Precision Recall Curve Decision Tree

- **Random Forest Classifier** : Random Forest is a machine learning algorithm for regression and classification tasks, based on the concept of decision tree. Achieved 0.8266 Validation Accuracy and Test accuracy of 0.8307.

Value	Validation	Test
Accuracy	0.8272	0.82925
Precision	0.8498	0.80630
Recall	0.7963	0.8667
F1-Score	0.8222	0.8354
Specificity	0.8583	0.8667
Confusion Matrix:	[[8553 1412] [2044 7991]]	[[8667 1333] [2082 7918]]
False Positive Rate	0.1416	0.1333

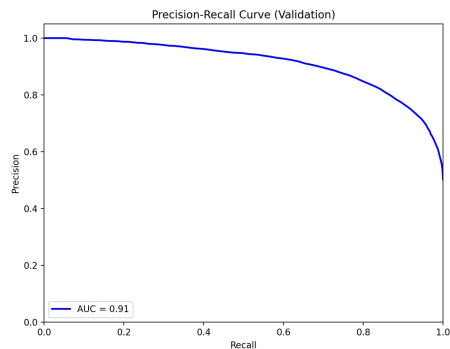


Figure 9. Precision Recall Curve Random Forest

- **Support Vector Machine** : Support vector Machine (SVM) is a supervised machine learning algorithm for regression and classification tasks, and especially effective in high-dimensional spaces. Achieved 0.814 Validation Accuracy and Test accuracy of 0.8108.

Value	Validation	Test
Accuracy	0.814	0.8108
Precision	0.81357	0.8144
Recall	0.8192	0.805
F1-Score	0.8163	0.8096
Specificity	0.8086	0.8166
Confusion Matrix:	[[1602 379] [ 365 1654]]	[[8166 1834] [1950 8050]]
False Positive Rate	0.1913	0.1834

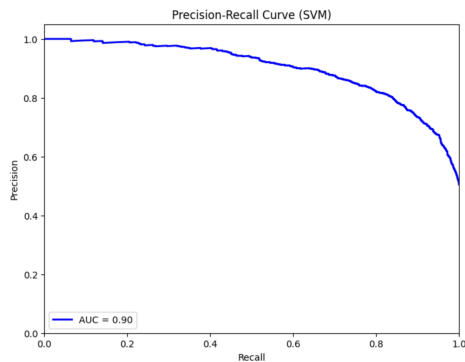


Figure 10. Precision Recall Curve SVM

- **Multi Layer Perceptron (MLP)** : Multi Layer Perceptron (MLP) is a supervised machine learning algorithm and is a type of artificial neural network. Achieved 0.7543 Validation Accuracy and Test accuracy of 0.7495.

Value	Validation	Test
Accuracy	0.7543	0.7495
Precision	0.7328	0.7304
Recall	0.7976	0.7909
F1-Score	0.7638	0.7594
Specificity	0.7112	0.7081
Confusion Matrix:	[[7137 2898] [2016 7949]]	[[7081 2919] [2091 7909]]
False Positive Rate	0.288	0.2919

- **Convolution Neural Network(CNN)** : Convolution Neural Network(CNN) is a type of deep neural network designed for processing structured grid data, such as images or videos. It has three Convolutional layers with increasing filter depths (32, 64, 128) and ReLU activation. It's compiled with Adam optimizer, binary crossentropy loss, and accuracy metric, making it suitable for binary image classification tasks.

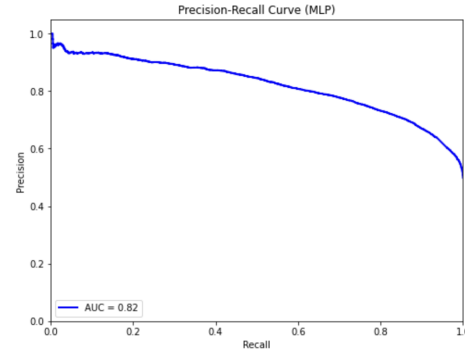


Figure 11. Precision Recall Curve MLP

Achieved 0.9263 Validation Accuracy and Test accuracy of 0.925.

Value	Validation	Test
Accuracy	0.9263	0.925
Precision	0.92022	0.9181
Recall	0.9330	0.9332
F1-Score	0.9155	0.9256
Specificity	0.8788	0.9168
Confusion Matrix:	[[8819 1216] [ 526 9439]]	[[9168 832] [ 668 9332]]
False Positive Rate	0.0803	0.0832

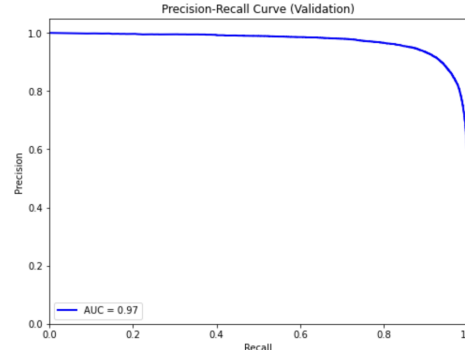


Figure 12. Precision Recall Curve CNN

## 5. Results and analysis

In the precision-recall curves (Figure 7-13), we observed and analysed the performance of our classifiers. These curves are graphical representations of how well a classification model distinguishes between positive and negative cases across different threshold settings. We observed that our best classifier, CNN has a curve that goes straight up the y-axis (recall) and then along the x-axis (precision). Also we observed that the area under the curve (AUC) for the precision-recall curve for the best performing model is maximum (0.97) and the least performing model is minimum (0.56).

In our evaluation of four machine learning models- Logistic Regression, Naïve Bayes, Decision Trees, Random Forests, SVM, MLP, and CNN-we observed varying levels of performance.

CNN achieved the highest test accuracy, reaching 91.47%, making it the top-performing model in our study. It is due to its ability to automatically learn hierarchical features and spatial hierarchies, enabling it to capture intricate patterns and relationships within images.

On the other hand, Naïve Bayes exhibited the poorest performance, with a test accuracy of 59.275%. This outcome is mainly due to its simplistic assumption of feature independence, which is inadequate for image data characterized by pixel correlations and complex spatial relationships.

We used various visualisation methods and found differences in pixel intensities, blur levels, and background defects between fake and real images.

Model	Validation Accuracy	Test Accuracy
Logistic Regression	0.6793	0.67725
Naïve Bayes	0.5893	0.59275
Decision Tree	0.69345	0.69765
Multi Layer Perceptron	0.7543	0.7495
Support Vector Machine	0.814	0.8108
Random Forest	0.8266	0.8307
Convolution Neural Network	0.9263	0.925

## 6. Conclusions

In conclusion, this project addresses the critical challenge of distinguishing real images from AI-generated ones, which have the potential to spread misinformation and manipulate public opinion. By utilizing the CIFAKE<sup>2</sup> dataset and employing various machine learning models, we have achieved promising results.

However, it's essential to acknowledge the project's limitations. Due to the novelty of topic, we couldn't find any suitable papers whose models that we could implement on our own. Since it is an image classification problem, we couldn't meet the hardware requirements of some of the pre-existing models and datasets (for example, the dataset introduced in GenImage<sup>3</sup> was 19 GBs, and our personal systems could not accommodate that). These include potential dataset diversity issues, and the absence of more complex deep learning models which require more computation using better hardware. The project mainly relies on accuracy and AUC as an evaluation metric, which may not fully capture model performance.

Despite these limitations, this project represents a valuable step in tackling the challenges posed by AI-generated images. Overall, this project underscores the ongoing need for research to ensure the authenticity of visual content in

the digital era.

## 7. References

1. Will Cukierski. (2013). CIFAR-10 - Object Recognition in Images. Kaggle. <https://kaggle.com/competitions/cifar-10>
2. Bird, J.J., Lotfi, A. (2023). CIFAKE: Image Classification and Explainable Identification of AI-Generated Synthetic Images. arXiv preprint arXiv:2303.14126.
3. M. Zhu, H. Chen, Q. Yan, X. Huang, G. Lin, W. Li, Z. Tu, H. Hu, J. Hu, and Y. Wang, "Genimage: A million-scale benchmark for detecting ai-generated image," arXiv preprint arXiv:2306.08571, 2023.
4. Haiwei Wu, Jiantao Zhou, Shile Zhang, "Generalizable Synthetic Image Detection via Language-guided Contrastive Learning", arXiv preprint arXiv:2305.13800, 2023.
5. Krizhevsky, A. (2009) Learning Multiple Layers of Features from Tiny Images. Technical Report TR-2009, University of Toronto, Toronto. <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>