
Project Report - ECE 285

Gayatri Kharche

Electrical and Computer Engineering
A69027290

Esha Sali

Electrical and Computer Engineering
A69027358

Abstract

Image inpainting, the process of reconstructing lost or deteriorated parts of an image, is a challenging problem in computer vision. In this project, we propose 3 different architectures for image inpainting. Inspired by Context Encoders, our method leverages unsupervised visual feature learning by predicting pixel context to reconstruct missing image regions. The network's encoder-decoder architecture compresses the input image into a latent representation and reconstructs it, ensuring coherence through a channel-wise fully connected layer. Our approach is validated on the human faces dataset, including, demonstrating its ability to learn meaningful features and perform high-fidelity reconstructions. Quantitative analysis shows the effectiveness of the learned features for tasks like classification, detection, and segmentation. This work underscores the potential of context-based feature learning for image inpainting, suggesting promising avenues for future research and application in various computer vision tasks.

1 Introduction

Image inpainting, the process of reconstructing missing parts of an image in a coherent and visually plausible manner, is a crucial task in the field of computer vision. This problem holds significant value in various applications, such as restoring old and damaged photographs, removing unwanted objects from images, and enhancing visual effects in multimedia.

The motivation behind solving this problem stems from the necessity to automate and improve the accuracy of image restoration techniques. Traditional methods often rely on manual editing, which is time-consuming and requires significant expertise. With the advent of deep learning, we have the opportunity to develop automated systems that can perform inpainting tasks efficiently and with high quality.

Understanding the problem involves recognizing the challenges associated with reconstructing missing parts of an image while maintaining the overall aesthetic and structural integrity. This requires models that can capture both the global context and the fine-grained details of the image.

To tackle this problem, we propose three distinct deep learning architectures:

- **Autoencoder1:** A basic convolutional autoencoder that employs a simple encoder-decoder structure. The encoder compresses the image into a lower-dimensional representation, and the decoder reconstructs it back to the original dimensions. This model serves as our baseline, providing insights into the fundamental capabilities of autoencoders for inpainting tasks.
- **Autoencoder2:** An enhanced version of the basic autoencoder, featuring additional layers and increased depth. This model aims to capture more complex features and provide a more detailed reconstruction. By incorporating more convolutional layers and using transposed convolutions for upsampling, Autoencoder2 seeks to improve the fidelity of the inpainted images.

- A more sophisticated architecture designed initially for biomedical image segmentation. U-Net includes a symmetrical encoder-decoder pathway with skip connections that link corresponding layers in the encoder and decoder. These skip connections help recover spatial information lost during downsampling, making U-Net particularly effective for tasks requiring precise localization.

Each of these models offers unique advantages, and their comparative analysis will shed light on their effectiveness in different scenarios of image inpainting. Preliminary results indicate that while Autoencoder1 provides a solid baseline, Autoencoder2 and U-Net show significant improvements in capturing details and maintaining the structural integrity of the images.

By the end of this project, we aim to demonstrate the strengths and weaknesses of each approach, providing a comprehensive understanding of their capabilities in solving the image inpainting problem.

2 Related Work

In the context of our image inpainting project, we have referred to several notable works that have contributed significantly to the image inpainting and are related to our method.

First, the work by Pathak et al. on "Context Encoders: Feature Learning by Inpainting" provides a foundational approach to unsupervised visual feature learning driven by context-based pixel prediction. The authors propose a method where a convolutional neural network, termed Context Encoders, is trained to generate the contents of an arbitrary image region conditioned on its surroundings. This task requires the network to understand the content of the entire image and produce a plausible hypothesis for the missing parts. They experimented with both standard pixel-wise reconstruction loss and a combination of reconstruction and adversarial loss, finding that the latter produced sharper results by handling multiple modes in the output more effectively. The learned features demonstrated effectiveness for CNN pre-training on classification, detection, and segmentation tasks, and the context encoders could be used for semantic inpainting tasks. Our method draws inspiration from this approach by utilizing autoencoders and UNet architectures to achieve effective image inpainting.

Second, the work on inpainting of irregular holes in old handwritten manuscripts by utilizing deep learning techniques, including UNet and partial convolution, is particularly relevant. This study highlights the challenges of restoring degraded manuscripts due to various factors like human mishandling and environmental conditions. The authors proposed a novel method to inpaint irregular holes with appropriate pixels, focusing on recovering words or letters in the manuscript. They leveraged deep learning's capability to solve complex problems, proposing a UNet architecture to generate masks of irregular holes and partial convolution for the inpainting task. This study informs our project by demonstrating the effectiveness of UNet in generating detailed and accurate inpainting results, particularly in challenging scenarios involving irregular holes and degraded images.

In our project, we utilize architectures inspired by these works, including autoencoders and UNet, to achieve robust image inpainting. The Autoencoder1 and Autoencoder2 architectures we implemented employ convolutional layers for feature extraction and upsampling layers for reconstruction, similar to the encoder-decoder approach used in Context Encoders. Additionally, the UNet architecture we used is directly influenced by the manuscript inpainting work, leveraging its powerful encoder-decoder structure with skip connections to achieve detailed and contextually accurate inpainting results.

By integrating these approaches, our method aims to enhance the inpainting performance, effectively handling various types of image degradation and missing regions, ensuring the restoration of high-quality and semantically coherent images.

3 Method

In this project, we utilize three distinct methods for image inpainting: Autoencoder1, Autoencoder2, and UNet. Each method comprises different network architectures and training algorithms to enhance the inpainting performance.

3.1 Autoencoder1 Methodology

3.1.1 Formulations

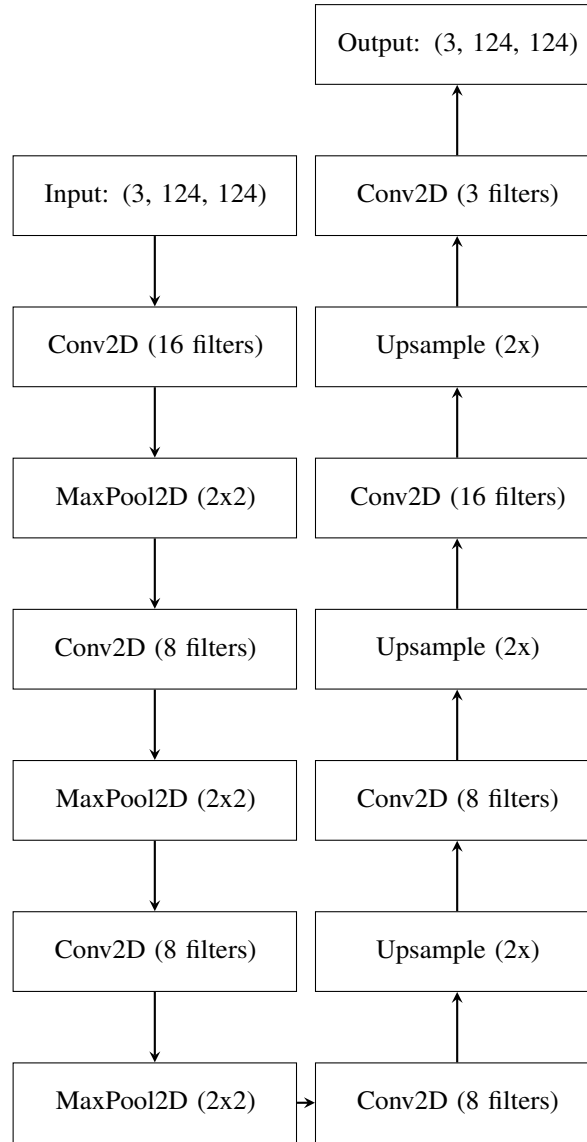
The image inpainting method utilizing Autoencoder1 aims to reconstruct missing or damaged parts of an input image I . Let $I_{\text{inpainted}}$ represent the inpainted image.

The autoencoder architecture comprises an encoder E and a decoder D . The encoder compresses the input image I into a latent representation z , while the decoder reconstructs the original image from this latent representation.

The formulation of the inpainting process using Autoencoder1 can be represented as follows:

$$I_{\text{inpainted}} = D(E(I))$$

3.1.2 Network Architecture



3.1.3 Training Algorithm

The training algorithm for Autoencoder1 involves optimizing the parameters of the network to minimize the reconstruction error between the inpainted images and the ground truth images.

The algorithm can be summarized as follows:

1. Initialize the parameters of Autoencoder1.
2. Iterate over the training dataset:
 - (a) Forward pass: Compute the reconstructed image $I_{\text{inpainted}}$ using Autoencoder1.
 - (b) Compute the reconstruction loss between $I_{\text{inpainted}}$ and the ground truth image I .
 - (c) Backpropagate the gradients and update the parameters using an optimization algorithm (e.g., Adam).

3.1.4 Testing Algorithm

The testing algorithm evaluates the performance of the trained Autoencoder1 on unseen data by inpainting missing patches in input images.

The algorithm can be summarized as follows:

1. Forward pass: Pass the input images through the trained Autoencoder1 to generate inpainted images.
2. Evaluate the quality of the inpainted images using appropriate metrics. Mean Squared Error loss is used here.

3.1.5 Proposed Techniques

The proposed technique in this project is the utilization of Autoencoder1 for image inpainting. Compared to previous works, autoencoders have demonstrated effectiveness in capturing meaningful representations of images and generating plausible completions for missing regions.

The strength of choosing Autoencoder1 lies in its architecture, which learns to compress input images into a lower-dimensional latent space representation and then reconstructs the original images from this representation.

3.2 Autoencoder2 Architecture

3.2.1 Formulation

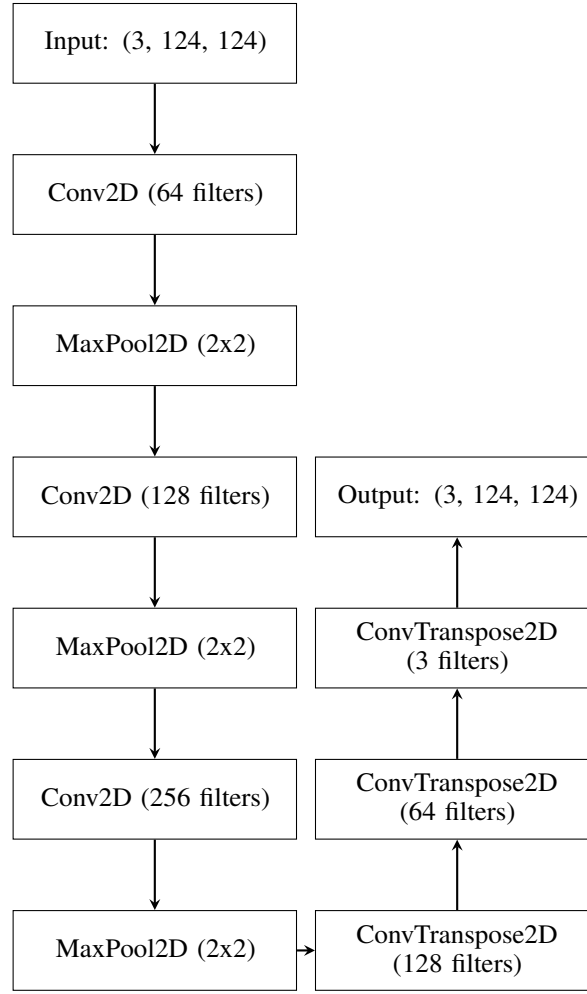
The image inpainting process using the Autoencoder2 model follows a similar formulation as described for Autoencoder1. Given an input image I , the goal is to reconstruct the missing or damaged parts of the image. Let $I_{\text{inpainted}}$ represent the inpainted image.

The autoencoder consists of an encoder E and a decoder D . The encoder compresses the input image I into a latent representation z , while the decoder reconstructs the original image from this latent representation.

The formulation of the inpainting process can be represented as follows:

$$I_{\text{inpainted}} = D(E(I))$$

3.2.2 Network Architecture



3.2.3 Training Algorithm

The training algorithm for Autoencoder2 follows a similar procedure as described earlier for Autoencoder1.

3.2.4 Testing Algorithm

Similarly, the testing algorithm for Autoencoder2 evaluates the performance of the trained model on unseen data by inpainting missing patches in input images.

3.2.5 Proposed Techniques

The proposed technique in this project, utilizing Autoencoder2, aims to explore a deeper and more complex neural network architecture compared to Autoencoder1. By increasing the depth and complexity of the network, we aim to capture more intricate features and patterns in the input images, potentially leading to improved inpainting performance.

3.3 UNet

3.3.1 Formulation

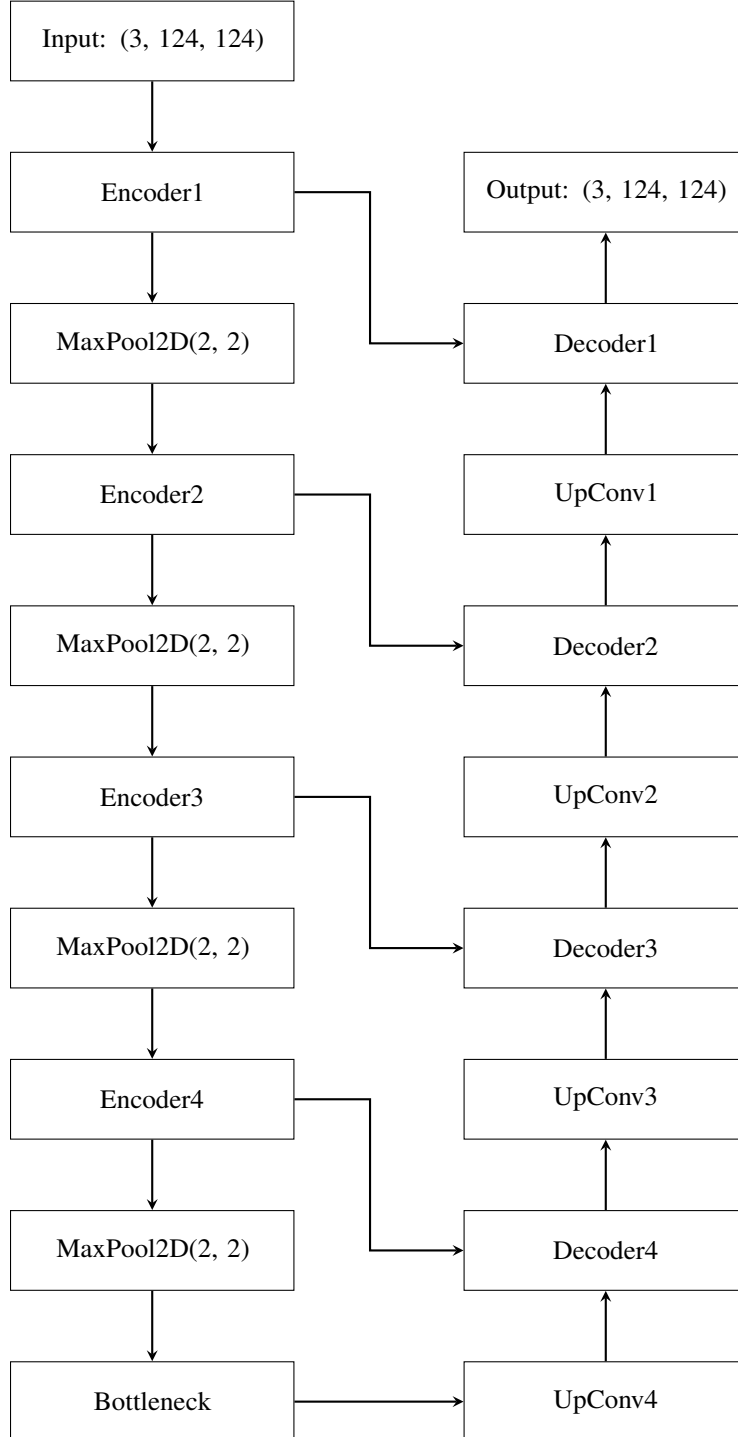
The U-Net architecture is designed for image segmentation tasks, where the goal is to predict a segmentation mask for a given input image. Let I represent the input image and M represent the

predicted segmentation mask. The formulation of the U-Net model can be represented as follows:

$$M = \text{U-Net}(I)$$

The U-Net model consists of an encoder-decoder architecture with skip connections between corresponding encoder and decoder layers. This allows the model to capture both low-level and high-level features from the input image and refine the segmentation predictions.

3.3.2 Network Architecture



3.3.3 Training Algorithm

The training algorithm for the U-Net model involves optimizing the parameters of the network using a suitable optimization algorithm such as Adam. The loss function, typically mean squared error (MSE) or a similar metric, is minimized between the inpainted images generated by the model and the ground truth images.

3.3.4 Testing Algorithm

During testing, the trained U-Net model is used to inpaint missing patches in input images. The model takes the incomplete images as input and generates the corresponding inpainted images using the learned features and parameters.

3.3.5 Proposed Techniques

The U-Net architecture, with its skip connections, offers a unique advantage for image inpainting tasks. The skip connections facilitate the flow of detailed information from early encoder layers to later decoder layers, enabling the model to preserve spatial information and capture fine details during the inpainting process. This approach often leads to more visually appealing and realistic inpainted images compared to traditional autoencoder architectures.

4 Experiments

4.1 Datasets Used

For this project, we utilized the Human Faces dataset, obtained from Kaggle. This dataset contains a diverse collection of human face images sourced from various online repositories. Each image in the dataset represents a human face captured under different lighting conditions, poses, and facial expressions. The dataset serves as a valuable resource for training and evaluating image inpainting models due to its rich and varied content.

The Human Faces dataset is organized in a structured manner, allowing for easy access and integration into machine learning pipelines. Each image is stored in a common image format such as JPEG or PNG, with dimensions typically ranging from a few hundred pixels to a few thousand pixels. The images are accompanied by metadata, which may include information about the identity of the person depicted, the source of the image, and any relevant annotations.

To prepare the dataset for training, we performed preprocessing steps such as resizing all images to a consistent size, applying data augmentation techniques like random rotation and horizontal flipping, and converting the images into a suitable format for input to our neural network models.

Additionally, we implemented a custom dataset class, CustomDataset, to load the images from the dataset directory. This class enables easy loading and preprocessing of images for training and testing purposes. The CustomDataset class takes the root directory of the dataset as input, along with optional transformations such as resizing, normalization, and data augmentation. It allows for efficient handling of large datasets and facilitates seamless integration with PyTorch DataLoader for batch processing during training.

4.2 Results

The performance of our image inpainting models was evaluated on the Human Faces dataset. We implemented three distinct architectures: Autoencoder1, Autoencoder2, and UNet. The results indicate that the UNet model outperformed the other architectures, followed by Autoencoder2, with Autoencoder1 providing a solid baseline.

4.2.1 Autoencoder1:

Loss: The model exhibited moderate reconstruction loss, successfully capturing basic features but lacking fine details. The model demonstrates a significant reduction in the loss over the 30 epochs of training, with an average loss of 0.0134.

Visual Quality: The inpainted images were generally coherent, but some structural integrity and finer details were not well-preserved.

Performance: Autoencoder1 served as an effective baseline, demonstrating the fundamental capabilities of autoencoders for inpainting tasks.

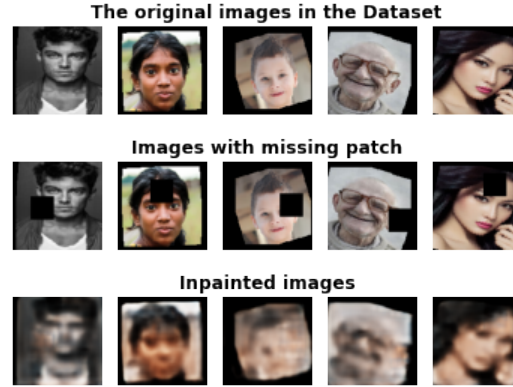


Figure 1: Output for Autoencoder1

4.2.2 Autoencoder2:

Loss: The enhanced architecture resulted in a lower reconstruction loss compared to Autoencoder1, indicating improved learning of complex features. The model demonstrates a significant reduction in the loss over the 30 epochs of training, with an average loss of 0.0091.

Visual Quality: The inpainted images showed better detail and structural coherence, effectively capturing more intricate features.

Performance: Autoencoder2 provided a significant improvement over Autoencoder1, particularly in terms of capturing and reconstructing finer image details.

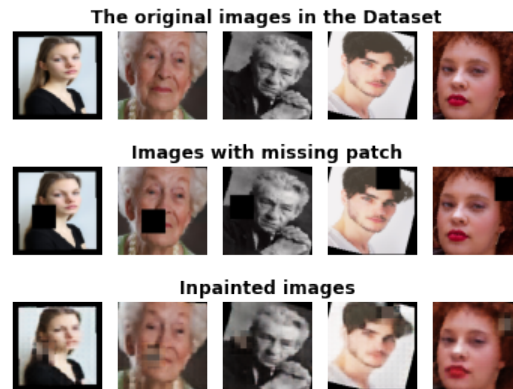


Figure 2: Output for Autoencoder2

4.2.3 Unet:

Loss: The UNet architecture achieved the lowest reconstruction loss among the three models, demonstrating superior feature learning and reconstruction capabilities. The model demonstrates a significant reduction in the loss over the 30 epochs of training, with an average loss of 0.0017.

Visual Quality: The inpainted images were highly detailed and visually appealing, with excellent preservation of spatial information and structural integrity.

Performance: UNet outperformed both Autoencoder1 and Autoencoder2, leveraging its skip connections to achieve high-fidelity reconstructions and maintaining the overall aesthetic and structural quality of the images.

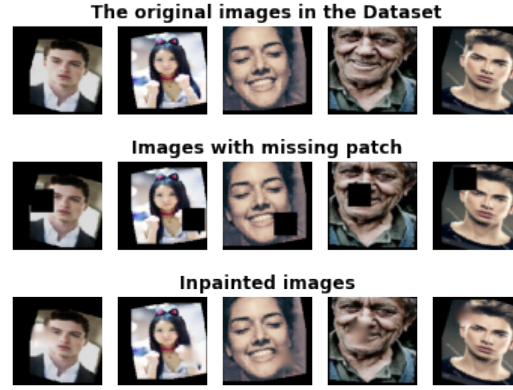


Figure 3: Output for UNet

Overall, the UNet model showed the best performance in our experiments, followed by Autoencoder2 and Autoencoder1. The results highlight the importance of model architecture in image inpainting tasks, with more sophisticated models like UNet providing substantial improvements in both quantitative and qualitative metrics.

References

1. Pathak, Deepak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. "Context encoders: Feature learning by inpainting." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2536-2544. 2016
2. Krizhevsky, Alex, and Geoffrey Hinton. "Learning multiple layers of features from tiny images." (2009): 7.
3. Kulkarni, Amey, Harsh Patel, Shivam Sahni, and Udit Vyas. "Image Inpainting using Partial Convolution." arXiv preprint arXiv:2108.08791 (2021).
4. Yu, Jiahui, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. "Generative image inpainting with contextual attention." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 5505-5514. 2018.
5. Kaur, Amreen, Ankit Raj, N. Jayanthi, and S. Indu. "Inpainting of irregular holes in a manuscript using unet and partial convolution." In 2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA), pp. 778-784. IEEE, 2020.

5 Implementation

In this project, several components were implemented from scratch to ensure a robust and comprehensive evaluation of the proposed image inpainting models. Both Autoencoder1 and Autoencoder2 architectures were designed and implemented from the ground up, including the detailed network layers, encoder-decoder pathways, and the training and testing algorithms. The CustomDataset class was also developed to facilitate the loading and preprocessing of the Human Faces dataset, integrating seamlessly with PyTorch DataLoader for efficient batch processing.

The destroy patches function, which randomly destroys patches in input images by setting pixel values to zero, was custom-built for this project. Additionally, the training and testing algorithms, incorporating loss functions, backpropagation, optimization, and evaluation metrics, were implemented specifically for our needs. A function for plotting reconstructed images was created to visually compare original, destroyed, and inpainted images, further aiding in the qualitative analysis of the results.

Some components were adapted from existing codes to leverage established methods. The UNet architecture was borrowed and modified for image inpainting tasks, utilizing its powerful encoder-decoder structure with skip connections. Standard data augmentation techniques, such as random rotation and horizontal flipping, were implemented using PyTorch functionalities. Moreover, existing PyTorch functionalities were used for optimization (e.g., Adam) and loss functions (e.g., MSE), ensuring efficient and effective training processes.

This blend of custom implementations and adapted components allowed for a robust and thorough examination of the different models, highlighting the strengths and weaknesses of each approach in solving the image inpainting problem.