

Automated Teeth Disease Detection Using Deep Learning

Ahmed Mostafa Gamal Eddin

Computer Vision intern

Cellula Technologies

July 2, 2025

Contents

1	Introduction	2
2	Problem Statement	2
3	Dataset Description	2
4	procedures and steps	3
5	Load and explore the data	3
6	plot the distribution of the data	4
7	Plotting some photos for each class	4
7.1	some random images	5
8	Preprocessing and Augmentation	5
8.0.1	Image Augmentation (Training Set Only)	5
8.0.2	Image Normalization (Validation Set)	6
8.0.3	Data Generator	6
8.1	Code Snippet	6
8.2	Augmented vs original	7
9	CNN Model Architecture	8
10	Training Configuration	9
11	train and Evaluation	9
12	Conclusion	12

1 Introduction

Teeth diseases can lead to severe health problems if left undiagnosed. However, manual classification by professionals is time-consuming and error-prone. With the rise of deep learning, especially Convolutional Neural Networks (CNNs), automated image classification can aid medical professionals in faster diagnosis.

2 Problem Statement

The goal of this project is to develop a CNN model that can accurately classify teeth diseases from images into one of seven categories:

- OC (Oral Cancer)
- CaS (Caries)
- OT (Other Types)
- CoS (Cosmetic Issues)
- Gum (Gum Disease)
- MC (Molar Conditions)
- OLP (Oral Lichen Planus)

3 Dataset Description

The dataset is divided into three subsets:

1. **Training:** Used to train the model
2. **Validation:** Used to tune hyperparameters and monitor performance
3. **Testing:** Final evaluation

Dataset Hierarchy

```
Teeth_Dataset/  
  Training/  
  Validation/  
  Testing/
```

4 procedures and steps

- Load and Explore the data by showing images for each class
- Apply data augmentation to enhance generalization.
- Apply data augmentation to enhance generalization.
- train CNN model to process the data
- evaluate the model using various metrics like accuracy classification report .

5 Load and explore the data

The dataset is organized into three main directories: Training, Validation, and Testing, each containing subfolders named after specific disease classes. These subfolders contain the corresponding images. The loading process reads these directories and automatically assigns labels based on folder names, ensuring that the data is structured appropriately for supervised learning.

Code Snippet

```
1 base_path = "/kaggle/input/teeth-dataset/Teeth_Dataset"
2
3 train_data = tf.keras.utils.image_dataset_from_directory(
4     base_path + "/Training",
5     image_size=(224, 224),
6     batch_size=32
7 )
8
9 val_data = tf.keras.utils.image_dataset_from_directory(
10     base_path + "/Validation",
11     image_size=(224, 224),
12     batch_size=32
13 )
14
15 test_data = tf.keras.utils.image_dataset_from_directory(
16     base_path + "/Testing",
17     image_size=(224, 224),
18     batch_size=32,
19     shuffle=False
20 )
21 )
```

Listing 1: Load the data

6 plot the distribution of the data

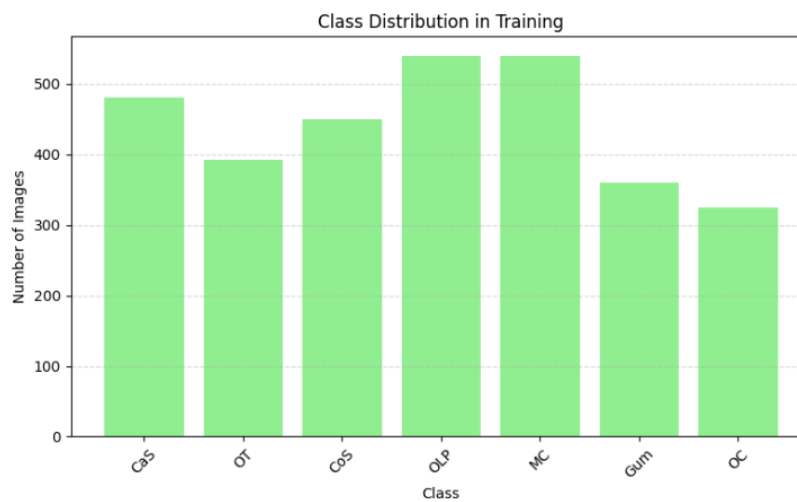


Figure 1: Distribution of dataset

7 Plotting some photos for each class

Code Snippet

```

1 Illness = ['OC', 'CaS', 'OT', 'CoS', 'Gum', 'MC', 'OLP']
2 Types = ['Validation', 'Training', 'Testing']
3 def plot_one_image_per_illness(category, illness_list):
4
5     fig, axes = plt.subplots(1, len(illness_list), figsize=(10, 10))
6     if len(illness_list) == 1:
7         axes = [axes]
8
9     for ax, illness in zip(axes, illness_list):
10         # illness_dir = "/kaggle/input/teeth-dataset/Teeth_Dataset/OC"
11         illness_dir = os.path.join(base_path, category, illness)
12         valid_extensions = ('.jpeg', '.jpg', '.bmp', '.png')
13         #illness_image = list of the path of each image
14         illness_images = list(filter(lambda x: x.endswith(valid_extensions), os.listdir(illness_dir)))
15
16         if illness_images:
17             img_path = os.path.join(illness_dir, random.choice(
illness_images))
18             img = plt.imread(img_path)
19             ax.imshow(img)
20             ax.set_title(illness)
21             ax.axis('off')
22         else:
23             ax.set_title(f'No images for {illness}')

```

```

24         ax.axis('off')
25
26     plt.tight_layout()
27     plt.show()
28
29
30 plot_one_image_per_illness(Types[0], Illness)
31 plt.suptitle('Random Images [Validation Dataset]:', y=1.00)
32
33 plot_one_image_per_illness(Types[1], Illness)
34 plt.suptitle('Random Images [Training Dataset]:', y=1.00)
35
36 plot_one_image_per_illness(Types[2], Illness)
37 plt.suptitle('Random Images [Test Dataset]:', y=1.00)

```

Listing 2: Load the data

7.1 some random images



Figure 2: Distribution of dataset

8 Preprocessing and Augmentation

The `Load_and_preprocessing()` function is responsible for preparing the training and validation datasets before feeding them into the Convolutional Neural Network (CNN). It performs two essential tasks:

8.0.1 Image Augmentation (Training Set Only)

To improve the generalization of the model and prevent overfitting, the training images undergo a series of random transformations:

- **Rescaling:** Pixel values are normalized from the range $[0, 255]$ to $[0, 1]$.

- **Data Augmentation Techniques:**

- Random rotation
- Width and height shifts
- Shearing transformation
- Zooming in and out
- Horizontal flipping
- Filling of missing pixels using the nearest strategy

8.0.2 Image Normalization (Validation Set)

For the validation set, only rescaling is applied:

- **Rescaling:** Pixel values are normalized to the $[0, 1]$ range.
- **No augmentation:** Validation images remain unchanged to ensure an accurate evaluation of model performance.

8.0.3 Data Generator

Both training and validation sets are passed through the `flow_from_directory()` method which:

- Automatically assigns labels based on folder names.
- Resizes images to the specified `target_size` (e.g., 150x150).
- Returns batches of image-label pairs using one-hot encoding.

8.1 Code Snippet

```

1 def Load_and_preprocessing(train_data, validation_data, target_size,
   batch_size):
2     train_datagen = ImageDataGenerator(
3         rescale=1.0/255,
4         rotation_range=20,
5         width_shift_range=0.2,
6         height_shift_range = 0.2,
7         shear_range= 0.2,
8         zoom_range=0.2,
9         horizontal_flip=True,
10        fill_mode='nearest'
11    )
12
13    validation_datagen = ImageDataGenerator(rescale=1./255)
14
15    train_generator = train_datagen.flow_from_directory(
16        train_data,

```

```

17     target_size = (150,150),
18     batch_size=batch_size,
19     class_mode = 'categorical'
20 )
21
22 validation_generator = validation_datagen.flow_from_directory(
23     validation_data,
24     target_size=(150, 150),
25     batch_size=batch_size,
26     class_mode='categorical'
27 )
28
29
30 return train_generator, validation_generator

```

Listing 3: Image Preprocessing and Augmentation

8.2 Augmented vs original



Figure 3: Augmented vs original

9 CNN Model Architecture

The Convolutional Neural Network (CNN) model used in this project is designed to classify input dental images into 7 disease categories. It follows a sequential architecture consisting of multiple convolutional and pooling layers, followed by dense layers for classification.

Layer-by-Layer Breakdown

- **Conv2D (32 filters, 3×3 kernel, ReLU activation):**
Extracts 32 local features from the input image using small 3×3 filters with non-linearity introduced via ReLU.
- **MaxPooling2D (2×2 pool size):**
Downsamples the feature maps by taking the maximum value in each 2×2 region, reducing dimensionality and computation.
- **Conv2D (64 filters, 3×3 kernel, ReLU activation):**
Increases the depth of the network to learn more complex patterns and features.
- **MaxPooling2D (2×2 pool size):**
Further reduces the spatial dimensions to retain only the most prominent features.
- **Conv2D (128 filters, 3×3 kernel, ReLU activation):**
Adds more abstraction by learning high-level features, useful for distinguishing fine-grained disease types.
- **MaxPooling2D (2×2 pool size):**
Reduces spatial dimensions again to prepare for the dense layers.
- **Flatten:**
Converts the 2D output of the last pooling layer into a 1D vector for input into the fully connected layers.
- **Dense (512 units, ReLU activation):**
A fully connected layer acting as the decision-making layer with ReLU activation to introduce non-linearity.
- **Dropout (rate = 0.5):**
Randomly drops 50% of neurons during training to prevent overfitting and improve generalization.
- **Dense (7 units, Softmax activation):**
Final output layer that produces a probability distribution over the 7 disease classes for classification.

Code Snippet

```
1 model = Sequential([
2     Conv2D(32,(3,3),activation='relu', input_shape=(150,150,3)),
3     MaxPooling2D((2,2)),
4     Conv2D(64,(3,3),activation='relu'),
5     MaxPooling2D((2,2)),
6     Conv2D(128,(3,3),activation='relu'),
7     MaxPooling2D((2,2)),
8     Flatten(),
9     Dense(512, activation='relu'),
10    Dropout(0.5),
11    Dense(7, activation='softmax')
12 ])
```

Listing 4: CNN Architecture

10 Training Configuration

Loss Function: Categorical Crossentropy

Optimizer: Adam

Metrics: Accuracy

Callbacks: EarlyStopping and ReduceLROnPlateau

```
1 model.compile(
2     optimizer=Adam(learning_rate=0.001),
3     loss='categorical_crossentropy',
4     metrics=['accuracy']
5 )
6 early_stopping = EarlyStopping(monitor='val_loss', patience=4,
7     restore_best_weights=True)
8 reduce_lr = ReduceLROnPlateau(monitor='val_loss', factor=0.2, patience=5,
9     min_lr=0.00001)
```

Listing 5: Callbacks for Training

11 train and Evaluation

After 100 epochs the model achieved the following:

- Training Accuracy: accuracy: 0.9112 - loss: 0.2765
- Validation Accuracy: accuracy: 0.9786 - loss: 0.0863
- Test Accuracy: accuracy: 0.9838 - loss: 0.0679
- Classification Report and Confusion Matrix generated using `sklearn.metrics`

Visualizations

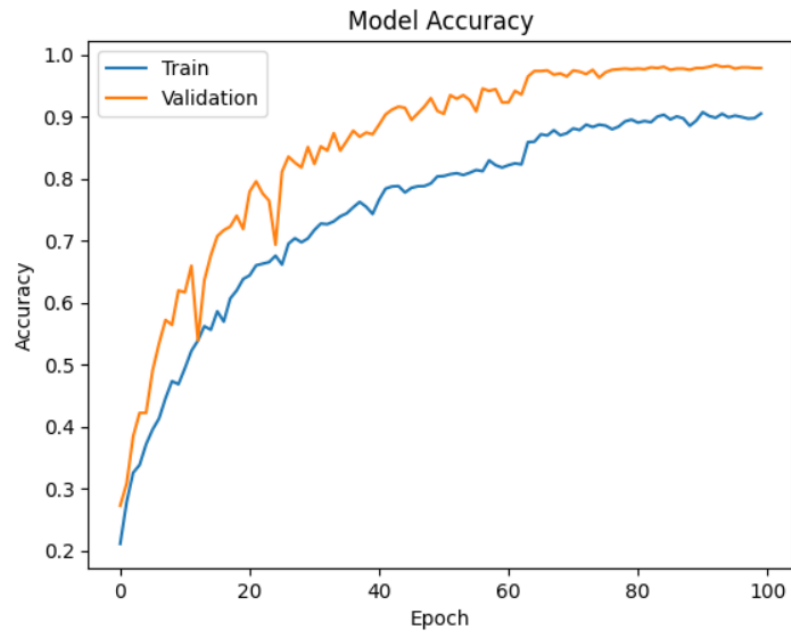


Figure 4: Training and Validation accuracy over Epochs

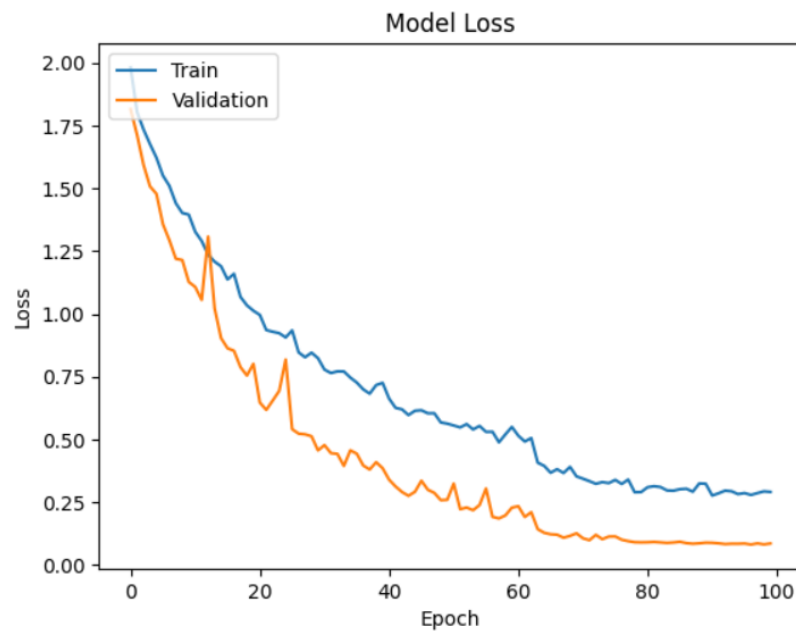


Figure 5: Training and Validation loss over Epochs

Classification Report

	precision	recall	f1-score	support
CaS	0.99	0.99	0.99	160
CoS	1.00	0.99	0.99	149
Gum	0.99	0.98	0.99	120
MC	0.97	0.96	0.97	180
OC	0.95	0.98	0.96	108
OLP	0.98	0.98	0.98	180
OT	0.98	0.99	0.98	131
accuracy			0.98	1028
macro avg	0.98	0.98	0.98	1028
weighted avg	0.98	0.98	0.98	1028

Figure 6: classification report

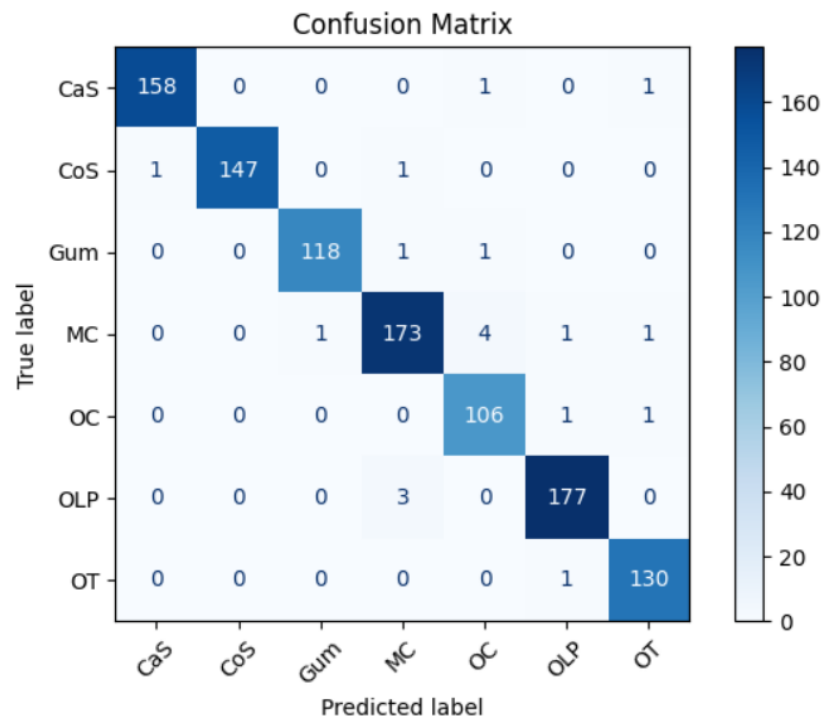


Figure 7: confusion matrix

12 Conclusion

In this project, we built a deep learning pipeline for classifying teeth diseases with decent accuracy. The CNN model learned to differentiate between 7 types of diseases using real-world image data. With further tuning, it could support clinical diagnosis as a decision support tool.