

Data Collection and Preprocessing Phase

Date	14 June 2025
Team ID	SWTID1749876754
Project Title	SynapseScan- AI Driven Classification of Ovarian Cancer Variants
Maximum Marks	6 Marks

Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers in ovarian cancer medical imaging data, with Python employed for preprocessing tasks like image normalization, data augmentation, and feature engineering. Data cleaning will address missing values and image quality issues, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for accurate cancer variant classification and predictions.

Section	Description
Data Overview	<p>Medical imaging dataset containing histopathological images of ovarian cancer variants</p> <p>Both training set and test set were divided into 5 subfolders: CC, EC, HGSC, LGSC, MC, which are the 5 main types of ovarian cancer.</p> <p>Each subfolder in training set had between 3000-12000 images.</p> <p>Each subfolder in test set had about 800 images.</p>

Data Preprocessing Code Screenshots

Preprocessing Step	Implementation
Loading Data	<pre>datadownload.py 1 import os 2 os.environ['KAGGLE_CONFIG_DIR'] = os.path.join(os.getcwd(), ".kaggle") 3 import subprocess 4 subprocess.run([5 "kaggle", "datasets", "download", 6 "-d", "sunilthite/ovarian-cancer-classification-dataset", 7 "--unzip" 8], check=True) 9 print('Data downloaded successfully') 10</pre>
Handling Missing Data	<pre>48 # Oversampling to fix discrepancy between classes in train set 49 ros = RandomOverSampler(random_state=42) 50 X_resampled, y_resampled = ros.fit_resample(df[['file_path']], df['category_encoded']) 51 df_resampled = pd.DataFrame(X_resampled, columns=['file_path']) 52 df_resampled['category_encoded'] = y_resampled 53</pre>
Data Transformation	<pre>79 tr_gen = ImageDataGenerator(rescale=1./255) 80 ts_gen = ImageDataGenerator(rescale=1./255) 81 82 train_gen_new = tr_gen.flow_from_dataframe(83 train_df_new, 84 x_col='file_path', 85 y_col='category_encoded', 86 target_size=img_size, 87 class_mode='sparse', 88 color_mode='rgb', 89 shuffle=True, 90 batch_size=batch_size 91)</pre>

	<pre> 43 df = pd.DataFrame(data, columns=['file_path', 'label']) 44 label_encoder = LabelEncoder() 45 df['category_encoded'] = label_encoder.fit_transform(df['label']) 46 df = df[['file_path', 'category_encoded']] 47 48 # Oversampling to fix discrepancy between classes in train set 49 ros = RandomOverSampler(random_state=42) 50 X_resampled, y_resampled = ros.fit_resample(df[['file_path']], df['category_encoded']) 51 df_resampled = pd.DataFrame(X_resampled, columns=['file_path']) 52 df_resampled['category_encoded'] = y_resampled 53 54 # Split dataset 55 df_resampled['category_encoded'] = df_resampled['category_encoded'].astype(str) 56 57 train_df_new, temp_df_new = train_test_split(58 df_resampled, 59 train_size=0.8, 60 shuffle=True, 61 random_state=42, 62 stratify=df_resampled['category_encoded'] 63) 64 65 valid_df_new, test_df_new = train_test_split(66 temp_df_new, 67 test_size=0.5, 68 shuffle=True, 69 random_state=42, 70 stratify=temp_df_new['category_encoded'] 71) </pre>
Feature Engineering	<pre> 154 def create_inception_model(input_shape): 155 inputs = Input(shape=input_shape) 156 base_model = InceptionV3(weights='imagenet', input_tensor=inputs, include_top=False) 157 for layer in base_model.layers: 158 layer.trainable = False 159 160 x = base_model.output 161 height, width, channels = 5, 5, 2048 162 x = Reshape((height * width, channels))(x) 163 attention_output = DifferentialAttention(num_heads=8, key_dim=channels)(x) 164 attention_output = Reshape((height, width, channels))(attention_output) </pre>
Save Processed Data	<pre> 33 dataset_path = "Train_Images" 34 data = [] 35 36 for label in os.listdir(dataset_path): 37 sub_dir = os.path.join(dataset_path, label) 38 if os.path.isdir(sub_dir): 39 for file_name in os.listdir(sub_dir): 40 file_path = os.path.join(sub_dir, file_name) 41 data.append([file_path, label]) </pre> <pre> 200 cnn_model.save(model_path) 201 print(f"Model saved to {model_path}") </pre>