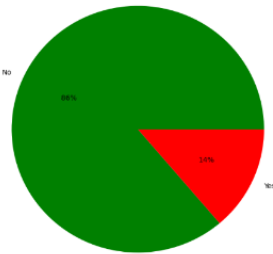
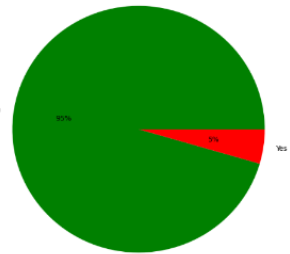


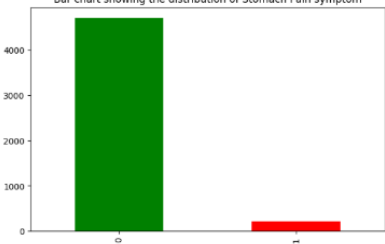
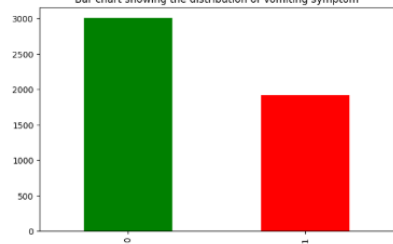
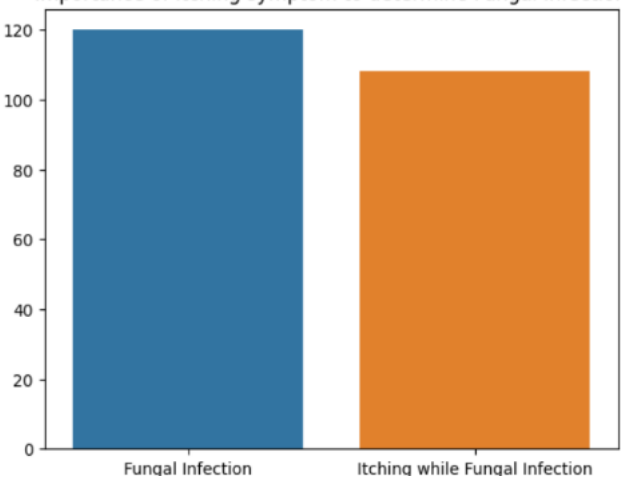
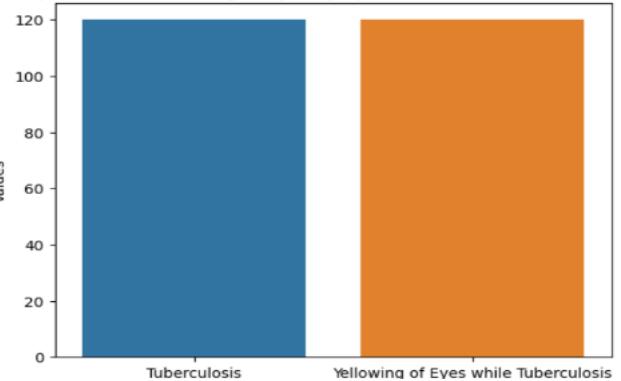
Data Collection and Preprocessing Phase

Date	15 June 2024
Team ID	740003
Project Title	Disease prediction using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description																																																																																																			
Data Overview	<div>Dimension: 8 rows x 131 columns</div> <div>Descriptive statistics:</div> <div><pre>train_data.describe()</pre><table><thead><tr><th></th><th>itching</th><th>skin_rash</th><th>nodal_skin_eruptions</th><th>continuous_sneezing</th><th>shivering</th><th>chills</th><th>joint_pain</th><th>stomach_pain</th><th>acidity</th><th>ulcers_on_t</th></tr></thead><tbody><tr><td>count</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td><td>4920.000000</td></tr><tr><td>mean</td><td>0.137805</td><td>0.159756</td><td>0.021951</td><td>0.045122</td><td>0.021951</td><td>0.162195</td><td>0.139024</td><td>0.045122</td><td>0.045122</td><td>0.021951</td></tr><tr><td>std</td><td>0.344730</td><td>0.366417</td><td>0.146539</td><td>0.207593</td><td>0.146539</td><td>0.368667</td><td>0.346007</td><td>0.207593</td><td>0.207593</td><td>0.146539</td></tr><tr><td>min</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr><tr><td>25%</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr><tr><td>50%</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr><tr><td>75%</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td></tr><tr><td>max</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td></tr></tbody></table><div>8 rows x 131 columns</div></div>		itching	skin_rash	nodal_skin_eruptions	continuous_sneezing	shivering	chills	joint_pain	stomach_pain	acidity	ulcers_on_t	count	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	mean	0.137805	0.159756	0.021951	0.045122	0.021951	0.162195	0.139024	0.045122	0.045122	0.021951	std	0.344730	0.366417	0.146539	0.207593	0.146539	0.368667	0.346007	0.207593	0.207593	0.146539	min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	max	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
		itching	skin_rash	nodal_skin_eruptions	continuous_sneezing	shivering	chills	joint_pain	stomach_pain	acidity	ulcers_on_t																																																																																									
count	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000	4920.000000																																																																																										
mean	0.137805	0.159756	0.021951	0.045122	0.021951	0.162195	0.139024	0.045122	0.045122	0.021951																																																																																										
std	0.344730	0.366417	0.146539	0.207593	0.146539	0.368667	0.346007	0.207593	0.207593	0.146539																																																																																										
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000																																																																																										
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000																																																																																										
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000																																																																																										
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000																																																																																										
max	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000																																																																																										
Univariate Analysis	<div><div>Pie chart showing the distribution of Itching symptom into number of Yes/No</div><div></div></div> <div><div>Pie Chart showing the distribution of Continuous Sneezing symptom into number of Yes/No</div><div></div></div>																																																																																																			

	<div> <div> <p>Bar chart showing the distribution of Stomach Pain symptom</p>  </div> <div> <p>Bar chart showing the distribution of Vomiting symptom</p>  </div> </div>
Bivariate Analysis	<p>Importance of Itching symptom to determine Fungal Infection</p> 
	<p>Importance of Yellowing of Eyes symptom to determine Tuberculosis</p> 

[illegible]

Handling Missing Data In train and test

```
[ ] train_data.isnull().sum()
```

```

↳ itching            0
   skin_rash         0
   nodal_skin_eruptions  0
   continuous_sneezing  0
   shivering          0
   ...
   blister            0
   red_sore_around_nose  0
   yellow_crust_ooze   0
   prognosis          0
   Unnamed: 133        4920
   Length: 134, dtype: int64

```

```
[ ] train_data.isna().sum().sum()
```

```
↳ 4920
```

REMOVING NULL COLUMNS IN TRAINING DATA

```
[ ] train_data['Unnamed: 133'].value_counts()
```

```
↳ Series([], Name: count, dtype: int64)
```

```
[ ] train_data.drop("Unnamed: 133",axis = 1,inplace=True)
   train_data.drop("fluid_overload",axis = 1,inplace=True)
```

```
[ ] train_data.shape
```

```
↳ (4920, 132)
```

```
[ ] test_data.isnull().sum()
```

```

↳ itching            0
   skin_rash         0
   nodal_skin_eruptions  0
   continuous_sneezing  0
   shivering          0
   ..
   inflammatory_nails  0
   blister            0
   red_sore_around_nose  0
   yellow_crust_ooze   0
   prognosis          0
   Length: 133, dtype: int64

```

```
test_data.drop("fluid_overload",axis = 1,inplace=True)
```

Data Transformation	<pre> from sklearn.preprocessing import LabelEncoder label_encoder = LabelEncoder() train_data['prognosis'] = label_encoder.fit_transform(train_data['prognosis']) train_data['prognosis'].unique() array([15, 4, 16, 9, 14, 33, 1, 12, 17, 6, 23, 30, 7, 32, 28, 29, 8, 11, 37, 40, 19, 20, 21, 22, 3, 36, 10, 34, 13, 18, 39, 26, 24, 25, 31, 5, 0, 2, 38, 35, 27]) [] label_encoder = LabelEncoder() test_data['prognosis'] = label_encoder.fit_transform(test_data['prognosis']) test_data['prognosis'].unique() array([15, 4, 16, 9, 14, 33, 1, 12, 17, 6, 23, 30, 7, 32, 28, 29, 8, 11, 37, 40, 19, 20, 21, 22, 3, 36, 10, 34, 13, 18, 39, 26, 24, 25, 31, 5, 0, 2, 38, 35, 27]) </pre>
Feature Engineering	-
Save Processed Data	-