

# Disease Prediction Using Machine Learning

This article aims to implement a robust machine-learning model that can efficiently predict the disease of a human, based on the symptoms that he/she possesses. Let us look into how we can approach this machine-learning problem

## Data Cleaning

```
In [13]: ▶ 1 import pandas as pd
2 import numpy as np
3 import pandas as pd
4 from scipy.stats import mode
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 from sklearn.preprocessing import LabelEncoder
8 from sklearn.model_selection import train_test_split, cross_val_score
9 from sklearn.svm import SVC
10 from sklearn.naive_bayes import GaussianNB
11 from sklearn.ensemble import RandomForestClassifier
12 from sklearn.metrics import accuracy_score, confusion_matrix
13
14 %matplotlib inline
```

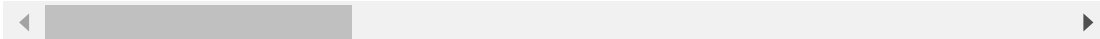
```
In [29]: ▶ 1 #disease = pd.read_csv("Training_disease.csv")
2 DATA_PATH = "Training_disease.csv"
3 data = pd.read_csv(DATA_PATH).dropna(axis = 1)
```

In [30]: 1 data.head(10)

Out[30]:

	itching	skin_rash	nodal_skin_eruptions	continuous_sneezing	shivering	chills	joint
0	1	1	1	0	0	0	
1	0	1	1	0	0	0	
2	1	0	1	0	0	0	
3	1	1	0	0	0	0	
4	1	1	1	0	0	0	
5	0	1	1	0	0	0	
6	1	0	1	0	0	0	
7	1	1	0	0	0	0	
8	1	1	1	0	0	0	
9	1	1	1	0	0	0	

10 rows × 133 columns



In [31]: 1 data.columns

Out[31]: Index(['itching', 'skin\_rash', 'nodal\_skin\_eruptions', 'continuous\_sneezing', 'shivering', 'chills', 'joint\_pain', 'stomach\_pain', 'acidity', 'ulcers\_on\_tongue', ..., 'blackheads', 'scurring', 'skin\_peeling', 'silver\_like\_dusting', 'small\_dents\_in\_nails', 'inflammatory\_nails', 'blister', 'red\_sore\_around\_nose', 'yellow\_crust\_ooze', 'prognosis'], dtype='object', length=133)

In [32]:

1 data.isnull()

Out[32]:

	itching	skin_rash	nodal_skin_eruptions	continuous_sneezing	shivering	chills	j
0	False	False	False	False	False	False	
1	False	False	False	False	False	False	
2	False	False	False	False	False	False	
3	False	False	False	False	False	False	
4	False	False	False	False	False	False	
...	...	...	...	...	...	...	
4915	False	False	False	False	False	False	
4916	False	False	False	False	False	False	
4917	False	False	False	False	False	False	
4918	False	False	False	False	False	False	
4919	False	False	False	False	False	False	

4920 rows × 133 columns

In [33]: 1 `sum(data["Unnamed: 133"].isnull())`

```
-----
-----
KeyError                                Traceback (most recent call
last)
File C:\ProgramData\anaconda3\lib\site-packages\pandas\core\indexes\ba
se.py:3802, in Index.get_loc(self, key, method, tolerance)
    3801 try:
-> 3802     return self._engine.get_loc(casted_key)
    3803 except KeyError as err:

File C:\ProgramData\anaconda3\lib\site-packages\pandas\_libs\index.py
x:138, in pandas._libs.index.IndexEngine.get_loc()

File C:\ProgramData\anaconda3\lib\site-packages\pandas\_libs\index.py
x:165, in pandas._libs.index.IndexEngine.get_loc()

File pandas\_libs\hashtable_class_helper.pxi:5745, in pandas._libs.has
htable.PyObjectHashTable.get_item()

File pandas\_libs\hashtable_class_helper.pxi:5753, in pandas._libs.has
htable.PyObjectHashTable.get_item()

KeyError: 'Unnamed: 133'
```

The above exception was the direct cause of the following exception:

```
KeyError                                Traceback (most recent call
last)
Cell In[33], line 1
----> 1 sum(data["Unnamed: 133"].isnull())

File C:\ProgramData\anaconda3\lib\site-packages\pandas\core\frame.py:3
807, in DataFrame.__getitem__(self, key)
    3805 if self.columns.nlevels > 1:
    3806     return self._getitem_multilevel(key)
-> 3807 indexer = self.columns.get_loc(key)
    3808 if is_integer(indexer):
    3809     indexer = [indexer]

File C:\ProgramData\anaconda3\lib\site-packages\pandas\core\indexes\ba
se.py:3804, in Index.get_loc(self, key, method, tolerance)
    3802     return self._engine.get_loc(casted_key)
    3803 except KeyError as err:
-> 3804     raise KeyError(key) from err
    3805 except TypeError:
    3806     # If we have a listlike key, _check_indexing_error will ra
ise
    3807     # InvalidIndexError. Otherwise we fall through and re-rai
se
    3808     # the TypeError.
    3809     self._check_indexing_error(key)

KeyError: 'Unnamed: 133'
```

In [9]: 1 `print(disease.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
```

In [34]: 1 `data = data.dropna()` *# Dropping the missing values.*  
2 `data.count()`

Out[34]:

```
itching            4920
skin_rash          4920
nodal_skin_eruptions  4920
continuous_sneezing  4920
shivering          4920
...
inflammatory_nails  4920
blister            4920
red_sore_around_nose  4920
yellow_crust_ooze   4920
prognosis          4920
Length: 133, dtype: int64
```

In [35]: 1 `print(data.isnull().sum())`

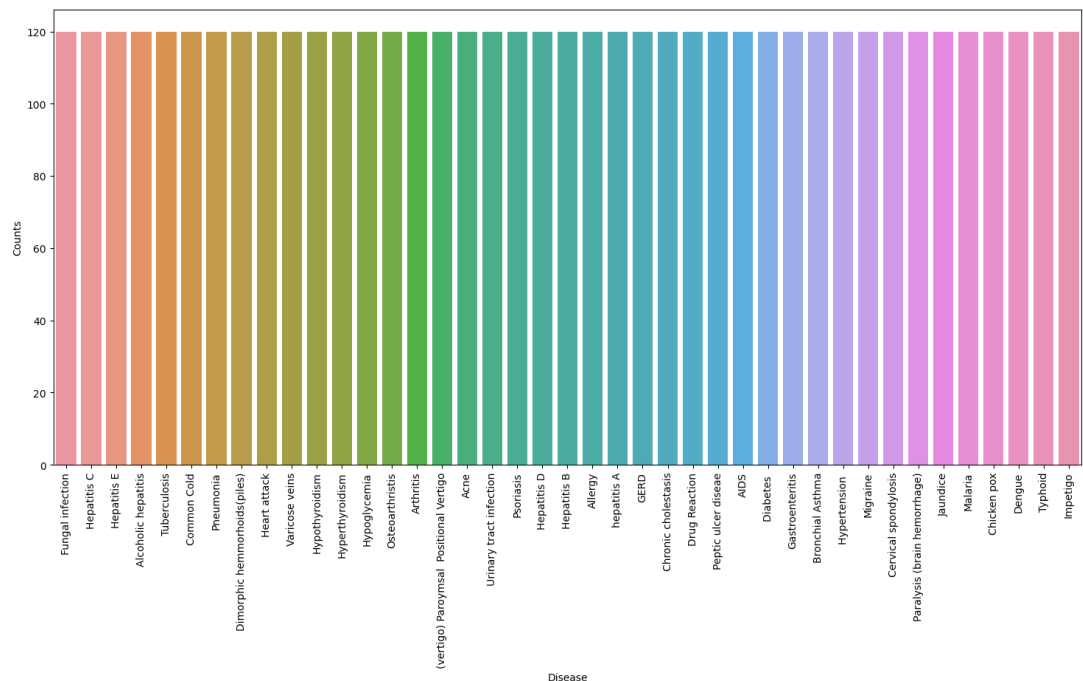
```
itching            0
skin_rash          0
nodal_skin_eruptions  0
continuous_sneezing  0
shivering          0
..
inflammatory_nails  0
blister            0
red_sore_around_nose  0
yellow_crust_ooze   0
prognosis          0
Length: 133, dtype: int64
```

In [36]: 1 `duplicate_rows_data = data[data.duplicated()]`  
2 `print("number of duplicate rows: ", duplicate_rows_data.shape)`

number of duplicate rows: (4616, 133)

```
In [37]: 1 # Checking whether the dataset is balanced or not
2 disease_counts = data["prognosis"].value_counts()
3 temp_df = pd.DataFrame({
4     "Disease": disease_counts.index,
5     "Counts": disease_counts.values
6 })
7
```

```
In [38]: 1 plt.figure(figsize = (18,8))
2 sns.barplot(x = "Disease", y = "Counts", data = temp_df)
3 plt.xticks(rotation=90)
4 plt.show()
```



## Encoding the target value into numerical

### value using LabelEncoder

```
In [39]: 1 encoder = LabelEncoder()
2 data["prognosis"] = encoder.fit_transform(data["prognosis"])
```

## Splitting the data for training and testing the model

```
In [40]: 1 X = data.iloc[:, :-1]
2 y = data.iloc[:, -1]
3 X_train, X_test, y_train, y_test = train_test_split(
4     X, y, test_size = 0.2, random_state = 24)
5
6 print(f"Train: {X_train.shape}, {y_train.shape}")
7 print(f"Test: {X_test.shape}, {y_test.shape}")
8
```

Train: (3936, 132), (3936,)

Test: (984, 132), (984,)

## Using K-Fold Cross-Validation for model selection

```
In [41]: 1 # Defining scoring metric for k-fold cross validation
2 def cv_scoring(estimator, X, y):
3     return accuracy_score(y, estimator.predict(X))
4
5 # Initializing Models
6 models = {
7     "SVC": SVC(),
8     "Gaussian NB": GaussianNB(),
9     "Random Forest": RandomForestClassifier(random_state=18)
10 }
11
12 # Producing cross validation score for the models
13 for model_name in models:
14     model = models[model_name]
15     scores = cross_val_score(model, X, y, cv = 10,
16                             n_jobs = -1,
17                             scoring = cv_scoring)
18     print("=="*30)
19     print(model_name)
20     print(f"Scores: {scores}")
21     print(f"Mean Score: {np.mean(scores)}")
```

```
=====
SVC
Scores: [1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
Mean Score: 1.0
=====
Gaussian NB
Scores: [1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
Mean Score: 1.0
=====
Random Forest
Scores: [1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
Mean Score: 1.0
```

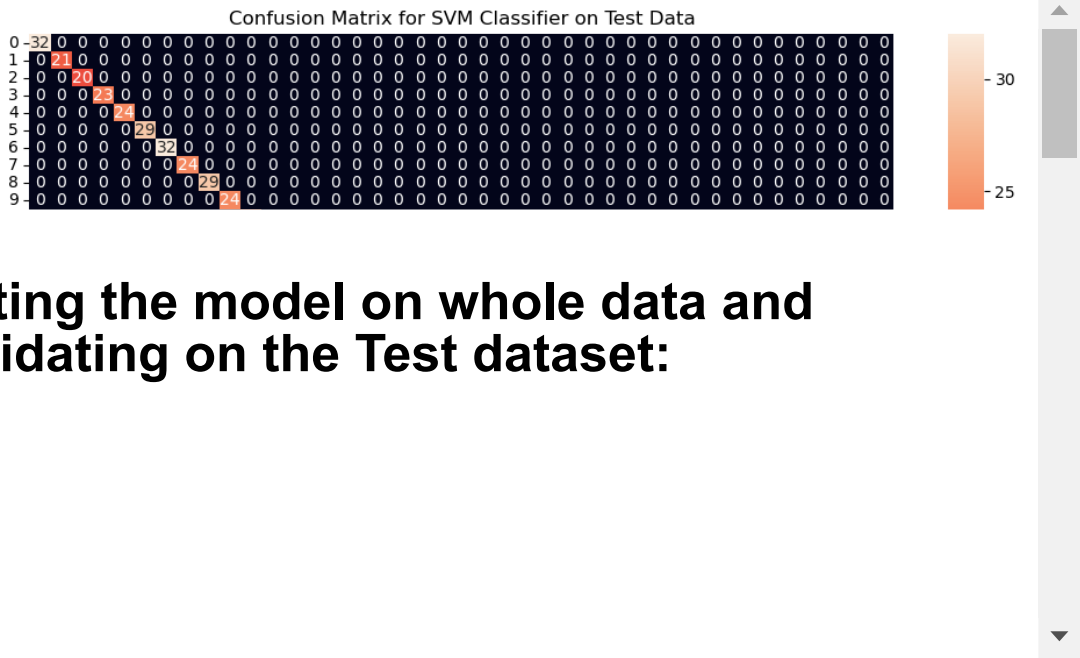
## Building robust classifier by combining all models:

```
In [42]: 1 # Training and testing SVM Classifier
2 svm_model = SVC()
3 svm_model.fit(X_train, y_train)
4 preds = svm_model.predict(X_test)
5
6 print(f"Accuracy on train data by SVM Classifier\
7 : {accuracy_score(y_train, svm_model.predict(X_train))*100}")
8
9 print(f"Accuracy on test data by SVM Classifier\
10 : {accuracy_score(y_test, preds)*100}")
11 cf_matrix = confusion_matrix(y_test, preds)
12 plt.figure(figsize=(12,8))
13 sns.heatmap(cf_matrix, annot=True)
14 plt.title("Confusion Matrix for SVM Classifier on Test Data")
15 plt.show()
16
17 # Training and testing Naive Bayes Classifier
18 nb_model = GaussianNB()
19 nb_model.fit(X_train, y_train)
20 preds = nb_model.predict(X_test)
21 print(f"Accuracy on train data by Naive Bayes Classifier\
22 : {accuracy_score(y_train, nb_model.predict(X_train))*100}")
23
24 print(f"Accuracy on test data by Naive Bayes Classifier\
25 : {accuracy_score(y_test, preds)*100}")
26 cf_matrix = confusion_matrix(y_test, preds)
27 plt.figure(figsize=(12,8))
28 sns.heatmap(cf_matrix, annot=True)
29 plt.title("Confusion Matrix for Naive Bayes Classifier on Test Data")
30 plt.show()
31
32 # Training and testing Random Forest Classifier
33 rf_model = RandomForestClassifier(random_state=18)
34 rf_model.fit(X_train, y_train)
35 preds = rf_model.predict(X_test)
36 print(f"Accuracy on train data by Random Forest Classifier\
37 : {accuracy_score(y_train, rf_model.predict(X_train))*100}")
38
39 print(f"Accuracy on test data by Random Forest Classifier\
40 : {accuracy_score(y_test, preds)*100}")
41
42 cf_matrix = confusion_matrix(y_test, preds)
43 plt.figure(figsize=(12,8))
44 sns.heatmap(cf_matrix, annot=True)
45 plt.title("Confusion Matrix for Random Forest Classifier on Test Data")
46 plt.show()
47
```

Accuracy on train data by SVM Classifier: 100.0

Accuracy on test data by SVM Classifier: 100.0





**Fitting the model on whole data and validating on the Test dataset:**

```

In [44]: ▶ 1 # Training the models on whole data
2 final_svm_model = SVC()
3 final_nb_model = GaussianNB()
4 final_rf_model = RandomForestClassifier(random_state=18)
5 final_svm_model.fit(X, y)
6 final_nb_model.fit(X, y)
7 final_rf_model.fit(X, y)
8
9 # Reading the test data
10 test_data = pd.read_csv("Testing_disease.csv").dropna(axis=1)
11
12 test_X = test_data.iloc[:, :-1]
13 test_Y = encoder.transform(test_data.iloc[:, -1])
14
15 # Making prediction by take mode of predictions
16 # made by all the classifiers
17 svm_preds = final_svm_model.predict(test_X)
18 nb_preds = final_nb_model.predict(test_X)
19 rf_preds = final_rf_model.predict(test_X)
20
21 final_preds = [mode([i,j,k])[0][0] for i,j,
22                 k in zip(svm_preds, nb_preds, rf_preds)]
23
24 print(f"Accuracy on Test dataset by the combined model\
25 : {accuracy_score(test_Y, final_preds)*100}")
26
27 cf_matrix = confusion_matrix(test_Y, final_preds)
28 plt.figure(figsize=(12,8))
29
30 sns.heatmap(cf_matrix, annot = True)
31 plt.title("Confusion Matrix for Combined Model on Test Dataset")
32 plt.show()
33

```

C:\Users\hp\AppData\Local\Temp\ipykernel\_19396\408761218.py:21: Future Warning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
final_preds = [mode([i,j,k])[0][0] for i,j,
```

Accuracy on Test dataset by the combined model: 100.0



## Creating a function that can take symptoms as input and generate predictions for disease

```
In [45]: 1 symptoms = X.columns.values
2
3 # Creating a symptom index dictionary to encode the
4 # input symptoms into numerical form
5 symptom_index = {}
6 for index, value in enumerate(symptoms):
7     symptom = " ".join([i.capitalize() for i in value.split("_")])
8     symptom_index[symptom] = index
9
10 data_dict = {
11     "symptom_index":symptom_index,
12     "predictions_classes":encoder.classes_
13 }
14
15 # Defining the Function
16 # Input: string containing symptoms separated by commas
17 # Output: Generated predictions by models
18 def predictDisease(symptoms):
19     symptoms = symptoms.split(",")
20
21     # creating input data for the models
22     input_data = [0] * len(data_dict["symptom_index"])
23     for symptom in symptoms:
24         index = data_dict["symptom_index"][symptom]
25         input_data[index] = 1
26
27     # reshaping the input data and converting it
28     # into suitable format for model predictions
29     input_data = np.array(input_data).reshape(1,-1)
30
31     # generating individual outputs
32     rf_prediction = data_dict["predictions_classes"][final_rf_model]
33     nb_prediction = data_dict["predictions_classes"][final_nb_model]
34     svm_prediction = data_dict["predictions_classes"][final_svm_model]
35
36     # making final prediction by taking mode of all predictions
37     final_prediction = mode([rf_prediction, nb_prediction, svm_prediction])
38     predictions = {
39         "rf_model_prediction": rf_prediction,
40         "naive_bayes_prediction": nb_prediction,
41         "svm_model_prediction": svm_prediction,
42         "final_prediction":final_prediction
43     }
44     return predictions
45
46 # Testing the function
47 print(predictDisease("Itching,Skin Rash,Nodal Skin Eruptions"))
```

```
{'rf_model_prediction': 'Fungal infection', 'naive_bayes_prediction':
'Fungal infection', 'svm_model_prediction': 'Fungal infection', 'final
_prediction': 'Fungal infection'}
```

```

C:\ProgramData\anaconda3\lib\site-packages\sklearn\base.py:420: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names
  warnings.warn(
C:\ProgramData\anaconda3\lib\site-packages\sklearn\base.py:420: UserWarning: X does not have valid feature names, but GaussianNB was fitted with feature names
  warnings.warn(
C:\ProgramData\anaconda3\lib\site-packages\sklearn\base.py:420: UserWarning: X does not have valid feature names, but SVC was fitted with feature names
  warnings.warn(
C:\Users\hp\AppData\Local\Temp\ipykernel_19396\2615708464.py:37: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.
  final_prediction = mode([rf_prediction, nb_prediction, svm_prediction])[0][0]
C:\Users\hp\AppData\Local\Temp\ipykernel_19396\2615708464.py:37: DeprecationWarning: Support for non-numeric arrays has been deprecated as of SciPy 1.9.0 and will be removed in 1.11.0. `pandas.DataFrame.mode` can be used instead, see https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.mode.html. (https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.mode.html.)
  final_prediction = mode([rf_prediction, nb_prediction, svm_prediction])[0][0]

```

**Observation :- The symptoms that are given as input to the function should be exactly the same among the 132 symptoms in the dataset.**

In [ ]: ▶

1