

```
import numpy as np #linear algebra library of Python
import pandas as pd # build on top of numpy for data analysis, data manipulation and d
import matplotlib.pyplot as plt #plotting library of Python
```

Now let's mount Google drive so that we can upload the diabetes.csv file. You can find the code in the 'Code s

```
from google.colab import drive
drive.mount('/content/gdrive')
```



First thing that we do is take a look at the shape of the dataframe (df.shape) and take a look at first 5 lines th

```
df=pd.read_csv('/content/gdrive/My Drive/Colab Notebooks/diabetes.csv') #import file f
df.head() #shows first 5 lines including column namesdf.shape # number of rows and col
```



```
df.shape # provides # rows and # columns of the dataframe df - 768 rows and 9 columns
```



Now we will assess if the dataset has the same proportion of diabetes vs. non-diabetes cases. At the same t
dataset we note that woman #2 has a skin thickness of zero and this is not realistic. It leads us to believe tha
was available. This does not apply to columns columns 1 and 9 for obvious reasons.

We use a trick to count the non-zero values of the columns. We convert the data type of the dataframe df to t
values to false=0 and all other entries to true=1 . We subsequently add up all True entries per column.

```
df.astype(bool).sum(axis=0) # counts the number of non-zeros for each column while act
```



The dataframe is unbalanced as we have 268 ones (diabetes) and thus 500 zeros (no diabetes).

The easiest option could be to eliminate all those patients with zero values, but in this way we would eliminat

Another option is to calculate the median value for a specific column and substitute the zero values for the co

```
median_BMI=df['BMI'].median()
df['BMI']=df['BMI'].replace(to_replace=0, value=median_BMI)

median_BloodPressure=df['BloodPressure'].median()
df['BloodPressure']=df['BloodPressure'].replace(to_replace=0, value=median_BloodPressu

median_Glucose=df['Glucose'].median()
df['Glucose']=df['Glucose'].replace(to_replace=0, value=median_Glucose)

median_SkinThickness=df['SkinThickness'].median()
df['SkinThickness']=df['SkinThickness'].replace(to_replace=0, value=median_SkinThickne
```

```
median_Insulin=df['Insulin'].median()  
df['Insulin']=df['Insulin'].replace(to_replace=0, value=median_Insulin)
```

```
df.head() #shows first 5 lines including column names
```



The skin thickness of woman #2 is now 23 (median of that column)

Let's create numpy arrays, one for the features (X) and one for the label (y)

```
X=df.drop('Outcome', 1).values #drop 'Outcome' column but you keep the index column  
y=df['Outcome'].values
```

We import the train_test_split function from sklearn to split the arrays or matrices into random train and test sets

Parameters:

test_size : in our case 20% (default=0.25)

random_state: is basically used for reproducing your problem the same every time it is run. If you do not use random_state, the split you might get a different set of train and test data points and will not help you in debugging in case your results number does not matter

stratify : array-like or None (default=None) If the number of values belonging to each class are unbalanced, use stratify, basically asking the model to take the training and test set such that the class proportion is same as of the whole dataset

```
from sklearn.model_selection import train_test_split #method to split training and test data  
X_train, X_test, y_train, y_test=train_test_split(X, y, test_size=0.2, random_state=42)
```

```
print(X_train)
```



The last preprocessing step is feature normalization transforming the data to have mean=0 and standard deviation=1. We will use the StandardScaler class from sklearn.preprocessing to execute the following code:

```
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
X_train=sc.fit_transform(X_train)
X_test=sc.transform(X_test)
```

```
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
```

```
svm_model = GridSearchCV(SVC(), params_grid, cv=5)
svm_model.fit(X_train, y_train)
```

```
params_grid = [{'kernel': ['rbf'], 'gamma': [1e-3, 1e-4], 'nu': [1, 10, 100, 1000]}], {'C': [1, 10, 100, 1000]}
```

```
print('Best score for training data:', svm_model.best_score_, "\n")

# View the best parameters for the model found using grid search
print('Best C:', svm_model.best_estimator_.C, "\n")
print('Best Kernel:', svm_model.best_estimator_.kernel, "\n")
print('Best Gamma:', svm_model.best_estimator_.gamma, "\n")
```



```
model = SVC(C=100, kernel='rbf', gamma=0.001)
model.fit(X_train, y_train)
```



```
y_pred = model.predict(X_test)
```

```
from sklearn import metrics
```

```
print("Accuracy:", metrics.accuracy_score(y_test, y_pred))
```



Now we will look at other classification KPIs that we discussed in our lessons: Confusion Matrix, ROC, AUC, F

```
from sklearn.metrics import confusion_matrix  
y_pred=model.predict(X_test)  
confusion_matrix(y_test,y_pred)
```



Classifier not so good: true positives=27, true negatives=82, false positives=18 and false negatives=27. We c
this case as we would tell a women that she is not diabetic whereas she actually is diabetic. One option to re
but this will increase the amount of False Positives (FPs) as we have seen in lesson 3. Recall in our case is TP
 $TP/(TP+FP)=60\%$