

We will repeat the Pima Indians exercise with PCA. As this is a low dimensional problem we think that applying PCA will not change a lot. It is mainly to get familiar with PCA and the associated code.

We will add sklearn's PCA library

```
import numpy as np #linear algebra library of Python
import pandas as pd # build on top of numpy for data analysis, data manipulation
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 snippets' tab of Colab

```
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 through df.head()

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



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



Now we will assess if the dataset has the same proportion of diabetes vs. non-diabetes cases. At the same time we will look if there are missing values. In our dataset we note that woman #2 has a skin thickness of zero and this is not realistic. It leads us to believe that there are a few zero entries that signal that no data 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 Boolean using df.astype(bool) converting all zero 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 whi
```



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 eliminate a lot of important data.

Another option is to calculate the median value for a specific column and substitute the zero values for the columns by that median value.

```
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_BloodP  
  
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_SkinT  
  
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 subsets>

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 a random_state in train_test_split, every time you make the split you might get a different set of train and test data points and will not help you in debugging in case you get an issue. We used random_state=42 but number does not matter

stratify : array-like or None (default=None) If the number of values belonging to each class are unbalanced, using stratified sampling is a good thing. You are basically asking the model to take the training and test set such that the class proportion is same as of the whole dataset, which is the right thing to do.

```
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)
```



```
print(X_test)
```



Now we are ready to use Logistic Regression

```
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
model.fit(X_train, y_train)
score = model.score(X_test, y_test)
print("Accuracy:", score)
```



Now we will look at other classification KPIs that we discussed in our lessons: Confusion Matrix, ROC, AUC, F1-Score

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



Classifier not so good: true positives=25, true negatives=86, false positives=14 and false negatives=29. We cannot accept the numerous False Negatives (FNs) in this case as we would tell a woman that she is not diabetic whereas she actually is diabetic. One option to reduce the FNs is to change the threshold of the classifier but this will increase the amount of False Positives (FPs) as we have seen in lesson 3. Recall in our case is $TP/(TP+FN) = 46\%$ and this is too high. Precision is $TP/(TP+FP)=64\%$

ROC (Receiver Operating Characteristic) curve

It is a plot of Recall vs. False Positive Rate (FPR) for the different possible thresholds of the classifier. It shows the tradeoff between Recall and Precision. The closer the curve follows the left-hand border and then the top border of the ROC space, the more accurate the test. The closer the curve comes to the 45-degree diagonal of the ROC space, the less accurate the test. The area under the curve is a measure of test accuracy.

```
from sklearn.metrics import roc_curve
y_pred_proba=model.predict_proba(X_test)[:,-1]
fpr, tpr, thresholds=roc_curve(y_test, y_pred_proba)
```

```
plt.plot([0,1], [0,1], '')
plt.plot(fpr, tpr, label='')
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title('Logistic Regression ROC curve')
plt.show()
```



```
from sklearn.metrics import roc_auc_score #area under the ROC curve
roc_auc_score(y_test, y_pred_proba)
```

