

The Pima are a group of Native Americans living in Arizona. A genetic predisposition allowed this group to survive in the desert. In the recent years, because of a sudden shift from traditional agricultural crops to processed foods, together with the highest prevalence of type 2 diabetes and for this reason they have been subject of many studies.

The dataset includes data from 768 women with 8 features:

Number of times pregnant

Plasma glucose concentration a 2 hours in an oral glucose tolerance test

Diastolic blood pressure (mm Hg)

Triceps skin fold thickness (mm)

2-Hour serum insulin (mu U/ml)

Body mass index (weight in kg/(height in m)^2)

Diabetes pedigree function

Age (years)

The last column of the dataset indicates if the person has been diagnosed with diabetes (1) or not (0)

We will start with importing a few key libraries

```
import numpy as np #linear algebra library of Python
import pandas as pd # build on top of numpy for data analysis, data manipulation and data I/O
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' section

```
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 of the dataframe (df.head())

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



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

```
↳ Drive already mounted at /content/gdrive; to attempt to forcibly remount, call
```

Now we will assess if the dataset has the same proportion of diabetes vs. non-diabetes cases. At the same time, in the dataset we note that woman #2 has a skin thickness of zero and this is not realistic. It leads us to believe that this 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, where 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 while acting as a boolean
```

```
↳
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabetes
0	6	148	72	35	0	33.6	
1	1	85	66	29	0	26.6	
2	8	183	64	0	0	23.3	
3	1	89	66	23	94	28.1	
4	0	137	40	35	168	43.1	

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 significant part of the dataset.
Another option is to calculate the median value for a specific column and substitute the zero values for the calculated median.

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

```
↳ (768, 9)
```

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', axis=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 are not as expected. random_state number does not matter

stratify : array-like or None (default=None) If the number of values belonging to each class are unbalanced, use stratify to make sure that the training and test sets have the same class proportion. 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)
```

The last preprocessing step is feature normalization transforming the data to have mean=0 and standard deviation=1. As the similarity measure in KNN we should not forget this step.

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

```
print(X_train)
```

➞ Pregnancies	657
Glucose	763
BloodPressure	733
SkinThickness	541
Insulin	394
BMI	757
DiabetesPedigreeFunction	768

Age	768
Outcome	268
dtype: int64	

Now we are ready to use the KNN algorithm

```
from sklearn.neighbors import KNeighborsClassifier # we import the K-Nearest Neighbor
neighbors=np.arange(1,20) #we will try different k - default step size is 1 - returns
train_accuracy=np.empty(len(neighbors)) # creates an array that will be used for storing
test_accuracy=np.empty(len(neighbors)) # creates an array that will be used for storing

print(neighbors)
```

➞ **Pregnancies Glucose BloodPressure SkinThickness Insulin BMI Diabetes**

A lot of times when dealing with iterators, we also get a need to keep a count of iterations. Python eases the pain by providing a built-in function called enumerate() for this task. Enumerate() method adds a counter to an iterable and returns it in a form of enumerate objects that can be iterated directly in 'for loops' or be converted into a list of tuples using list() method.

```
for i,k in enumerate(neighbors): #k goes from 1 to 19 en i is de counter
    knn=KNeighborsClassifier(n_neighbors=k)
    knn.fit(X_train, y_train)
    train_accuracy[i]=knn.score(X_train, y_train)
    test_accuracy[i]=knn.score(X_test, y_test)
```

```
plt.title('k-NN Varying number of neighbors')
plt.plot(neighbors, test_accuracy, label='Testing Accuracy')
plt.plot(neighbors, train_accuracy, label='Training Accuracy')
plt.legend()
plt.xlabel('Number of neighbors')
plt.ylabel('Accuracy')
plt.show()
```

➞ [[-0.85135507 -0.98013068 -0.40478372 ... -0.60767846 0.31079384
 -0.79216928]

```
[ 0.35657564  0.16144422  0.46536842 ... -0.30213902 -0.11643851
 0.56103382]
[-0.5493724  -0.50447447 -0.62232176 ...  0.3725939  -0.76486207
-0.70759409]
...
[-0.85135507 -0.75815778  0.03029235 ...  0.77997981 -0.78607218
-0.28471812]
[ 1.86648903 -0.31421198  0.03029235 ... -0.56948603 -1.01938346
 0.56103382]
[ 0.05459296  0.73223168 -0.62232176 ... -0.31486983 -0.57700104
 0.30730824]]
```

We get maximum testing accuracy for k=8, so we will setup a KNN classifier with hyperparameter k=8 (we will use this classifier to predict a new test point)

```
knn=KNeighborsClassifier(n_neighbors=8)
knn.fit(X_train, y_train)
knn.score(X_test, y_test) #score method represents accuracy
```



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=knn.predict(X_test)
confusion_matrix(y_test,y_pred)
```



Classifier not so good: true positives=26, true negatives=92, false positives=8 and false negatives=28. We can see that we have a high number of false positives, which means we are predicting a woman is diabetic when she is not. This is a bad case as we would tell a woman that she is not diabetic whereas she actually is diabetic. One option to reduce false positives is to increase the threshold for predicting a woman is diabetic. This will increase the amount of False Positives (FPs) as we have seen in lesson 3. Recall in our case is $TP/(TP+FN)=77\%$

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 how well the classifier performs. The closer the curve follows the left-hand border and then the top border of the ROC space, the more accurate the classifier. The closer the curve follows the 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=knn.predict_proba(X_test)[:,-1]
fpr, tpr, thresholds=roc_curve(y_test, y_pred_proba)
```

```
plt.plot([0,1], [0,1], 'k--')
plt.plot(fpr, tpr, label='knn')
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title('Knn(neighbors=8) ROC curve')
plt.show()
```



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



Cross Validation

Now before getting into the details of Hyperparameter tuning, let us understand the concept of Cross validation

The trained model's performance is dependent on way the data is split. It might not be representative of the model's performance on new data.

The solution is cross validation.

Cross-validation is a technique to evaluate predictive models by partitioning the original sample into a training set and a validation set.

In k-fold cross-validation, the original sample is randomly partitioned into k equal size subsamples. Of the k subsamples, one subsample is used as the validation data for testing the model, and the remaining k-1 subsamples are used as training data. The process is repeated k times, with each of the k subsamples used exactly once as the validation data. The cross-validated performance is the average of the k performance scores.

validation data for testing the model, and the remaining k-1 subsamples are used as training data. The cross-validation process is repeated k times, with each of the k subsamples used exactly once as the validation data. The k results from the folds can then be averaged to produce a single estimation. The advantage of this method is that all observations are used for both training and validation exactly once

(Kilian Weinberger explains this)

Hyperparameter tuning

The value of k (i.e 7) we selected above was selected by observing the curve of accuracy vs number of neighbors

There is a better way of doing it which involves:

- 1) Trying a bunch of different hyperparameter values
- 2) Fitting all of them separately
- 3) Checking how well each performs
- 4) Choosing the best performing one
- 5) Using cross-validation every time

Scikit-learn provides a simple way of achieving this using GridSearchCV i.e Grid Search cross-validation.

Let's see what accuracy we get using logistics regression and naive Bayes

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



Accuracy is only 74%, same as KNN

Now SVM

```
from sklearn import svm
```

```
clf = svm.SVC(kernel='linear') # Linear Kernel
```

```
clf.fit(X_train, y_train)
score=clf.score(X_test, y_test)
print("Accuracy:", score)
```



Also 74% accuracy, which is not that good. If we use a different kernel, we can get a better accuracy.

Also 74%...a bit strange that all classifiers have the same accuracy