# Pulsars prediction

•••
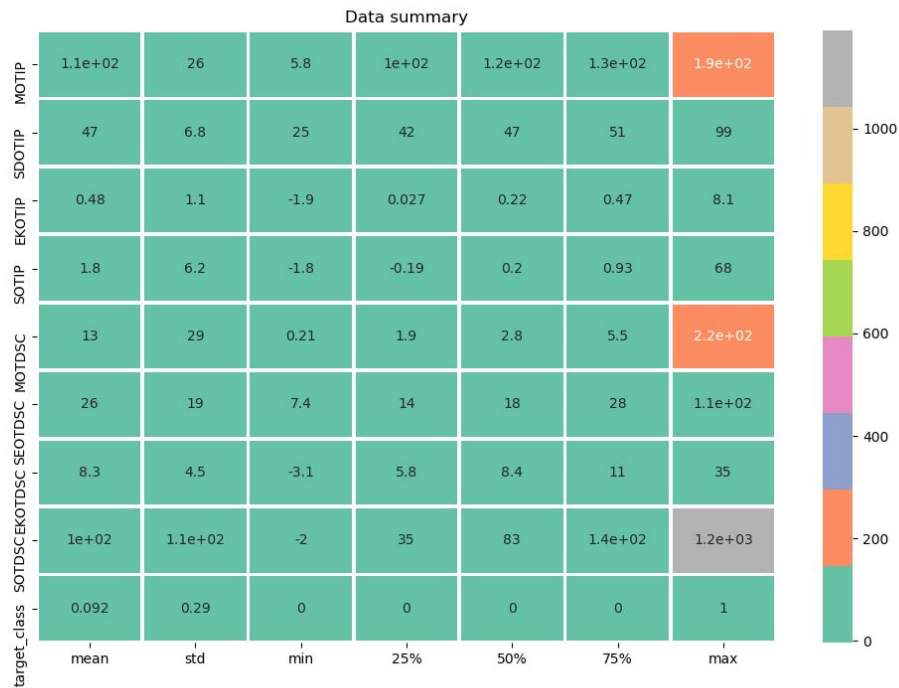
# 1. Checking the dataset

We will need to scale it.

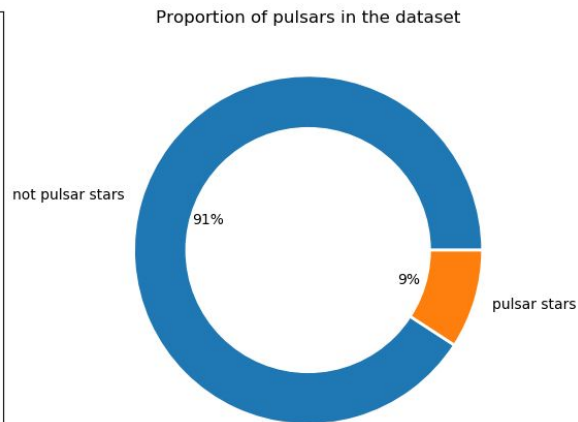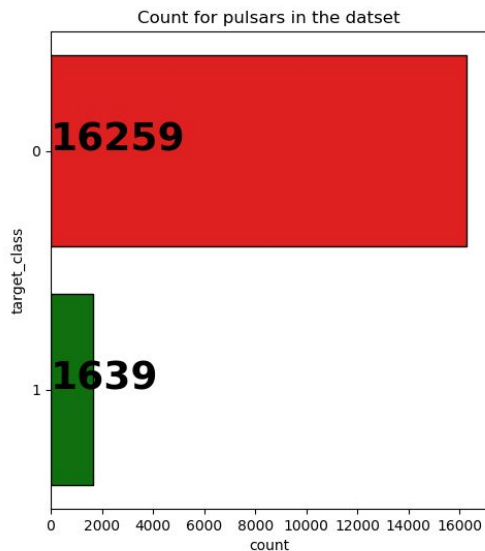Using **StandardScaler**



Data summary

# 1. Checking the dataset

There is a strong correlation between **SOTIP** and **EKOTIP** and also between **SOTDSC** and **EKOTDSC**, so we will drop **SOTIP** and **SOTDSC**
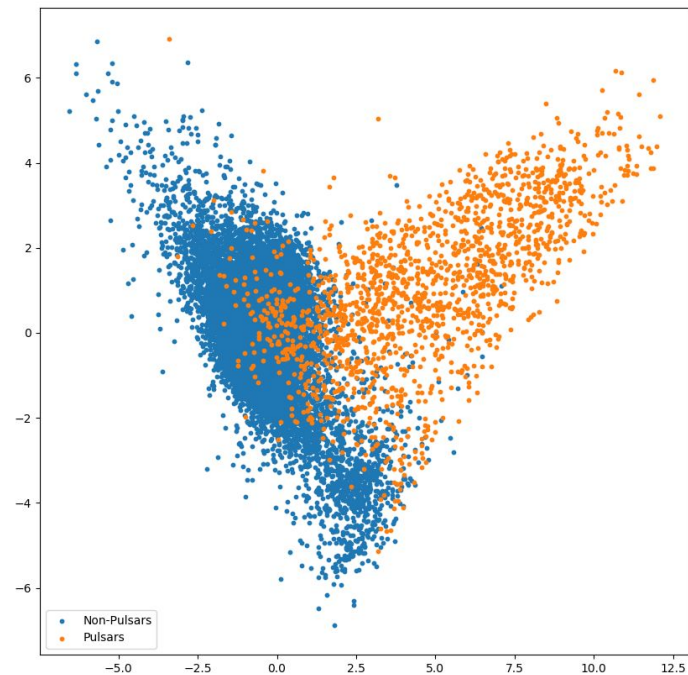


Variables correlation

# 1. Checking the dataset

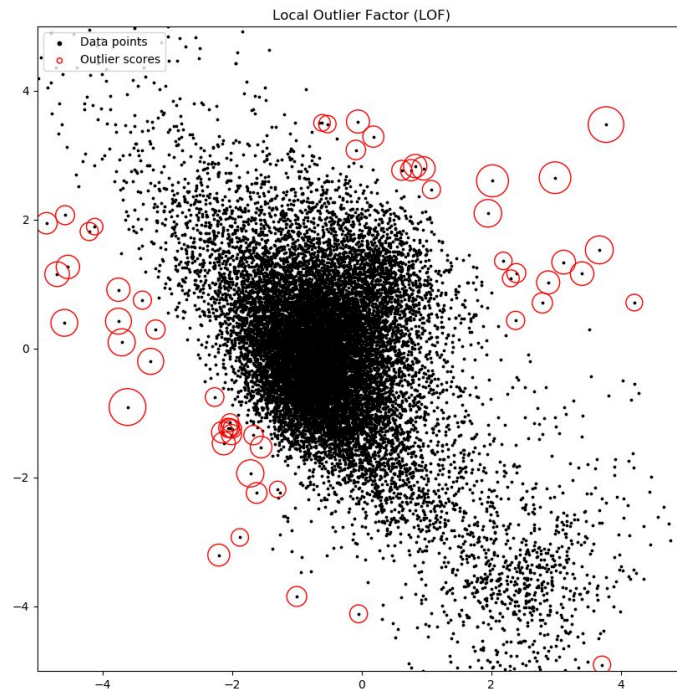Dataset is imbalanced - we will need to use stratification

# 1. Checking the dataset

Let's search for anomalies using Principal Component Analysis and Local Outlier Factor

# 1. Checking the dataset

We can see some outliers for non-pulsars...



Local Outlier Factor (LOF)

# 1. Checking the dataset

...and for pulsars too.
Let's remove them.

# 1. Checking the dataset

Splitting the dataset

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, stratify=y,
random_state=0)
```

And then go on testing different models

# 2. Testing models

```
RandomForestClassifier:

Classification Report:

              precision    recall  f1-score   support

           0       0.99      0.99      0.99      4855
           1       0.94      0.86      0.90       487

    accuracy                           0.98      5342
   macro avg       0.97      0.93      0.95      5342
weighted avg       0.98      0.98      0.98      5342


Confusion Matrix:

[[4830   25]
 [  66  421]]

Cross validation:

Recall: 98.09%
[[4825   30]
 [  72  415]]
```

```
LinearSVC:

//anaconda3/lib/python3.7/site-packages/sklearn/svm/base
  "the number of iterations.", ConvergenceWarning)
Classification Report:

              precision    recall  f1-score   support

           0       0.98      1.00      0.99      4855
           1       0.95      0.84      0.89       487

    accuracy                           0.98      5342
   macro avg       0.97      0.92      0.94      5342
weighted avg       0.98      0.98      0.98      5342


Confusion Matrix:

[[4832   23]
 [  77  410]]

[[4835    20]
 [  80   407]]
```

# 2. Testing models

```
GradientBoostingClassifier:

Classification Report:

              precision    recall  f1-score   support

           0       0.99      0.99      0.99      4855
           1       0.93      0.87      0.90       487

    accuracy                           0.98      5342
   macro avg       0.96      0.93      0.94      5342
weighted avg       0.98      0.98      0.98      5342


Confusion Matrix:

[[4823   32]
 [  65  422]]

Cross validation:

Recall: 98.0%
[[4816   39]
 [  68  419]]
```

# 3. Tuning RF using **GridSearchCV**

```python
forest = RandomForestClassifier(bootstrap=True, class_weight='balanced_subsample',
                        criterion='gini', max_depth=15, max_features=4,
                        max_leaf_nodes=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=1,
                        min_samples_split=40, min_weight_fraction_leaf=0.0,
                        n_estimators=120, n_jobs=None, oob_score=False,
                        random_state=0, verbose=0, warm_start=False)
```

# 3. Tuning RF using **GridSearchCV**

## Tuning to make FN minimal

```
Classification Report:

              precision    recall  f1-score   support

           0       0.99      0.98      0.99      4855
           1       0.86      0.92      0.89       487

    accuracy                           0.98      5342
   macro avg       0.93      0.95      0.94      5342
weighted avg       0.98      0.98      0.98      5342


Confusion Matrix:

[[4781    74]
 [  39   448]]
Recall: 97.58%
```
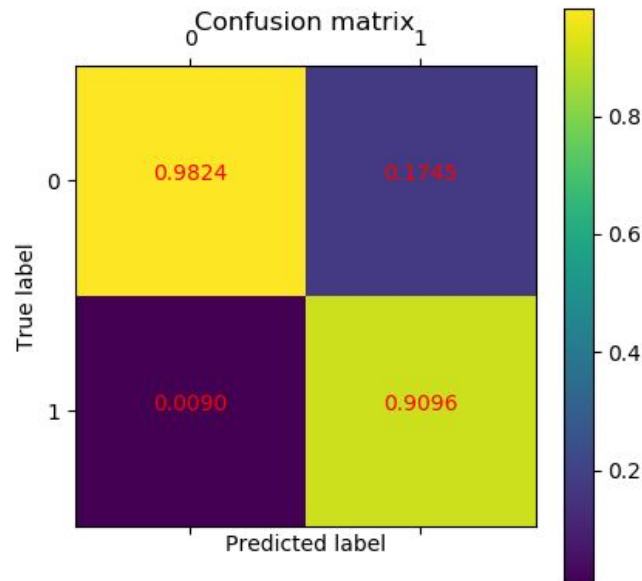


Confusion matrix

# Pneumonia prediction
...

# Using CNN

```python
# Build the CNN
classifier = Sequential()
# Convolution
classifier.add(Conv2D(32, (3, 3), activation="relu", input_shape=(64, 64, 3)))
# Pooling
classifier.add(MaxPooling2D(pool_size = (2, 2)))
# Pooling is made with a 2x2 array
# Add 2nd convolutional layer with the same structure as the 1st to improve predictions
classifier.add(Conv2D(32, (3, 3), activation="relu"))
classifier.add(MaxPooling2D(pool_size = (2, 2)))
# Flattening
classifier.add(Flatten())
# Full Connection
classifier.add(Dense(activation = 'relu', units = 128))
classifier.add(Dense(activation = 'sigmoid', units = 1))
# Compile the CNN
classifier.compile(optimizer = 'adam', loss = 'binary_crossentropy', metrics = ['accuracy'])
```

# Preparing Image Generator

```python
train_datagen = ImageDataGenerator(rescale = 1./255,
                                   shear_range = 0.2,
                                   zoom_range = 0.2,
                                   horizontal_flip = True)

test_datagen = ImageDataGenerator(rescale = 1./255)
training_set = train_datagen.flow_from_directory('./chest_xray/train',
                                                 target_size = (64, 64),
                                                 batch_size = 32,
                                                 class_mode = 'binary')
test_set = test_datagen.flow_from_directory('./chest_xray/test',
                                            target_size = (64, 64),
                                            batch_size = 32,
                                            class_mode = 'binary')
```

# Running on 20 epoch

# We can even continue training, accuracy growing