

## ▼ Configuration

```
#Parameters
PROJECT_NAME = 'ML1010_Weekly'
ENABLE_COLAB = True

#Root Machine Learning Directory. Projects appear underneath
GOOGLE_DRIVE_MOUNT = '/content/gdrive'
COLAB_ROOT_DIR = GOOGLE_DRIVE_MOUNT + '/MyDrive/Colab Notebooks'
COLAB_INIT_DIR = COLAB_ROOT_DIR + '/utility_files'

LOCAL_ROOT_DIR = '/home/magni/Documents/ML_Projects'
LOCAL_INIT_DIR = LOCAL_ROOT_DIR + '/utility_files'
```

## ▼ Bootstrap Environment

```
#add in support for utility file directory and importing
import sys
import os

if ENABLE_COLAB:
    #Need access to drive
    from google.colab import drive
    drive.mount(GOOGLE_DRIVE_MOUNT, force_remount=True)

    #add in utility directory to syspath to import
    INIT_DIR = COLAB_INIT_DIR
    sys.path.append(os.path.abspath(INIT_DIR))

    #Config environment variables
    ROOT_DIR = COLAB_ROOT_DIR

else:
    #add in utility directory to syspath to import
    INIT_DIR = LOCAL_INIT_DIR
    sys.path.append(os.path.abspath(INIT_DIR))

    #Config environment variables
    ROOT_DIR = LOCAL_ROOT_DIR

#Import Utility Support
from jarvis import Jarvis
jarvis = Jarvis(ROOT_DIR, PROJECT_NAME)

import mv python utils as mvutils
```

```
Mounted at /content/gdrive
Wha...where am I?
I am awake now.
```

```
I have set your current working directory to /content/gdrive/MyDrive/Colab Notebooks/ML1
The current time is 08:16
Hello sir. An early morning I see.
```

## ▼ Setup Runtime Environment

```
if ENABLE_COLAB:
    #!pip install scipy -q
    #!pip install scikit-learn -q
    #!pip install pycaret -q
    #!pip install matplotlib -q
    #!pip install joblib -q
    #!pip install pandasql -q
    display('Google Colab has been enabled')

else:
    display('Google Colab not enabled')

#Common imports
#import json
#import gzip
import pandas as pd
import numpy as np
import matplotlib
#import re
import nltk
import matplotlib.pyplot as plt

pd.set_option('mode.chained_assignment', None)
nltk.download('stopwords')
%matplotlib inline

'Google Colab has been enabled'
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
```

## ▼ Load Data

```
jarvis.showProjectDataFiles()
```

Here are all your project data files

```
[D] /content/gdrive/MyDrive/Colab Notebooks/data/ML1010_Weekly
```

```
---[ gz][ csv]--> complaints.csv.gz (370.67 MB)
```

```
[*][ csv]-----> movie_reviews_cleaned.csv (38.37 MB)
```

```
[*][ csv]-----> pima-indians-diabetes.csv (22.73 KB)
```

```
---[ gz][ tsv]--> rspct.tsv.gz (347.13 MB)
```

```
---[ gz][ csv]--> subreddit_info.csv.gz (37.80 KB)
```

```
[D] /content/gdrive/MyDrive/Colab Notebooks/data/ML1010_Weekly/01_original
```

```
----->** No files **
```

```
[D] /content/gdrive/MyDrive/Colab Notebooks/data/ML1010_Weekly/02_working
```

```
----->** No files **
```

```
[D] /content/gdrive/MyDrive/Colab Notebooks/data/ML1010_Weekly/03_train
```

```
----->** No files **
```

```
[D] /content/gdrive/MyDrive/Colab Notebooks/data/ML1010_Weekly/04_test
```

```
----->** No files **
```

# Feature Selection with Univariate Statistical Tests

```
from numpy import set_printoptions
```

```
from sklearn.feature_selection import SelectKBest
```

```
from sklearn.feature_selection import f_classif
```

# load data

```
filename = 'pima-indians-diabetes.csv'
```

```
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
```

```
dataframe = pd.read_csv(jarvis.DATA_DIR + "/" + filename, names=names)
```

```
array = dataframe.values
```

```
X = array[:,0:8]
```

```
Y = array[:,8]
```

# feature extraction

```
test = SelectKBest(score_func=f_classif, k=4)
```

```
fit = test.fit(X, Y)
```

# summarize scores

```
set_printoptions(precision=3)
```

```
print(fit.scores_)
```

```
features = fit.transform(X)
```

# summarize selected features

```
print(features[0:5,:])
```

```
[ 39.67  213.162   3.257   4.304  13.281  71.772  23.871  46.141]
```

```
[[ 6.  148.   33.6  50. ]
```

```
 [ 1.   85.   26.6  31. ]
```

```
 [ 8.  183.   23.3  32. ]
```

```
 [ 1.   89.   28.1  21. ]
```

```
 [ 0.  137.   43.1  33. ]]
```

```
# Feature Extraction with RFE
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression

# load data
filename = 'pima-indians-diabetes.csv'
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
dataframe = pd.read_csv(jarvis.DATA_DIR + "/" + filename, names=names)

array = dataframe.values
X = array[:,0:8]
Y = array[:,8]
# feature extraction
model = LogisticRegression(solver='lbfgs')
rfe = RFE(model, n_features_to_select=3)
fit = rfe.fit(X, Y)
print("Num Features: %d" % fit.n_features_)
print("Selected Features: %s" % fit.support_)
print("Feature Ranking: %s" % fit.ranking_)

Num Features: 3
Selected Features: [ True False False False False  True  True False]
Feature Ranking: [1 2 4 5 6 1 1 3]
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: ConvergenceWarning:
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max\_iter) or scale the data as shown in:  
<https://scikit-learn.org/stable/modules/preprocessing.html>  
Please also refer to the documentation for alternative solver options:  
[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)  
extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG,

```
# Feature Extraction with PCA
import numpy
from sklearn.decomposition import PCA

# load data
filename = 'pima-indians-diabetes.csv'
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
dataframe = pd.read_csv(jarvis.DATA_DIR + "/" + filename, names=names)

array = dataframe.values
X = array[:,0:8]
Y = array[:,8]
# feature extraction
pca = PCA(n_components=3)
fit = pca.fit(X)
# summarize components
print("Explained Variance: %s" % fit.explained_variance_ratio_)
print(fit.components_)
```

```

Explained Variance: [0.889 0.062 0.026]
[[-2.022e-03  9.781e-02  1.609e-02  6.076e-02  9.931e-01  1.401e-02
   5.372e-04 -3.565e-03]
 [-2.265e-02 -9.722e-01 -1.419e-01  5.786e-02  9.463e-02 -4.697e-02
  -8.168e-04 -1.402e-01]
 [-2.246e-02  1.434e-01 -9.225e-01 -3.070e-01  2.098e-02 -1.324e-01
  -6.400e-04 -1.255e-01]]

```

```

# Feature Importance with Extra Trees Classifier
from sklearn.ensemble import ExtraTreesClassifier
# load data
filename = 'pima-indians-diabetes.csv'
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
dataframe = pd.read_csv(jarvis.DATA_DIR + "/" + filename, names=names)

array = dataframe.values
X = array[:,0:8]
Y = array[:,8]
# feature extraction
model = ExtraTreesClassifier(n_estimators=10)
model.fit(X, Y)
print(model.feature_importances_)

[0.104 0.22  0.108 0.089 0.078 0.142 0.115 0.144]

```