# Supervised Learning with scikit-learn

## 4). Preprocessing & Pipelining:

**a). Exploring Categorical Features**

**# Import pandas**

**import pandas as pd**

**# Read 'gapminder.csv' into a DataFrame: df**

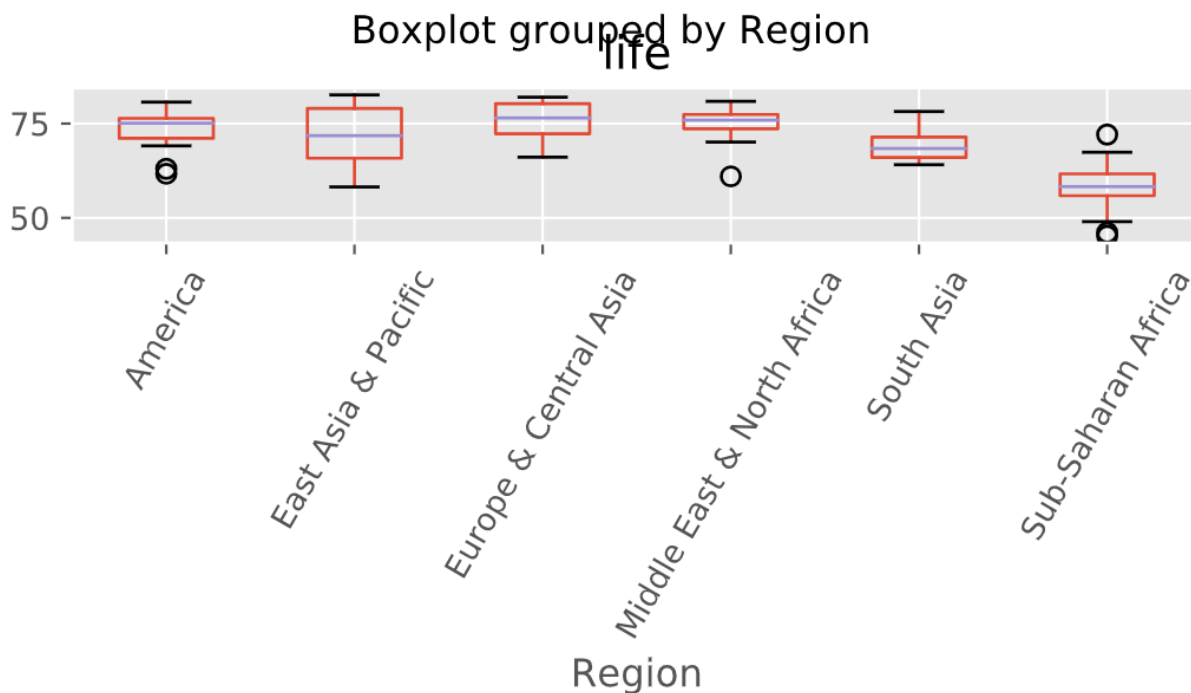**df = pd.read_csv('gapminder.csv')**

**# Create a boxplot of life expectancy per region**

**df.boxplot('life', 'Region', rot=60)**

**# Show the plot**

**plt.show()**



Boxplot grouped by Region

**b). <u>Creating Dummy Variables</u>**

**# Create dummy variables: df_region**

**df_region = pd.get_dummies(df)**


**# Print the columns of df_region**

**print(df_region.columns)**


**# Drop 'Region_America' from df_region**

**df_region = pd.get_dummies(df, drop_first=True)**


**# Print the new columns of df_region**

**print(df_region.columns)**


```
<script.py> output:
    Index(['population', 'fertility', 'HIV', 'CO2', 'BMI_male', 'GDP',
        'BMI_female', 'life', 'child_mortality', 'Region_America',
        'Region_East Asia & Pacific', 'Region_Europe & Central Asia',
        'Region_Middle East & North Africa', 'Region_South Asia',
        'Region_Sub-Saharan Africa'],
        dtype='object')
    Index(['population', 'fertility', 'HIV', 'CO2', 'BMI_male', 'GDP',
        'BMI_female', 'life', 'child_mortality', 'Region_East Asia & Pacific',
        'Region_Europe & Central Asia', 'Region_Middle East & North Africa',
        'Region_South Asia', 'Region_Sub-Saharan Africa'],
        dtype='object')
```

**c). Regression with categorical features**

**# Import necessary modules**

**from sklearn.linear_model import Ridge**

**from sklearn.model_selection import cross_val_score**


**# Instantiate a ridge regressor: ridge**

**ridge = Ridge(alpha=0.5, normalize=True)**


**# Perform 5-fold cross-validation: ridge_cv**

**ridge_cv = cross_val_score(ridge,X,y,cv=5)**


**# Print the cross-validated scores**

**print(ridge_cv)**


&lt;script.py&gt; output:

   [ 0.86808336  0.80623545  0.84004203  0.7754344   0.87503712]

**d). <u>Dropping Missing Data:</u>**

**# Convert '?' to NaN**

**df[df == '?'] = np.nan**

**# Print the number of NaNs**

**print(df.isnull().sum())**

**# Print shape of original DataFrame**

**print("Shape of Original DataFrame: {}".format(df.shape))**

**# Drop missing values and print shape of new DataFrame**

**df = df.dropna()**

**# Print shape of new DataFrame**

**print("Shape of DataFrame After Dropping All Rows with Missing Values: {}".format(df.shape))**

&lt;script.py&gt; output:

```
party               0
infants            12
water              48
budget             11
physician          11
salvador           15
religious          11
satellite          14
aid                15
missile            22
immigration         7
synfuels           21
education          31
superfund          25
crime              17
duty_free_exports  28
eaa_rsa           104
dtype: int64
```

Shape of Original DataFrame: (435, 17)

Shape of DataFrame After Dropping All Rows with Missing Values: (232, 17)

**e). <u>Imputing missing data in a ML Pipeline</u>**

**# Import necessary modules**

**from sklearn.preprocessing import Imputer**

**from sklearn.pipeline import Pipeline**

**from sklearn.svm import SVC**

**# Setup the pipeline steps: steps**

**steps = [('imputation', Imputer(missing_values='NaN', strategy='most_frequent', axis=0)),**

**('SVM', SVC())]**

**# Create the pipeline: pipeline**

**pipeline = Pipeline(steps)**

**# Create training and test sets**

**X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.3,random_state=42)**

**# Fit the pipeline to the train set**

**pipeline.fit(X_train,y_train)**

**# Predict the labels of the test set**

**y_pred = pipeline.predict(X_test)**

**# Compute metrics**

**print(classification_report(y_test, y_pred))**

<script.py> output:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| democrat | 0.99 | 0.96 | 0.98 | 85 |
| republican | 0.94 | 0.98 | 0.96 | 46 |
| avg / total | 0.97 | 0.97 | 0.97 | 131 |

**f). <u>Centring and Scaling your Data</u>**

**# Import scale**

**from sklearn.preprocessing import scale**

**# Scale the features: X_scaled**

**X_scaled = scale(X)**

**# Print the mean and standard deviation of the unscaled features**

**print("Mean of Unscaled Features: {}".format(np.mean(X)))**

**print("Standard Deviation of Unscaled Features: {}".format(np.std(X)))**

**# Print the mean and standard deviation of the scaled features**

**print("Mean of Scaled Features: {}".format(np.mean(X_scaled)))**

**print("Standard Deviation of Scaled Features: {}".format(np.std(X_scaled)))**

&lt;script.py&gt; output:

    Mean of Unscaled Features: 18.432687072460002

    Standard Deviation of Unscaled Features: 41.54494764094571

    Mean of Scaled Features: 2.7314972981668206e-15

    Standard Deviation of Scaled Features: 0.9999999999999999

**g). <u>Centering and Scaling in a Pipeline</u>**

**# Import the necessary modules**

**from sklearn.preprocessing import StandardScaler**

**from sklearn.pipeline import Pipeline**

**# Setup the pipeline steps: steps**

**steps = [('scaler', StandardScaler()),**

**        ('knn', KNeighborsClassifier())]**

**# Create the pipeline: pipeline**

**pipeline = Pipeline(steps)**

**# Create train and test sets**

**X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.3, random_state=42)**

**# Fit the pipeline to the training set: knn_scaled**

**knn_scaled = pipeline.fit(X_train,y_train)**

**# Instantiate and fit a k-NN classifier to the unscaled data**

**knn_unscaled = KNeighborsClassifier().fit(X_train, y_train)**

**# Compute and print metrics**

**print('Accuracy with Scaling: {}'.format(knn_scaled.score(X_test, y_test)))**

**print('Accuracy without Scaling: {}'.format(knn_unscaled.score(X_test, y_test)))**

<script.py> output:

   Accuracy with Scaling: 0.7700680272108843

   Accuracy without Scaling: 0.6979591836734694

**h). <u>Pipeline for Classification</u>**

**# Setup the pipeline**

**steps = [('scaler', StandardScaler()),**

**    ('SVM', SVC())]**

**pipeline = Pipeline(steps)**

**# Specify the hyperparameter space**

**parameters = {'SVM__C':[1, 10, 100],**

**        'SVM__gamma':[0.1, 0.01]}**

**# Create train and test sets**

**X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state=21)**

**# Instantiate the GridSearchCV object: cv**

**cv = GridSearchCV(pipeline,param_grid=parameters)**

**# Fit to the training set**

**cv.fit(X_train,y_train)**

**# Predict the labels of the test set: y_pred**

**y_pred = cv.predict(X_test)**

**# Compute and print metrics**

**print("Accuracy: {}".format(cv.score(X_test, y_test)))**

**print(classification_report(y_test, y_pred))**

**print("Tuned Model Parameters: {}".format(cv.best_params_))**

&lt;script.py&gt; output:

  Accuracy: 0.7795918367346939

          precision   recall  f1-score   support

    False      0.83     0.85     0.84      662

    True      0.67     0.63     0.65      318

  avg / total    0.78     0.78     0.78      980

  Tuned Model Parameters: {'SVM__gamma': 0.1, 'SVM__C': 10}

**i). <u>Pipelining for Regression</u>**

**# Setup the pipeline steps: steps**

**steps = [('imputation', Imputer(missing_values='NaN', strategy='mean', axis=0)),**

**('scaler', StandardScaler()),**

**('elasticnet', ElasticNet())]**


**# Create the pipeline: pipeline**

**pipeline = Pipeline(steps)**


**# Specify the hyperparameter space**

**parameters = {'elasticnet__l1_ratio':np.linspace(0,1,30)}**


**# Create train and test sets**

**X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.4, random_state=42)**


**# Create the GridSearchCV object: gm_cv**

**gm_cv = GridSearchCV(pipeline, param_grid=parameters)**


**# Fit to the training set**

**gm_cv.fit(X_train,y_train)**


**#Compute and print the metrics**

**r2 = gm_cv.score(X_test, y_test)**

**print("Tuned ElasticNet Alpha: {}".format(gm_cv.best_params_))**

**print("Tuned ElasticNet R squared: {}".format(r2))**


<script.py> output:

   Tuned ElasticNet Alpha: {'elasticnet__l1_ratio': 1.0}

   Tuned ElasticNet R squared: 0.8862016570888217