# Chapter 2 - Ex2: Iris - Multiple Linear Regression

## Cho dữ liệu Iris.xls

## Yêu cầu: Thực hiện linenear regression để từ sepallength, sepalwidth, petallength => dự đoán petalwidth

1. Đọc dữ liệu, trực quan hóa dữ liệu.
2. Tạo X_train, X_test, y_train, y_test từ dữ liệu đọc được là sepallength, sepalwidth, petallength (inputs) và petalwidth (outputs) với tỷ lệ dữ liệu test là 0.2
3. Áp dụng linrear regression
4. Vẽ hình. Nhận xét kết quả
5. Nếu sepallength, sepalwidth, petallength là 4.5, 3.1, 1.6 => petalwidth là bao nhiêu?
6. Áp dụng lựa chọn thuộc tính quan trọng cho model. Xây dựng lại model sau khi lựa chọn các thuộc tính quan trọng.

In [1]:
```python
# from google.colab import drive
# drive.mount("/content/gdrive", force_remount=True)

# %cd '/content/gdrive/My Drive/LDS6_MachineLearning/practice/Chapter2_Linear_Reg
```

In [3]:
```python
import pandas as pd
iris = pd.read_excel("Iris.xls")
iris.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
sepallength    150 non-null float64
sepalwidth     150 non-null float64
petallength    150 non-null float64
petalwidth     150 non-null float64
iris           150 non-null object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```
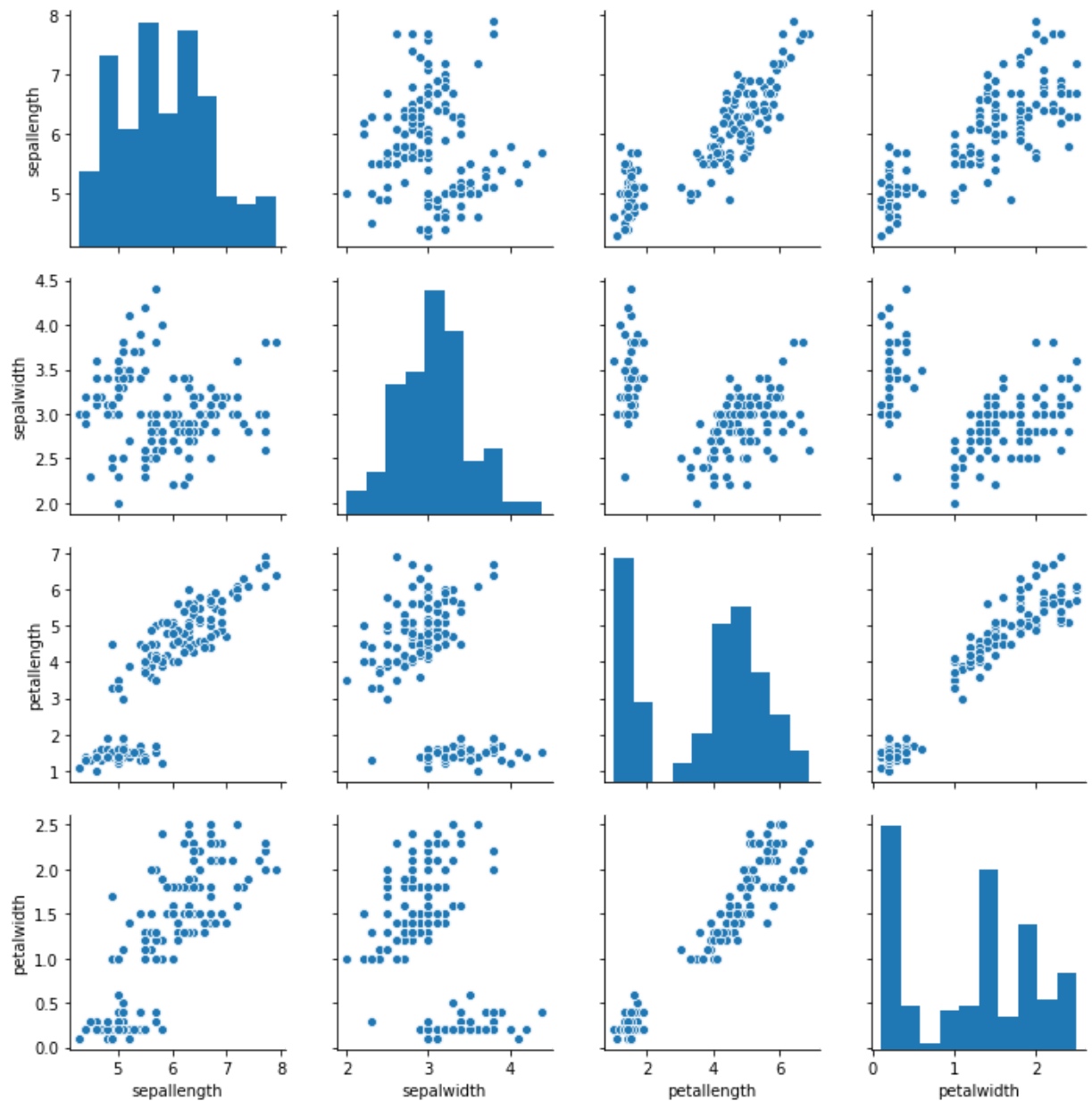
In [4]:
```python
iris.head()
```

Out[4]:

|   | sepallength | sepalwidth | petallength | petalwidth | iris |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

In [5]:
```python
import matplotlib.pyplot as plt
import seaborn as sns
```

In [6]:
```python
sns.pairplot(iris)
plt.show()
```

In [7]:
```python
inputs = iris[['sepallength','sepalwidth', 'petallength']]
inputs.head()
```

Out[7]:

|   | sepallength | sepalwidth | petallength |
|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 |
| 1 | 4.9 | 3.0 | 1.4 |
| 2 | 4.7 | 3.2 | 1.3 |
| 3 | 4.6 | 3.1 | 1.5 |
| 4 | 5.0 | 3.6 | 1.4 |

In [8]:
```python
outputs = iris[['petalwidth']]
outputs.head()
```

Out[8]:

|   | petalwidth |
|---|---|
| 0 | 0.2 |
| 1 | 0.2 |
| 2 | 0.2 |
| 3 | 0.2 |
| 4 | 0.2 |

In [9]:
```python
import numpy as np
from sklearn import datasets, linear_model
from sklearn.metrics import mean_squared_error, r2_score
```

In [10]:
```python
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(inputs, outputs,
                                                    test_size=0.20)

regr1 = linear_model.LinearRegression()
regr1 = regr1.fit(X_train, y_train)
```

In [11]:
```python
y_pred = regr1.predict(X_test)
```

In [12]:
```python
# The mean squared error
print("Mean squared error: %.2f"
      % mean_squared_error(outputs, regr1.predict(inputs)))
# Explained variance score: 1 is perfect prediction
print('Variance score: %.2f' % regr1.score(inputs, outputs))
```

```
Mean squared error: 0.04
Variance score: 0.94
```

In [13]:
```python
# Score = 94% => model fits with ~ 94% data => This is suitable model.
```

In [14]:
```python
print('Variance score: %.2f' % r2_score(y_test, y_pred)) # y real, y predict
```

```
Variance score: 0.92
```

In [15]: 
```python
# Check the score of train and test
```

In [16]: 
```python
regr1.score(X_train, y_train)
```

Out[16]: 0.9412513438876041

In [17]: 
```python
regr1.score(X_test, y_test)
```
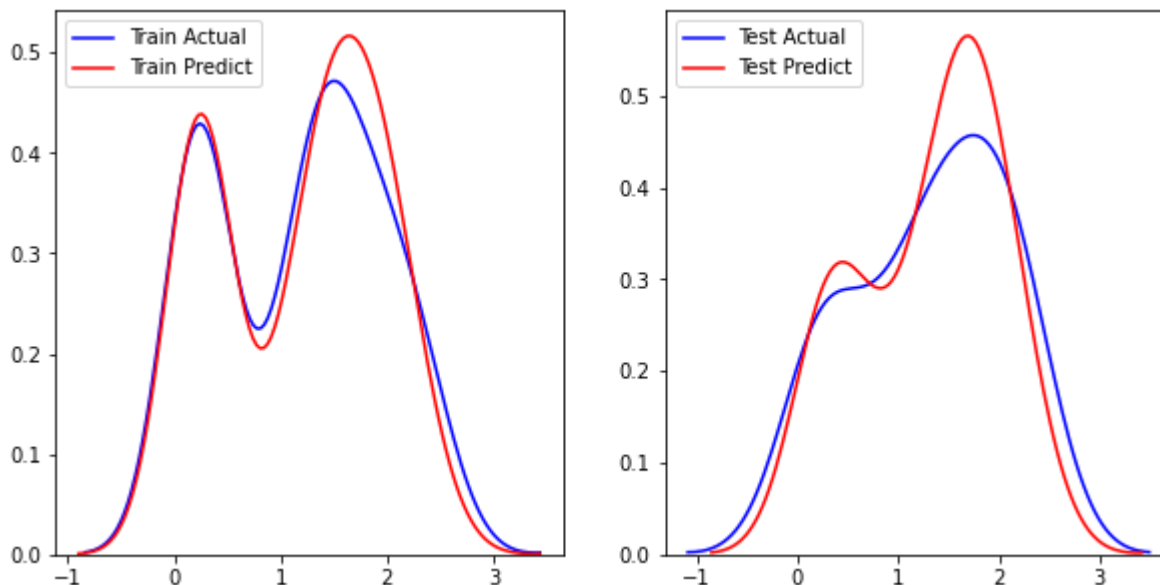
Out[17]: 0.9205639872829582

In [18]: 
```python
# Both training data and testing data have high score.
# => Choose this model.
```

In [19]: 
```python
# The coefficients
m=regr1.coef_
b=regr1.intercept_
print('Coefficients: \n', m)
print('Interceft: \n', b)
```

```
Coefficients:
 [[-0.19423398  0.23716335  0.51604647]]
Interceft:
 [-0.33344587]
```

In [20]: 
```python
# Visualization
y_train_hat = regr1.predict(X_train)
y_test_hat = regr1.predict(X_test)
```

```python
In [21]: plt.figure(figsize=(10,5))
         plt.subplot(1, 2, 1)
         ax1 = sns.distplot(y_train, hist=False, color="b", label='Train Actual')
         sns.distplot(y_train_hat, hist=False, color="r", label='Train Predict', ax=ax1)
         plt.subplot(1,2,2)
         ax2 = sns.distplot(y_test, hist=False, color="b", label='Test Actual')
         sns.distplot(y_test_hat, hist=False, color="r", label='Test Predict', ax=ax2)
         plt.show()
```



```python
In [22]: # Make new prediction
         x_now = [[4.5, 3.1, 1.6]]
         y_now = regr1.predict(x_now)
         print(y_now)
```

```
[[0.35338191]]
```

# Select important features

### Solution 1: SelectKBest

```python
In [23]: # Univariate Selection
         from sklearn.feature_selection import SelectKBest
         from sklearn.feature_selection import f_regression
```

In [24]:
```python
# Apply SelectKBest class to extract all best features
bestfeatures = SelectKBest(score_func=f_regression, k='all')
fit = bestfeatures.fit(inputs,outputs)
dfscores = pd.DataFrame(fit.scores_)
dfcolumns = pd.DataFrame(inputs.columns)
```

```
c:\program files\python36\lib\site-packages\sklearn\utils\validation.py:724: Da
taConversionWarning: A column-vector y was passed when a 1d array was expected.
Please change the shape of y to (n_samples, ), for example using ravel().
  y = column_or_1d(y, warn=True)
```

In [25]:
```python
#concat two dataframes for better visualization
featureScores = pd.concat([dfcolumns,dfscores],axis=1)
featureScores.columns = ['Specs','Score']  # naming the dataframe columns
print(featureScores.nlargest(3,'Score'))  # print 3 best features
```

```
        Specs       Score
2  petallength  1876.657813
0  sepallength   299.194957
1   sepalwidth    21.554378
```

In [26]:
```python
# 2 features have highest scores
X_now = inputs[['petallength', 'sepallength']]
```

In [27]:
```python
X_train_n, X_test_n, y_train_n, y_test_n = train_test_split(X_now, outputs,
                                                          test_size=0.20)

regr_n = linear_model.LinearRegression()
regr_n = regr1.fit(X_train_n, y_train_n)
```

In [28]:
```python
# The mean squared error
print("Mean squared error: %.2f"
      % mean_squared_error(outputs, regr_n.predict(X_now)))
# Explained variance score: 1 is perfect prediction
print('Variance score: %.2f' % regr_n.score(X_now, outputs))
```

```
Mean squared error: 0.04
Variance score: 0.93
```

In [29]:
```python
print("Train's score:", regr_n.score(X_train_n, y_train_n))
```

```
Train's score: 0.9282284238669857
```

In [30]:
```python
print("Test's score:", regr_n.score(X_test_n, y_test_n))
```

```
Test's score: 0.9254702078887744
```

## Solution 2: ExtraTreesRegressor

- https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesRegressor.html (https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesRegressor.html)

In [31]:
```python
from sklearn.ensemble import ExtraTreesRegressor
```
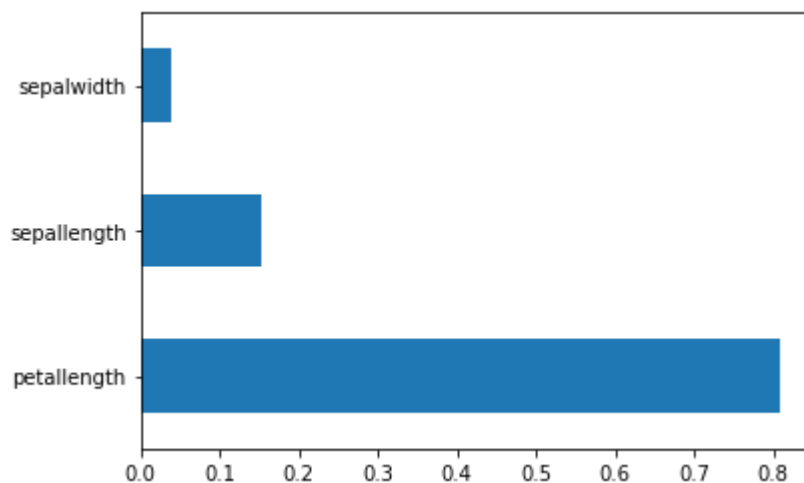
In [32]:
```python
model = ExtraTreesRegressor()
model.fit(inputs,outputs)
```

c:\program files\python36\lib\site-packages\ipykernel_launcher.py:2: DataConver sionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

Out[32]:
```
ExtraTreesRegressor(bootstrap=False, criterion='mse', max_depth=None,
                    max_features='auto', max_leaf_nodes=None,
                    min_impurity_decrease=0.0, min_impurity_split=None,
                    min_samples_leaf=1, min_samples_split=2,
                    min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=None,
                    oob_score=False, random_state=None, verbose=0,
                    warm_start=False)
```

In [33]:
```python
print(model.feature_importances_)
# use inbuilt class feature_importances of tree based regressor
# plot graph of feature importances for better visualization
feat_importances = pd.Series(model.feature_importances_, index=inputs.columns)
feat_importances.nlargest(3).plot(kind='barh')
plt.show()
```

[0.15356028 0.03893405 0.80750567]



In [34]:
```python
# Tương tự: 2 thuộc tính quan trọng nhất vẫn là 'petallength', 'sepallength'
```

## Solution 3: Correlation Matrix with Heatmap

In [35]:
```python
#get correlations of each features in dataset
data_sub = iris.iloc[:,0:4]
corrmat = data_sub.corr()
top_corr_features = corrmat.index
```
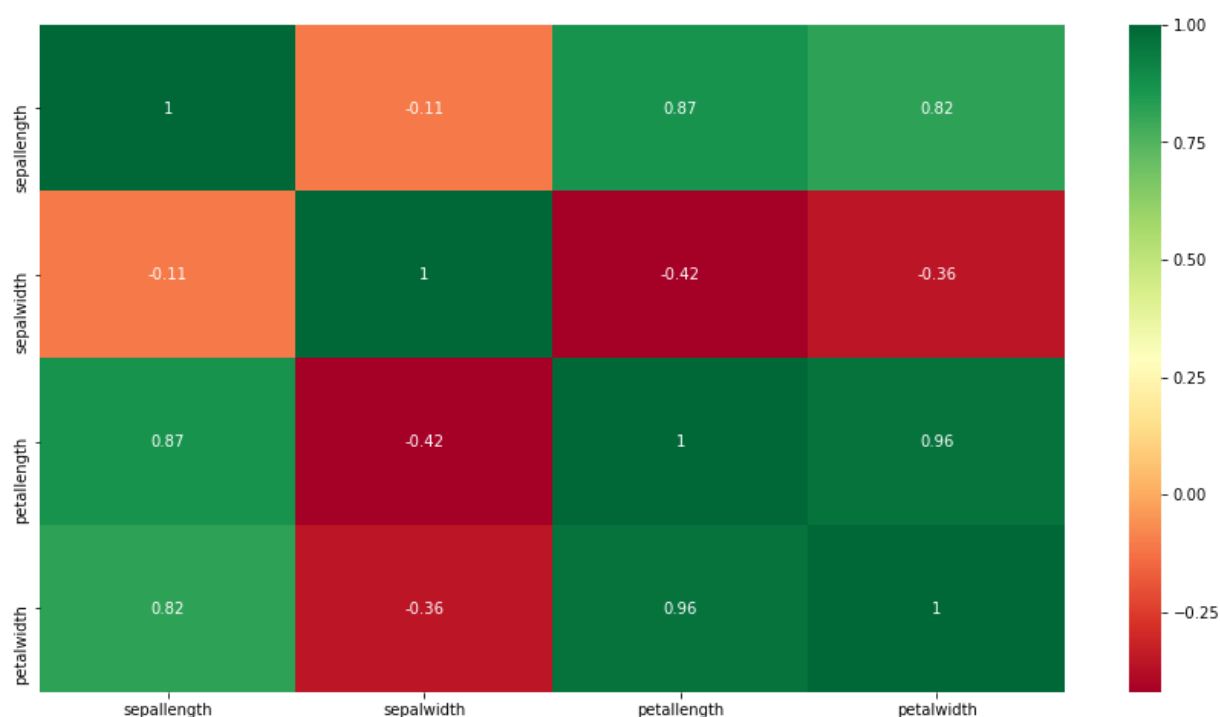
In [36]: 
```python
data_sub.corr()
```

Out[36]:

|  | sepallength | sepalwidth | petallength | petalwidth |
|---|---|---|---|---|
| sepallength | 1.000000 | -0.109369 | 0.871754 | 0.817954 |
| sepalwidth | -0.109369 | 1.000000 | -0.420516 | -0.356544 |
| petallength | 0.871754 | -0.420516 | 1.000000 | 0.962757 |
| petalwidth | 0.817954 | -0.356544 | 0.962757 | 1.000000 |

In [37]: 
```python
plt.figure(figsize=(15,8))
#plot heat map
g=sns.heatmap(data_sub[top_corr_features].corr(),cmap="RdYlGn", annot=True) # an
```



In [ ]: