

# Supervised Learning with scikit-learn

## 2). Regression

### a). Importing Data for Supervised Learning

# Import numpy and pandas

```
import numpy as np
```

```
import pandas as pd
```

# Read the CSV file into a DataFrame: df

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

# Create arrays for features and target variable

```
y = df['life'].values
```

```
X = df['fertility'].values
```

# Print the dimensions of X and y before reshaping

```
print("Dimensions of y before reshaping: {}".format(y.shape))
```

```
print("Dimensions of X before reshaping: {}".format(X.shape))
```

# Reshape X and y

```
y = y.reshape(-1,1)
```

```
X = X.reshape(-1,1)
```

# Print the dimensions of X and y after reshaping

```
print("Dimensions of y after reshaping: {}".format(y.shape))
```

```
print("Dimensions of X after reshaping: {}".format(X.shape))
```

<script.py> output:

Dimensions of y before reshaping: (139,)

Dimensions of X before reshaping: (139,)

Dimensions of y after reshaping: (139, 1)

Dimensions of X after reshaping: (139, 1)

**b). Fit & predict for regression**

```
# Import LinearRegression
```

```
from sklearn.linear_model import LinearRegression
```

```
# Create the regressor: reg
```

```
reg = LinearRegression()
```

```
# Create the prediction space
```

```
prediction_space = np.linspace(min(X_fertility), max(X_fertility)).reshape(-1,1)
```

```
# Fit the model to the data
```

```
reg.fit(X_fertility, y)
```

```
# Compute predictions over the prediction space: y_pred
```

```
y_pred = reg.predict(prediction_space)
```

```
# Print R^2
```

```
print(reg.score(X_fertility, y))
```

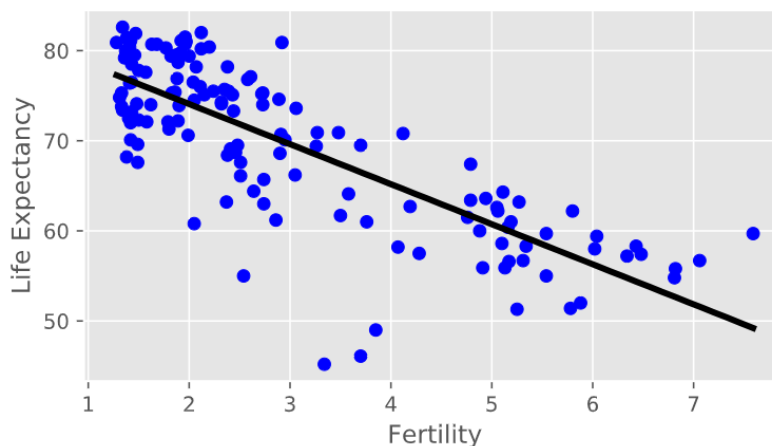
```
# Plot regression line
```

```
plt.plot(prediction_space, y_pred, color='black', linewidth=3)
```

```
plt.show()
```

<script.py> output:

0.619244216774



**c). Train/Test Split for regression****# Import necessary modules****from sklearn.linear\_model import LinearRegression****from sklearn.metrics import mean\_squared\_error****from sklearn.model\_selection import train\_test\_split****# 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)****# Create the regressor: reg\_all****reg\_all = LinearRegression()****# Fit the regressor to the training data****reg\_all.fit(X\_train, y\_train)****# Predict on the test data: y\_pred****y\_pred = reg\_all.predict(X\_test)****# Compute and print R^2 and RMSE****print("R^2: {}".format(reg\_all.score(X\_test, y\_test)))****rmse = np.sqrt(mean\_squared\_error(y\_test, y\_pred))****print("Root Mean Squared Error: {}".format(rmse))**

&lt;script.py&gt; output:

R^2: 0.838046873142936

Root Mean Squared Error: 3.2476010800377213

**d). 5-fold cross-validation**

**# Import the necessary modules**

**from sklearn.linear\_model import LinearRegression**

**from sklearn.model\_selection import cross\_val\_score**

**# Create a linear regression object: reg**

**reg = LinearRegression()**

**# Compute 5-fold cross-validation scores: cv\_scores**

**cv\_scores = cross\_val\_score(reg, X, y, cv=5)**

**# Print the 5-fold cross-validation scores**

**print(cv\_scores)**

**print("Average 5-Fold CV Score: {}".format(np.mean(cv\_scores)))**

<script.py> output:

[ 0.81720569 0.82917058 0.90214134 0.80633989 0.94495637]

Average 5-Fold CV Score: 0.8599627722793232

**e). K Fold CV Comparison**

**# Import necessary modules**

**from sklearn.linear\_model import LinearRegression**

**from sklearn.model\_selection import cross\_val\_score**

**# Create a linear regression object: reg**

**reg = LinearRegression()**

**# Perform 3-fold CV**

**cvscores\_3 = cross\_val\_score(reg,X,y, cv=3)**

**print(np.mean(cvscores\_3))**

**# Perform 10-fold CV**

**cvscores\_10 = cross\_val\_score(reg,X,y,cv=10)**

**print(np.mean(cvscores\_10))**

<script.py> output:

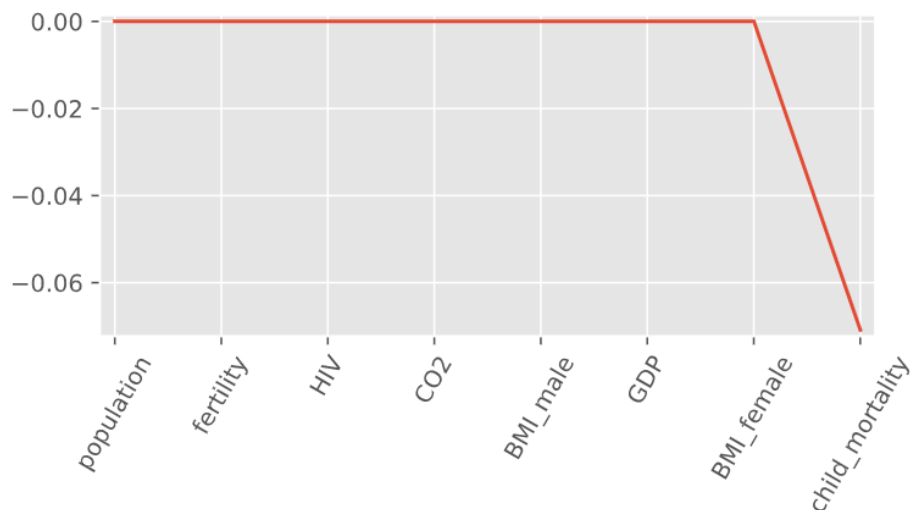
0.871871278262

0.843612862013

**f). Regularization I: Lasso****# Import Lasso****from sklearn.linear\_model import Lasso****# Instantiate a lasso regressor: lasso****lasso = Lasso(alpha=0.4, normalize=True)****# Fit the regressor to the data****lasso.fit(X, y)****# Compute and print the coefficients****lasso\_coef = lasso.coef\_****print(lasso\_coef)****# Plot the coefficients****plt.plot(range(len(df\_columns)), lasso\_coef)****plt.xticks(range(len(df\_columns)), df\_columns.values, rotation=60)****plt.margins(0.02)****plt.show()**

&lt;script.py&gt; output:

```
[-0.    -0.    -0.    0.    0.    0.   -0.
 -0.07087587]
```



**g). Regularization II: Ridge**

```
# Import necessary modules
from sklearn.linear_model import Ridge
from sklearn.model_selection import cross_val_score

# Setup the array of alphas and lists to store scores
alpha_space = np.logspace(-4, 0, 50)
ridge_scores = []
ridge_scores_std = []

# Create a ridge regressor: ridge
ridge = Ridge(normalize=True)

# Compute scores over range of alphas
for alpha in alpha_space:

    # Specify the alpha value to use: ridge.alpha
    ridge.alpha = alpha

    # Perform 10-fold CV: ridge_cv_scores
    ridge_cv_scores = cross_val_score(ridge,X,y,cv=10)

    # Append the mean of ridge_cv_scores to ridge_scores
    ridge_scores.append(np.mean(ridge_cv_scores))

    # Append the std of ridge_cv_scores to ridge_scores_std
    ridge_scores_std.append(np.std(ridge_cv_scores))

# Display the plot
display_plot(ridge_scores, ridge_scores_std)
```

