

# Chapter 1

## 1)

- Import `KNeighborsClassifier` from `sklearn.neighbors`.
- Create arrays `x` and `y` for the features and the target variable. Here this has been done for you. Note the use of `.drop()` to drop the target variable `'party'` from the feature array `x` as well as the use of the `.values` attribute to ensure `x` and `y` are NumPy arrays. Without using `.values`, `x` and `y` are a DataFrame and Series respectively; the scikit-learn API will accept them in this form also as long as they are of the right shape.
- Instantiate a `KNeighborsClassifier` called `knn` with 6 neighbors by specifying the `n_neighbors` parameter.
- Fit the classifier to the data using the `.fit()` method.

```
# Import KNeighborsClassifier from sklearn.neighbors  
from sklearn.neighbors import KNeighborsClassifier
```

```
# Create arrays for the features and the response variable  
y = df['party'].values  
X = df.drop('party', axis=1).values
```

```
# Create a k-NN classifier with 6 neighbors  
knn = KNeighborsClassifier(n_neighbors=6)
```

```
# Fit the classifier to the data  
knn.fit(X,y)
```

## 2)

```
# Predict the labels for the training data X  
y_pred = knn.predict(X)
```

```
# Predict and print the label for the new data point X_new  
new_prediction = knn.predict(X_new)  
print("Prediction: {}".format(new_prediction))
```

## 3)

```
# Import necessary modules  
from sklearn import datasets
```

```
import matplotlib.pyplot as plt

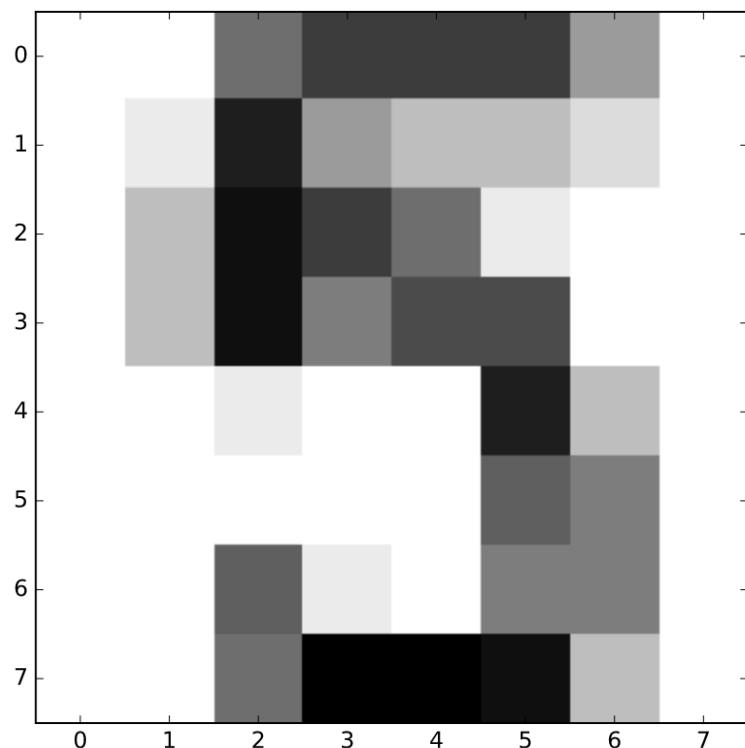
# Load the digits dataset: digits
digits = datasets.load_digits()

# Print the keys and DESCR of the dataset
print(digits.keys())
dict_keys(['DESCR', 'data', 'target_names', 'images', 'target'])

# Print the shape of the images and data keys
print(digits.images.shape)
(1797, 8, 8)

print(digits.data.shape)
(1797, 64)

# Display digit 1010
plt.imshow(digits.images[1010], cmap=plt.cm.gray_r, interpolation='nearest')
plt.show()
```



## 4)

- Import `KNeighborsClassifier` from `sklearn.neighbors` and `train_test_split` from `sklearn.model_selection`.
- Create stratified training and test sets using `0.2` for the size of the test set. Use a random state of `42`. **Stratify the split according to the labels so that they are distributed in the training and test sets as they are in the original dataset.**

```
# Import necessary modules
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split

# Create feature and target arrays
X = digits.data
y = digits.target

# Split into training and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2,
random_state=42, stratify=y)

# Create a k-NN classifier with 7 neighbors: knn
knn = KNeighborsClassifier(n_neighbors=7)

# Fit the classifier to the training data
knn.fit(X_train,y_train)

# Print the accuracy
print(knn.score(X_test, y_test))
```

## 5)

```
# Setup arrays to store train and test accuracies
neighbors = np.arange(1, 9)
train_accuracy = np.empty(len(neighbors))
test_accuracy = np.empty(len(neighbors))

# Loop over different values of k
for i, k in enumerate(neighbors):
    # Setup a k-NN Classifier with k neighbors: knn
    knn = KNeighborsClassifier(n_neighbors=k)
```

```

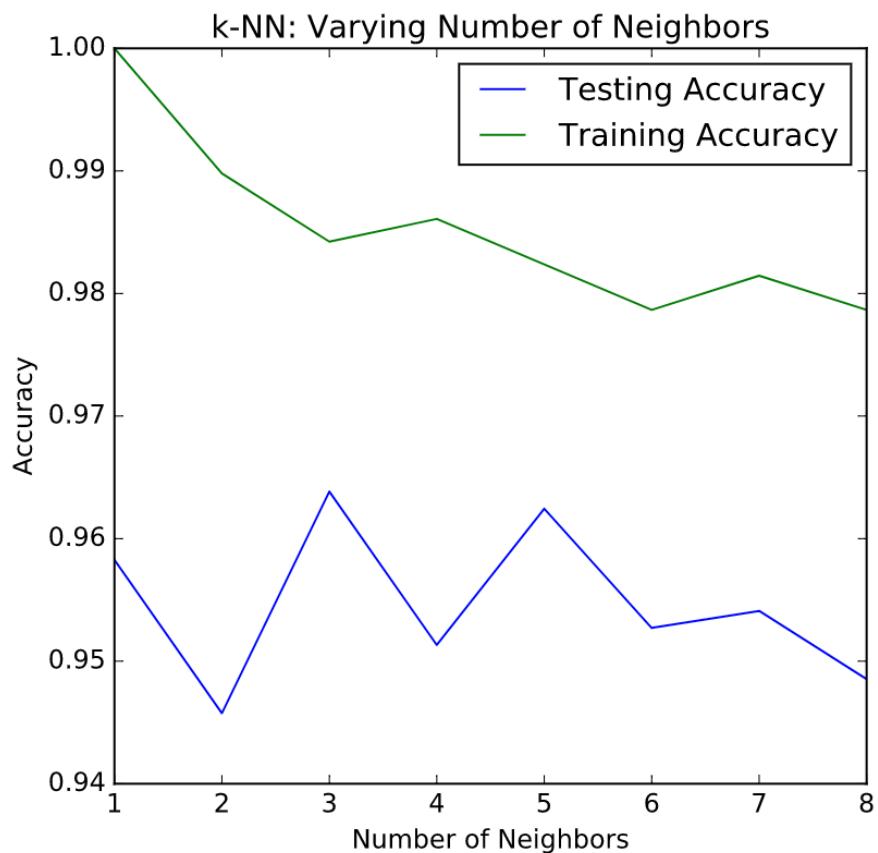
# Fit the classifier to the training data
knn.fit(X_train,y_train)

#Compute accuracy on the training set
train_accuracy[i] = knn.score(X_train, y_train)

#Compute accuracy on the testing set
test_accuracy[i] = knn.score(X_test, y_test)

# Generate plot
plt.title('k-NN: Varying Number of Neighbors')
plt.plot(neighbors, test_accuracy, label = 'Testing Accuracy')
plt.plot(neighbors, train_accuracy, label = 'Training Accuracy')
plt.legend()
plt.xlabel('Number of Neighbors')
plt.ylabel('Accuracy')
plt.show()

```



## Chapter 2

1)

```
# 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']
X = df['fertility']

# 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))
```

Dimensions of y before reshaping: (139,)

Dimensions of X before reshaping: (139,)

```
# 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))
```

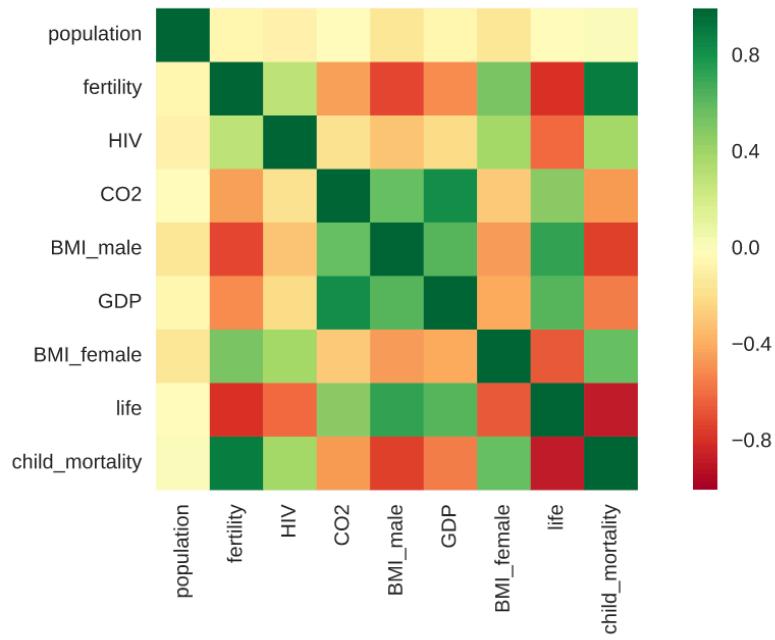
```
print("Dimensions of X after reshaping: {}".format(X.shape))
```

Dimensions of y after reshaping: (139, 1)

Dimensions of X after reshaping: (139, 1)

2)

```
import seaborn as sns
plt.figure()
sns.heatmap(df.corr(), square=True, cmap='RdYlGn')
plt.show()
```



**3)**

```
# 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))

R^2: 0.838046873142936
Root Mean Squared Error: 3.2476010800377213
```

## 4)

### Cross-validation sử dụng $R^2$ để đánh giá model

```
# 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(cv_scores.mean()))
[ 0.81720569  0.82917058  0.90214134  0.80633989  0.94495637]
Average 5-Fold CV Score: 0.8599627722793232
```

## 5)

### Lasso

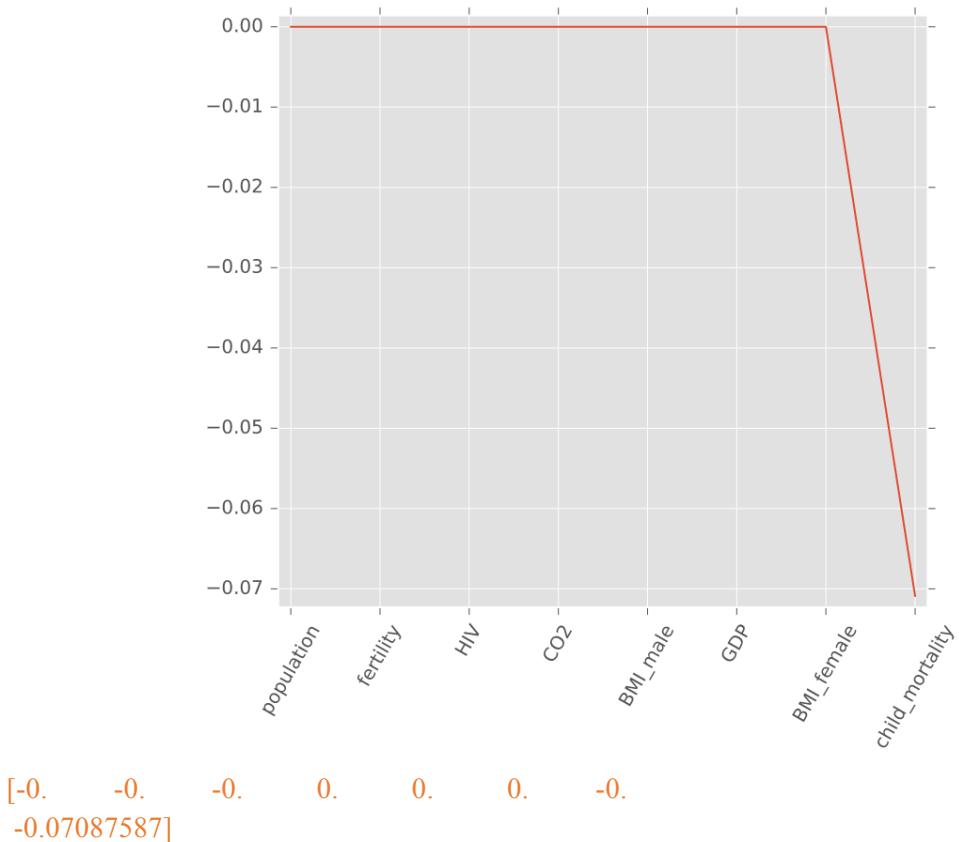
```
# Import Lasso
from sklearn.linear_model import Lasso

# 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()
```



## 6) Ridge

```
def display_plot(cv_scores, cv_scores_std):
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)
    ax.plot(alpha_space, cv_scores)

    std_error = cv_scores_std / np.sqrt(10)

    ax.fill_between(alpha_space, cv_scores + std_error, cv_scores -
    std_error, alpha=0.2)
    ax.set_ylabel('CV Score +/- Std Error')
    ax.set_xlabel('Alpha')
    ax.axhline(np.max(cv_scores), linestyle='--', color='.5')
    ax.set_xlim([alpha_space[0], alpha_space[-1]])
    ax.set_xscale('log')
    plt.show()
```

```
# 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

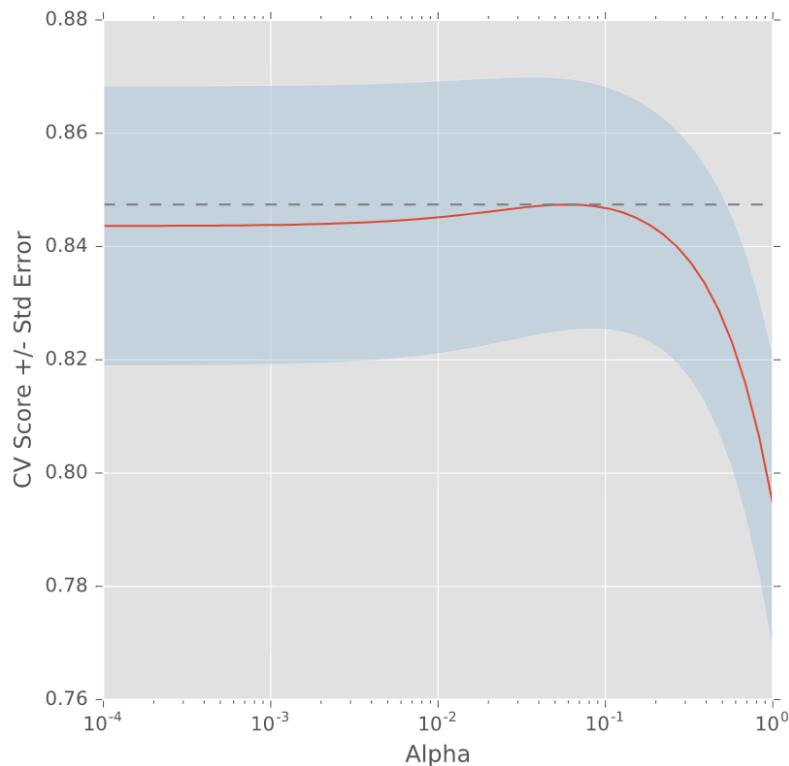
    # Perform 10-fold CV: ridge_cv_scores
    ridge_cv_scores = cross_val_score(ridge,X,y,cv=10)

    # 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)

```



## 7

```

# 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)

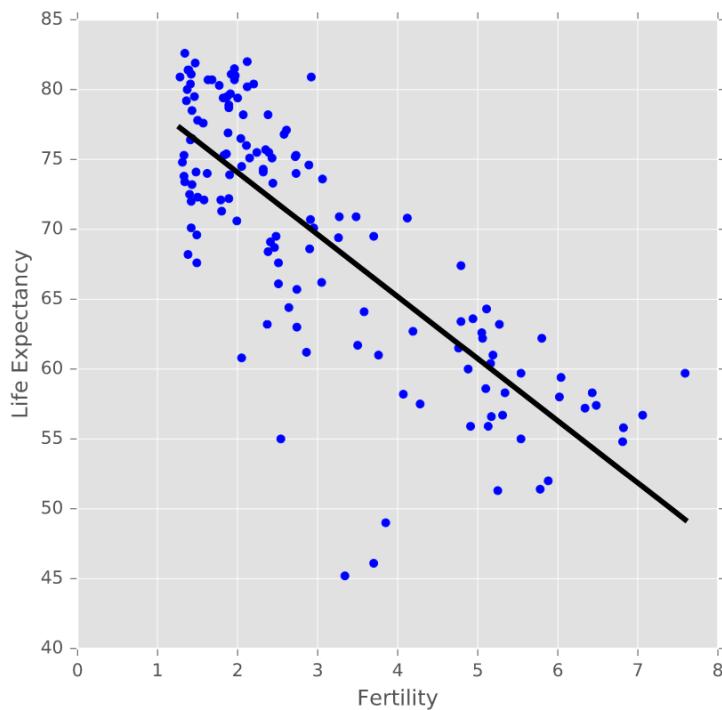
# Fit the model to the data
reg.fit(X_fertility,y)

# 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()

```



## Chapter 3

Note)

Confusion matrix

	Predicted: Spam Email	Predicted: Real Email
Actual: Spam Email	True Positive	False Negative
Actual: Real Email	False Positive	True Negative

## Metrics from the confusion matrix

- Precision :  $\frac{tp}{tp + fp}$
- Recall :  $\frac{tp}{tp + fn}$
- F1 score :  $2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$
- High precision: Not many real emails predicted as spam
- High recall: Predicted most spam emails correctly



1)

```
# Import necessary modules
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
```

```
# Create training and test set
X_train, X_test, y_train, y_test =
train_test_split(X,y,test_size=0.4,random_state=42)
```

```
# Instantiate a k-NN classifier: knn
knn = KNeighborsClassifier(n_neighbors=6)
```

```

# Fit the classifier to the training data
knn.fit(X_train,y_train)

# Predict the labels of the test data: y_pred
y_pred = knn.predict(X_test)

# Generate the confusion matrix and classification report
print(confusion_matrix(y_test, y_pred))
[[176 30]
 [ 52 50]]

print(classification_report(y_test, y_pred))
precision    recall  f1-score   support

          0       0.77      0.85      0.81     206
          1       0.62      0.49      0.55     102

avg / total       0.72      0.73      0.72     308

```

**Note)**

## Plotting the ROC curve

```

In [1]: from sklearn.metrics import roc_curve

In [2]: y_pred_prob = logreg.predict_proba(X_test)[:,1]

In [3]: fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)

In [4]: plt.plot([0, 1], [0, 1], 'k--')

In [5]: plt.plot(fpr, tpr, label='Logistic Regression')

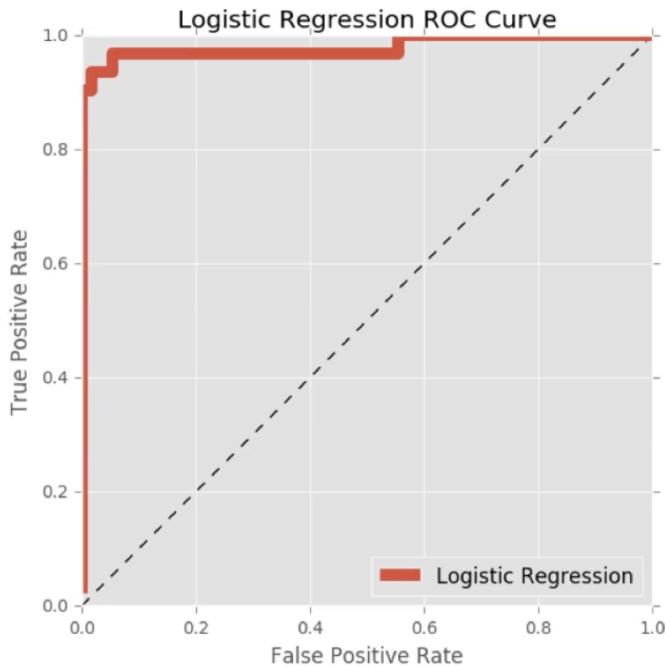
In [6]: plt.xlabel('False Positive Rate')

In [7]: plt.ylabel('True Positive Rate')

In [8]: plt.title('Logistic Regression ROC Curve')

In [9]: plt.show();

```



```
logreg.predict_proba(X_test)[:, 1]
```

## 2)

```
# Import the necessary modules
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix, classification_report

# Create training 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 classifier: logreg
logreg = LogisticRegression()

# Fit the classifier to the training data
logreg.fit(X_train,y_train)

# Predict the labels of the test set: y_pred
y_pred = logreg.predict(X_test)

# Compute and print the confusion matrix and classification report
print(confusion_matrix(y_test, y_pred))
print(classification_report(y_test, y_pred))
```

```

[[176 30]
 [35 67]]
precision    recall   f1-score   support

          0       0.83      0.85      0.84     206
          1       0.69      0.66      0.67     102

avg / total   0.79      0.79      0.79     308

```

### 3)

```

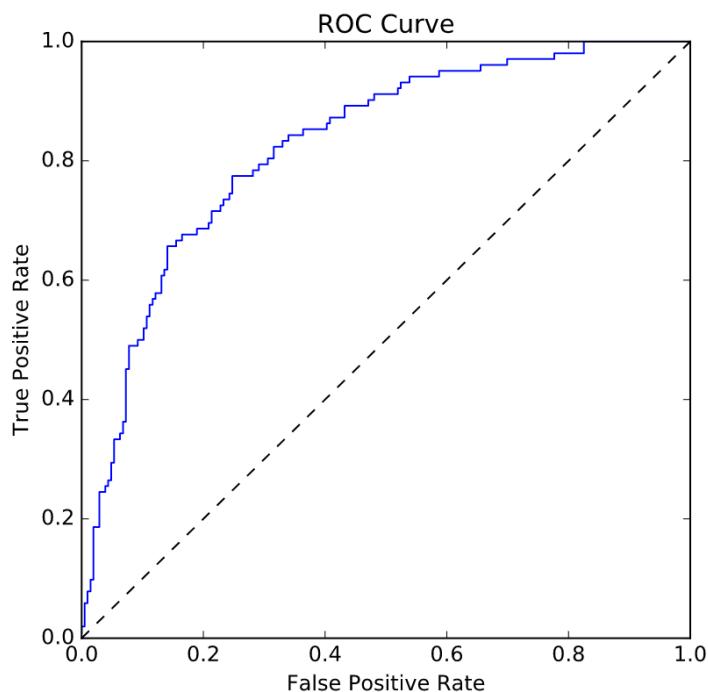
# Import necessary modules
from sklearn.metrics import roc_curve

# Compute predicted probabilities: y_pred_prob
y_pred_prob = logreg.predict_proba(X_test)[:,1]

# Generate ROC curve values: fpr, tpr, thresholds
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)

# Plot ROC curve
plt.plot([0, 1], [0, 1], 'k--')
plt.plot(fpr,tpr)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve')
plt.show()

```



## 4)

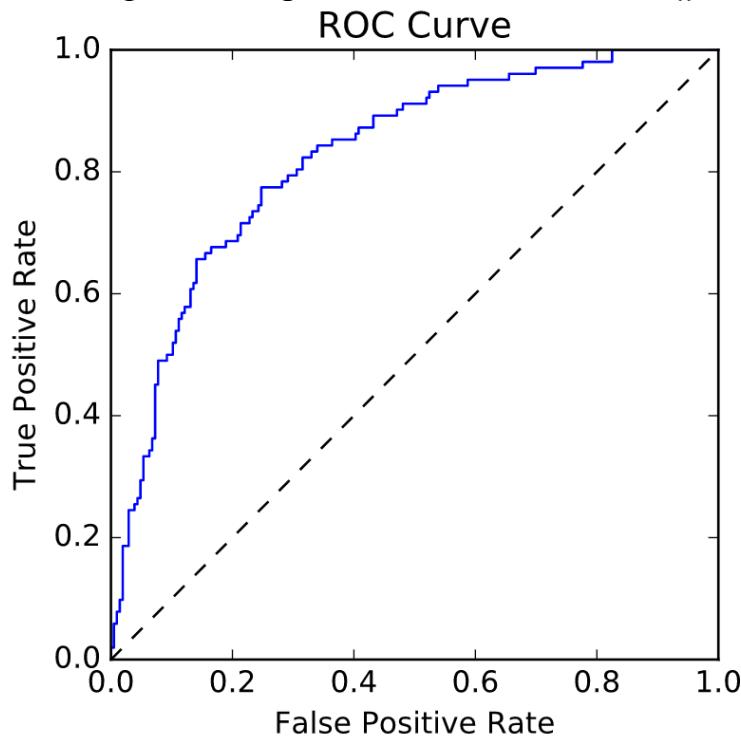
```
# Import necessary modules
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import cross_val_score

# Compute predicted probabilities: y_pred_prob
y_pred_prob = logreg.predict_proba(X_test)[:,1]

# Compute and print AUC score
print("AUC: {}".format(roc_auc_score(y_test, y_pred_prob)))
AUC: 0.825480677707976

# Compute cross-validated AUC scores: cv_auc
cv_auc = cross_val_score(logreg,X,y,cv=5,scoring='roc_auc')

# Print list of AUC scores
print("AUC scores computed using 5-fold cross-validation: {}".format(cv_auc))
```



AUC scores computed using 5-fold cross-validation: [ 0.80148148 0.8062963 0.81481481 0.86245283 0.8554717 ]

Note)

## Grid search cross-validation

C	0.5	0.701	0.703	0.697	0.696
	0.4	0.699	0.702	0.698	0.702
	0.3	0.721	0.726	0.713	0.703
	0.2	0.706	0.705	0.704	0.701
	0.1	0.698	0.692	0.688	0.675
		0.1	0.2	0.3	0.4

Alpha

## GridSearchCV in scikit-learn

```
In [1]: from sklearn.model_selection import GridSearchCV  
  
In [2]: param_grid = {'n_neighbors': np.arange(1, 50)}  
  
In [3]: knn = KNeighborsClassifier()  
  
In [4]: knn_cv = GridSearchCV(knn, param_grid, cv=5)  
  
In [5]: knn_cv.fit(X, y)  
  
In [6]: knn_cv.best_params_  
Out[6]: {'n_neighbors': 12}  
  
In [7]: knn_cv.best_score_  
Out[7]: 0.933216168717
```

5)

## Hyperparameter tuning with GridSearchCV

Hugo demonstrated how to use to tune the `n_neighbors` parameter of the `KNeighborsClassifier()` using GridSearchCV on the voting dataset. You will now practice this yourself, but by using logistic regression on the diabetes dataset instead!

Like the alpha parameter of lasso and ridge regularization that you saw earlier, logistic regression also has a regularization parameter: C. C controls the *inverse* of the regularization strength, and this is what you will tune in this exercise. A large C can lead to an *overfit* model, while a small C can lead to an *underfit* model.

The hyperparameter space for C has been setup for you. Your job is to use GridSearchCV and logistic regression to find the optimal C in this hyperparameter space. The feature array is available as `x` and target variable array is available as `y`.

```
# Import necessary modules
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV

# Setup the hyperparameter grid
c_space = np.logspace(-5, 8, 15)
param_grid = {'C': c_space}

# Instantiate a logistic regression classifier: logreg
logreg = LogisticRegression()

# Instantiate the GridSearchCV object: logreg_cv
logreg_cv = GridSearchCV(logreg, param_grid, cv=5)

# Fit it to the data

logreg_cv.fit(X,y)
# Print the tuned parameters and score
print("Tuned Logistic Regression Parameters:
{}".format(logreg_cv.best_params_))
print("Best score is {}".format(logreg_cv.best_score_))

Tuned Logistic Regression Parameters: {'C': 3.7275937203149381}
Best score is 0.7708333333333334
```

## 6)

## Hyperparameter tuning with RandomizedSearchCV

GridSearchCV can be computationally expensive, especially if you are searching over a large hyperparameter space and dealing with multiple

hyperparameters. A solution to this is to use `RandomizedSearchCV`, in which not all hyperparameter values are tried out. Instead, a fixed number of hyperparameter settings is sampled from specified probability distributions. You'll practice using `RandomizedSearchCV` in this exercise and see how this works.

Here, you'll also be introduced to a new model: the Decision Tree. Don't worry about the specifics of how this model works. Just like k-NN, linear regression, and logistic regression, decision trees in scikit-learn

have `.fit()` and `.predict()` methods that you can use in exactly the same way as before. Decision trees have many parameters that can be tuned, such as `max_features`, `max_depth`, and `min_samples_leaf`: This makes it an ideal use case for `RandomizedSearchCV`.

```
# Import necessary modules
from scipy.stats import randint
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import RandomizedSearchCV

# Setup the parameters and distributions to sample from: param_dist
param_dist = {"max_depth": [3, None],
              "max_features": randint(1, 9),
              "min_samples_leaf": randint(1, 9),
              "criterion": ["gini", "entropy"]}

# Instantiate a Decision Tree classifier: tree
tree = DecisionTreeClassifier()

# Instantiate the RandomizedSearchCV object: tree_cv
tree_cv = RandomizedSearchCV(tree, param_dist, cv=5)

# Fit it to the data
tree_cv.fit(X,y)

# Print the tuned parameters and score
print("Tuned Decision Tree Parameters: {}".format(tree_cv.best_params_))
print("Best score is {}".format(tree_cv.best_score_))

Tuned Decision Tree Parameters: {'criterion': 'gini', 'max_depth': 3, 'min_samples_leaf': 2,
'max_features': 8}
Best score is 0.7213541666666666
```

## 7)

```
# Import necessary modules
from sklearn.model_selection import train_test_split
```

```

from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV

# Create the hyperparameter grid
c_space = np.logspace(-5, 8, 15)
param_grid = { 'C': c_space, 'penalty': ['l1', 'l2']}

# Instantiate the logistic regression classifier: logreg
logreg = LogisticRegression()

# 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)

# Instantiate the GridSearchCV object: logreg_cv
logreg_cv = GridSearchCV(logreg,param_grid,cv=5)

# Fit it to the training data

logreg_cv.fit(X_train,y_train)
# Print the optimal parameters and best score
print("Tuned Logistic Regression Parameter:
{} ".format(logreg_cv.best_params_))
print("Tuned Logistic Regression Accuracy:
{} ".format(logreg_cv.best_score_))

Tuned Logistic Regression Parameter: {'C': 0.43939705607607948, 'penalty': 'l1'}
Tuned Logistic Regression Accuracy: 0.7652173913043478

```

## 8)

```

# Import necessary modules
from sklearn.linear_model import ElasticNet
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import train_test_split

# 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 hyperparameter grid
l1_space = np.linspace(0, 1, 30)
param_grid = {'l1_ratio': l1_space}

```

```
# Instantiate the ElasticNet regressor: elastic_net
elastic_net = ElasticNet()

# Setup the GridSearchCV object: gm_cv
gm_cv = GridSearchCV(elastic_net, param_grid, cv=5)

# Fit it to the training data
gm_cv.fit(X_train,y_train)

# Predict on the test set and compute metrics
y_pred = gm_cv.predict(X_test)
r2 = gm_cv.score(X_test, y_test)
mse = mean_squared_error(y_test, y_pred)
print("Tuned ElasticNet l1 ratio: {}".format(gm_cv.best_params_))
print("Tuned ElasticNet R squared: {}".format(r2))
print("Tuned ElasticNet MSE: {}".format(mse))

Tuned ElasticNet l1 ratio: {'l1_ratio': 0.20689655172413793}
Tuned ElasticNet R squared: 0.8668305372460283
Tuned ElasticNet MSE: 10.05791413339844
```

## Chapter 4

### Dealing with categorical features

- Scikit-learn will not accept categorical features by default
- Need to encode categorical features numerically
- Convert to ‘dummy variables’
  - 0: Observation was NOT that category
  - 1: Observation was that category

### Dealing with categorical features in Python

- scikit-learn: OneHotEncoder()
- pandas: get\_dummies()

### Encoding dummy variables

```
In [1]: import pandas as pd  
  
In [2]: df = pd.read_csv('auto.csv')  
  
In [3]: df_origin = pd.get_dummies(df)  
  
In [4]: print(df_origin.head())  
      mpg  displ   hp  weight  accel  size  origin_Asia  origin_Europe  \\  
0  18.0  250.0   88    3139   14.5  15.0          0              0  
1   9.0  304.0  193    4732   18.5  20.0          0              0  
2  36.1   91.0   60    1800   16.4  10.0          1              0  
3  18.5  250.0   98    3525   19.0  15.0          0              0  
4  34.3   97.0   78    2188   15.8  10.0          0              1  
  
      origin_US  
0            1  
1            1  
2            0  
3            1  
4            0
```

### Encoding dummy variables

```
In [5]: df_origin = df_origin.drop('origin_Asia', axis=1)  
  
In [6]: print(df_origin.head())  
      mpg  displ   hp  weight  accel  size  origin_Europe  origin_US  
0  18.0  250.0   88    3139   14.5  15.0          0             1  
1   9.0  304.0  193    4732   18.5  20.0          0             1  
2  36.1   91.0   60    1800   16.4  10.0          0             0  
3  18.5  250.0   98    3525   19.0  15.0          0             1  
4  34.3   97.0   78    2188   15.8  10.0          1             0
```

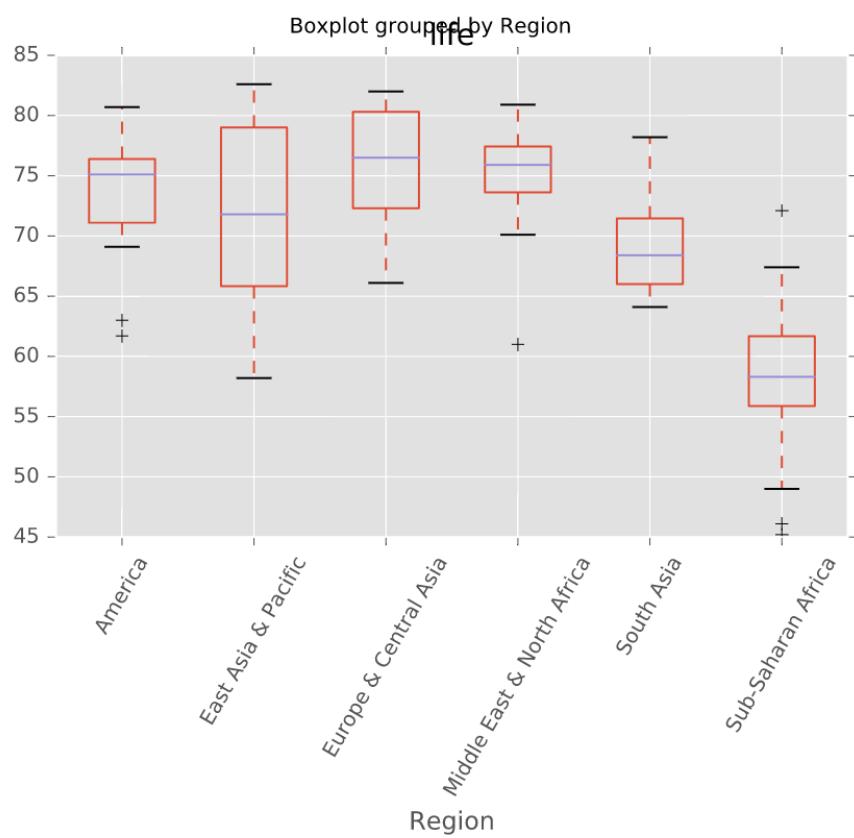
**1)**

```
# 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()
```



**2)**

```
# Create dummy variables: df_region
df_region = pd.get_dummies(df)

# Print the columns of df_region
print (df_region.columns)
```

```

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')

# Create dummy variables with drop_first=True: df_region
df_region = pd.get_dummies(df, drop_first=True)

# Print the new columns of df_region
print(df_region.columns)

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')

```

### 3)

```

# 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)
[ 0.86808336  0.80623545  0.84004203  0.7754344   0.87503712]

```

## Note)

# Imputing missing data

- Making an educated guess about the missing values
- Example: Using the mean of the non-missing entries

```

In [1]: from sklearn.preprocessing import Imputer
In [2]: imp = Imputer(missing_values='NaN', strategy='mean', axis=0)

```

# Imputing within a pipeline

```
In [1]: from sklearn.pipeline import Pipeline  
  
In [2]: from sklearn.preprocessing import Imputer  
  
In [3]: imp = Imputer(missing_values='NaN', strategy='mean', axis=0)  
  
In [4]: logreg = LogisticRegression()  
  
In [5]: steps = [('imputation', imp),  
...:             ('logistic_regression', logreg)]  
  
In [6]: pipeline = Pipeline(steps)
```

# Imputing within a pipeline

```
In [8]: pipeline.fit(X_train, y_train)  
  
In [9]: y_pred = pipeline.predict(X_test)  
  
In [10]: pipeline.score(X_test, y_test)  
Out[10]: 0.75324675324675328
```

4)

```
# 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))
```

Shape of Original DataFrame: (435, 17)

```
# 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))
```

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

## 5)

```
# Import the Imputer module
from sklearn.preprocessing import Imputer
from sklearn.svm import SVC

# Setup the Imputation transformer: imp
imp = Imputer(missing_values='NaN', strategy='most_frequent', axis=0)

# Instantiate the SVC classifier: clf
clf = SVC()

# Setup the pipeline with the required steps: steps
steps = [('imputation', imp),
          ('SVM', clf)]
```

## 6)

```
# 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))
```

```

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

```

**Note)**

## Scaling in a pipeline

```

In [6]: from sklearn.preprocessing import StandardScaler

In [7]: steps = [('scaler', StandardScaler()),
   ....      ('knn', KNeighborsClassifier())]

In [8]: pipeline = Pipeline(steps)

In [9]: X_train, X_test, y_train, y_test = train_test_split(X, y,
   ....: test_size=0.2, random_state=21)

In [10]: knn_scaled = pipeline.fit(X_train, y_train)

In [11]: y_pred = pipeline.predict(X_test)

In [12]: accuracy_score(y_test, y_pred)
Out[12]: 0.956

In [13]: knn_unscaled = KNeighborsClassifier().fit(X_train, y_train)

In [14]: knn_unscaled.score(X_test, y_test)
Out[14]: 0.928

```

7)

```
# 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)))
```

Mean of Unscaled Features: 18.432687072460002

Standard Deviation of Unscaled Features: 41.54494764094571

```
# 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)))
```

```
Mean of Scaled Features: 2.7314972981668206e-15  
Standard Deviation of Scaled Features: 0.9999999999999999
```

## 8)

```
# 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)))  
Accuracy with Scaling: 0.7700680272108843  
Accuracy without Scaling: 0.6979591836734694
```

## 9)

```
# 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,parameters, cv=3)

# 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_))

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}

## 10)

```

# 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,parameters,cv=3)

# 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))

Tuned ElasticNet Alpha: {'elasticnet__l1_ratio': 1.0}
Tuned ElasticNet R squared: 0.8862016570888217
```