# Data Science Track: Course 17

# Machine Learning with the Experts: School Budgets

## Chap 1: Exploring the raw data

In [1]:
```python
# Import plotting modules
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
import numpy as np
plt.style.use('ggplot')
```

### Exploring the data

In [4]:
```python
# A column for each possible value

df = pd.DataFrame({'Eyes':['Brown','Brown','Blue','Blue'],
                   'Hair':['Curly','Straight','Wavy','Straight']},
                  index=['Jamal','Luisa','Jenny','Max'])
df
```

Out[4]:

|  | Eyes | Hair |
|---|---|---|
| **Jamal** | Brown | Curly |
| **Luisa** | Brown | Straight |
| **Jenny** | Blue | Wavy |
| **Max** | Blue | Straight |

In [6]:
```python
df_dummies = pd.get_dummies(df)
df_dummies
```

Out[6]:

|  | Eyes_Blue | Eyes_Brown | Hair_Curly | Hair_Straight | Hair_Wavy |
|---|---|---|---|---|---|
| **Jamal** | 0 | 1 | 1 | 0 | 0 |
| **Luisa** | 0 | 1 | 0 | 1 | 0 |
| **Jenny** | 1 | 0 | 0 | 0 | 1 |
| **Max** | 1 | 0 | 0 | 1 | 0 |

In [7]:
```python
# Load and preview the data
sample_df = pd.read_csv('datasets/drivendata/sample_data.csv')
sample_df.head()
```

Out[7]:

|  | label | numeric | text | with_missing |
|---|---|---|---|---|
| **0** | a | -4.167578 | bar | -4.084883 |
| **1** | b | -0.562668 | NaN | 2.043464 |
| **2** | a | -21.361961 | NaN | -33.315334 |
| **3** | a | 16.402708 | foo bar | 30.884604 |
| **4** | a | -17.934356 | foo | -27.488405 |

In [9]:
```python
# Summarize the data
sample_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 4 columns):
label           5 non-null object
numeric         5 non-null float64
text            3 non-null object
with_missing    5 non-null float64
dtypes: float64(2), object(2)
memory usage: 240.0+ bytes
```

```
In [10]: sample_df.describe()
```

Out[10]:

|       | numeric    | with_missing |
|-------|------------|--------------|
| count | 5.000000   | 5.000000     |
| mean  | -5.524771  | -6.392111    |
| std   | 15.100440  | 25.670748    |
| min   | -21.361961 | -33.315334   |
| 25%   | -17.934356 | -27.488405   |
| 50%   | -4.167578  | -4.084883    |
| 75%   | -0.562668  | 2.043464     |
| max   | 16.402708  | 30.884604    |

```
In [12]: # EXERCISES
```

```
In [13]: # Loading the data
         df = pd.read_csv('datasets/drivendata/TrainingData.csv',index_col=0)
         df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 400277 entries, 134338 to 415831
Data columns (total 25 columns):
Function                400277 non-null object
Use                     400277 non-null object
Sharing                 400277 non-null object
Reporting               400277 non-null object
Student_Type            400277 non-null object
Position_Type           400277 non-null object
Object_Type             400277 non-null object
Pre_K                   400277 non-null object
Operating_Status        400277 non-null object
Object_Description       375493 non-null object
Text_2                  88217 non-null object
SubFund_Description     306855 non-null object
Job_Title_Description   292743 non-null object
Text_3                  179964 non-null object
Text_4                  53746 non-null object
Sub_Object_Description  91603 non-null object
Location_Description    162054 non-null object
FTE                     126071 non-null float64
Function_Description    342195 non-null object
Facility_or_Department  53886 non-null object
Position_Extra          264764 non-null object
Total                   395722 non-null float64
Program_Description     304660 non-null object
Fund_Description        202877 non-null object
Text_1                  292285 non-null object
dtypes: float64(2), object(23)
memory usage: 79.4+ MB
```

```python
In [16]:  # Summarizing the data

          # Print the summary statistics
          print(df.describe())

          # Import matplotlib.pyplot as plt
          import matplotlib.pyplot as plt

          # Create the histogram
          plt.hist(df['FTE'].dropna())

          # Add title and labels
          plt.title('Distribution of %full-time \n employee works')
          plt.xlabel('% of full-time')
          plt.ylabel('num employees')
          plt.xlim([0,6])

          # Display the histogram
          plt.show()
```
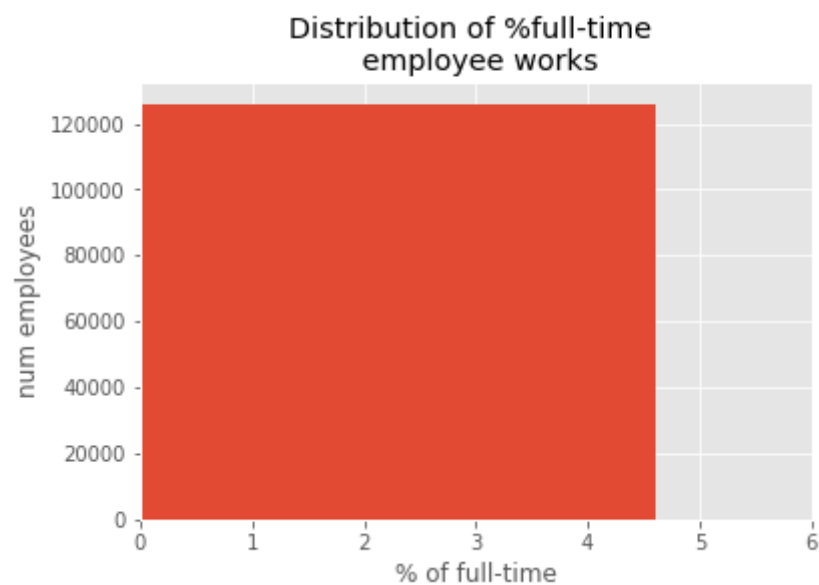
```
                   FTE          Total
count   126071.000000   3.957220e+05
mean         0.426794   1.310586e+04
std          0.573576   3.682254e+05
min         -0.087551  -8.746631e+07
25%          0.000792   7.379770e+01
50%          0.130927   4.612300e+02
75%          1.000000   3.652662e+03
max         46.800000   1.297000e+08
```



## Looking at the datatypes

```python
In [18]:  sample_df['label'].head()
```

```
Out[18]:  0    a
          1    b
          2    a
          3    a
          4    a
          Name: label, dtype: object
```

Why encode labels as categories?

- ML algorithms work on numbers, not strings
- Strings can be slow compared to numbers (take more space)

```python
In [19]:  # Encode labels as categories (sample data)

          sample_df.label.head(2)
```

```
Out[19]:  0    a
          1    b
          Name: label, dtype: object
```

```python
In [20]:  sample_df.label = sample_df.label.astype('category')
          sample_df.label.head(2)
```

```
Out[20]:  0    a
          1    b
          Name: label, dtype: category
          Categories (2, object): [a, b]
```

```
In [22]:  # Dummy variable encoding

          dummies = pd.get_dummies(sample_df[['label']], prefix_sep='_')
          dummies.head(2)
```

Out[22]:

|   | label_a | label_b |
|---|---------|---------|
| **0** | 1 | 0 |
| **1** | 0 | 1 |

```
In [24]:  # Lambda functions

          square = lambda x: x*x
          square(2)
```

Out[24]: 4

```
In [25]:  # Encode labels as categories

          categorize_label = lambda x: x.astype('category')
          sample_df[['label']] = sample_df[['label']].apply(categorize_label,axis=0)
          sample_df.info()
```

```
          <class 'pandas.core.frame.DataFrame'>
          RangeIndex: 5 entries, 0 to 4
          Data columns (total 4 columns):
          label           5 non-null category
          numeric         5 non-null float64
          text            3 non-null object
          with_missing    5 non-null float64
          dtypes: category(1), float64(2), object(1)
          memory usage: 301.0+ bytes
```

```
In [ ]:   # EXERCISES
```

```
In [27]:  # Exploring datatypes in pandas
          df.dtypes.value_counts()
```

Out[27]:
```
          object     23
          float64     2
          dtype: int64
```

```
In [28]:  # Encode the labels as categorical variables
          LABELS = ['Function','Use','Sharing','Reporting','Student_Type',
                    'Position_Type','Object_Type','Pre_K','Operating_Status']
```

```
In [29]:  df[LABELS].dtypes
```

Out[29]:
```
          Function          object
          Use               object
          Sharing           object
          Reporting         object
          Student_Type      object
          Position_Type     object
          Object_Type       object
          Pre_K             object
          Operating_Status  object
          dtype: object
```

```
In [30]:  # Define the lambda function: categorize_label
          categorize_label = lambda x: x.astype('category')

          # Convert df[LABELS] to a categorical type
          df[LABELS] = df[LABELS].apply(categorize_label,axis=0)

          # Print the converted dtypes
          print(df[LABELS].dtypes)
```

```
          Function          category
          Use               category
          Sharing           category
          Reporting         category
          Student_Type      category
          Position_Type     category
          Object_Type       category
          Pre_K             category
          Operating_Status  category
          dtype: object
```

```
In [31]:  # Counting unique labels

          # Import matplotlib.pyplot
          import matplotlib.pyplot as plt

          # Calculate number of unique values for each label: num_unique_labels
          num_unique_labels = df[LABELS].apply(lambda x: pd.Series.nunique(x))

          # Plot number of unique values for each label
          num_unique_labels.plot(kind='bar')

          # Label the axes
          plt.xlabel('Labels')
          plt.ylabel('Number of unique values')

          # Display the plot
          plt.show()
```
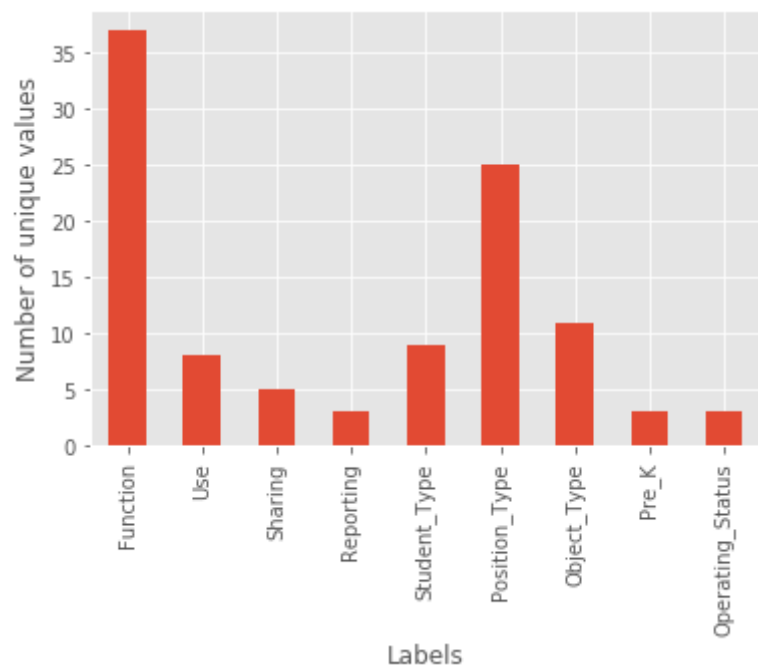


## How do we measure success?

Log loss binary classification
logloss(N=1) = y log(p) + (1-y) log(1-p)

```
In [33]:  # Computing log loss with NumPy

          def compute_log_loss(predicted, actual, eps=1e-14):
              """ Computes the logarithmic loss between predicted and
                  actual when these are 1D arrays.

              :param predicted: The predicted probabilities as floats between 0-1
              :param actual: The actual binary labels. Either 0 or 1.
              :param eps (optional): log(0) is inf, so we need to offset our
                                     predicted values slightly by eps from 0 or 1.
              """
              predicted = np.clip(predicted, eps, 1 - eps)
              loss = -1 * np.mean(actual * np.log(predicted)
                                  + (1 - actual) * np.log(1 - predicted))

              return loss
```

```
In [34]:  compute_log_loss(predicted=0.9, actual=0)
```

```
Out[34]:  2.3025850929940459
```

```
In [35]:  compute_log_loss(predicted=0.5, actual=1)
```

```
Out[35]:  0.69314718055994529
```

```
In [36]:  # EXERCISES
```

```
In [39]:  # Penalizing highly confident wrong answers
          print('A: {}'.format(compute_log_loss(predicted=0.85, actual=1)))
          print('B: {}'.format(compute_log_loss(predicted=0.99, actual=0)))
          print('C: {}'.format(compute_log_loss(predicted=0.51, actual=0)))

          A: 0.16251892949777494
          B: 4.605170185988091
          C: 0.7133498878774648
```

```
In [40]:  # Computing log loss with NumPy
          # 5 one-dimensional numeric arrays simulating different types of predictions

          actual_labels = np.array([ 1.,  1.,  1.,  1.,  1.,
                                     0.,  0.,  0.,  0.,  0.])
          correct_confident = np.array([ 0.95,  0.95,  0.95,  0.95,  0.95,
                                         0.05,  0.05,  0.05,  0.05,  0.05])
          correct_not_confident = np.array([ 0.65,  0.65,  0.65,  0.65,  0.65,
                                             0.35,  0.35,  0.35,  0.35,  0.35])
          wrong_not_confident = np.array([ 0.35,  0.35,  0.35,  0.35,  0.35,
                                           0.65,  0.65,  0.65,  0.65,  0.65])
          wrong_confident = np.array([ 0.05,  0.05,  0.05,  0.05,  0.05,
                                       0.95,  0.95,  0.95,  0.95,  0.95])
```

```
In [41]:  # Compute and print log loss for 1st case
          correct_confident = compute_log_loss(correct_confident, actual_labels)
          print("Log loss, correct and confident: {}".format(correct_confident))

          # Compute log loss for 2nd case
          correct_not_confident = compute_log_loss(correct_not_confident, actual_labels)
          print("Log loss, correct and not confident: {}".format(correct_not_confident))

          # Compute and print log loss for 3rd case
          wrong_not_confident = compute_log_loss(wrong_not_confident, actual_labels)
          print("Log loss, wrong and not confident: {}".format(wrong_not_confident))

          # Compute and print log loss for 4th case
          wrong_confident = compute_log_loss(wrong_confident, actual_labels)
          print("Log loss, wrong and confident: {}".format(wrong_confident))

          # Compute and print log loss for actual labels
          actual_labels = compute_log_loss(actual_labels, actual_labels)
          print("Log loss, actual labels: {}".format(actual_labels))
```

```
Log loss, correct and confident: 0.05129329438755058
Log loss, correct and not confident: 0.4307829160924542
Log loss, wrong and not confident: 1.049822124498678
Log loss, wrong and confident: 2.9957322735539904
Log loss, actual labels: 9.99200722162646e-15
```

# Chap 2: Creating a simple first model

```
In [65]:  # Import plotting modules
          import matplotlib.pyplot as plt
          import seaborn as sns
          import pandas as pd
          import numpy as np
          plt.style.use('ggplot')

          from warnings import warn
          from sklearn.linear_model import LogisticRegression
          from sklearn.multiclass import OneVsRestClassifier
```

In this chapter:

- Build first-pass model based only on numeric data
- Multi-class logistic regression
- Format predictions and save to csv
- Compute log-loss score

### It's time to build a model

**multilabel_train_test_split**

```
In [62]: def multilabel_sample(y, size=1000, min_count=5, seed=None):
             """ Takes a matrix of binary labels `y` and returns
                 the indices for a sample of size `size` if
                 `size` > 1 or `size` * len(y) if size =< 1.
                 The sample is guaranteed to have > `min_count` of
                 each label.
             """
             try:
                 if (np.unique(y).astype(int) != np.array([0, 1])).all():
                     raise ValueError()
             except (TypeError, ValueError):
                 raise ValueError('multilabel_sample only works with binary indicator matrices')

             if (y.sum(axis=0) < min_count).any():
                 raise ValueError('Some classes do not have enough examples. Change min_count if necessary.')

             if size <= 1:
                 size = np.floor(y.shape[0] * size)

             if y.shape[1] * min_count > size:
                 msg = "Size less than number of columns * min_count, returning {} items instead of {}."
                 warn(msg.format(y.shape[1] * min_count, size))
                 size = y.shape[1] * min_count

             rng = np.random.RandomState(seed if seed is not None else np.random.randint(1))

             if isinstance(y, pd.DataFrame):
                 choices = y.index
                 y = y.values
             else:
                 choices = np.arange(y.shape[0])

             sample_idxs = np.array([], dtype=choices.dtype)

             # first, guarantee > min_count of each label
             for j in range(y.shape[1]):
                 label_choices = choices[y[:, j] == 1]
                 label_idxs_sampled = rng.choice(label_choices, size=min_count, replace=False)
                 sample_idxs = np.concatenate([label_idxs_sampled, sample_idxs])

             sample_idxs = np.unique(sample_idxs)

             # now that we have at least min_count of each, we can just random sample
             sample_count = int(size - sample_idxs.shape[0])

             # get sample_count indices from remaining choices
             remaining_choices = np.setdiff1d(choices, sample_idxs)
             remaining_sampled = rng.choice(remaining_choices,
                                            size=sample_count,
                                            replace=False)
             return np.concatenate([sample_idxs, remaining_sampled])

         def multilabel_sample_dataframe(df, labels, size, min_count=5, seed=None):
             """ Takes a dataframe `df` and returns a sample of size `size` where all
                 classes in the binary matrix `labels` are represented at
                 least `min_count` times.
             """
             idxs = multilabel_sample(labels, size=size, min_count=min_count, seed=seed)
             return df.loc[idxs]


         def multilabel_train_test_split(X, Y, size, min_count=5, seed=None):
             """ Takes a features matrix `X` and a label matrix `Y` and
                 returns (X_train, X_test, Y_train, Y_test) where all
                 classes in Y are represented at least `min_count` times.
             """
             index = Y.index if isinstance(Y, pd.DataFrame) else np.arange(Y.shape[0])

             test_set_idxs = multilabel_sample(Y, size=size, min_count=min_count, seed=seed)
             train_set_idxs = np.setdiff1d(index, test_set_idxs)

             test_set_mask = index.isin(test_set_idxs)
             train_set_mask = ~test_set_mask

             return (X[train_set_mask], X[test_set_mask], Y[train_set_mask], Y[test_set_mask])
```

```
In [63]: # Splitting the multi-class dataset

         NUMERIC_COLUMNS = ['FTE','Total']
         LABELS = ['Function','Use','Sharing','Reporting','Student_Type',
                   'Position_Type','Object_Type','Pre_K','Operating_Status']
         data_to_train = df[NUMERIC_COLUMNS].fillna(-1000)
         labels_to_use = pd.get_dummies(df[LABELS])
```

```
In [64]:  X_train, X_test, y_train, y_test = multilabel_train_test_split(
              data_to_train, labels_to_use,size=0.2, seed=123)
```
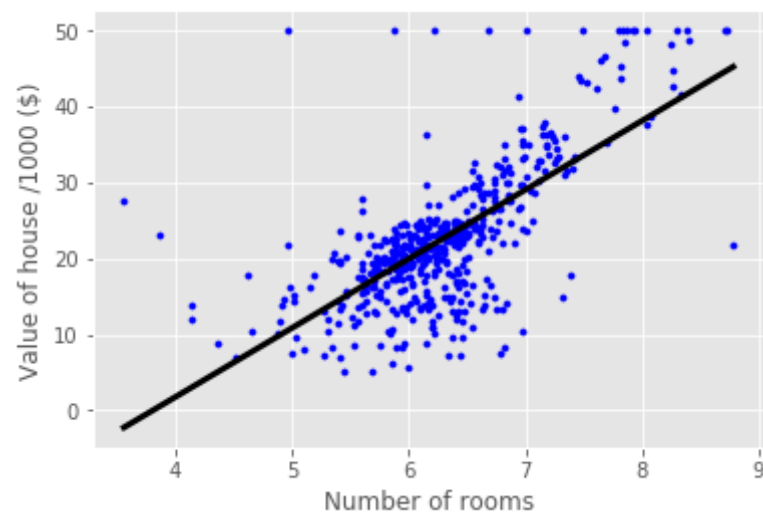
```
In [66]:  # Training the model

          clf = OneVsRestClassifier(LogisticRegression())
          clf.fit(X_train, y_train)
```

```
Out[66]:  OneVsRestClassifier(estimator=LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                    penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
                    verbose=0, warm_start=False),
                    n_jobs=1)
```

```
In [67]:  # EXERCISES
```

```
In [272]:  # Setting up a train-test split in scikit-learn
```



```
In [273]:  # EXERCISES
```

```
In [276]:  # Importing data for supervised learning
           # Gapminder Countries GDP data

           # Read the CSV file into a DataFrame: df
           df = pd.read_csv('datasets/gm_2008_region.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))
```
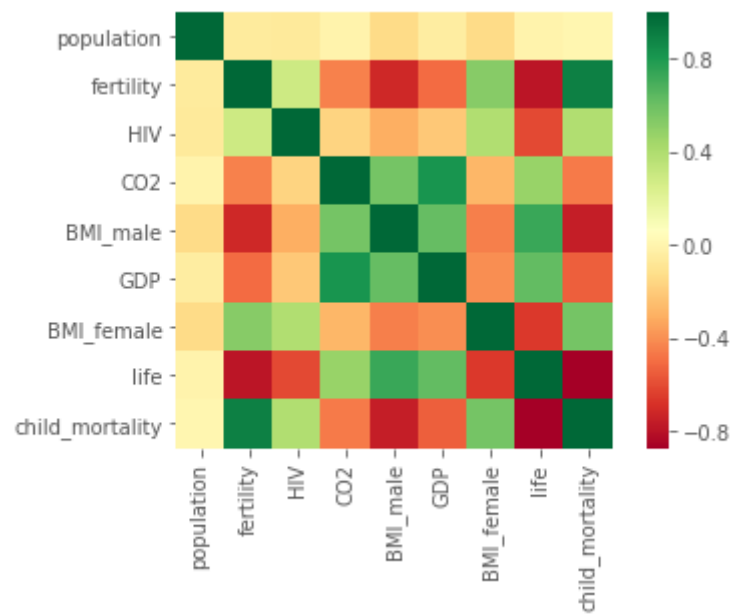
```
           Dimensions of y before reshaping: (139,)
           Dimensions of X before reshaping: (139,)
           Dimensions of y after reshaping: (139, 1)
           Dimensions of X after reshaping: (139, 1)
```

```
In [280]:  # Exploring the Gapminder data

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



## Making predictions

```
In [282]:  # Linear regression on all features

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

           reg_all = linear_model.LinearRegression()
           reg_all.fit(X_train, y_train)

           y_pred = reg_all.predict(X_test)
           reg_all.score(X_test, y_test)
```

Out[282]:  0.72989873609074984

```
In [285]:  # EXERCISE
```

```
In [299]:  # Fit & predict for regression

           # Read the CSV file into a DataFrame: df
           df = pd.read_csv('datasets/gm_2008_region.csv')

           # Create arrays for features and target variable
           y = df['life'].values.reshape(-1,1)
           X_fertility = df['fertility'].values.reshape(-1,1)
```

```
In [304]:  plt.scatter(X_fertility,y,c='blue',s=10);
```

```
In [305]:  # 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()
```
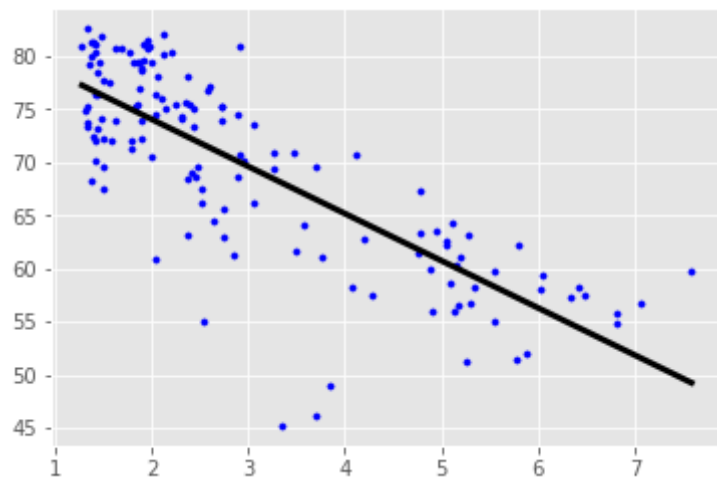
0.619244216774



```
In [306]:  # 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))
```

R^2: 0.7298987360907498
Root Mean Squared Error: 4.194027914110239

## A very brief introduction to NLP

```
In [322]:  from sklearn.model_selection import cross_val_score

           reg = linear_model.LinearRegression()
           cv_results = cross_val_score(reg, X, y, cv=5)
           print(cv_results)
           np.mean(cv_results)
```

[ 0.71001079  0.75007717  0.55271526  0.547501    0.52410561]

Out[322]:  0.61688196444251187

```
In [325]:  # EXERCISE
```

```python
In [326]: # 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)))
```

```
[ 0.71001079  0.75007717  0.55271526  0.547501    0.52410561]
Average 5-Fold CV Score: 0.6168819644425119
```

## Representing text numerically

**Ridge regression** takes the sum of the squared values of the coefficients multiplied by some alpha, this is also known as the L2 regularization.

```python
In [328]: # Ridge regression in scikit-learn

          from sklearn.linear_model import Ridge
          X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=42)
          ridge = Ridge(alpha=0.1, normalize=True)
          ridge.fit(X_train, y_train)
          ridge_pred = ridge.predict(X_test)
          ridge.score(X_test, y_test)
```

Out[328]: 0.74001557383978234

**Lasso regression** performs regularization by adding to the loss function a penalty term of the absolute value of each coefficient multiplied by some alpha. This is also known as L1 regularization because the regularization term is the L1 norm of the coefficients.

```python
In [329]: # Lasso regression in scikit-learn

          from sklearn.linear_model import Lasso
          X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state=42)
          lasso = Lasso(alpha=0.1, normalize=True)
          lasso.fit(X_train, y_train)
          lasso_pred = lasso.predict(X_test)
          lasso.score(X_test, y_test)
```
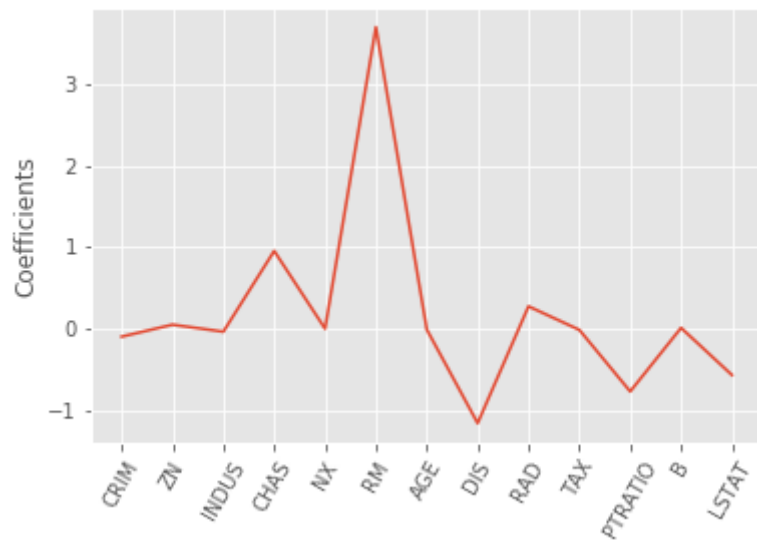
Out[329]: 0.73913024600881294

```python
In [334]: # Lasso for feature selection in scikit-learn

          X = boston.drop('MEDV', axis=1).values
          y = boston['MEDV'].values
```

In [339]:
```python
from sklearn.linear_model import Lasso

names = boston.drop('MEDV', axis=1).columns
lasso = Lasso(alpha=0.1)
lasso_coef = lasso.fit(X, y).coef_

_ = plt.plot(range(len(names)), lasso_coef)
_ = plt.xticks(range(len(names)), names, rotation=60)
_ = plt.ylabel('Coefficients')
plt.show()
```



In [343]:
```python
# EXERCISES
```

In [369]:
```python
# Regularization I: Lasso

df_columns = df.drop(['life','Region'],axis=1).columns
X = df.drop(['life','Region'],axis=1).values
y = df['life'].values
df_columns
```

Out[369]:
```
Index(['population', 'fertility', 'HIV', 'CO2', 'BMI_male', 'GDP',
       'BMI_female', 'child_mortality'],
      dtype='object')
```

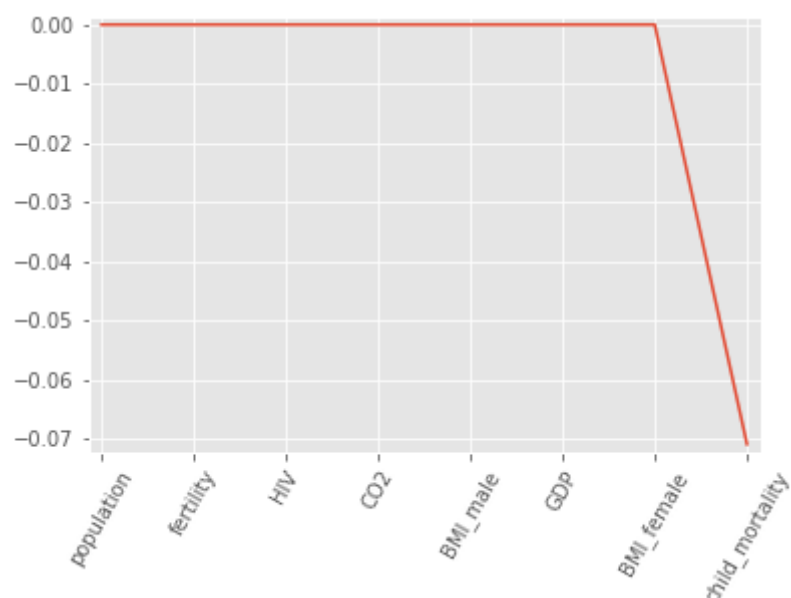In [371]:
```python
from sklearn.linear_model import Lasso

# 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()
```

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



In [385]:
```python
# Regularization II: Ridge
# fitting ridge regression models over a range of different alphas, and plot cross-validated R^2 scores for each.
```

```
In [386]: # function to visualize the scores and standard deviations
          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()
```

```
In [387]: # 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)
```
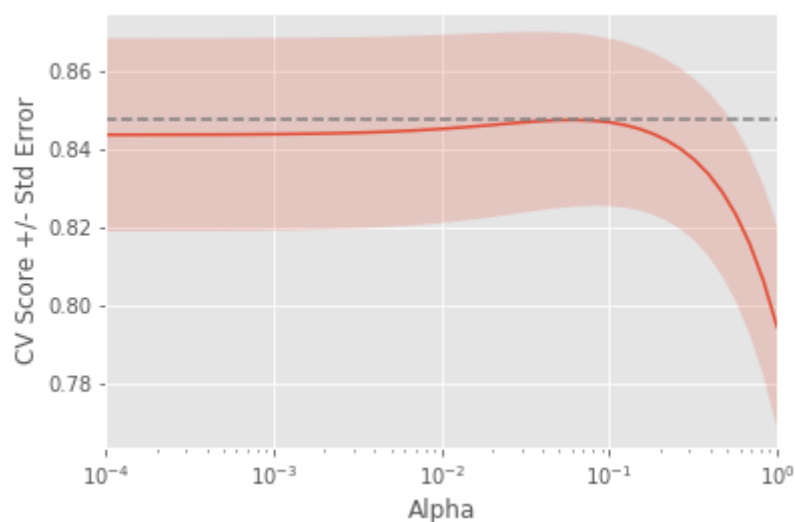


Notice how the cross-validation scores change with different alphas. Which alpha should you pick? How can you fine-tune your model?

# Chap 3: Improving your model

```
In [391]: # Import plotting modules
          import matplotlib.pyplot as plt
          import seaborn as sns
          import pandas as pd
          import numpy as np
          from sklearn import datasets
          plt.style.use('ggplot')
```

### Pipelines, feature & text preprocessing

```
In [398]: # Confusion matrix in scikit-learn

          df = pd.read_csv('datasets/house-votes-84.csv',header=None)
          df.columns = ['party','infants','water','budget','physician',
                        'salvador','religious','satellite','aid','missile',
                        'immigration','synfuels','education','superfund',
                        'crime','duty_free_exports','eaa_rsa']
          df.replace({'n':0,'y':1,'?':0},inplace=True)
```

```
In [399]: y = df['party'].values
          X = df.drop('party', axis=1).values
```

```
In [400]: from sklearn.metrics import classification_report
          from sklearn.metrics import confusion_matrix

          knn = KNeighborsClassifier(n_neighbors=8)
          X_train, X_test, y_train, y_test = train_test_split(X, y,test_size=0.4, random_state=42)

          knn.fit(X_train, y_train)
          y_pred = knn.predict(X_test)
```

```
In [401]: print(confusion_matrix(y_test, y_pred))

          print(classification_report(y_test, y_pred))
```

```
[[108   7]
 [  6  53]]
             precision    recall  f1-score   support

   democrat       0.95      0.94      0.94       115
 republican       0.88      0.90      0.89        59

avg / total       0.93      0.93      0.93       174
```

```
In [411]: # EXERCISES
```

```
In [412]: # Metrics for classification
          # computing a confusion matrix and generating a classification report

          df = pd.read_csv('datasets/diabetes.csv')
          y = df['diabetes'].values
          X = df.drop('diabetes',axis=1).values
```

```
In [413]: # Import necessary modules
          from sklearn.model_selection import train_test_split
          from sklearn.neighbors import KNeighborsClassifier
          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))
          print(classification_report(y_test, y_pred))
```

```
[[176  30]
 [ 56  46]]
             precision    recall  f1-score   support

          0       0.76      0.85      0.80       206
          1       0.61      0.45      0.52       102

avg / total       0.71      0.72      0.71       308
```

## Text features and feature unions

```
In [429]:  # Logistic regression in scikit-learn

           df = pd.read_csv('datasets/house-votes-84.csv',header=None)
           df.columns = ['party','infants','water','budget','physician',
                         'salvador','religious','satellite','aid','missile',
                         'immigration','synfuels','education','superfund',
                         'crime','duty_free_exports','eaa_rsa']
           df.replace({'n':0,'y':1,'?':0},inplace=True)

           y = df['party'].map({'republican':0,'democrat':1}).values
           X = df.drop('party', axis=1).values
```

```
In [430]:  from sklearn.linear_model import LogisticRegression
           from sklearn.model_selection import train_test_split

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

           logreg.fit(X_train, y_train)
           y_pred = logreg.predict(X_test)
```

```
In [431]:  # Plotting the ROC curve

           from sklearn.metrics import roc_curve

           y_pred_prob = logreg.predict_proba(X_test)[:,1]
           fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob)

           plt.plot([0, 1], [0, 1], 'k--')
           plt.plot(fpr, tpr, label='Logistic Regression')
           plt.xlabel('False Positive Rate')
           plt.ylabel('True Positive Rate')
           plt.title('Logistic Regression ROC Curve')
           plt.show();
```



```
In [443]:  # EXERCISES
```

```
In [444]:  # Building a logistic regression model

           df = pd.read_csv('datasets/diabetes.csv')
           y = df['diabetes'].values
           X = df.drop('diabetes',axis=1).values
```

In [445]:
```python
# 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))
```

```
[[174  32]
 [ 36  66]]
            precision    recall  f1-score   support

         0       0.83      0.84      0.84       206
         1       0.67      0.65      0.66       102

avg / total       0.78      0.78      0.78       308
```
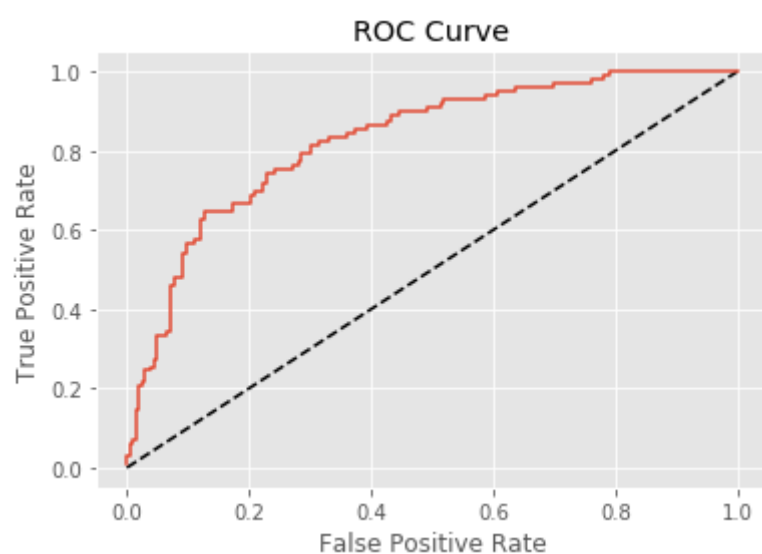
In [447]:
```python
# Plotting an ROC curve

# 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()
```
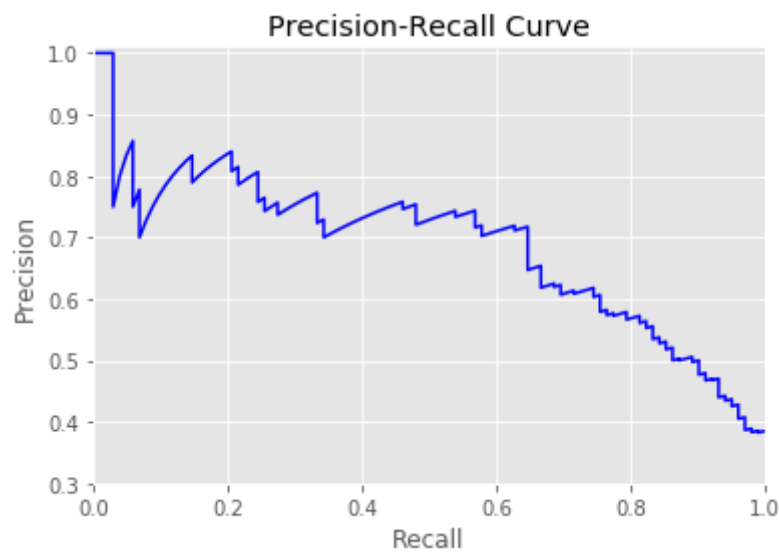
```
In [462]:  # Precision-recall Curve

           from sklearn.metrics import recall_score
           from sklearn.metrics import precision_score
           from sklearn.metrics import precision_recall_curve

           pr, rc, thld = precision_recall_curve(y_test, y_pred_prob)

           # Plot precision-recall curve
           plt.plot(rc, pr,color='blue')
           plt.xlabel('Recall')
           plt.ylabel('Precision')
           plt.xlim([0.0,1.0])
           plt.ylim([0.3,1.01])
           plt.title('Precision-Recall Curve')
           plt.show()
```



## Choosing a classification model

```
In [464]:  # AUC in scikit-learn

           from sklearn.metrics import roc_auc_score

           logreg = LogisticRegression()
           X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=42)
           logreg.fit(X_train, y_train)

           y_pred_prob = logreg.predict_proba(X_test)[:,1]
           roc_auc_score(y_test, y_pred_prob)
```

Out[464]:  0.82686084142394833

```
In [468]:  # AUC using cross-validation

           from sklearn.model_selection import cross_val_score

           cv_scores = cross_val_score(logreg, X, y, cv=5, scoring='roc_auc')
           print(cv_scores)
```

```
[ 0.7987037   0.80759259  0.81944444  0.86622642  0.85056604]
```

```
In [469]:  # EXERCISES
```

```
In [470]:  # AUC computation

           # 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)))

           # 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: 0.826608414239483
AUC scores computed using 5-fold cross-validation: [ 0.7987037   0.80759259  0.81944444  0.86622642  0.85056604]
```

```
In [471]:  # GridSearchCV in scikit-learn

           from sklearn.model_selection import GridSearchCV

           param_grid = {'n_neighbors': np.arange(1, 50)}
           knn = KNeighborsClassifier()
           knn_cv = GridSearchCV(knn, param_grid, cv=5)

           knn_cv.fit(X, y)
           print(knn_cv.best_params_)
           print(knn_cv.best_score_)
```

```
{'n_neighbors': 14}
0.7578125
```

```
In [472]:  # EXERCISES
```

```
In [477]:  # Hyperparameter tuning with GridSearchCV

           df = pd.read_csv('datasets/diabetes.csv')
           y = df['diabetes']
           X = df.drop('diabetes',axis=1)
```

```
In [480]:  # Import necessary modules
           from sklearn.model_selection import GridSearchCV
           from sklearn.linear_model import LogisticRegression

           # 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': 163789.3706954068}
Best score is 0.7721354166666666
```

```
In [483]:  # Hyperparameter tuning with RandomizedSearchCV

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

           # 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': 'entropy', 'max_depth': 3, 'max_features': 8, 'min_samples_leaf': 6}
Best score is 0.7395833333333334
```

```
In [ ]:  # EXERCISES
```

```
In [486]:  # Hold-out set in practice I: Classification

           # 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_))
```
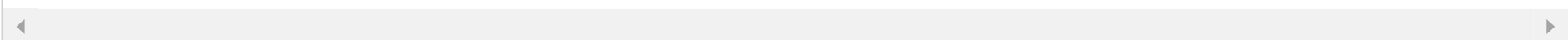
```
Tuned Logistic Regression Parameter: {'C': 31.622776601683793, 'penalty': 'l2'}
Tuned Logistic Regression Accuracy: 0.7673913043478261
```

```
In [497]:  # Hold-out set in practice II: Regression

           # Gapminder Countries GDP data
           # Read the CSV file into a DataFrame: df
           df = pd.read_csv('datasets/gm_2008_region.csv')

           # Create arrays for features and target variable
           y = df['life'].values
           X = df.drop(['life','Region'],axis=1).values
```

```
In [502]:  import warnings; warnings.filterwarnings('ignore')

           # 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))
```
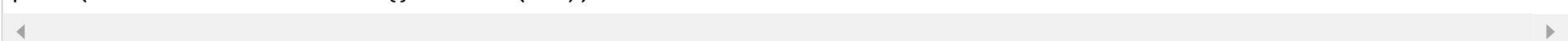
```
Tuned ElasticNet l1 ratio: {'l1_ratio': 0.20689655172413793}
Tuned ElasticNet R squared: 0.8668305372460283
Tuned ElasticNet MSE: 10.057914133398445
```

## Chap 4: Learning from the experts

```
In [84]:  # Import plotting modules
          import matplotlib.pyplot as plt
          import seaborn as sns
          import pandas as pd
          import numpy as np
          from sklearn import datasets
          plt.style.use('ggplot')
```
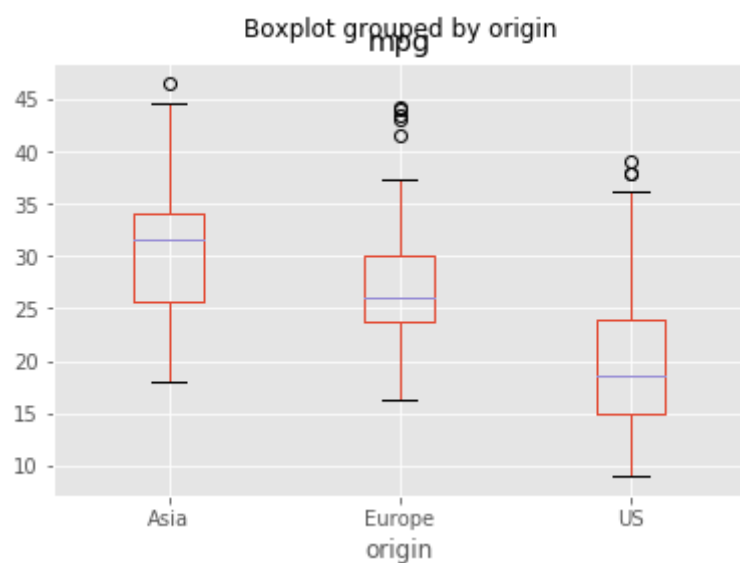
### Learning from the expert: processing

```
In [507]:  df = pd.read_csv('datasets/auto.csv')
           df.head()
```

Out[507]:

|   | mpg | displ | hp | weight | accel | origin | size |
|---|-----|-------|-----|--------|-------|--------|------|
| 0 | 18.0 | 250.0 | 88 | 3139 | 14.5 | US | 15.0 |
| 1 | 9.0 | 304.0 | 193 | 4732 | 18.5 | US | 20.0 |
| 2 | 36.1 | 91.0 | 60 | 1800 | 16.4 | Asia | 10.0 |
| 3 | 18.5 | 250.0 | 98 | 3525 | 19.0 | US | 15.0 |
| 4 | 34.3 | 97.0 | 78 | 2188 | 15.8 | Europe | 10.0 |

```
In [537]:  # EDA w/ categorical feature
           df.boxplot(column='mpg',by='origin');
           plt.show()
```



```
In [544]:  # Encoding dummy variables

           df_origin = pd.get_dummies(df,drop_first=True)
           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
```

```
In [550]:  X=df_origin.drop(['origin_Europe','origin_US'],axis=1)
           y=df_origin[['origin_Europe','origin_US']]
```

```
In [555]:  # Linear regression with dummy variables

           from sklearn.model_selection import train_test_split
           from sklearn.linear_model import Ridge

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

           ridge = Ridge(alpha=0.5, normalize=True).fit(X_train,y_train)
           ridge.score(X_test, y_test)
```

Out[555]:  0.28837060656093705
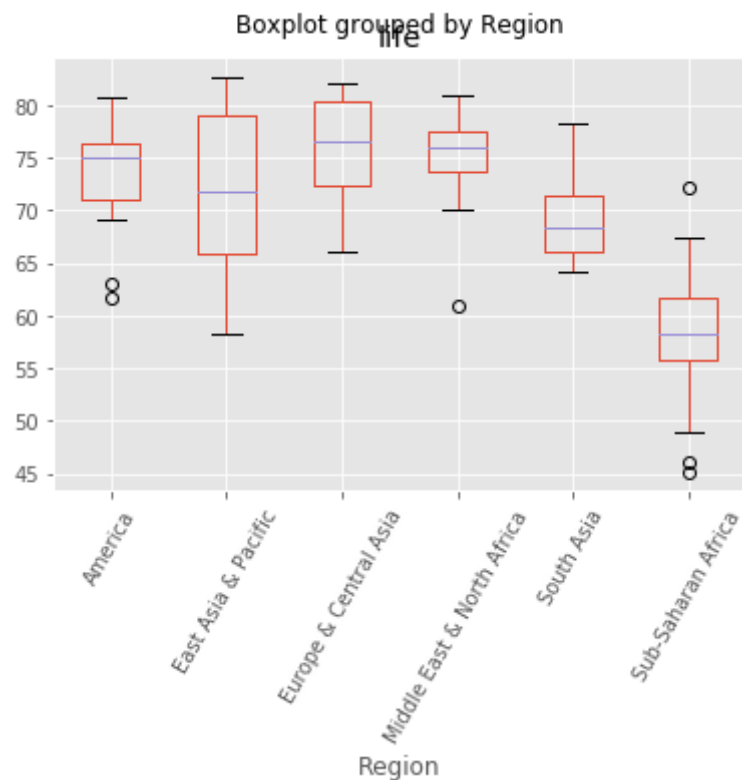
```
In [552]:  # EXERCISES
```

```
In [560]:  # Exploring categorical features

           # Gapminder Countries GDP data
           # Read the CSV file into a DataFrame: df
           df = pd.read_csv('datasets/gm_2008_region.csv')

           # Create a boxplot of life expectancy per region
           df.boxplot('life', 'Region', rot=60)

           # Show the plot
           plt.show()
```


Boxplot grouped by Region

```
In [564]:  # Creating dummy variables

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

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

           # 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_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')
```

```
In [565]:  # Create arrays for features and target variable
           y = df_region['life']
           X = df_region.drop(['life'],axis=1)
           X.shape, y.shape
```

```
Out[565]:  ((139, 13), (139,))
```

```
In [567]:  # 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)
```

```
           [ 0.86808336  0.80623545  0.84004203  0.7754344   0.87503712]
```

# Learning from the expert: a stats trick

```
In [570]:  # PIMA Indians dataset
           df = pd.read_csv('datasets/diabetes.csv')
           df.info()

           <class 'pandas.core.frame.DataFrame'>
           RangeIndex: 768 entries, 0 to 767
           Data columns (total 9 columns):
           pregnancies     768 non-null int64
           glucose         768 non-null int64
           diastolic       768 non-null int64
           triceps         768 non-null int64
           insulin         768 non-null int64
           bmi             768 non-null float64
           dpf             768 non-null float64
           age             768 non-null int64
           diabetes        768 non-null int64
           dtypes: float64(2), int64(7)
           memory usage: 54.1 KB
```

```
In [571]:  print(df.head())

              pregnancies  glucose  diastolic  triceps  insulin   bmi    dpf  age  \
           0            6      148         72       35        0  33.6  0.627   50
           1            1       85         66       29        0  26.6  0.351   31
           2            8      183         64        0        0  23.3  0.672   32
           3            1       89         66       23       94  28.1  0.167   21
           4            0      137         40       35      168  43.1  2.288   33

              diabetes
           0         1
           1         0
           2         1
           3         0
           4         1
```

```
In [573]:  df.insulin.replace(0, np.nan, inplace=True)
           df.triceps.replace(0, np.nan, inplace=True)
           df.bmi.replace(0, np.nan, inplace=True)
           df.info()

           <class 'pandas.core.frame.DataFrame'>
           RangeIndex: 768 entries, 0 to 767
           Data columns (total 9 columns):
           pregnancies     768 non-null int64
           glucose         768 non-null int64
           diastolic       768 non-null int64
           triceps         541 non-null float64
           insulin         394 non-null float64
           bmi             757 non-null float64
           dpf             768 non-null float64
           age             768 non-null int64
           diabetes        768 non-null int64
           dtypes: float64(4), int64(5)
           memory usage: 54.1 KB
```

```
In [574]:  # Dropping missing data
           df.dropna().shape
```

```
Out[574]:  (393, 9)
```

```
In [592]:  # Imputing missing data
           y = df['diabetes']
           X = df.drop('diabetes',axis=1)
```

```
In [593]:  from sklearn.preprocessing import Imputer
           imp = Imputer(missing_values='NaN', strategy='mean', axis=0)
           imp.fit(X)
           X = imp.transform(X)
```

```
In [597]:  # Imputing within a pipeline
           # Pipeline: All steps before last must be transformers (like impute)
           # Pipeline: Last step can be transformer or estimator (like classifier/regressor)

           from sklearn.pipeline import Pipeline
           from sklearn.preprocessing import Imputer

           imp = Imputer(missing_values='NaN', strategy='mean', axis=0)
           logreg = LogisticRegression()

           steps = [('imputation', imp),('logistic_regression', logreg)]
           pipeline = Pipeline(steps)

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

           pipeline.fit(X_train, y_train)
           y_pred = pipeline.predict(X_test)
           pipeline.score(X_test, y_test)
```

Out[597]:  0.76190476190476186

```
In [ ]:  # EXERCISES
```

```
In [598]:  # Dropping missing data

           df = pd.read_csv('datasets/house-votes-84.csv',header=None)
           df.columns = ['party','infants','water','budget','physician',
                         'salvador','religious','satellite','aid','missile',
                         'immigration','synfuels','education','superfund',
                         'crime','duty_free_exports','eaa_rsa']
           df.replace({'n':0,'y':1},inplace=True)
```

```
In [601]:  # 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))
```

```
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)
```

```
In [602]:   # Imputing missing data in a ML Pipeline I

            # 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)]
```

```
In [603]:   # Imputing missing data in a ML Pipeline II

            # 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

                  0       0.65      1.00      0.79       151
                  1       0.00      0.00      0.00        80

        avg / total       0.43      0.65      0.52       231
```

### Learning from the expert: a computational trick and the winning model

```
In [663]:   # Why scale your data?

            df = pd.read_csv('datasets/winequality-red.csv',header=0,sep=';')
            X = df.drop('quality',axis=1).values
            y = df['quality'].values
            df.head()
```

Out[663]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |
| 1 | 7.8 | 0.88 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.9968 | 3.20 | 0.68 | 9.8 | 5 |
| 2 | 7.8 | 0.76 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.9970 | 3.26 | 0.65 | 9.8 | 5 |
| 3 | 11.2 | 0.28 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.9980 | 3.16 | 0.58 | 9.8 | 6 |
| 4 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |

```
In [664]:   # Scaling in scikit-learn

            from sklearn.preprocessing import scale

            X_scaled = scale(X)
            print(np.mean(X), np.std(X))
            print(np.mean(X_scaled), np.std(X_scaled))
```

```
8.13421922452 16.7265339794
2.54662653149e-15 1.0
```

```
In [665]:  # Scaling in a pipeline

           from sklearn.preprocessing import StandardScaler
           from sklearn.metrics import accuracy_score

           steps = [('scaler', StandardScaler()),('knn', KNeighborsClassifier())]
           pipeline = Pipeline(steps)

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

           knn_scaled = pipeline.fit(X_train, y_train)
           y_pred = pipeline.predict(X_test)
           print(accuracy_score(y_test, y_pred))

           knn_unscaled = KNeighborsClassifier().fit(X_train, y_train)
           print(knn_unscaled.score(X_test, y_test))
```

```
0.615625
0.49375
```

```
In [673]:  # CV and scaling in a pipeline

           steps = [('scaler', StandardScaler()),(('knn', KNeighborsClassifier()))]
           pipeline = Pipeline(steps)
           parameters = {'knn__n_neighbors':np.arange(1, 50)}

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

           cv = GridSearchCV(pipeline, param_grid=parameters)
           cv.fit(X_train, y_train)
           y_pred = cv.predict(X_test)
```

```
In [674]:  print(cv.best_params_)
           print(cv.score(X_test, y_test))
           print(classification_report(y_test, y_pred))
```

```
{'knn__n_neighbors': 1}
0.634375
             precision    recall  f1-score   support

          3       0.00      0.00      0.00         1
          4       0.18      0.12      0.15        16
          5       0.66      0.72      0.69       127
          6       0.68      0.60      0.64       131
          7       0.63      0.69      0.66        42
          8       0.25      0.33      0.29         3

avg / total       0.63      0.63      0.63       320
```

```
In [675]:  # EXERCISES
```

```
In [684]:  # Centering and scaling your data

           # White wine quality dataset.
           df = pd.read_csv('datasets/white-wine.csv')
           y = df.quality < 5
           X = df.drop('quality',axis=1).values
```

```
In [685]:  # 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)))
```

```
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
```

In [687]:
```python
# Centering and scaling in a pipeline

# 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: 0.964625850340136
Accuracy without Scaling: 0.9666666666666667
```

In [689]:
```python
# Bringing it all together I: Pipeline for classification

from sklearn.pipeline import Pipeline
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.grid_search import GridSearchCV
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
```

In [690]:
```python
# 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_))
```

```
Accuracy: 0.9693877551020408
             precision    recall  f1-score   support

      False       0.97      1.00      0.98       951
       True       0.43      0.10      0.17        29

avg / total       0.96      0.97      0.96       980

Tuned Model Parameters: {'SVM__C': 100, 'SVM__gamma': 0.01}
```

In [702]: 
```python
# Bringing it all together II: Pipeline for regression

df = pd.read_csv('datasets/gm_2008_region.csv')
y = df['life'].values
X = df.drop(['life','Region'],axis=1).values
```

In [705]: 
```python
# 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))
```

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

### Next steps and the social impact of your work

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]: