# Using Python for Research Homework: Week 5, Case Study Part 2

The movie dataset on which this case study is based is a database of 5000 movies catalogued by The Movie Database (TMDb). The information available about each movie is its budget, revenue, rating, actors and actresses, etc. In this case study, we will use this dataset to determine whether any information about a movie can predict the total revenue of a movie. We will also attempt to predict whether a movie's revenue will exceed its budget.

In Part 2, we will use the dataset prepared in Part 1 for an applied analysis.

```
In [26]:  # DO NOT EDIT THIS CODE
          import pandas as pd
          import numpy as np

          from sklearn.model_selection import cross_val_score
          from sklearn.linear_model import LinearRegression
          from sklearn.linear_model import LogisticRegression
          from sklearn.ensemble import RandomForestRegressor
          from sklearn.ensemble import RandomForestClassifier

          from sklearn.metrics import accuracy_score
          from sklearn.metrics import r2_score

          import matplotlib.pyplot as plt

          import warnings
          warnings.filterwarnings("ignore")

          # EDIT THIS CODE TO LOAD THE SAVED DF FROM THE LAST HOMEWORK
          df = pd.read_csv('movies_clean.csv')
```

## Exercise 1

In Part 2 of this case study, we will primarily use the two models we recently discussed: linear/logistic regression and random forests to perform prediction and classification. We will use these methods to predict revenue, and we will use logistic regression to classify whether a movie was profitable.

In this exercise, we will instantiate regression and classification models. Code is provided that prepares the covariates and outcomes we will use for data analysis.

### Instructions

- Instantiate `LinearRegression()`, `LogisticRegression()`, `RandomForestRegressor()`, and `RandomForestClassifier()` objects, and assign them to `linear_regression`, `logistic_regression`, `forest_regression`, and `forest_classifier`, respectively.
- For the random forests models, specify `max_depth=4` and `random_state=0`.

```
In [38]:   # Define all covariates and outcomes from `df`.
           regression_target = 'revenue'
           classification_target = 'profitable'
           all_covariates = ['budget','popularity','runtime','vote_count','vote_average','Acti
                             'Science Fiction','Crime','Drama','Thriller','Animation','Family'
                             'Horror','Mystery','War','History','Music','Documentary','TV Movie'

           regression_outcome = df[regression_target]
           classification_outcome = df[classification_target]
           covariates = df[all_covariates]

           # Instantiate all regression models and classifiers.
           linear_regression = LinearRegression(fit_intercept=True)
           logistic_regression = LogisticRegression()
           forest_regression = RandomForestRegressor(max_depth=4, random_state=0)
           forest_classifier = RandomForestClassifier(max_depth=4, random_state=0)
```

## Exercise 2

In this exercise, we will create two functions that compute a model's score. For regression models, we will use correlation as the score. For classification models, we will use accuracy as the score.

### Instructions

- Define a function called `correlation` with arguments `estimator`, `X`, and `y`. The function should compute the correlation between the observed outcome `y` and the outcome predicted by the model.
    - To obtain predictions, the function should first use the `fit` method of `estimator` and then use the `predict` method from the fitted object.
    - The function should return the first argument from `r2_score` comparing `predictions` and `y`.
- Define a function called `accuracy` with the same arguments and code, substituting `accuracy_score` for `r2_score`.

```
In [39]:   def correlation(estimator, X, y):
               predictions = estimator.fit(X, y).predict(X)
               return r2.score(predictions, y)

           def accuracy(estimator, X, y):
               predictions = estimator.fit(X, y).predict(X)
               return accuracy.score(predictions, y)
```

## Exercise 3

In this exercise, we will compute the cross-validated performance for the linear and random forest regression models.

### Instructions

- Call `cross_val_score` using `linear_regression` and `forest regression` as models. Store the output as `linear_regression_scores` and `forest_regression_scores`, respectively.
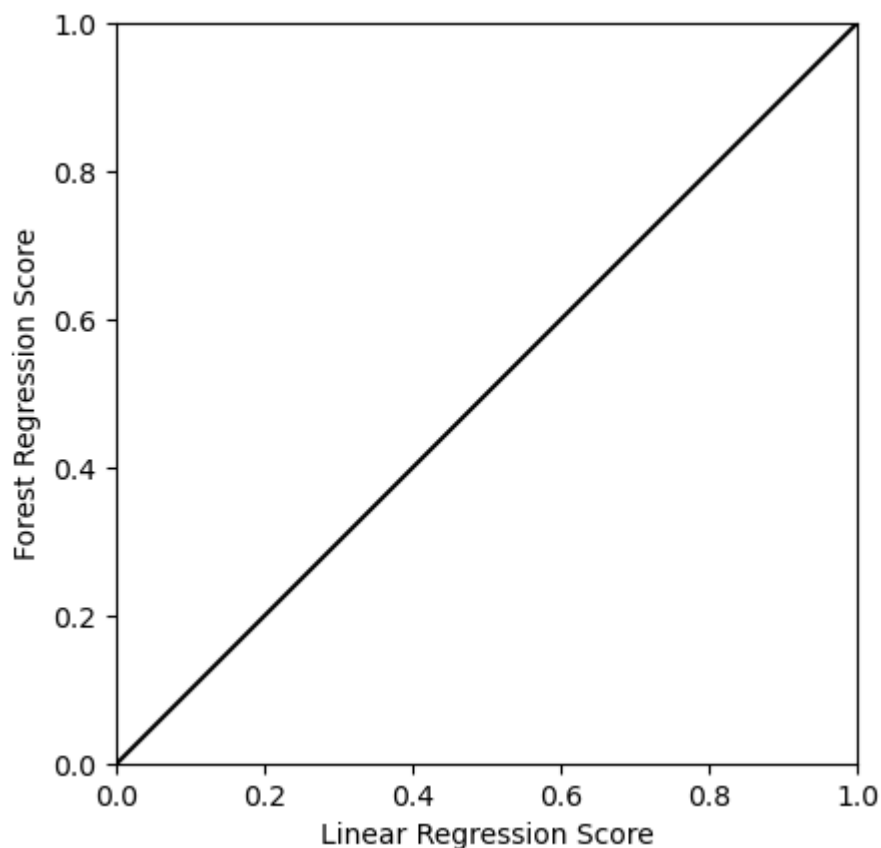
  - Set the parameters `cv=10` to use 10-fold cross-validation and `scoring=correlation` to use our `correlation` function defined in the previous exercise.
- Plotting code has been provided to compare the performance of the two models. Use `plt.show()` to plot the correlation between actual and predicted revenue for each cross-validation fold using the linear and random forest regression models.
- Which of the two models exhibits a better fit?

In [40]:
```python
# Determine the cross-validated correlation for linear and random forest models.
linear_regression_scores = cross_val_score(linear_regression, covariates, regressio
forest_regression_scores = cross_val_score(forest_regression, covariates, regressio

# Plot Results
plt.axes().set_aspect('equal', 'box')
plt.scatter(linear_regression_scores, forest_regression_scores)
plt.plot((0, 1), (0, 1), 'k-')

plt.xlim(0, 1)
plt.ylim(0, 1)
plt.xlabel("Linear Regression Score")
plt.ylabel("Forest Regression Score")

plt.show()
```



## Exercise 4

In this exercise, we will compute cross-validated performance for the linear and random forest classification models.

### Instructions

- Call `cross_val_score` using `logistic_regression` and `forest_classifier` as models. Store the output as `logistic_regression_scores` and `forest_classification_scores`, respectively.
  - Set the parameters `cv=10` to use 10-fold cross-validation and `scoring=accuracy` to use our accuracy function defined in the previous exercise.
- Plotting code has been provided to compare the performance of the two models. Use `plt.show()` to plot the accuracy of predicted profitability for each cross-validation fold using the logistic and random forest classification models.
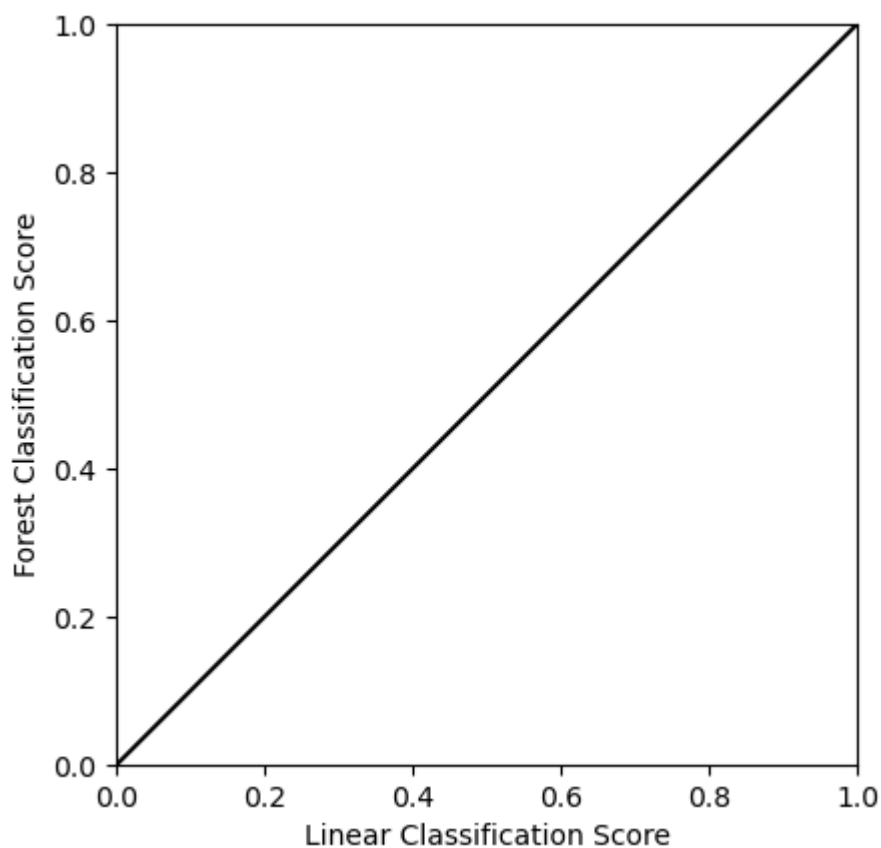- Which of the two models exhibits a better fit?

```
In [41]: logistic_regression_scores = cross_val_score(logistic_regression, covariates, class
         forest_classification_scores = cross_val_score(forest_classifier, covariates, class

         # Determine the cross-validated accuracy for logistic and random forest models.

         # Plot Results
         plt.axes().set_aspect('equal', 'box')
         plt.scatter(logistic_regression_scores, forest_classification_scores)
         plt.plot((0, 1), (0, 1), 'k-')

         plt.xlim(0, 1)
         plt.ylim(0, 1)
         plt.xlabel("Linear Classification Score")
         plt.ylabel("Forest Classification Score")

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



## Exercise 5

In Exercise 3, we saw that predicting revenue was only moderately successful. It might be the case that predicting movies that generated precisely no revenue is difficult. In the next three exercises, we will exclude these movies, and rerun the analyses to determine if the fits improve. In this exercise, we will rerun the regression analysis for this subsetted dataset.

## Instructions

- Define `positive_revenue_df` as the subset of movies in `df` with `revenue` greater than zero.
- Code is provided below that creates new instances of model objects. Replace all instances of `df` with `positive_revenue_df`, and run the given code.

```
In [47]:   positive_revenue_df = df[df["revenue"] > 0]

           regression_outcome = positive_revenue_df[regression_target]
           classification_outcome = positive_revenue_df[classification_target]
           covariates = positive_revenue_df[all_covariates]

           linear_regression = LinearRegression()
           logistic_regression = LogisticRegression()
           forest_regression = RandomForestRegressor(max_depth=4, random_state=0)
           forest_classifier = RandomForestClassifier(max_depth=4, random_state=0)
           linear_regression_scores = cross_val_score(linear_regression, covariates, regressic
           forest_regression_scores = cross_val_score(forest_regression, covariates, regressic
           logistic_regression_scores = cross_val_score(logistic_regression, covariates, class
           forest_classification_scores = cross_val_score(forest_classifier, covariates, class

           np.nanmean(forest_regression_scores)
```

Out[47]:   nan

# Exercise 6

In this exercise, we will compute the cross-validated performance for the linear and random forest regression models for positive revenue movies only.

## Instructions

- Call `cross_val_score` using `linear_regression` and `forest regression` as models. Store the output as `linear_regression_scores` and `forest_regression_scores`, respectively.
  - Set the parameters `cv=10` to use 10-fold cross-validation and `scoring=correlation` to use our `correlation` function defined in the previous exercise.
- Plotting code has been provided to compare the performance of the two models. Use `plt.show()` to plot the correlation between actual and predicted revenue for each cross-validation fold using the linear and random forest regression models.
- Which of the two models exhibits a better fit? Is this result different from what we observed when considering all movies?
- Code is provided for you that prints the importance of each covariate in predicting revenue using the random forests classifier.
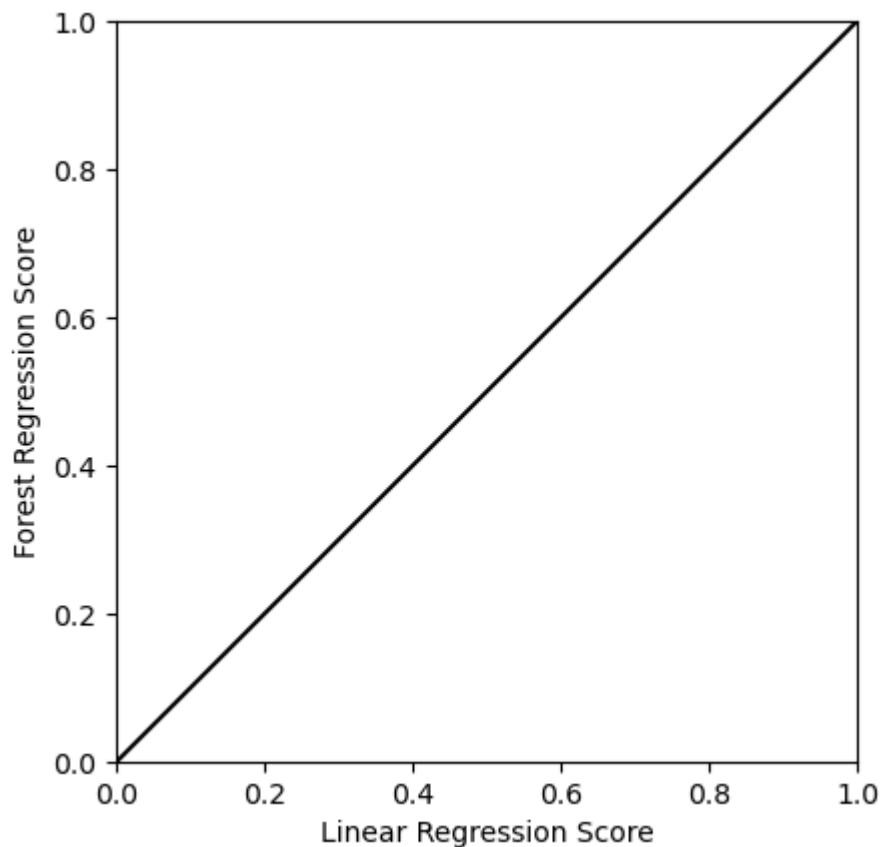  - Which variables are most important?

In [43]:
```python
# Determine the cross-validated correlation for linear and random forest models.
logistic_regression_scores = cross_val_score(logistic_regression, covariates, class
forest_classification_scores = cross_val_score(forest_classifier, covariates, class

# Plot Results
plt.axes().set_aspect('equal', 'box')
plt.scatter(linear_regression_scores, forest_regression_scores)
plt.plot((0, 1), (0, 1), 'k-')

plt.xlim(0, 1)
plt.ylim(0, 1)
plt.xlabel("Linear Regression Score")
plt.ylabel("Forest Regression Score")

# Show the plot.
plt.show()

# Print the importance of each covariate in the random forest regression.
forest_regression.fit(positive_revenue_df[all_covariates], positive_revenue_df[regr
sorted(list(zip(all_covariates, forest_regression.feature_importances_)), key=lambd
```

Out[43]:
```
[('Fantasy', 0.0),
 ('Western', 0.0),
 ('Mystery', 0.0),
 ('Music', 0.0),
 ('TV Movie', 0.0),
 ('Documentary', 3.948347422319195e-06),
 ('Foreign', 1.213773760536117e-05),
 ('Romance', 6.302030482832048e-05),
 ('Family', 0.00017893326581542565),
 ('Animation', 0.00033957110631423735),
 ('War', 0.00039666244687010813),
 ('Horror', 0.0005407257444202018),
 ('History', 0.0007927009679097137),
 ('Crime', 0.0010242988124237993),
 ('Thriller', 0.0011393486577492154),
 ('Action', 0.001266502656791017),
 ('Drama', 0.0014597518446305246),
 ('Comedy', 0.0023190793717307763),
 ('Adventure', 0.0023284914850058283),
 ('Science Fiction', 0.00233390552560341),
 ('vote_average', 0.014225606673103894),
 ('runtime', 0.01619249422581429),
 ('popularity', 0.07872384599889447),
 ('budget', 0.3136921377219631),
 ('vote_count', 0.5629668371051042)]
```

## Exercise 7

In this exercise, we will compute cross-validated performance for the linear and random forest classification models for positive revenue movies only.

### Instructions

- Call `cross_val_score` using `logistic_regression` and `forest classifer` as models. Store the output as `logistic_regression_scores` and `forest_classification_scores`, respectively.
  - Set the parameters `cv=10` to use 10-fold cross-validation and `scoring=accuracy` to use our `accuracy` function defined in the previous exercise.
- Plotting code has been provided to compare the performance of the two models. Use `plt.show()` to plot the correlation between actual and predicted revenue for each cross-validation fold using the linear and random forest regression models.
- Which of the two models exhibits a better fit? Is this result different from what we observed when considering all movies?
- Code is provided for you that prints the importance of each covariate in predicting profitabilitiy using the random forests classifier.
  - Which variables are most important?

In [45]:
```
# Determine the cross-validated accuracy for logistic and random forest models.


# Plot Results
plt.axes().set_aspect('equal', 'box')
plt.scatter(logistic_regression_scores, forest_classification_scores)
plt.plot((0, 1), (0, 1), 'k-')

plt.xlim(0, 1)
```
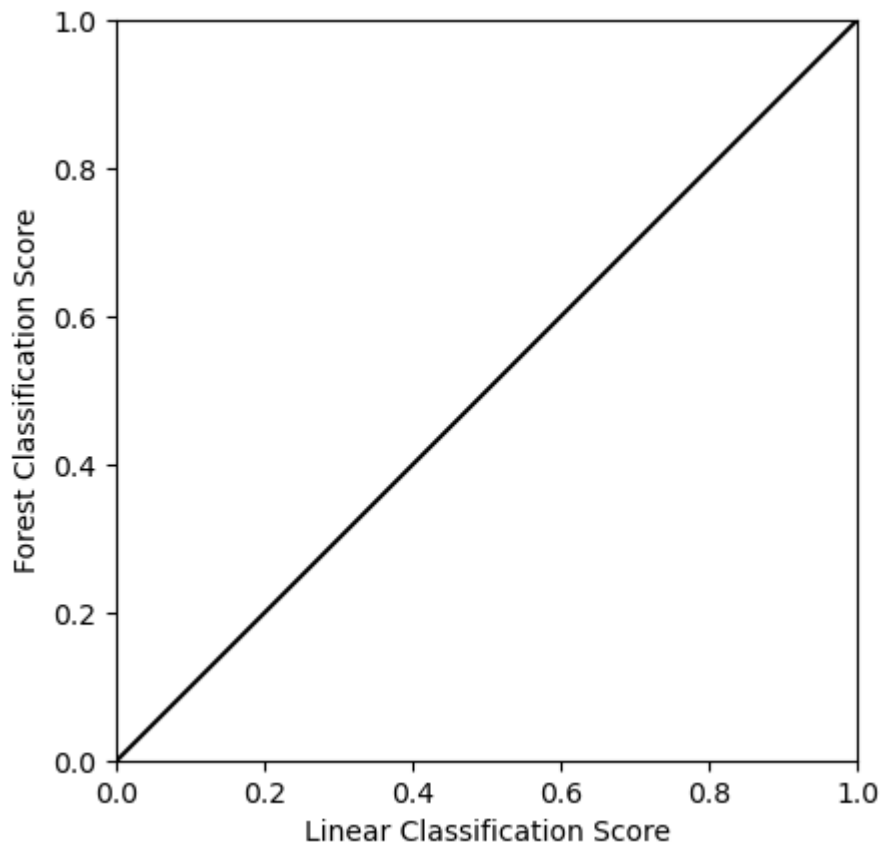
```python
plt.ylim(0, 1)
plt.xlabel("Linear Classification Score")
plt.ylabel("Forest Classification Score")

# Show the plot.

# Print the importance of each covariate in the random forest classification.
forest_classifier.fit(positive_revenue_df[all_covariates], positive_revenue_df[clas
sorted(list(zip(all_covariates, forest_classifier.feature_importances_)), key=lambd
```

Out[45]:
```
[('TV Movie', 0.0),
 ('Horror', 0.001715202327676785),
 ('Animation', 0.0019388197444951466),
 ('Comedy', 0.0022574689899296065),
 ('Foreign', 0.0022801352325337114),
 ('Documentary', 0.002846458591904433),
 ('Romance', 0.0031608732977368944),
 ('Thriller', 0.0035569898966812397),
 ('Mystery', 0.004282452349394276),
 ('Music', 0.004308655018573079),
 ('Fantasy', 0.0051937079152913745),
 ('Western', 0.005480591973153852),
 ('Family', 0.0066609392542522055),
 ('Crime', 0.006772395781754328),
 ('History', 0.006793172805113654),
 ('Action', 0.0073412694021133835),
 ('Adventure', 0.007596959755592538),
 ('Science Fiction', 0.010816587516514861),
 ('War', 0.011275947022575308),
 ('Drama', 0.023093574562804687),
 ('runtime', 0.04154729351420867),
 ('budget', 0.08765680648089587),
 ('vote_average', 0.10261105225795153),
 ('popularity', 0.2811360280003983),
 ('vote_count', 0.36967661830845444)]
```

In [ ]: