

## Predicting future environmental intensity (Time Series) and Collect Pilot Stocks Companies Descriptions

In this notebook, we will predict future environmental intensity for all the companies in the 'Excel data'. We will be using data from years to predict the future environmental intensity.

First, we will create the following columns:

- 1) Industry Indicator

- 1 if above the industry average in current year
- 0 if at industry average in current year

- (-1) if below the Environmental

At the end, we collected

- ```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import missingno as msno
import warnings
from sklearn import linear_model
from sklearn import metrics
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
import requests
```

```
from bs4 import BeautifulSoup
warnings.filterwarnings('ignore')
```

```
# df = pd.read_csv('Users/MarshallTorres/Documents/STAMP/Predicting-Environmental-and-Social-Actions/16_07_19/The dataset has (df.shape[0]) rows and (df.shape[1]) columns')
df.head(3)
```

The dataset has 14515 rows and 39 columns

|   | ISIN         | Year | Company Name    | Country        | Industry(Xlobase)                                | EnvironmentalIntensity(Sales) | EnvironmentalIntensity(OpInc) |
|---|--------------|------|-----------------|----------------|--------------------------------------------------|-------------------------------|-------------------------------|
| 0 | DE000545503  | 2016 | 1&1 DRILLISH AG | Germany        | Post and telecommunications (64)                 | -0.07%                        | -0.82%                        |
| 1 | GB00B1YW4409 | 2010 | 3i GROUP PLC    | United Kingdom | Financial intermediation except insurance and... | -0.12%                        | -0.11%                        |
| 2 | GB00B1YW4409 | 2011 | 3i GROUP PLC    | United Kingdom | Financial intermediation except insurance and... | -0.16%                        | -0.16%                        |

3 rows x 39 columns

Now, we are going to subset for the columns that we are interested in. As discussed in the notebook 'Environmental-Impact Data Cleaning', we transformed the EnvironmentalIntensity(Sales) to decimals and removed spaces in columns CompanyName and Country.

|   | ISIN         | Year | CompanyName      | Country        | Industry(Exiobase)                                | EnvironmentalIntensity(Sales) | Env. Int. |
|---|--------------|------|------------------|----------------|---------------------------------------------------|-------------------------------|-----------|
| 0 | DE0005455503 | 2016 | 1&1 DRILLISCH AG | Germany        | Post and telecommunications (54)                  |                               | -0.07%    |
| 1 | GB00B1YW4409 | 2010 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance and... |                               | -0.12%    |
| 2 | GB00B1YW4409 | 2011 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance and... |                               | -0.16%    |
| 3 | GB00B1YW4409 | 2012 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance and... |                               | -0.15%    |
| 4 | US88579Y1010 | 2010 | 3M COMPANY       | United States  | Activities of membership organisation n.e.c. (91) |                               | -7.90%    |

### Creating industry indicator

```
industry_avg = df.groupby('Industry(Exiobase)')['Env_intensity'].mean().reset_index()
df['industry_avg'] = df['Env_intensity'].groupby(df['Industry(Exiobase)']).transform('mean')

def create_ind(df):
    if df['Env_intensity'] > df['industry_avg']:
        return 1
    elif df['Env_intensity'] == df['industry_avg']:
        return 0
    elif df['Env_intensity'] < df['industry_avg']:
        return -1

df['industry_indicator'] = df.apply(create_ind, axis = 1)
df.head()
```

| ISIN         | Year | CompanyName      | Country        | Industry(Exiobase)                                | EnvironmentalIntensity(Sales) | Env. Intensity | Industry Avg | Indicator |
|--------------|------|------------------|----------------|---------------------------------------------------|-------------------------------|----------------|--------------|-----------|
| DE0005455503 | 2016 | 1&1 DRILLISCH AG | Germany        | Post and telecommunications (54)                  |                               | -0.07%         | -0.07%       | 0         |
| GB00B1YW4409 | 2010 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance and... |                               | -0.12%         | -0.12%       | 0         |
| GB00B1YW4409 | 2011 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance and... |                               | -0.16%         | -0.16%       | 0         |
| GB00B1YW4409 | 2012 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance and... |                               | -0.15%         | -0.15%       | 0         |
| US88579Y1010 | 2010 | 3M COMPANY       | United States  | Activities of membership organisation n.e.c. (91) |                               | -7.90%         | -7.90%       | 0         |

| ISIN         |
|--------------|
| DE0005545503 |

|   |              |      |              |                | (64)                                                 |        |         |           |
|---|--------------|------|--------------|----------------|------------------------------------------------------|--------|---------|-----------|
| 1 | G80081YW4409 | 2010 | 3i GROUP PLC | United Kingdom | Financial intermediation, except insurance and...    | -0.12% | -0.0012 | -0.028537 |
| 2 | G80081YW4409 | 2011 | 3i GROUP PLC | United Kingdom | Financial intermediation, except insurance and...    | -0.16% | -0.0016 | -0.028537 |
| 3 | G80081YW4409 | 2012 | 3i GROUP PLC | United Kingdom | Financial intermediation, except insurance and...    | -0.15% | -0.0015 | -0.028537 |
| 4 | US85879Y1010 | 2010 | 3M COMPANY   | United States  | Activities of membership organisation n.e.c.<br>(31) | -7.90% | -0.0790 | -0.175838 |

```
def create_ind_year(df):
    if df['Env_Intensity'] > df['industry_avg_year']:
        return 1
    elif df['Env_Intensity'] == df['industry_avg_year']:
        return 0
    elif df['industry_avg_year'] < df['industry_avg_year+1']:
        return 1
```

```
def get_industry_avg_year(industry):
    return -1

df['industry_avg_year'] =
```

```
df["Industry_indicator_year"] = df.apply([create_ind_year, axis = 1])
df.head()
```

|   | ISIN          | Year | CompanyName      | Country        | Industry(Sic6)                            | EnvironmentalIntensity(Sales) | Env_intensity | industry_avg | ln |
|---|---------------|------|------------------|----------------|-------------------------------------------|-------------------------------|---------------|--------------|----|
| 0 | DE0005454503  | 2015 | 1&1 DRILLISCH AG | Germany        | Post and telecommunications (64)          | -0.07%                        | -0.0007       | -0.020506    |    |
| 1 | GB00B11YW4409 | 2010 | 3i GROUP PLC     | United Kingdom | Financial Intermediation except insurance | -0.12%                        | -0.0012       | -0.028537    |    |

| 2                                                                                                        | G800B1YW4409 | 2011 | 3i GROUP PLC    | United Kingdom | Financial intermediation, except insurance and... | -0.16%                        | -0.0016       | -0.028537   |  |
|----------------------------------------------------------------------------------------------------------|--------------|------|-----------------|----------------|---------------------------------------------------|-------------------------------|---------------|-------------|--|
| 3                                                                                                        | G800B1YW4409 | 2012 | 3i GROUP PLC    | United Kingdom | Financial intermediation, except insurance and... | -0.15%                        | -0.0015       | -0.028537   |  |
| 4                                                                                                        | US88579Y1010 | 2010 | 3M COMPANY      | United States  | Activities of membership organisation n.e.c. (91) | -7.90%                        | -0.0790       | -0.175838   |  |
| df.loc[(df['Industry(Exibase)'] == 'Activities auxiliary to financial intermediation (67)')[:,1].sort... |              |      |                 |                |                                                   |                               |               |             |  |
|                                                                                                          | ISIN         | Year | CompanyName     | Country        | Industry(Exibase)                                 | EnvironmentalIntensity(Sales) | Env_intensity | industry_av |  |
| 190                                                                                                      | CH0012138605 | 2010 | ADECCO GROUP AG | Switzerland    | Activities auxiliary to financial intermediat...  | -0.14%                        | -0.0014       | -0.0045     |  |
| 2052                                                                                                     | G800B23K0M20 | 2010 | CAPITA PLC      | United Kingdom | Activities auxiliary to financial...              | -0.40%                        | -0.0040       | -0.0045     |  |

|      |              |      |                                     |                |                                                  | intermediate... |         |         |  |
|------|--------------|------|-------------------------------------|----------------|--------------------------------------------------|-----------------|---------|---------|--|
| 1919 | FR0006174348 | 2010 | BUREAU VERITAS SA                   | France         | Activities auxiliary to financial intermediat... | -0.52%          | -0.0052 | -0.0045 |  |
| 3524 | DE0005810055 | 2010 | DEUTSCHE BOERSE AG                  | Germany        | Activities auxiliary to financial intermediat... | -0.28%          | -0.0028 | -0.0045 |  |
| 1718 | ES0115056139 | 2010 | BOLSA Y MERCADOS ESPAÑOLAS SIMOT SA | Spain          | Activities auxiliary to financial intermediat... | -0.26%          | -0.0026 | -0.0045 |  |
| ...  | ...          | ...  | ...                                 | ...            | ...                                              | ...             | ...     | ...     |  |
| 4428 | GB00B19NVL48 | 2019 | EXPERIAN PLC                        | United Kingdom | Activities auxiliary to financial intermediat... | -0.20%          | -0.0020 | -0.0045 |  |
| 979  | FR0000074148 | 2019 | ASSYSTEM SA                         | France         | Activities auxiliary to financial intermediat... | -0.31%          | -0.0031 | -0.0045 |  |
| 3529 | DE0005810055 | 2019 | DEUTSCHE BOERSE AG                  | Germany        | Activities auxiliary to financial intermediat... | 1.88%           | 0.0188  | -0.0045 |  |
| ...  | ...          | ...  | PAGEGROUP                           | United         | Activities auxiliary                             | ...             | ...     | ...     |  |

| ISIN | Year         | Company Name     | Country        | Industry (Exiobase)                        | Environmental Intensity (Sales) | Env_intensity | industry_avg | in |
|------|--------------|------------------|----------------|--------------------------------------------|---------------------------------|---------------|--------------|----|
| 0    | DE000545503  | 1&1 DRILLISCH AG | Germany        | Post and telecommunications (64)           | -0.07%                          | -0.0007       | -0.020596    |    |
| 1    | GBO0011W4449 | 3i GROUP PLC     | United Kingdom | Financial intermediation, except insurance | -0.12%                          | -0.0012       | -0.028537    |    |

|   |              |      |              |                |                                                   |        |         |           |
|---|--------------|------|--------------|----------------|---------------------------------------------------|--------|---------|-----------|
| 2 | G800B1YW4409 | 2011 | 3i GROUP PLC | United Kingdom | Financial intermediation, except insurance and... | -0.16% | -0.0016 | -0.028537 |
| 3 | G800B1YW4409 | 2012 | 3i GROUP PLC | United Kingdom | Financial intermediation, except insurance and... | -0.15% | -0.0015 | -0.028537 |
| 4 | US88579Y1010 | 2010 | 3M COMPANY   | United States  | Activities of membership organisation n.e.c. (91) | -7.90% | -0.0790 | -0.175838 |

```
df[["Environmental_Growth"] = df.groupby(["CompanyName"])["Env_intensity"].apply(lambda x: x.pct_change(df.head(1))
```

| ISIN         | Year | CompanyName     | Country        | Industry(Exioabase)              | EnvironmentalIntensity(Sales) | Env_intensity | industry_avg | in |
|--------------|------|-----------------|----------------|----------------------------------|-------------------------------|---------------|--------------|----|
| DE0000545503 | 2016 | 1&1 DRLLISCH AG | Germany        | Post and telecommunications (64) | -0.07%                        | -0.0007       | -0.020506    |    |
|              |      |                 | United Kingdom | Financial intermediation         |                               |               |              |    |

|   |              |      |              |                |                                                   |        |         |           |
|---|--------------|------|--------------|----------------|---------------------------------------------------|--------|---------|-----------|
| 1 | G800B1YW4409 | 2010 | 3I GROUP PLC | United Kingdom | Financial intermediation, except insurance and... | -0.12% | -0.0012 | -0.028537 |
| 2 | G800B1YW4409 | 2011 | 3I GROUP PLC | United Kingdom | Financial intermediation, except insurance and... | -0.16% | -0.0016 | -0.028537 |
| 3 | G800B1YW4409 | 2012 | 3I GROUP PLC | United Kingdom | Financial intermediation, except insurance and... | -0.15% | -0.0015 | -0.028537 |
| 4 | US88579Y1010 | 2010 | 3M COMPANY   | United States  | Activities of membership organisation n.e.c. (91) | -7.90% | -0.0790 | -0.175838 |

## Model 1 - Past years Environmental Intensity

```
df1 = df.copy()

def p2f(x):
    return float(x.strip('%')/100)
df1['EnvironmentalIntensity(Sales)'] = df1['EnvironmentalIntensity(Sales)'].apply(p2f)
```

```
companies_2018 = list(dfl[dfl['Year']
```

```
companies_2019 = list(dfl[(dfl['Year'] == 2019) & (dfl['CompanyName'].isin(companies_2018))])

#Getting companies that are in both years
list2018 as set = set(companies_2018)
```

```

intersection = list2018_as_set.intersection(companies_2019)

X = df1[(df1['Year'] == 2018) & (df1['CompanyYear'].isin(intersection))][['Env_intensity']]
y = df1[(df1['Year'] == 2019) & (df1['CompanyYear'].isin(intersection))][['Env_intensity']]

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42)

reg = LinearRegression().fit(X_train, y_train)

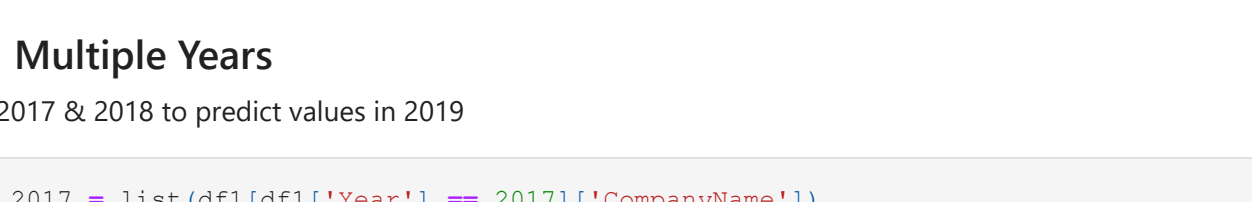
y_pred = reg.predict(X_test)

print('R2 score', metrics.r2_score(y_test, y_pred))
print('MSE', metrics.mean_squared_error(y_test, y_pred))

R2 score: 0.8533746128153976
MSE: 0.0152514501102281

fig, ax = plt.subplots(figsize=(12,8))
ax.scatter(y_test, y_pred, c='green')
ax.plot([y_test.min(), y_test.max()], [y.min(), y.max()], 'k--', lw=4)
ax.set_xlabel('Observed 2019 Environmental Intensity')
ax.set_ylabel('Predicted 2019 Environmental Intensity')
ax.set_title('Linear Regression - Predicting 2019 Environmental intensity using 2018');

```



## Including Multiple Years

Starting with 2017 & 2018 to predict values in 2019

```
companies_2017 = list(df1[df1['Year'] == 2017])['CompanyName']
list2017_as_set = set(companies_2017)
intersection2 = list2017_as_set.intersection(intersection)

X = df1[df1['Year'].isin([2017, 2018])] & df1['CompanyName'].isin(intersection2)][['CompanyName', 'log2019_sales']]
X.groupby('CompanyName').mean()[['EnvironmentalIntensity(Sales)']]

y = df1[df1['Year'] == 2019] & df1['CompanyName'].isin(intersection2)][['EnvironmentalIntensity(Sales)']]

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42)

reg = LinearRegression().fit(X_train, y_train)

y_pred = reg.predict(X_test)

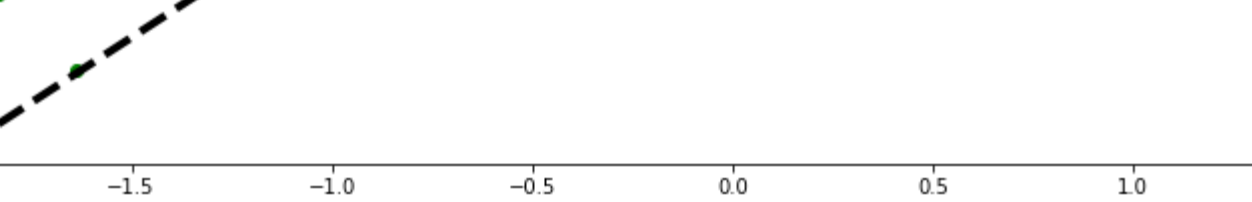
print('R2 score:', metrics.r2_score(y_test, y_pred))
```

```
print('Adjusted R^2: ', 1-'0.8729')*(len(intersection2)-1)/(len(intersection2)-2-1))

R2 score: 0.9640133945498272
MSE: 0.002836465211312934
Adjusted R2: 0.837016996642953

fig, ax = plt.subplots(figsize=(12,8))
ax.scatter(y_test, y_pred, c='green')
ax.plot([y_test.min(), y_test.max()], [y.min(), y.max()], 'k--', lw=4)
ax.set_xlabel('Observed 2019 Environmental Intensity')
ax.set_ylabel('Predicted 2019 Environmental Intensity')
ax.set_title('Linear Regression- Predicting 2019 Environmental intensity using 2017 and 2018'))
```

19 Environ  
-0.5



We see that the results of this model are worse than our first model which predicted 2019 numbers with 2018. We go on to include

```

Companies_2016 = list(df1[df1['Year'] == 2016]['CompanyName'])
list2016_as_set = set(Companies_2016)
intersection3= list2016_as_set.intersection(intersection2)

X = df1[df1['Year'].isin([2016, 2017, 2018])] & (df1['CompanyName'].isin(intersection3))
X = X.groupby('CompanyName').mean()[['EnvironmentalIntensity(Sales)']]

y = df1[df1['Year'] == 2019] & (df1['CompanyName'].isin(intersection3))
y = y.groupby('CompanyName').mean()[['EnvironmentalIntensity(Sales)']]

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42)

reg = LinearRegression().fit(X_train, y_train)

```

```

y_pred = reg.predict(X_test)

print('R2 score:', metrics.r2_score(y_test, y_pred))
print('MSE:', metrics.mean_squared_error(y_test, y_pred))
print('Adjusted R2:', 1-((1-.92741)/(len(intersection3)-1))/(len(intersection3)-3-1))

R2 score: 0.916454447105587
MSE: 0.00421395537780859
Adjusted R2: 0.9272047502336268

When we include 2016 both our R score and MSE improved.

fig, ax = plt.subplots(figsize=(12,8))
ax.scatter(y_test, y_pred, c='green')
ax.plot(y_test.min(), y_test.max(), [y.min(), y.max()], 'k--', lw=4)
ax.set_ylabel('Observed 2019 Environmental Intensity')
ax.set_xlabel('Predicted 2019 Environmental Intensity')
ax.set_title('Linear Regression- Predicting 2019 Environmental intensity using 2016 to 2018');

```

```
df2 = df2.copy()
EI_2016 = df2[df2['Year'] == 2016]
EI_2017 = df2[df2['Year'] == 2017]
EI_2018 = df2[df2['Year'] == 2018]
EI_2019 = df2[df2['Year'] == 2019]

df_mod1 = EI_2016.merge(EI_2017, how='inner', on='CompanyName', suffixes=('_2016', '_2017'))
df_mod2 = df_mod1.merge(EI_2018, how='inner', on='CompanyName', suffixes=('_2016', '_2018'))
df_mod3 = df_mod2.merge(EI_2019, how='inner', on='CompanyName', suffixes=('_2018', '_2019'))

df_mod3 = df_mod3[['CompanyName', 'Year_2016', 'Env_intensity_2016', 'Year_2017', 'Env_intensity_2017', 'Year_2018', 'Env_intensity_2018', 'Year_2019', 'Env_intensity_2019']]

df_mod3.head()
```

|   | CompanyName       | Year_2016 | Env_intensity_2016 | Year_2017 | Env_intensity_2017 | Year_2018 | Env_intensity_2018 | Year_2019 | Env_intensity_2019 |
|---|-------------------|-----------|--------------------|-----------|--------------------|-----------|--------------------|-----------|--------------------|
| 0 | 3M COMPANY        | 2016      | -0.0705            | 2017      | -0.0660            | 2018      | -0.0710            | 2019      | -0.0710            |
| 1 | AGV PRODUCTS CORP | 2016      | -0.0170            | 2017      | -0.0190            | 2018      | -0.0140            | 2019      | -0.0140            |
| 2 | AA PIP            | 2016      | -0.0130            | 2017      | -0.0124            | 2018      | -0.0117            | 2019      | -0.0117            |

|   |                           |      |         |      |         |
|---|---------------------------|------|---------|------|---------|
| 3 | AAC TECHNOLOGIES HOLDINGS | 2016 | -0.0389 | 2017 | -0.0591 |
|---|---------------------------|------|---------|------|---------|

|   |                |      |         |      |         |      |         |
|---|----------------|------|---------|------|---------|------|---------|
| 4 | AAREAL BANK AG | 2016 | -0.0024 | 2017 | -0.0021 | 2018 | -0.0019 |
|---|----------------|------|---------|------|---------|------|---------|

```
df_melt2 = pd.melt(df_mod2, id_vars=['CompanyName'], value_vars=['Env_intensity_2016', 'Env_intensity_2017', 'Env_intensity_2018', 'Env_intensity_2019'], var_name='myVarName', value_name='Environmental_intensity')
```

```
df_melt2.head()
```

|   | CompanyName                             | myVarName          | Environmental_intensity |
|---|-----------------------------------------|--------------------|-------------------------|
| 0 | 3M COMPANY                              | Env_intensity_2016 | -0.0705                 |
| 1 | A.G.V. PRODUCTS CORP                    | Env_intensity_2016 | -0.0170                 |
| 2 | AA PLC                                  | Env_intensity_2016 | -0.0130                 |
| 3 | AAC TECHNOLOGIES HOLDINGS INCORPORATION | Env_intensity_2016 | -0.0389                 |
| 4 | AAREAL BANK AG                          | Env_intensity_2016 | -0.0024                 |

```
xi_avg = df_melt2.groupby('CompanyName').mean().reset_index()
X = xi_avg[['Environmental_intensity']]
y = df_mod2[['Env_intensity_2019']]

print(X.shape)
```

```
print(y.shape)

(1065, 1)
(1065, 1)

x_train, x_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
# train Linear Regression
LTrainer = LinearRegression()
LTrainer.fit(x_train, y_train)

# mse for linear regression
y_pred_lr = LTrainer.predict(x_test)
print(metrics.mean_squared_error(y_pred_lr, y_test))

0.006269966737857914

# display the parameters
print('Model intercept: ', LTrainer.intercept_)
print('Model coefficients: ', LTrainer.coef_)

Model intercept: [-0.0051724]
Model coefficients: [[0.80432622]]

print('R2 score:', metrics.r2_score(y_test, y_pred_lr))

R2 score: 0.8756921251462437
```

```
fig, ax = plt.subplots(figsize=(12,8))
ax.scatter(y_test, y_pred, s=100, c='green')
ax.plot([y_test.min(), y_test.max()], [y_min(), y_max()], 'k--', lw=4)
ax.set_xlabel('Observed 2019 Environmental Intensity')
ax.set_ylabel('Predicted 2019 Environmental Intensity')
```

[illegible]

```

df_industry = df.groupby('Industry(Exiabase)').count()[['CompanyName']].reset_index()

industries = df_industry[df_industry['CompanyName'] > 3]['Industry(Exiabase)']

df_industry_count4 = df[df['Industry(Exiabase)'].isin(industries)]

df_2018 = df_industry_count4.loc[df.Year == 2018, ]
df_2019 = df_industry_count4.loc[df.Year == 2019, ]
df_mod1 = pd.merge(df_2018, df_2019, on='CompanyName', how='inner')
df_mod3 = df_mod3[['Year_x', 'CompanyName', 'industry_avg_year_x', 'Year_y', 'Env_intensity_y']]
df_mod3.head()

```

|   | Year_x | CompanyName          | industry_avg_year_x | Year_y | Env_intensity_y |
|---|--------|----------------------|---------------------|--------|-----------------|
| 0 | 2018   | 3M COMPANY           | -0.229308           | 2019   | -0.0641         |
| 1 | 2018   | 3SBIO INC            | -0.027793           | 2019   | -0.0240         |
| 2 | 2018   | A.G.V. PRODUCTS CORP | -0.072254           | 2019   | -0.0172         |
| 3 | 2018   | AA PLC               | -0.073777           | 2019   | -0.0070         |

```
4 2018 AAC TECHNOLOGIES HOLDINGS INCORPORATION -0.023555 2019 -0.1080
```

```
X = df_mod3[['industry_avg_year_x']]
y = df_mod3.iloc[:,4]
print(X.shape)
print(y.shape)

(1336, 1)
(1336,)

x_train, x_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
# Train Linear Regression
lRtrainer = LinearRegression()
lRtrainer.fit(x_train, y_train)

# use for linear regression
y_pred_lr = lRtrainer.predict(x_test)
print(metrics.mean_squared_error(y_pred_lr, y_test))

0.06598743409385226

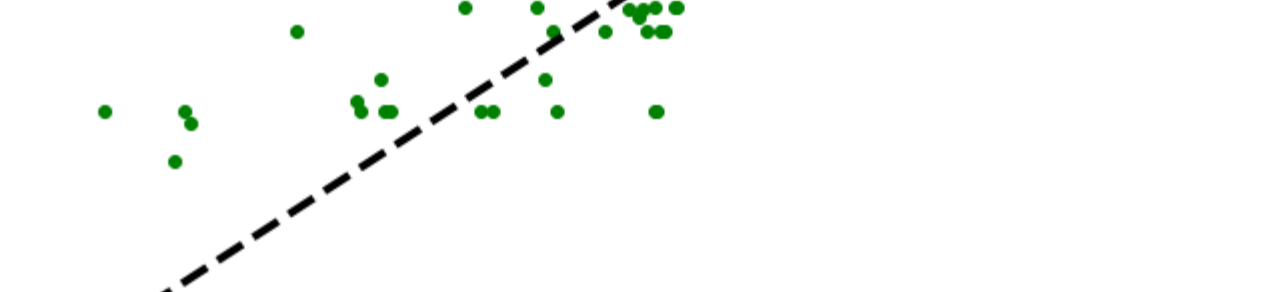
# display the parameters
print('Model intercept: ', lRtrainer.intercept_)
print('Model coefficients: ', lRtrainer.coef_)
```

```
Model intercept: 0.000269631900767966
Model coefficients: [0.91303819]
```

```
print('R2 score', metrics.r2_score(y_test, y_pred_lr))
```

R2 score: 0.36647102462935643

```
fig, ax = plt.subplots(figsize=(12,8))
ax.scatter(y_test, y_pred_lr, c='green')
ax.plot(y_test.min(), y_test.max(), 'r--', lw=4)
ax.set_xlabel('Observed 2019 Environmental Intensity')
ax.set_ylabel('Predicted 2019 Environmental Intensity')
ax.set_title('Linear Regression- Predicting 2019 Environmental Intensity using 2018 industry average e
```



A scatter plot showing the relationship between Observed 2019 Environmental Intensity (x-axis) and Predicted 2019 Environmental Intensity (y-axis). The x-axis ranges from -2.0 to 1.5, and the y-axis ranges from -2.0 to 0.0. A dashed black line represents the linear regression fit. The data points are green dots, showing a positive correlation between the observed and predicted values.

Now, we are going to consider 2017 - 2018 to predict 2019

```
df_2019=df_industry_count4(df_industry_count4['Year'] == 2019)
df_2018=df_industry_count4(df_industry_count4['Year'] == 2018)
df_2017=df_industry_count4(df_industry_count4['Year'] == 2017)
```

```
df2019=df2019[['Year', 'CompanyName', 'Env_intensity', 'Ind_Yearavg2019', inplace=True]]
df2018.rename(columns={'Env_intensity': 'Env_intensity2019', inplace=True})
df2018=df2018[['CompanyName', 'Industry_avg_year']]
df2018.rename(columns={'Industry_avg_year': 'Ind_Yearavg2018', inplace=True})
df2018.rename(columns={'Industry_avg_year': 'Industry_avg_year', inplace=True})
df2017.rename(columns={'Industry_avg_year': 'Ind_Yearavg2017', inplace=True})
mdl = pd.merge(df2019, df2018, on=['CompanyName'])
mdl1 = pd.merge(mdl, df2017, on=['CompanyName'])
mdl1
```

| Year | CompanyName | Env_intensity2019                       | Ind_Yearavg2018 | Ind_Yearavg2017 |           |
|------|-------------|-----------------------------------------|-----------------|-----------------|-----------|
| 0    | 2019        | 3M COMPANY                              | -0.0641         | -0.229308       | -0.225496 |
| 1    | 2019        | 3SBIO INC                               | -0.0340         | -0.027793       | -0.031771 |
| 2    | 2019        | A.G.V. PRODUCTS CORP                    | -0.0172         | -0.072254       | -0.063618 |
| 3    | 2019        | AA PLC                                  | -0.0079         | -0.073777       | -0.096157 |
| 4    | 2019        | AAC TECHNOLOGIES HOLDINGS INCORPORATION | -0.1080         | -0.023555       | -0.027254 |
| ...  | ...         | ...                                     | ...             | ...             | ...       |
| 1190 | 2019        | ZEON CORPORATION                        | -0.0730         | -0.218707       | -0.257385 |
| 1191 | 2019        | ZHEN DING TECHNOLOGY HOLDING LIMITED    | -0.0602         | -0.076812       | -0.068294 |
| 1192 | 2019        | ZIG SHENG INDUSTRIAL COMPANY LIMITED    | -0.1615         | -0.124750       | -0.117456 |

[illegible]

```

Model train: 1.064, test: 0.038

# display the parameters
print('Model intercept: ', regr.intercept_)
print('Model coefficients: ', regr.coef_)

Model intercept: 0.0005763983280114432
Model coefficients: [ 0.94555899 -0.06948928]

print('R2 score: ', metrics.r2_score(y_test, y_pred))

R2 score: 0.26873239153605133

fig, ax = plt.subplots(figsize=(12,8))
ax.scatter(y_test, y_pred, c='green')
ax.plot([y_test.min(), y_test.max()], [y_min(), y_max()], 'k--', lw=4)
ax.set_ylabel('Observed 2019 Environmental Intensity')
ax.set_xlabel('Predicted 2019 Environmental Intensity')
ax.set_title('Linear Regression- Predicting 2019 Environmental intensity using 2017- 2018 industry average')

Linear Regression- Predicting 2019 Environmental intensity using 2017-2018 industry average environmental intensity

```







