- **Scikit-learn**

    It  is  a free  software machine  learning library for  the Python programming language. The framework is built on top of several popular Python packages, namely NumPy, SciPy, and matplotlib. Scikit-Learn is characterized by a clean, uniform, and streamlined API, as well as by very useful and complete online documentation. A benefit of this uniformity is that once you understand the basic use and syntax of Scikit-Learn for one type of model, switching to a new model or algorithm is very straightforward.

## 3.5 Configuration

### 3.5.1 Scraping

```python
driver.maximize_window()
driver.get('https://www.boerse-frankfurt.de/?lang=en')
driver.find_element_by_xpath('//button[contains(text(),"Accept")]').click()
search_bar = driver.find_element_by_xpath('//input')
search_bar.clear()
search_bar.send_keys("DAX")
sleep(3)
search_bar.send_keys(Keys.ENTER)
print(driver.current_url)
dax_page = driver.find_element_by_xpath('//a[contains(text(),"DAX")]').get_attribute('href')
driver.get(dax_page)
sleep(3)
driver.find_element_by_xpath('//button[contains(text(),"Constituents")]').click()
sleep(4)
driver.find_element_by_xpath('//button[contains(@class,"page-bar-type-button") and contains(text(),"100")]').click()
sleep(3)
lst = driver.find_elements_by_xpath('//table/tbody/tr/td/div/a')
df = pd.DataFrame()
for i in lst:
    df = df.append({'constituent_name' : i.text, 'ISIN' : i.get_attribute('href').split('/')[-1]},ignore_index=True)
```

Figure 3.3 Collect all the constituents' name and WKN from website

The Figure 3.3 the code for collecting the constituents and its details provided by the website is collected. This data is stored in a sqlite3 database.

```python
for index,row in df.iterrows():
    final_data = {}
    driver.get(url.format(row['constituent_ISIN']))
    sleep(4)
    tabs = driver.find_elements_by_xpath('//button[contains(@class,"data-menue-button")]')
    for tab in tabs:
        if tab.text in ['Charts','News','Company Details']:
            continue
        tab.click()
        sleep(4)
        tables = driver.find_elements_by_xpath('//table')
        table_names = []
        for table in tables:
            try:
                table_name = table.find_element_by_xpath(
                    './..//preceding-sibling::h2[contains(@class,"widget-table-headline")]').text
            except:
                table_name = ''
            if table_name.find(row['constituent_name']) != -1:
                table_name = table_name[ :table_name.find(row['constituent_name'])].strip()
            table_names.append(table_name)
        data = pd.read_html(driver.page_source)
        for each_df,table_name in zip(data,table_names):
            if not table_name:
                continue
            final_data[table_name] = each_df
    cleaner = Cleaner(final_data)
```

Figure 3.4 Collect stock related data for each constituent

The Figure 3.4 shows code snippet which collects stock related data from each constituent.

### 3.5.2 Cleaning

```python
if i == 'Historical key data':
    self.data[i] = self.data[i].T.reset_index()
    self.data[i]['index'].iloc[0] = 'Year'
    self.data[i].columns = self.data[i].iloc[0]
    self.data[i] = self.data[i].iloc[1:]
if i != 'Historical key data'and i != 'Historical prices and volumes':
    self.data[i] = self.data[i].T
    self.data[i].columns = self.data[i].iloc[0]
    self.data[i] = self.data[i].iloc[1:]
try:
    self.data[i].columns = self.data[i].columns.str.replace(r'[^a-zA-Z0-9\s]*', '', regex=True).str.strip()
    self.data[i].columns = self.data[i].columns.str.strip().str.replace(r'[\s]+', '_', regex=True).str.strip()
    self.data[i].columns = self.data[i].columns.str.lower().str.strip()
    self.data[i].columns = self.data[i].columns.str.replace(r'_in$','',regex=True).str.strip()
except:
    pass
finally:
    for j in self.data[i].columns:
        try:
            if j == 'year':
                raise
            self.data[i][j] = self.data[i][j].apply(float)
        except:
            pass
        try:
            self.data[i][j] = self.data[i][j].str.replace(r'm$|bn$','',regex = True).str.strip()
            self.data[i][j] = self.data[i][j].apply(lambda x : datetime.strptime(x[:-2] + '20' + x[-2:],'%d/%m/%Y'))
        except:
            pass
```

Figure 3.5 Clean and pre-processing tabular data collected by collection script

Figure 5.5 shows code snippet to preprocess the tabular stock data collected. The data collected will not conform with the types in database, it may contain unwanted text characters which will is noise according to the model and hence these discrepancies will be removed.

### 3.5.3 Storing

```python
def run_query(self,query,params=None):
    return (self.conn.execute(query))

def insert_bulk(self,table_name,data):
    #insertion
    data.to_sql(table_name,self.conn,if_exists='append',index=False)
    print('Data inserted')
    # print(data.head())

def get_date(self,table_name):
    """Get latest collection data"""
    try:
        query = "SELECT max(collection_date) from {}".format(table_name)
        date = list(self.conn.execute(query))[0][0]
        return date
    except:
        return None
```

Figure 3.6 Storage utility which includes database function

The storage utility handles all data manipulation process on the database end. The code in Figure 3.6 shows the functions used for this manipulation which includes process like running a query to retrieve data and inserting data

### 3.5.4 Model

```python
@app.route('/yearly')
def yearly():
    conn = sqlite3.connect('constituents.db')
    df = pd.read_sql_query('select * from yearly',conn).groupby('wkn')['number_of_employees','sales_in_mio'].sum()#.to_dict(orient=
    df.columns = ['x','y']
    model = LinearRegression()
    model.fit(df.drop('y',axis=1),df['y'])
    wkn = ' '.join(df.index.tolist())
    test = df.to_dict(orient='records')
    preds = model.predict(df.drop('y',axis=1))
    df['y'] = preds
    preds = df.to_dict(orient='records')
    conn.close()
    # return render_template('yearly.html',test=test,labels="-".join(labels),data=" ".join(map(str,data)))
    return render_template('yearly.html',test=test,wkn=wkn,preds=preds)

if __name__ == '__main__':
    app.run(debug = True)
```

Figure 3.7 Linear regression model

Figure 3.7 shows the code which predicts the stock trends using linear regression.

## 4.2 Results

Daily Historical



Figure 4.1 Yearly data trend

The figure 4.1 shows the linear model fitted for number of employees vs sales in million.
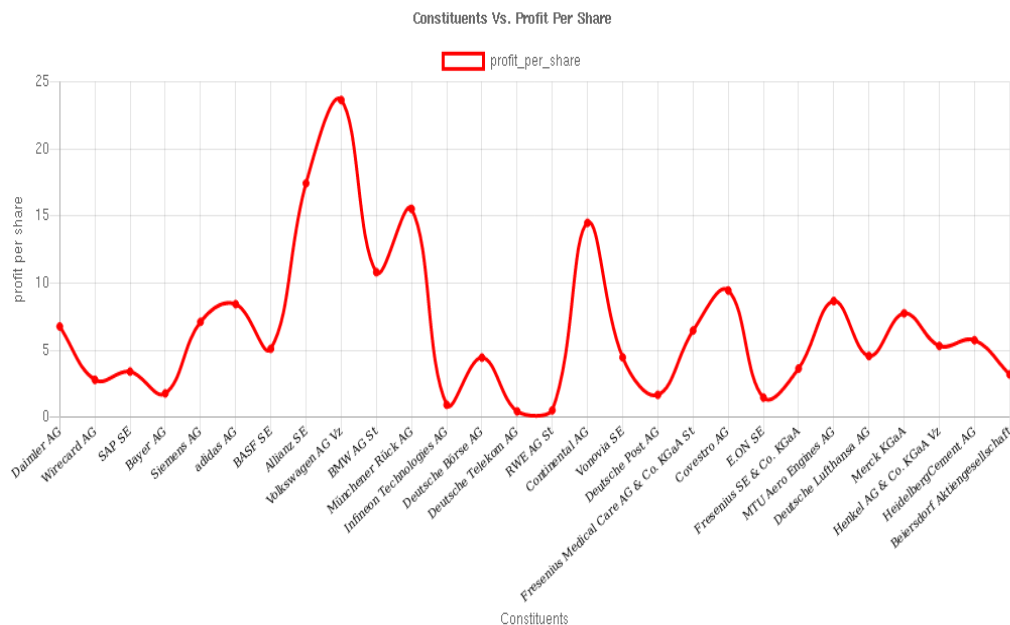
Yearly Historical



Figure 4.2 Daily data trend

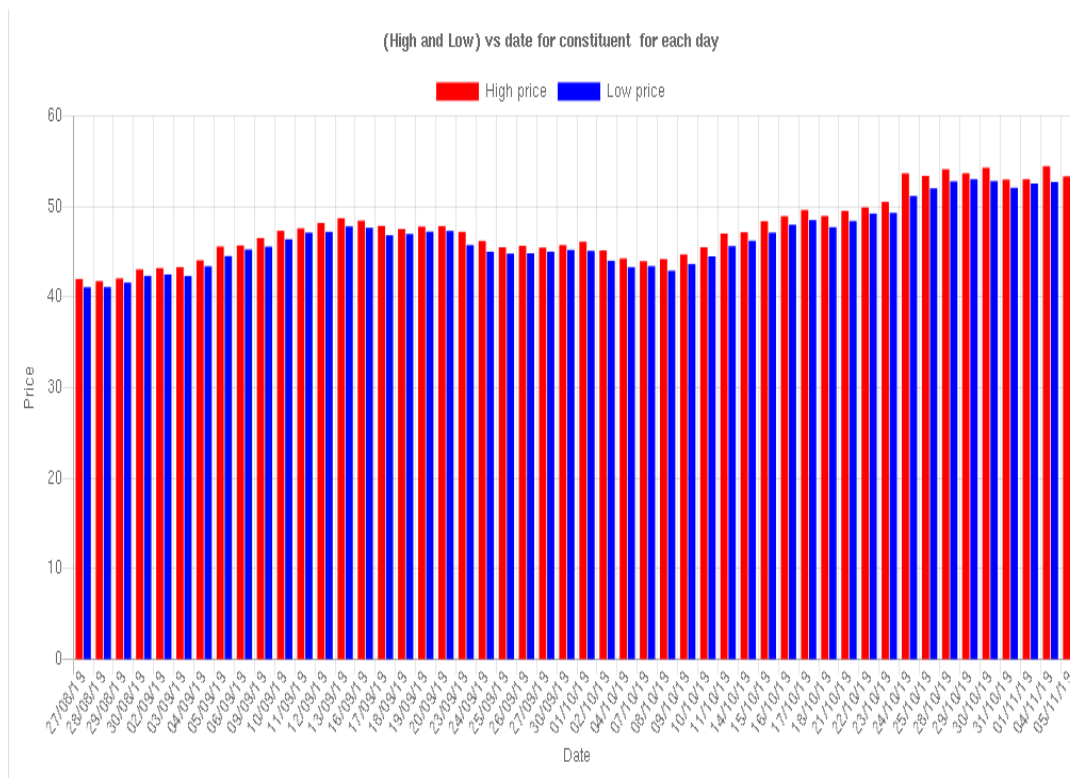Figure 4.2 shows profit per share for each constituent to help make decisions on which stock to buy or sell.

Figure 4.3 Historical data Bar chart

Figure 4.3 shows the high and low prices of the stock of a constituent on each day for the past month.
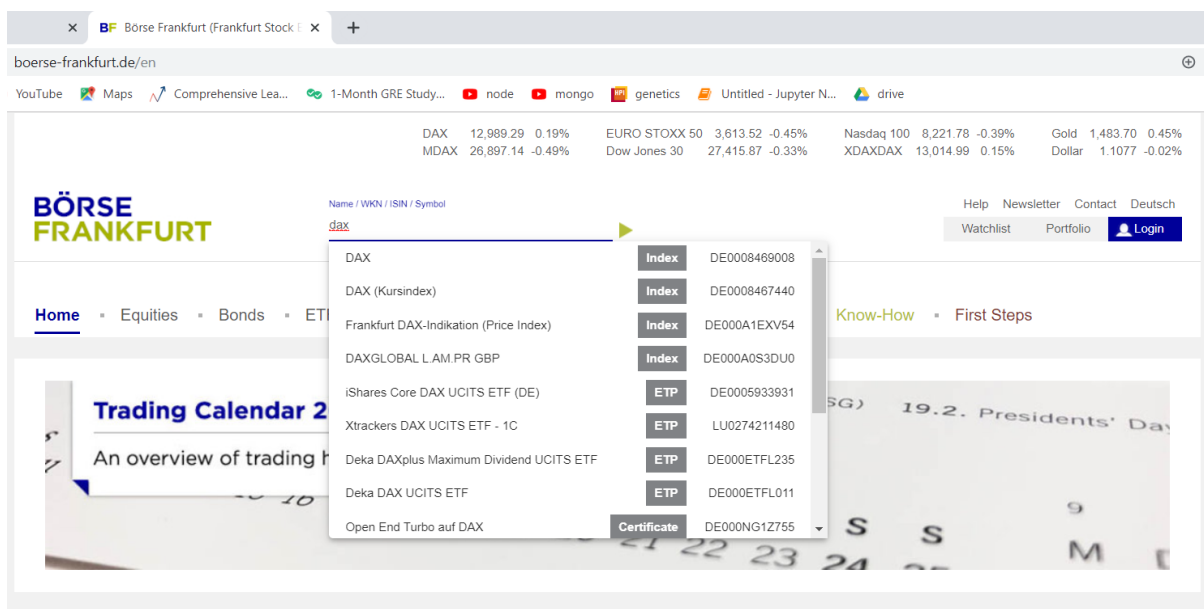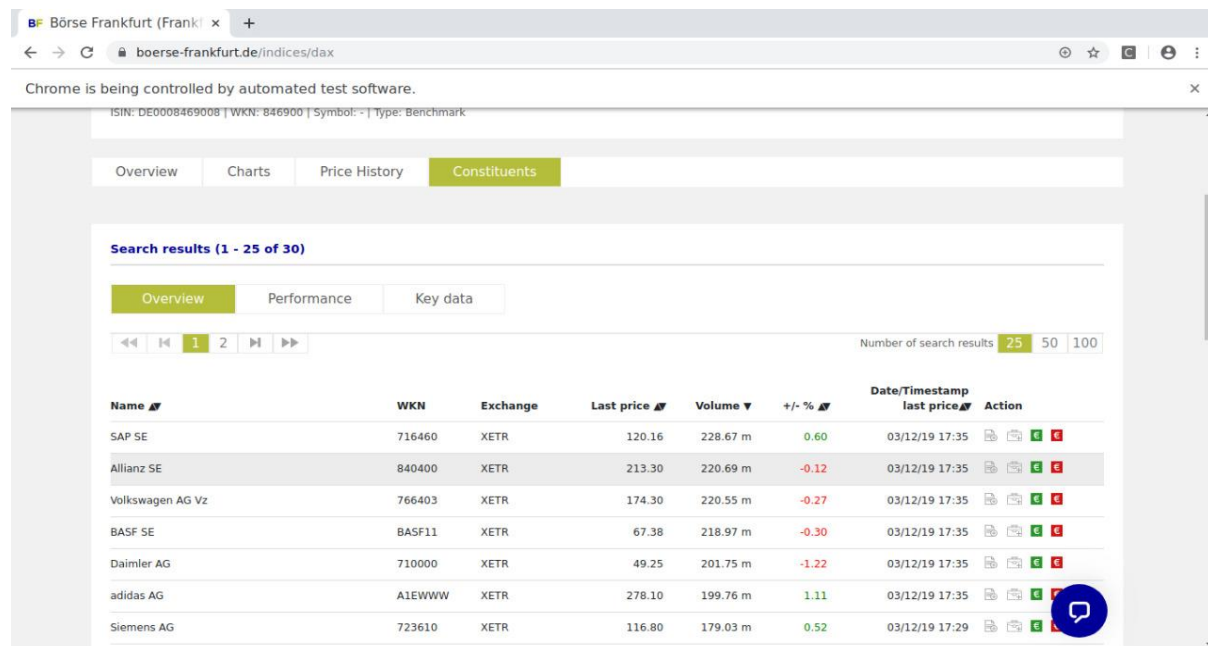
## 4.3 Snapshots



Figure 4.4 Sending keys to search box

The Figure 4.5 shows the selenium tool being used to enter keys into a search box.

Figure 4.5 Collecting tabular data

The Figure 4.5 shows collecting of an entire table using the selenium tool



Figure 4.6 Collecting multiple tables

The Figure 4.6 shows how the selenium tool collects multiple tables of different format and stores in a pandas dataframe.