**Text Analysis - Companies description** In this notebook, we will do some text analysis for the companies existing in the WRDS and Harvard Business School, Impact-Weighted Accounts Project report. We will apply some Nature Language Processing (NLP) using the pre-trained DistilBERT. First, we need to merge the datasets and obtain the companies description from Yahoo Finance. 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') ## Useful for the DistilBERT model import torch import transformers as ppb # pytorch transformers from sklearn.linear model import LogisticRegression from sklearn.model selection import cross val score from sklearn.model selection import train test split df ishares = pd.read csv('/Users/maralinetorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/De df\_ei = pd.read\_csv('/Users/maralinetorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/Dataset stocks = pd.read csv("/Users/maralinetorres/Documents/GitHub/Predicting-Environmental-and-Social-Actions/Datase For the ishares, we don't need the tickers in the sector 'Cash and devirates' and we are also going to filter by the necessary columns df ishares = df ishares.loc[df ishares.Sector != 'Cash and/or Derivatives',['Ticker','Name','Sector','CUSIP','] df ishares.drop duplicates(inplace=True) df ishares.head() Ticker Name Sector CUSIP ISIN Location AAPL APPLE INC Information Technology 37833100 US0378331005 United States MSFT MICROSOFT CORP Information Technology 594918104 US5949181045 United States 2 AMZN AMAZON COM INC Consumer Discretionary US0231351067 United States 23135106 3 TSLA TESLA INC Consumer Discretionary 88160R101 US88160R1014 United States FACEBOOK CLASS A INC Communication 30303M102 US30303M1027 United States In [4]: len(df ishares.Ticker.unique()) Out[4]: 855 df ei2 = df ei.copy()df ei2 = df ei2.loc[:,['ISIN','CompanyName','Country']] df\_ei2.drop\_duplicates(inplace=True) df ei2.head() ISIN CompanyName Country 1&1 DRILLISCH AG DE0005545503 Germany 3I GROUP PLC United Kingdom 1 GB00B1YW4409 US88579Y1010 3M COMPANY **United States 3SBIO INC** KYG8875G1029 China GI000A0F6407 888 HOLDINGS PLC Gibraltar len(df ei2.ISIN.unique()) Out[6]: 2628 stocks = stocks.iloc[:,0:2] stocks.head() Year Ticker **0** 2005 **AEE** 2006 AEE **2** 2007 AEE 2008 AEE **4** 2009 AEE Let's start merging the datasets. First, we are going to merge iShares and Environmental impact by ISIN. df = pd.merge(df\_ishares, df\_ei2, on='ISIN', how='outer') df.head() **Ticker CUSIP** ISIN Location Name Sector CompanyName Country 0 AAPL APPLE INC Information Technology US0378331005 **United States** NaN 37833100 NaN **MSFT** MICROSOFT CORP 594918104 US5949181045 **United States** MICROSOFT CORPORATION Information Technology **United States** 2 **AMZN** AMAZON COM INC Consumer Discretionary US0231351067 **United States** AMAZON.COM, INC. 23135106 **United States** TESLA INC Consumer Discretionary 3 **TSLA** 88160R101 US88160R1014 United States NaN NaN FACEBOOK CLASS A INC 30303M102 US30303M1027 United States FACEBOOK INCORPORATION United States FB Communication print(df.shape) print(len(df.Ticker.unique())) (3005, 8)856 mising Ticker = df.loc[df.Ticker.isna(),'ISIN'] df\_missing = df\_ei.loc[df\_ei.ISIN.isin(mising\_Ticker), ['ISIN','CompanyName','Country']] df missing.drop duplicates(inplace=True) print(f'We are unable to match {df\_missing.shape[0]} ISIN') df missing.head() We are unable to match 2136 ISIN ISIN CompanyName Country DE0005545503 1&1 DRILLISCH AG Germany KYG8875G1029 **3SBIO INC** China GI000A0F6407 888 HOLDINGS PLC Gibraltar PHY000221069 A Brown Company, Inc. 20 Philippines GB00B6XZKY75 A.G. BARR P.L.C. United Kingdom With this merge, we were able to find a 492 tickers for the 2,628 companies in the Environmental Impact dataset. We need to figure out how to map their ISIN to a Ticker. df\_missing.to\_csv('Companies\_missing\_Tickers.csv', index=False) In the meantime, we are going to export the companies that we were unable to match and continue working with the ones that have the Ticker. From the pilot stocks, we have 52 unique tickers. Let's see if all these tickers appear in the merged dataframe. tickers = stocks.Ticker.unique() x = df.loc[(df.Ticker.isin(tickers)), ['Ticker']] print(x.shape) print(f'{len(x.Ticker.unique())} Tickers') x[x.duplicated()] (55, 1)52 Tickers Ticker 540 DTE 726 MRO 802 ATO We noticed that these three tickers are duplicated. Let's see why? df ishares[df ishares.Ticker == 'ATO'] **Ticker** Name Sector **CUSIP** ISIN Location 49560105 US0495601058 United States ATO ATMOS ENERGY CORP Utilities 809 S56547813 FR0000051732 ATO ATOS Information Technology France It seems the ticker is not globally unique. We are going to filter by United States because our pilot stocks belongs to the SP500 US Stocks. In [14]: df us = df.loc[(df.Ticker.isin(tickers)) & (df.Location == 'United States'), ['Ticker']] print(f'Now, we have {len(df us.Ticker.unique())} tickers.') Now, we have 52 tickers. We know that we weren't able to match the ISIN for all the companies in the Environmental Intensity dataset. However, let's see the ones that we were able too and validate they matched correctly. found = df.loc[(df.Ticker.notna()) & (df.CompanyName.notna()),] print(len(found.Ticker.unique())) print(found.shape) 487 (495, 8)Could it be possible to have duplicate values? Let's subset for the records that have the same ticker. ticker agg = found.groupby('Ticker')[['Name']].count().sort values(by='Name', ascending=False).reset index() tickers2 = ticker agg[ticker agg.Name >= 2]['Ticker'] df3 = found[found.Ticker.isin(tickers2)].sort values(by='Ticker') df3.head(6) **Ticker** Name Sector **CUSIP** ISIN Location CompanyName Country AMERICAN AIRLINES GROUP **AMERICAN AIRLINES** United United 428 AAL Industrials 02376R102 US02376R1023 **GROUP INC** States States United United 557 AAL ANGLO AMERICAN PLC Materials GB00B1XZS820 ANGLO AMERICAN PLC Kingdom Kingdom United **AUTOMATIC DATA** Information United **AUTOMATIC DATA** 53015103 US0530151036 99 ADP PROCESSING INC Technology States PROCESSING, INC. States 866 AEROPORTS DE PARIS SA Industrials FR0010340141 France **AEROPORTS DE PARIS** France United United 271 DTE DTE ENERGY Utilities 233331107 US2333311072 DTE ENERGY COMPANY States States 540 DTE DEUTSCHE TELEKOM N AG Communication S58423591 DE0005557508 Germany DEUTSCHE TELEKOM AG Germany It happened the same thing as before, it seems the Ticker is not globally unique. However, it seems it matches correctly by ISIN. We did a little research and found that AAL in UK actually is AAL.L in Yahoo Finance. Another example, DTE- DEUTSCHE TELEKOM ticker in Yahoo Finance is DTE.DE. We decided to re-format the Ticker for the companies that are not from US. To some of them, we need to add '.' + 'the first two letters from the ISIN' to the current ticker. Let's see. ticker agg = df.groupby('Ticker')[['Name']].count().sort values(by='Name', ascending=False).reset index() tickers2 = ticker agg[ticker agg.Name >= 2]['Ticker'] df3 = df[df.Ticker.isin(tickers2)].sort values(by='Ticker') df3.head(6) **Ticker** Sector **CUSIP** ISIN Location CompanyName Name Country AMERICAN AIRLINES GROUP United **AMERICAN AIRLINES** United 428 AAL Industrials 02376R102 US02376R1023 INC States **GROUP INC** States United United 557 AAL ANGLO AMERICAN PLC Materials GB00B1XZS820 ANGLO AMERICAN PLC Kingdom Kingdom United United 735 **ADM** ADMIRAL GROUP PLC **Financials** GB00B02J6398 ADMIRAL GROUP PLC Kingdom Kingdom United 234 **ADM** ARCHER DANIELS MIDLAND **Consumer Staples** 39483102 US0394831020 NaN NaN States 866 ADP AEROPORTS DE PARIS SA Industrials FR0010340141 **AEROPORTS DE PARIS** France France AUTOMATIC DATA **AUTOMATIC DATA** Information United United US0530151036 99 **ADP** 53015103 PROCESSING INC Technology States PROCESSING, INC. States df4 df3[df3.Location != 'United States'] df4 **Ticker** Sector **CUSIP** ISIN Location **CompanyName** Name Country United United GB00B1XZS820 557 AAL ANGLO AMERICAN PLC Materials ANGLO AMERICAN PLC Kingdom Kingdom United United 735 **ADM** ADMIRAL GROUP PLC **Financials** GB00B02J6398 ADMIRAL GROUP PLC Kingdom Kingdom AEROPORTS DE PARIS 866 ADP **AEROPORTS DE PARIS SA** Industrials FR0010340141 France France Information 802 **ATO ATOS** S56547813 FR0000051732 ATOS SE France France Technology **VINCI SA** 553 DG FR0000125486 France VINCI Industrials France DEUTSCHE TELEKOM N S58423591 Germany **DEUTSCHE TELEKOM AG** 540 DTE Communication DE0005557508 Germany AG Consumer 560 EL **ESSILORLUXOTTICA SA** S72124779 FR0000121667 France **ESSILORLUXOTTICA SA** France Discretionary 731 LEG LEG IMMOBILIEN N AG Real Estate DE000LEG1110 Germany LEG IMMOBILIEN AG Germany 631 MRK **MERCK** Health Care S47418447 DE0006599905 Germany MERCK KGAA Germany **MELROSE INDUSTRIES** United **MELROSE INDUSTRIES** United 726 MRO GB00BZ1G4322 Industrials **PLC** Kingdom **PLC** Kingdom United United 563 **PRU** PRUDENTIAL PLC **Financials** S07099542 GB0007099541 PRUDENTIAL PLC Kingdom Kingdom S57059461 547 SAN **BANCO SANTANDER SA Financials** ES0113900J37 **BANCO SANTANDER SA** Spain Spain 516 SAN SANOFI SA SANOFI S.A. Health Care S56717358 FR0000120578 France France **TELENOR GROUP ASA** 739 TEL **TELENOR** S47324959 NO0010063308 Communication Norway Norway United United 586 TSCO **TESCO PLC** Consumer Staples S08847097 GB0008847096 **TESCO PLC** Kingdom Kingdom ticker = df4.Ticker.tolist() yahoo ticker = ['AAL.L','ADM.L','ADP.PA','ATO.PA','DG.PA','DTE.DE','EL.PA','LEG.DE','MRK.DE','MRO.L','PRU.L','S df = df.sort values(by='Ticker') df.loc[(df.Ticker.isin(ticker)) & (df.Location != 'United States'), 'Ticker'] = yahoo ticker **Ticker CUSIP** ISIN Location CompanyName Country Name Sector 867 **IBERDROLA SA** Utilities ES0144583236 Spain NaN NaN 711 1COV **COVESTRO AG** Materials DE0006062144 Germany **COVESTRO AG** Germany Health AGILENT TECHNOLOGIES INC AGILENT TECHNOLOGIES, INC. 170 Α 00846U101 US00846U1016 United States **United States** Care AMERICAN AIRLINES GROUP AMERICAN AIRLINES GROUP 428 **AAL** Industrials 02376R102 US02376R1023 **United States United States** United United 557 AAL.L ANGLO AMERICAN PLC Materials GB00B1XZS820 ANGLO AMERICAN PLC Kingdom Kingdom print(len(df.Ticker.unique())) 870 Now, we have 870 unique tickers that we can send to yahoo finance and get the company description. Below, we will create a tickers list. Also, we want to export the merged dataset to do some analysis onwards. tickers = df.Ticker.unique().tolist() df.to csv('df merged.csv', index=False) Collecting companies data in Reuters website Reuters, the news and media division of Thomson Reuters, is the world's largest multimedia news provider, reaching billions of people worldwide every day. First, we are going to add '.OQ' to each ticker. T In [24]: tickers US = df.loc[(df.Location == 'United States') & (df.Ticker.notna()), 'Ticker'].unique().tolist() tickers nonUS = df.loc[(df.Location != 'United States') & (df.Ticker.notna()), 'Ticker'].unique().tolist() reuter\_ticker = [] for ticket in tickers\_US: ticket = str(ticket) + '.0Q' reuter\_ticker.append(ticket) # Create a loop to store URLs of all stocks' description page URL = []DES = []comp desc = {} for i in reuter ticker: url ='https://www.reuters.com/companies/'+i URL.append(url) page = requests.get(url) # visits the URL htmldata = BeautifulSoup(page.content, 'html.parser') Business Description = htmldata.find('p', {'class':'TextLabel text-label 3oCVw TextLabel black 2FN-Z 5 if Business Description is not None: DES.append (Business Description.text) comp desc[i] = [Business Description.text] comp desc[i] = 'Not exists' exists = 0do not = 0companies\_desc = [] for val in comp\_desc.values(): if val == 'Not exists': do\_not+=1 else: exists**+=**1 companies\_desc = [] print(do\_not) print(exists) 362 143 We were able to find 143 companies descriptions which we will use later on this notebook. Collect companies descriptions from Yahoo Finance In this section, we will grab the company description directly from the Yahoo Finance website for the 870 companies We were unable to send 870 tickers to Yahoo Finance. It seems their page crashed and wasn't allowing web scraping # Create a loop to store URLs of all stocks' description page URL = []DES = []comp\_desc = {} for i in tickers[:2]: i = str(i)url ='https://finance.yahoo.com/quote/'+i+'/profile' URL.append(url) page = requests.get(url) # visits the URL htmldata = BeautifulSoup(page.content, 'html.parser') Business\_Description = htmldata.find('p', {'class':'Mt(15px) Lh(1.6)'}) # finds the business description par DES.append(Business\_Description.text) comp\_desc[i] = [Business\_Description.text] df comp desc = pd.DataFrame.from dict(comp desc, orient='index', columns = ['Description']) df comp desc.reset index(inplace = True) df comp desc.rename(columns = {'index':'Ticker'}, inplace=True) df comp desc.head() df\_comp\_desc.to\_csv('Companydescription.csv',index=False) Start working with the DistilBERT First, we are going to create a new dataframe using the 202 companies\_desc.csv. This is a merged file from the companies descriptions that we collected in Sprint # 6 and the descriptions from Reuters collected earlier in the notebook. stock des=pd.read csv('202 companies des.csv') #Yahoo descriptions df ei 2019 = df ei[df ei.Year == 2019] df ei merged = pd.merge(df,df ei 2019, on='ISIN',how='inner') df ei merged = df ei merged[['Ticker', 'Env intensity']] df ei merged.dropna(inplace=True) df ei merged.head() Ticker Env\_intensity 1COV -0.14230.0037 2 AAL -0.2853 3 AAL.L -0.1605 AAP -0.0145 Now, we will merge this dataset with the descriptions so we have the necessary fields to start running the DistilBERT df\_final = pd.merge(df\_ei\_merged, stock\_des, on ='Ticker', how='inner') df\_final = df\_final.iloc[1:,:] #Removing AAL Dupplicate df2 = df\_final.copy().reset\_index() print(f'We are going to be able to send {df2.shape[0]} companies descriptions to the DistilBERT model') df2.head() We are going to be able to send 74 companies descriptions to the DistilBERT model index Ticker Env\_intensity description 0 AAL -0.2853 American Airlines Group Inc. is a holding comp... 2 AC-0.1211 Associated Capital Group, Inc. provides invest... 2 3 **ADSK** 0.0093 Autodesk, Inc. is a design software and servic... 3 AEE -1.3412 Ameren Corporation, together with its subsidia... 4 5 **AEP** -1.6381 American Electric Power Company, Inc., an elec... We created a binary variable that is 1 if the Environmental Intensity (Sales) is above its median and 0 otherwise. This is the dependent variable (label) that we'll try to predict. df2['HIGH\_EI'] = (df2['Env\_intensity'].gt(df2['Env\_intensity'].median())).astype(int) Load a pre-trained BERT model. model class, tokenizer class, pretrained weights = (ppb.DistilBertModel, ppb.DistilBertTokenizer, 'distilbert-k # Load pretrained model/tokenizer tokenizer = tokenizer class.from pretrained(pretrained weights) model = model class.from pretrained(pretrained weights) Some weights of the model checkpoint at distilbert-base-uncased were not used when initializing DistilBertMode 1: ['vocab\_projector.weight', 'vocab\_layer\_norm.weight', 'vocab\_projector.bias', 'vocab\_transform.bias', 'vocab \_layer\_norm.bias', 'vocab\_transform.weight'] - This IS expected if you are initializing DistilBertModel from the checkpoint of a model trained on another ta sk or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTrain - This IS NOT expected if you are initializing DistilBertModel from the checkpoint of a model that you expect t o be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model). Tokenize the textual data for DistilBERT which in our case would be the companies description tokenized = df2['description'].apply((lambda x: tokenizer.encode(x, add special tokens=True))) Pad all lists of tokenized values to the same size In [34]: max len = 0for i in tokenized.values: if len(i) > max\_len:  $\max$  len = len(i) padded = np.array([i + [0]\*(max len-len(i)) for i in tokenized.values]) np.array(padded).shape Out[35]: (74, 463) Create attention mask variable for BERT to ignore (mask) the padding when it's processing its input. attention\_mask = np.where(padded != 0, 1, 0) attention\_mask.shape Out[36]: (74, 463) We run the pretrained DistilBERT model on the prepared predictor and keep the result in last\_hidden\_states variable. input\_ids = torch.tensor(padded) attention\_mask = torch.tensor(attention\_mask) with torch.no\_grad(): last\_hidden\_states = model(input\_ids, attention\_mask=attention\_mask) Keep the first layer of the hidden states and assign the outcome variable to labels. features = last\_hidden\_states[0][:,0,:].numpy() labels = df2['HIGH\_EI'] Split the data in train and test subsets, train the Logistic Regression on train set and evaluate its accuracy on the test set. train\_features, test\_features, train\_labels, test\_labels = train\_test\_split(features, labels,test\_size=0.2,rand lr clf = LogisticRegression(max iter=5000) lr\_clf.fit(train\_features, train\_labels) print(lr\_clf.score(test\_features, test\_labels)) 0.8 Our model can 80% accuratly capture whether the company is high or low environmental intensity. In [40]: test labels Out[40]: 4 63 1 18 1 0 0 28 73 0 10 0 34 1 12 55 1 65 0 31 0 9 45 1 0 Name: HIGH EI, dtype: int64 Conclusion We had difficulties trying to map the 4,270 companies to their corresponding Ticker. Also, we discovered that sometimes 'Tickers' are not globally unique. Therefore, it can be difficult to find an accurate way to map the values. Other difficulty was the technical issue with Yahoo Finance page. Their page wasn't allowing web scraping. Hence, we were unable to gather the 870 companies descriptions from the companies Tickers that we were able to map. We decided to collect some descriptions from the Reuters.com and, together with other companies description collected in earlier sprints, we were able to run the DistilBERT model. The dependent variable, the one that we want to predict/classify, is High Environmental Intensity and after training the model with 80% of the observations, it was able to accurate capture whether the company is high or low environmental intensity 80% of the times.

This is the end of our analysis.