

今天主要进行的是**第三步**，第一、二步为将最近的代码工作串起来而写上

## 1. 获取要处理的基本数据

执行程序名为 公司全称和公司缩写建立上下游关系.py 程序，

```
import xlwt

# global firm names 和 us names 的 Global Company Key 和 Company Name 提取
global_workbk = xlrd.open_workbook(r'C:\Users\jc\Documents\Pydata'+
                                   '\Database Table\global firm names.xlsx')
```

上面最后一行在读取 *global firm names.xlsx* 文件，使用时请改成**对应文件保存的位置**

```
global_worksh = global_workbk.sheet_by_name('0x77igavdumz8vul')
global_cpnames = global_worksh.col_values(colx = 7, start_rowx = 1)
global_cpkey = global_worksh.col_values(colx = 0, start_rowx = 1)
# 组成列表
global_namekey = list(zip(global_cpnames, global_cpkey))

us_workbk = xlrd.open_workbook(r'C:\Users\jc\Documents\Pydata'+
                               +'\Database Table\us names.xlsx')
```

上面最后一行在读取 *us\_names.xlsx* 文件，使用时请改成**对应文件保存的位置**

```
us_worksh = us_workbk .sheet_by_name('76aqys7wh9axjpme')
us_cpnames = us_worksh.col_values(colx = 9, start_rowx = 1)
us_cpkey = us_worksh.col_values(colx = 0, start_rowx = 1)
# 组成列表
us_namekey = list(zip(us_cpnames,us_cpkey))
# 将customer表中的Customer Name列和Global Company Key 提取出
customer_workbk = xlrd.open_workbook(r'C:\Users\jc\Documents\Pydata'+
                                       '\Database Table\customer.xlsx')
```

上面最后一行在读取*customer.xlsx*文件位置，使用时请改成**对应文件保存的位置**

```

customer_worksh = customer_workbk.sheet_by_name('vozkv0ioajsw5wov')
customer_downstream = customer_worksh.col_values(colx = 2, start_rowx = 1)
customer_upkey = customer_worksh.col_values(colx = 0, start_rowx = 1)
# 组成列表
customer_up_down = list(zip(customer_upkey,customer_downstream))

# global firm names 和 us names 的Global Company Key 和 Company Name列表合并后去重
Allnamekey_lst = list(set( global_namekey+us_namekey))
#对customer_up_down列表也去重
customer_up_down = list(set(customer_up_down))
#按照Global Company Key进行排序
Allnamekey_lst.sort(key=lambda x:x[1])

# 将Global Company Key 和 Company Name 写入Allcompany表中
Allcompany = xlwt.Workbook()
Allcompany_sheet = Allcompany.add_sheet('sheet1')
name_list = ['Global Company Key','Company Name(upstream)','Downstream']
for i in name_list:
    Allcompany_sheet.write(0, name_list.index(i), i)

for namekey in Allnamekey_lst:
    Allcompany_sheet.write(Allnamekey_lst.index(namekey)+1,0,namekey[1])
    Allcompany_sheet.write(Allnamekey_lst.index(namekey)+1,1,namekey[0])

cust_num = [[i[1],0] for i in Allnamekey_lst]
#比较 upstream key(customer_up_down[i][0])和global company key(Allnamekey_lst[j][1])
for upkey_down in customer_up_down:
    for name_key in Allnamekey_lst:
        if upkey_down[0] == name_key[1]:
            Allcompany_sheet.write(Allnamekey_lst.index(name_key)+1
                                   ,2+cust_num[Allnamekey_lst.index(name_key)][1]
                                   ,upkey_down[1])
            cust_num[Allnamekey_lst.index(name_key)][1] += 1

# 保存文件
Allcompany.save('Allcompany.xlsx')

```

最后一行是保存*Allcompany.xlsx*文件的位置，可以不修改，默认保存在同py文件的文件夹下。



```

        +'\Allcompany.xlsx')
Allcompany_sheet = Allcompany.sheet_by_name('sheet1')
Allcpkey = Allcompany_sheet.col_values(colx = 0, start_rowx = 1)
Allcpname = Allcompany_sheet.col_values(colx = 1, start_rowx = 1)

```

创建表Allcpdata,全称All Company Data表格，用于存储全称公司的：

Global Company key, Full name, company type , country , Cleande Full Name

其中 Cleande Full Name 表示去除了公司后缀的全称,用于后续与缩写的匹配  
country 用于存储公司可能所在的国家，有些无法识别出就空着

### 1. 创建All company match表格的表头

```

# 1.创建All Company Data表格，写好表头
Allcpdata_wb = xlwt.Workbook()
Allcpdata = Allcpdata_wb.add_sheet('sheet1')
name_list = ['Global Company key', 'Full name', 'company type' , 'country'
, 'Cleaned Full Name']

for i in range(len(name_list)):
    Allcpdata.write(0, i , name_list[i])

```

### 2. 将Global Company key写进表中 第i+1行,第0列

```

for i in range(len(Allcpkey)):
    Allcpdata.write(i+1, 0 , Allcpkey[i])

```

### 3. 将Full name写进表中 第i+1行，第1列

```

for i in range(len(Allcpname)):
    Allcpdata.write(i+1, 1, Allcpname[i])

```

### 4. 将company type写入表中 第i+1行，第2列, 用 cleanco中typesources()和match()判断

```

for i in range(len(Allcpname)):
    Allcpdata.write(i+1, 2, matches(Allcpname[i], typesources()))

```

5. 将country 写入表中 第i+1行, 第3列, 用cleanco中match()和countrysources()判断

```
for i in range(len(Allcpname)):
    Allcpdata.write(i+1, 3, matches(Allcpname[i], countrysources()))
```

6. 清理不必要的后缀与符号, 并转换为小写字符 写入表中, 第i+1行, 第4列, 用clean\_name()判断。

```
for i in range(len(Allcpname)):
    Allcpdata.write(i+1, 4, (clean_name(Allcpname[i]).lower()))
```

关于 clean\_name 方法解释如下:

1. 用cleanco中的basename删除后缀  
(A.S./ N.V./ SA/SG S.A.与这一步共同进行单独处理)  
(-CL A,-CL B,-CL C ,-CL I)单独处理
2. 删除未识别出的后缀  
ETF , A/S , -ADR , INC-OLD,group,
3. 再用cleanno中的basename删除一次, 确保删除干净了

以下是clean\_name方法, 基本能将company name中所有后缀与国家名都删除, 并不丢失此外的其他字符。

```
def clean_name(string):
    """
    1. 用cleanco中的basename删除后缀
    (A.S./ N.V./ SA/SG S.A.与这一步共同进行单独处理)
    (-CL A,-CL B,-CL C ,-CL I)单独处理
    2. 删除未识别出的后缀
    ETF , A/S , -ADR , INC-OLD,group,
    3. 再用cleanno中的basename删除一次, 确保删除干净了
    """
    return 公司清理后缀后的名称
    """
    #1.
    # 截取最后5个判断是否有-CL_?类后缀, 如果有则删去
    str_suffix0 = string[:len(string)-6:-1][::-1].replace(' ','')
    if '-CL' in str_suffix0:
        string = string[:string.find('-')]
    string = re.sub(r'[\W]',' ',string) # 首先删去字符影响, 替换为空格
```

```
string = cleanco.basename(string) # 用cleanco里自带的basename做第一次后缀清除
clean_suffix0 = ['A/S','N.V.','SA/AG','S.A.']
for cs in clean_suffix0:
    if cs in string:
        string = string[:string.rfind(cs[0])]
#2.
# 进行第二次后缀清理
# 截取最后 5 个字符,并删除其中的空格 为什么是 5 个（因为根据观察，自动清除的后缀里，最长的是GROUP，所以取最后5个判断不会有漏）
str_suffix = string[:len(string)-6:-1][::-1].replace(' ','')
clean_suffix = ['CO','ETF','AS','ADR','SA','AG','OLD','GROUP'] # 判断是否包含这些后缀
for cs in clean_suffix:
    if cs in str_suffix:
        string = string[:string.rfind(cs[0])]
#3.
#用basename进行第三次后缀处理
string = cleanco.basename(string)
return string
```

最后进行文件的保存：

```
Allcpdata_wb.save('All Company Data.xlsx')
```

数据处理结果如下：

A	B	C	D	E	F
57186751335	IPD GROUP LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	vid	
57187751336	SHAPE AUSTRALIA CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	shape australia	
57188751340	LATENT VIEW ANALYTICS LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	latent view analytics	
57189751341	FINATEXT HOLDINGS LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	finatext holdings	
57190751343	KOHOKU KOGYO CO. LTD	CorporationLimited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	kohoku kogyo	
57191751346	POLITREX S.A.	CorporationLimited Liability Company	AlgeriaindonesiaChileColumbiaDominican RepublicEcuadorGuatemalaLuxembourgMexicoPeruPolandPortugalRomaniaSpain	politrex	
57192751345	MARKS ELECTRICAL GRP PLC	Limited Liability Company	NigeriaUnited Kingdom	marks electrical grp	
57193751350	XPON TECHNOLOGIES GROUP LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	xpon technologies	
57194751351	HANGZHOU ZHENGQIANG CORP LTD	CorporationLimited	PhilippinesUnited States of AmericaHong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	hangzhou zhengqiang	
57195751355	HEALTHLEAD PUBLIC COMPANY	Corporation	United States of America	healthlead public	
57196751356	DUELL OYJ	Limited Liability Company	Finland	duell	
57197751358	NICO RESOURCES LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	nico resources	
57198751362	RENEWABLE JAPAN CO LTD	Limited		renewable japan	
57199751366	GUANGDONG LIFESTRONG PHARMAC			guangdong lifestrong pharmac	
57200751367	IMM ROBOT & AUTOMATION CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	imm robot automation	
57201751368	CAM RANH PORT JOINT STOCK CO	Corporation	United States of America	cam ranh port joint stock	
57202751370	SICHUAN DISCOVERY DREAM			sichuan discovery dream	
57203751372	OMNIPOTENT INDUSTRIES LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	omnipotent industries	
57204751373	GO FASHION LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	go fashion	
57205751374	ARISTON HOLDING N.V		Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	ariston holding n v	
57206751384	ALFONSIÑO S P A	Corporation	ChileItaly	alfonsoño s p a	
57207751385	INNER MONGOLIA ANKUA			inner mongolia ankua	
57208751387	ISBIR SENTETIK DOKUMA SANAYI			isbir sentetik dokuma sanayi	
57209751389	T.R.V. RUBBER PRODUCTS			t r v rubber products	
57210751391	PG ILETISIM VE MEDYA			pg iletisim v e medya	
57211751392	NIVIKIA FASTIGHETER AB PUBL	Limited Liability Company	LithuaniaSwedenSwitzerland	nivikia fastigheter ab publ	
57212751393	SANWAY UXA INDUSTRY CORPORATI			sanway uxa industry corporati	
57213751398	EWINGCARDS INC	Corporation	PhilippinesUnited States of AmericaUnited States of America	ewingcards	
57214751399	GROOOBER CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	grooobar	
57215751401	HYBRID TECHNOLOGIES CO LTD	LimitedCorporation	PhilippinesUnited States of AmericaUnited States of America	hybrid technologies	
57216751404	BEIJING FJR OPTOELECTRONIC		Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	beijing fjr optoelectronic	
57217751406	CORLUS INC	Corporation	PhilippinesUnited States of AmericaUnited States of America	corlus	
57218751407	TURIN CLOUD TECHNOLOGY			turin cloud technology	
57219751408	ENERAQUA TECHNOLOGIES PLC	Limited Liability Company	NigeriaUnited Kingdom	eneraqua technologies	
57220751416	GRUPE BERHEM SA	Limited Liability Company	ChileDominican RepublicEcuadorGuatemalaItalyPeruSwitzerland	grupe berhem	
57221751417	CLOUD MUSIC INC	Limited	PhilippinesUnited States of AmericaUnited States of America	cloud music	
57222751420	NIFTY LIFESTYLE CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	nifty lifestyle	
57223751421	ASIAQUEST CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	asiaquest	
57224751426	UNIVERSAL NETWORK SYSTEMS			universal network systems	
57225751427	ADVANCED GRAPHENE PRODUCTS			advanced graphene products	
57226751434	ATTURIA LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	atturia	
57227751436	DOSILUCON CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	dosilucion	
57228751436	HSEND BUSINESS SOLUTIONS PL			hsend business solutions pl	
57229751437	SUNMOW HOLDING BERHAD			sunmow holding berhad	
57230751439	ANHUI PROVINCIAL ARCHITECTUR			anhui provincial architectur	
57231751443	TARSONS PRODUCTS LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	tarsons products ltd	
57232751446	KLARABO SVERIGE AB (PUBL)	Limited Liability Company	LithuaniaSwedenSwitzerland	klarabo sverige ab publ	
57233751447	SOLID FORSAKRINGSAKTIEBOLAG			solid forsakringsaktiebolag	
57234751448	DANWA HOUSE LOGISTICS TRUST			danwa house logistics trust	
57235751454	SHANGHAI YI-ZHONG PHARMA CO	Corporation	United States of America	shanghai yihong pharma	
57236751456	ZHANG XIAOQUAN INC	Corporation	PhilippinesUnited States of AmericaUnited States of America	zhang xiaoquan	
57237751456	BEFOREPAY GROUP LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	beforepay	
57238751457	SECURE INC	Corporation	PhilippinesUnited States of AmericaUnited States of America	secure	
57239751463	HEGEM AUTOLINERS ASA	Limited Liability Company	Norway	hegem autoliners	
57240751465	LAMOR CORPORATION OYJ	CorporationLimited Liability Company	United States of AmericaFinland	lamor	
57241751478	ASHTEAD TECHNOLOGY HOLDINGS			ashthead technology holdings	
57242751485	HANGSHU TONGLIN			hangshu tonglin	
57243751481	IVECO GROUP N.V			iveco group n v	
57244751482	BIRDODG TECHNOLOGY LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	birdodg technology	
57245751487	HEALTHRACON PLC	Limited Liability Company	NigeriaUnited Kingdom	healthracon	
57246751486	HANGZHOU BIO SINCERITY PHARM			hangzhou bio sincerity pharm	
57247751500	ZJME ENVIRONMENTAL ENERGY CO	Corporation	United States of America	zjme environmental energy	

### 3.对获取到的缩写进行预处理

执行`公司缩写的预处理.py`文件,基本步骤与第二步差不多

第一步,从`customer.xlsx`中读取数据,`company\_abbr`用于保存公司缩写,`company\_upkey`用于保存用于识别缩写对应的上游公司的特殊数字串

```
customer_wb = xlrd.open_workbook(r'C:\Users\jc\Documents\Pydata'+
                                  '\Database Table\customer.xlsx')
customer = customer_wb.sheet_by_name('vozkv0ioajsw5wov')
company_abbr = customer.col_values(colx = 2, start_rowx = 1)
company_upkey = customer.col_values(colx = 0, start_rowx = 1)
# 组成列表
abbr_upkey_zip = list(set((zip(company_abbr,company_upkey))))
# 第i组:  abbr_upkey_zip[i][0]公司缩写  abbr_upkey_zip[i][1]上游公司key
```

第二步,创建表Abbrdata,全称Abbreviation Data表格,用于存储缩写公司的:  
Global Company key(upstream), Abbreviation, company type , country , Cleaned Abbreviation  
其中 Cleaned Abbreviation 表示去除了公司后缀的缩写,用于与公司全称的匹配  
country 用于存储公司可能所在的国家,有些无法识别出就空着。

1.创建Abbreviation Data表格,写好表头

```
Abbrdata_wb = xlwt.Workbook()
Abbrdata = Abbrdata_wb.add_sheet('sheet1')
name_list = ['Global Company key(upstream)', 'Abbreviation',
             'company type' , 'country' , 'Cleaned Abbreviation']

for i in range(len(name_list)):
    Abbrdata.write(0, i , name_list[i])
```

2.将Global Company key(upstream)写进表中 第i+1行,第0列

```
for i in range(len(abbr_upkey_zip)):
    Abbrdata.write(i+1, 0 , abbr_upkey_zip[i][1])
```

3.将Abbreviation写进表中 第i+1行, 第1列

```
for i in range(len(abbr_upkey_zip)):
    Abbrdata.write(i+1, 1, abbr_upkey_zip[i][0])
```

4.将company type写入表中 第i+1行, 第2列, 用 cleanco中typesources()和match()判断

```
for i in range(len(abbr_upkey_zip)):
    Abbrdata.write(i+1, 2, matches(abbr_upkey_zip[i][0], typesources()))
```

5.将country 写入表中 第i+1行, 第3列, 用cleanco中match()和countrysources()判断

```
for i in range(len(abbr_upkey_zip)):
    Abbrdata.write(i+1, 3, matches(abbr_upkey_zip[i][0], countrysources()))
```

6.清理不必要的后缀与符号, 并转换为小写字符 写入表中, 第i+1行, 第4列, 用 cleanco.clean\_name()判断

```
for i in range(len(abbr_upkey_zip)):
    Abbrdata.write(i+1, 4, (clean_name(abbr_upkey_zip[i][0]).lower()))

Abbrdata_wb.save('Abbreviation Data.xlsx')
```

最后保存数据在 Abbreviation

---

第三步中对缩写的数据清洗有点瓶颈  
面临的问题:

1. 多个缩写配对同一个公司全称
2. 多个公司全称匹配一个缩写

以上进一步判断是否该通过人工进行, 讨论之后再决定

其他工作:

1. 相关的 Allcompany.xlsx , All Company Data.xlsx , Abbreviation Data.xlsx 文件已上传至OneDrive--Liuchenxin文件夹
2. 阅读《Using Natural Language Processing for Supply Chain Mapping》论文, 阅读部分论文的笔记 (NOTE) Schöpper, H., Kersten, W., Using Natural Language Processing for Supply Chain Mapping 已上传至OneDrive---Liuchenxin文件夹



