

今天主要进行的是第二步，第一步是为了将最近的代码工作串起来而写上的

## 1. 执行程序名为 公司名称识别并建立上下游关系.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。

*Allcomany.xlsx*文件里保存的有

第一列：Global Company Key 每个公司特有的一串数字

第二列：Company Name(upstream) 上游公司的全称

第三列以后：Downstream 与上有公司对应的，下游公司的缩写

对数据进行了哪些处理：

1. 里面上游公司的名称没有重复，都只出现一次
2. 数据都是按照 Global Company Key 公司特有的关键字 从小到大排列

在我的电脑上代码执行大约需要25min

执行后结果如下：

	A	B	C	D	E
1	Global Company Key	Company Name(upstream)	Downstream		
2	001004	AAR CORP	U.S. Government	Not Reported	North America
3	001013	ADC TELECOMMUNICATIONS INC			
4	001019	AFA PROTECTIVE SYSTEMS INC	Not Reported		
5	001045	AMERICAN AIRLINES GROUP INC			
6	001050	CECO ENVIRONMENTAL CORP	Foreign		
7	001082	ASA GOLD AND PRECIOUS METALS			
8	001072	AVX CORP	Not Reported	Electronic Distributors	
9	001075	PINNACLE WEST CAPITAL CORP	Wholesale energy sales	Transmission services for others	Retail residential electric service
10	001076	PROG HOLDINGS INC	Home Exercise and Home Improvement	Other	Furniture and Mattresses
11	001078	ABBOTT LABORATORIES	International	Other Emerging Markets	United States
12	001082	SERVIDYNE INC			
13	001084	WORLDS INC			
14	001094	ACETO CORP	McKesson Corp	Europe	AmerisourceBergen Corp
15	001096	MORGUARD CORP			
16	001097	ACMAT CORP -CL A			
17	001104	ACME UNITED CORP	E-commerce	International	Not Reported
18	001117	BK TECHNOLOGIES CORP	Public Safety	Business and Industrial	Industrial
19	001119	ADAMS DIVERSIFIED EQUITY FD			
20	001121	ADAMS RESOURCES & ENERGY INC	Not Reported		
21	001161	ADVANCED MICRO DEVICES	Sony Corp	HEWLETT-PACKARD CO	Not Reported
22	001166	ASM INTERNATIONAL NV	7 Customers	China	Taiwan
23	001173	AEROSONIC CORP			
24	001177	AETNA INC	U.S. Federal Government	Medicaid	Foreign
25	001186	AGNICO EAGLE MINES LTD	4 Customers	Not Reported	
26	001209	AIR PRODUCTS & CHEMICALS INC	Sale of equipment	Merchant	Outside the United States
27	001210	AIR T INC	International	Federal Express Corp	United States
28	001224	SPIRE ALABAMA INC	Residential	Commercial and Industrial	Other Customer
29	001225	ALABAMA POWER CO	Wholesale	Residential-Retail	Other
30	001228	ALANCO TECHNOLOGIES INC	3 Customers		
31	001230	ALASKA AIR GROUP INC	Direct to customer	Reservation Call Centers	Traditional Agencies
32	001234	ATRION CORP	Not Reported	Outside the United States	

## 2. 对获取到的全称进行预处理:

执行 公司全名的预处理.py 文件

从第一步中得到的Allcompany.xlsx中获取所有公司的Allcpkey(每个公司特有的一串数字) 和 Allcpname (公司全名)

注：第一行路径对应上面第一部分代码中生成的Allcompany文件位置

```
...
```

```
从Allcompany中读取数据
```

```
Allcpkey 用于保存公司缩写
```

```
Allcpname用于保存公司全称
```

```
...
```

```
Allcompany = xlrd.open_workbook(r'C:\Users\jc\Documents\大学'  
                                + '\\0大三其他\\_金融数据挖掘科研课题'  
                                + '\\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，并不丢失其他的其他字符。

```
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
51186751335	IPD GROUP LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	vid	
51187751336	SHAPE AUSTRALIA CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	shape australia	
51188751340	LATENT VIEW ANALYTICS LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	latent view analytics	
51189751341	FINATOTX HOLDINGS LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	finatext holdings	
51190751343	KOHOKU KODYO CO. LTD	CorporationLimited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	kohoku kogyo	
51191751346	POLITRES S.A.	CorporationLimited Liability Company	AlgeriaindonesiaColumbiaDominican RepublicEcuadorGuatemalaLuxembourgMexicoPeruPolandPortugalRomaniaSpain	politres	
51192751348	MARICO ELECTRICAL GRP PLC	Limited Liability Company	NigeriaUnited Kingdom	marico electrical grp	
51193751350	XPON TECHNOLOGIES GROUP LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	xpon technologies	
51194751351	HANGZHOU ZHENQIANG CORP LTD	CorporationLimited	PhilippinesUnited States of AmericaHong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	hangzhou zhenqiang	
51195751355	HEALTHLEAD PUBLIC COMPANY	Corporation	United States of America	healthlead public	
51196751356	DUELL OYJ	Limited Liability Company	Finland	duell	
51197751358	NICO RESOURCES LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	nico resources	
51198751362	RENEWABLE JAPAN CO LTD.	Limited		renewable japan	
51199751366	GUANGDONG LIFESTRONG PHARMAC	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	guangdong lifestrong pharmaco	
51200751367	IMM ROBOT & AUTOMATION CO LTD	Corporation	United States of America	imm robot automation	
51201751368	CAM RANH PORT JOINT STOCK CO	Corporation	United States of America	cam ranh port joint stock	
51202751370	SICHUAN DISCOVERY DREAM	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	sichuan discovery dream	
51203751372	OMNIPOTENT INDUSTRIES LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	omnipotent industries	
51204751373	GO FASHION LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	go fashion	
51205751374	ARISTON HOLDING N.V.	Corporation	ChileItaly	ariston holding n v	
51206751384	ALFONSIKO S.P.A.			alfonsiko s p a	
51207751385	INNER MONGOLIA KHUKHA			inner mongolia khukha	
51208751387	ISBIR SENTETIK DOKUMA SANAYI			isbir sentetik dokuma sanayi	
51209751389	T.R.V. RUBBER PRODUCTS			t r v rubber products	
51210751391	PO ILETISIA VE MEDIA			po iletsia v e media	
51211751392	NIVIKIA FASTIGHETER AB PUBL	Limited Liability Company	LithuaniaSwedenSwitzerland	nivikia fastigheter ab publ	
51212751393	SANWAY LUXA INDUSTRY CORPORATI	Corporation	PhilippinesUnited States of AmericaUnited States of America	sanway luxa industry corporati	
51213751398	ENAWICARDS INC	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	enawicards	
51214751399	OROOOBER CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	orooober	
51215751401	HYBRID TECHNOLOGIES CO LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	hybrid technologies	
51216751404	BEIJING FIR OPTOELECTRONIC	Corporation	PhilippinesUnited States of AmericaUnited States of America	beijing fir optoelectronic	
51217751406	COPULS INC	Limited	NigeriaUnited Kingdom	copuls	
51218751407	ENERAQUE TECHNOLOGY	Limited Liability Company	ChileDominican RepublicEcuadorGuatemalaItalyPeruSwitzerland	eneraque technologies	
51219751416	GROUPE BERHEM SA	Limited Liability Company	PhilippinesUnited States of AmericaUnited States of America	groupe berhem	
51220751417	CLOUD MUSIC INC	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	cloud music	
51222751420	NIFTY LIFESTYLE CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	nifty lifestyle	
51223751421	ASIAQUEST CO LTD	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	asiaquest	
51224751426	UNIVERSAL NETWORK SYSTEMS	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	universal network systems	
51225751427	ADVANCED GRAPHENE PRODUCTS	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	advanced graphene products	
51226751434	ATTURRA LIMITED	LimitedCorporation	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	atturra	
51227751435	DOSILICON CO LTD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	dosilicon	
51228751436	HSENO BUSINESS SOLUTIONS PL	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	hseno business solutions pl	
51229751437	SUNMOW HOLDING BERHAD	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	sunmow holding berhad	
51230751439	ANHUI PROVINCIAL ARCHITECTUR	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of AmericaUnited States of America	anhui prov incial architectur	
51231751443	TARSONS PRODUCTS LTD	Limited Liability Company	LithuaniaSwedenSwitzerland	tarsons products	
51232751445	KARABO SVERIGE AB (PUBL)	Limited	United States of America	karabo sverige ab publ	
51233751447	SOLID FORSAKRINGSAGTIKTEBOLAG	Corporation	United States of America	solid forsakringsaktiebol	
51234751448	DANWA HOUSE LOGISTICS TRUST	Corporation	United States of America	danwa house logistics trust	
51235751454	SHANGHAI YIZHONG PHARMA CO	Corporation	United States of America	shanghai yizhong pharma	
51236751455	ZHANG XIAOQUAN INC	Corporation	PhilippinesUnited States of AmericaUnited States of America	zhang xiaoquan	
51237751456	BEFOREMY GROUP LIMITED	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	beforemy	
51238751457	SECURE INC	Corporation	PhilippinesUnited States of AmericaUnited States of America	secure	
51239751463	HCEGEN AUTOMENERS ASA	Limited Liability Company	Norway	hcegen autolmers	
51240751465	LAMOR CORPORATION OYJ	CorporationLimited Liability Company	United States of AmericaFinland	lamor	
51241751478	ASHTEAD TECHNOLOGY HOLDINGS	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	ashtead technology holdings	
51242751480	HANGSHU TONGJIN	Limited	NigeriaUnited Kingdom	hangshu tongjin	
51243751481	IVECO GROUP N.V.	Limited	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	iveco group n v	
51244751482	BIRDROD TECHNOLOGY LIMITED	Limited Liability Company	Hong KongIsraelNew ZealandPakistanUnited KingdomUnited States of America	birdrod technology	
51245751487	HEALTHBRACON P.L.C.	Limited	NigeriaUnited Kingdom	healthbrac	
51246751488	HANGZHOU BIO SINCERITY PHARM	Corporation	United States of America	hangzhou bio sincerity pharm	
51247751500	ZJM ENVIRONMENTAL ENERGY CO	Corporation	United States of America	zjm environmental energy	

第三步是3. 对Customer name进行数据预处理

另外:

1. 相关的 Allcompany.xlsx 与 All Company Data.xlsx 文件都已上传至OneLive--liuchenxin文件夹
2. 方翔提出"识别customer name中是否含country" 的问题, 已向他建议pycountry中自带的模糊查询方法