

Investment Assignment

by Ayush Mandowara

Imports

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Loading Data to Dataframes

```
In [2]: companies = pd.read_csv('companies.csv')
rounds2 = pd.read_csv('rounds2.csv')
```

Understanding the Data

```
In [3]: companies.head()

Out[3]:
```

	permalink	name	homepage_url	category_list	status	country_code	state_code	region	city	founded_at
0	/Organization/-fame	fame	http://fifame.com	Media	operating	IND	16	Mumbai	Mumbai	2000-01-01
1	/Organization/-Quarter	Quarter	http://www.quarter.com	Application Platforms Real Time Social Network...	operating	USA	DE	DE - Other	Delaware City	2014-01-01
2	/Organization/-The One-Of-Them-Inc	THE ONE-Of-THem-INc	http://oneothem.jp	Apps Games Mobile	operating	NaN	NaN	NaN	NaN	2014-01-01
3	/Organization/O-6-com	O-6-com	http://www.O-6.com	Curated Web	operating	CHN	22	Beijing	Beijing	2014-01-01
4	/Organization/O04-technologies	O04 Technologies	http://O04gmhb.de/en/O04-interact	Software	operating	USA	IL	Springfield, Illinois	Champaign	2014-01-01

```
In [4]: rounds2.head()
```

```
Out[4]:
```

	company_permalink	funding_round_permalink	funding_round_type	funding_round_code	funded_at	raised_amount_usd
0	/organization/-fame	round/9a01d05181ba9f734eeff7acebf7f638	venture	B	05-01-2015	1000000.0
1	/ORGANIZATION/-QUINTER	round/22acff4f96eb7acb2b901dec1cfe5633	venture	A	14-10-2014	NaN
2	/organization/-quarter	round/b44fb941539fcdcf13085330bb48030	seed	NaN	01-03-2014	700000.0
3	/ORGANIZATION/-THE ONE-Of-THem-INc	round/650b870d4416801069b6b178a1418776b	venture	B	30-01-2014	3406878.0
4	/organization/O-6-com	round/5727acaea57461bd2a29d9d945382d	venture	A	19-03-2008	2000000.0

Feature Extraction

- We will choose the required columns for analysis and drop the rest

In companies dataframe, the columns of interest are

- permalink
- category_list
- country_code

Dropping the rest

```
In [5]: companies.drop(['name', 'homepage_url', 'status', 'state_code', 'region', 'city', 'founded_at'],
#companies_df.head()

In [6]: companies_df permalink.apply(lambda x: str(x).lower()).describe()

Out[6]:
```

	count	unique	top	freq	Name
permalink	66368	66368	/organization/solarflare	19	company_permalink, dtype: object

```
In [7]: companies_df.describe()
```

```
Out[7]:
```

	permalink	category_list	country_code
count	66368	63220	59410
unique	66368	27296	137
top	/Organization/Gn420-20	Software	USA
freq	1	3995	37601

In rounds2 dataframe, the columns of interest are

- company_permalink
- funding_round_type
- funding_round_code
- raised_amount_usd

Dropping the rest

```
In [8]: rounds2_df = rounds2.drop(['funding_round_permalink', 'funded_at'], axis=1)
#rounds2_df.head()

In [9]: rounds2_df company_permalink.apply(lambda x: str(x).lower()).describe()

Out[9]:
```

	count	unique	top	freq	Name
company_permalink	114949	66373	/organization/solarflare	19	company_permalink, dtype: object

Null Value Analysis - rounds2

```
In [10]: rounds2_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 114949 entries, 0 to 114948
Data columns (total 4 columns):
# Column Non-Null Count Dtype
---  ---
0 company_permalink 114949 non-null object
1 funding_round_type 114949 non-null object
2 funding_round_code 11440 non-null object
3 raised_amount_usd 94939 non-null float64
dtypes: float64(1), object(3)
memory usage: 3.5+ MB

In [11]: # percentage of missing value in raised amount usd
rounds2_df.raised_amount_usd.isnull().value_counts()

Out[11]:
```

	count
0	1739

```
In [12]: # percentage of missing value in funding round code
rounds2_df.funding_round_code.isnull().value_counts()

Out[12]:
```

	count
0	7291

As per the above analysis, we find that

- there are 114949 entries in the dataframe
- no missing values in company_permalink
- no missing values in funding_round_type
- funding_round_code has ~73% missing values
- raised_amount_usd has ~17% missing values

Treating missing values - rounds2

As more than 70% data is missing in 'funding_round_code', we will drop this column from the dataset

```
In [13]: rounds2_df_treated = rounds2_df.drop(['funding_round_code'], axis=1)
# rounds2_df_treated.head()

In [14]: rounds2_df_treated.head()

Out[14]:
```

	company_permalink	funding_round_type	raised_amount_usd
0	/organization/-fame	venture	1000000.0
1	/ORGANIZATION/-QUINTER	venture	NaN
2	/organization/-quarter	seed	700000.0
3	/ORGANIZATION/-THE ONE-Of-THem-INc	venture	3406878.0
4	/organization/O-6-com	venture	2000000.0

Raised amount USD is an important feature and has about 17% missing values, hence we should fill the missing values

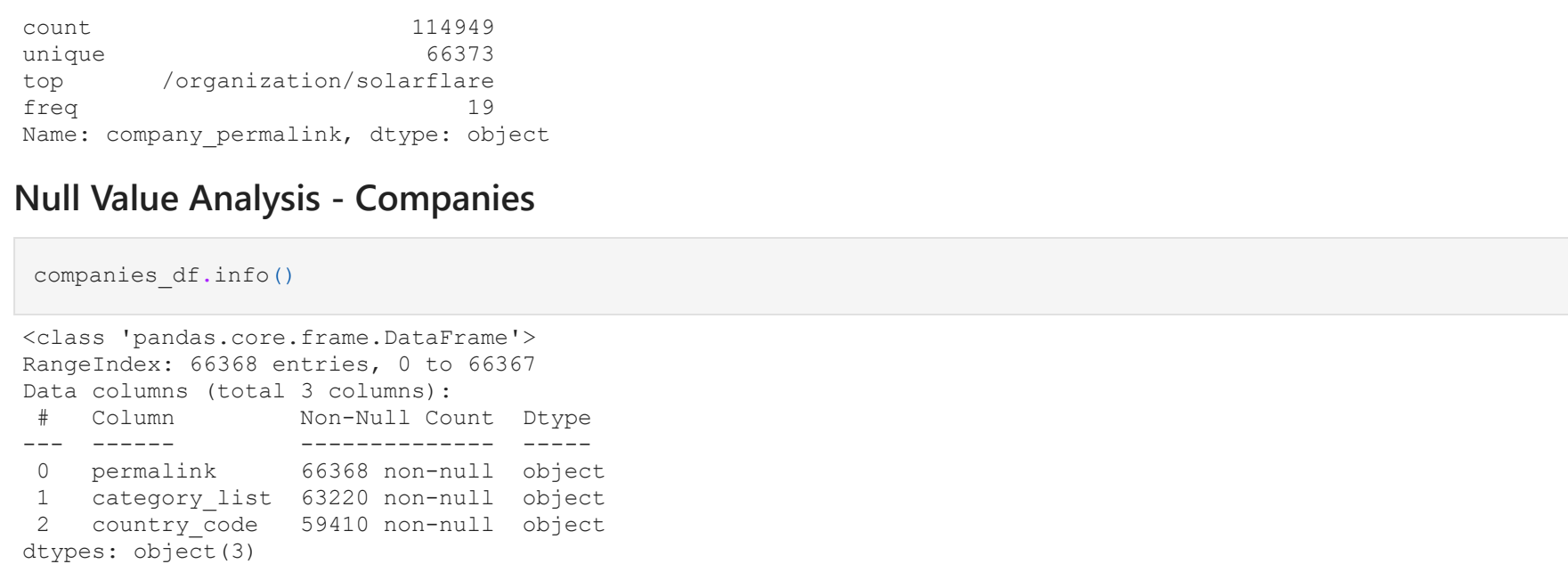
```
In [15]: rounds2_df_treated.raised_amount_usd.describe()

Out[15]:
```

	count	mean	std	min	25%	50%	75%	max	Name
raised_amount_usd	94939	9.495900e+04	1.042837e+07	0.000000e+00	1.482128e+08	2.225000e+00	1.690311e+06	7.000000e+06	raised_amount_usd, dtype: float64

Boxplot to perform outlier analysis for raised_amount_usd column

```
In [16]: sns.set_style('whitegrid')
amount_box_plot = sns.boxplot(x='raised_amount_usd', data=rounds2_df_treated)
amount_box_plot.set(xlabel='Amount in USD (x1000 Million)')
plt.show()
```



As we can see from the plot above,

- there are outliers in the dataset to the far right (top most values)
- moreover, a lot of values are NaN (i.e. missing)

We will remove these values for better visualization.

To remove, we will select all values below the 95th percentile

```
In [17]: rounds2_df_outlier_removed = rounds2_df_treated[rounds2_df_treated.raised_amount_usd < rounds2_df_treated.raised_amount_usd.quantile(0.95)]
amount_box_plot.set(xlabel='Amount in USD (x10 Million)')
plt.show()
```



```
In [18]: rounds2_df_outlier_removed.describe()

Out[18]:
```

	count	mean	std	min	25%	50%	75%	max
raised_amount_usd	9016400e+04	6.419239e+06	1.000000e+00	3.000000e+05	1.500000e+06	5.527561e+06	3293696e+07	

After removing outliers, (& NaN values, as that is default behavior of quantile selection in pandas see #1)

- The mean has come down significantly from 10 Million, to 4.3 Million
- The boxplot is plotted such that we can see that the median is somewhere in 1-2 million (near to 1.5 Million), which is what we saw earlier as well
- There are many values above the upper-fence, however, they concentrated and form a line, so these values are acc

Given this information, it is clear that median will be a better approximation for filling missing values

Note: We aren't removing outliers from the dataset, the above process was followed for visualization only

Further, we try to find median for each type of funding

```
In [19]: # Pivot Table to see median grouped by funding_round_type
rounds2_df_treated.pivot_table(
    index='funding_round_type',
    values='raised_amount_usd',
    aggfunc='median'
)

Out[19]:
```

	raised_amount_usd
angel	400000.0
convertible_note	272000.0
debt_financing	1100000.0
equity_crowdfunding	100000.0
grant	201684.0
non_equity_assistance	60000.0
post_ipo_debt	19950000.0
post_ipo_equity	12262852.5
private_equity	20000000.0
product_crowdfunding	183915.0
secondary_market	32600000.0
seed	275000.0
undisclosed	1018680.0
venture	5000000.0

We are choosing to fill missing values based on round_type, instead of the general median, since the raised amount is directly linked with what type of round it is.

```
In [20]: rounds2_df_treated.raised_amount_usd = rounds2_df_treated.raised_amount_usd.fillna(rounds2_df_treated.groupby('funding_round_type')['raised_amount_usd'].median())
# rounds2_df_treated.head()
```

rounds2 missing values treated

```
In [21]: rounds2_df_treated.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 114949 entries, 0 to 114948
Data columns (total 3 columns):
# Column Non-Null Count Dtype
---  ---
0 company_permalink 114949 non-null object
1 funding_round_type 114949 non-null object
2 raised_amount_usd 114949 non-null float64
dtypes: float64(1), object(2)
memory usage: 2.8+ MB

In [22]: # converting company_permalink to lowercase
rounds2_df_treated.company_permalink = rounds2_df_treated.company_permalink.apply(lambda x: str(x).lower())

In [23]: rounds2_df_treated.company_permalink.describe()

Out[23]:
```

	count	unique	top	freq	Name
company_permalink	114949	66373	/organization/solarflare	19	company_permalink, dtype: object

Null Value Analysis - Companies

```
In [24]: companies_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 66368 entries, 0 to 66367
Data columns (total 3 columns):
# Column Non-Null Count Dtype
---  ---
0 permalink 66368 non-null object
1 category_list 63220 non-null object
2 country_code 59410 non-null object
dtypes: object(3)
memory usage: 1.3+ MB

In [25]: # percentage of missing value in category_list
rounds(company_df.category_list.isnull()).value_counts()

Out[25]:
```

	count
0	474

```
In [26]: # percentage of missing value in country_code
round(sum(companies_df.country_code.isnull()))/len(companies_df)*100, 2)

Out[26]:
```

	count
0	10.48

As per the above analysis, we find that

- there are 66368 entries in the dataframe
- no missing values in permalink
- category_list has ~5% missing values
- country_code has ~10% missing values

Treating missing values - companies

As 5% data is missing in category_list, we can choose the most prominent category (mode) and fill the missing values with it.

- Filling Mode for NaN values see #2

```
In [27]: companies_df.category_list.mode()[0]

Out[27]: 'Software'

In [28]: companies_df_treated = companies_df.copy()

In [29]: companies_df_treated.category_list = companies_df_treated.category_list.fillna(companies_df_treated.category_list.mode()[0])
#companies_df_treated.head()
```

```
In [30]: #companies_df[companies_df.category_list.isnull()]

In [31]: companies_df_treated[companies_df_treated.category_list.isnull()]

Out[31]:
```

	permalink	category_list	country_code
0	/organization/-fame	Media	IND
1	/organization/-quarter	Application Platforms Real Time Social Network...	USA
2	/organization/-the one-of-them-inc	Apps Games Mobile	USA
3	/organization/O-6-com	Curated Web	CHN
4	/organization/O04-technologies	Software	USA

As 10% data is missing from country_code, we are choosing to fill it with this with the mode (most prominent value) of this column.

```
In [32]: companies_df_treated.country_code.mode()[0]

Out[32]: 'USA'

In [33]: companies_df_treated.country_code = companies_df_treated.country_code.fillna(companies_df_treated.country_code.mode()[0])
# companies_df_treated.head()
```

companies missing values treated

```
In [34]: companies_df_treated.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 66368 entries, 0 to 66367
Data columns (total 3 columns):
# Column Non-Null Count Dtype
---  ---
0 permalink 66368 non-null object
1 category_list 66368 non-null object
2 country_code 66368 non-null object
dtypes: object(3)
memory usage: 1.5+ MB

In [35]: # converting company_permalink to lowercase
companies_df_treated.company_permalink = companies_df_treated_permalink.apply(lambda x: str(x).lower())
companies_df_treated.head()
```

```
Out[35]:
```

	permalink	category_list	country_code
0	/organization/-fame	Media	IND
1	/organization/-quarter	Application Platforms Real Time Social Network...	USA
2	/organization/-the one-of-them-inc	Apps Games Mobile	USA
3	/organization/O-6-com	Curated Web	CHN
4	/organization/O04-technologies	Software	USA

Permalink Analysis - rounds2 vs companies

```
In [36]: rounds2_df_treated.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 114949 entries, 0 to 114948
Data columns (total 3 columns):
# Column Non-Null Count Dtype
---  ---
0 company_permalink 114949 non-null object
1 funding_round_type 114949 non-null object
2 raised_amount_usd 114949 non-null float64
dtypes: float64(1), object(2)
memory usage: 2.6+ MB

In [37]: rounds2_df_treated.company_permalink.describe()

Out[37]:
```

	count	unique	top	freq	Name
company_permalink	114949	66373	/organization/solarflare	19	company_permalink, dtype: object

```
In [38]: companies_df_treated_permalink.describe()

Out[38]:
```

	count	unique	top	freq	Name
company_permalink	66368	66368	/organization/solarflare	1	company_permalink, dtype: object

```
In [39]: set(rounds2_df_treated.company_permalink) - set(companies_df_treated_permalink)
```

```
Out[39]:
```

```
['/organization/affluent-attaché-club-2',
'/organization/asiamebook-v4',
'/organization/boreal-bikes-incorporated',
'/organization/briocreative-com',
'/organization/capptv',
'/organization/crème-ciseaux',
'/organization/e-ebicia',
'/organization/grafica-en-lima',
'/organization/huizhou-com',
'/organization/ignia-bienes-raices',
'/organization/irvengo-奥威诺',
'/organization/iproof-the-foundation-for-the-internet-of-things-v8',
'/organization/jean-pierre-produkte',
'/organization/jiwu-吉武',
'/organization/know-v8',
'/organization/lawpadi',
'/organization/magnet-tech',
'/organization/magnet-tech-利磁',
'/organization/monnier-frères',
'/organization/moonlight-2',
'/organization/patrofin',
'/organization/proditi-cs',
'/organization/presmetip',
'/organization/sailo-vip',
'/organization/talentign-v8',
'/organization/tid-pig-v8',
'/organization/tipcat-interactive-沙恩利利',
'/organization/tio-comes',
'/organization/va-de-lal',
'/organization/welche-tech-利威利',
'/organization/welche-tech-利威利',
'/organization/whodas-tech-利威利',
'/organization/zengame-v8',
'/organization/zeon',
'/organization/zeon-tiff-reklam-ve-tant-hizmetleri-tic']
```

Remove Non Ascii Characters from permalink - see #4

```
In [40]: companies_df_treated_permalink = companies_df_treated['permalink'].str.encode('ascii', 'ignore').str.decode('ascii')
companies_df_treated.head()
```

```
Out[40]:
```

	permalink	category_list	country_code
0	/organization/-fame	Media	IND
1	/organization/-quarter	Application Platforms Real Time Social Network...	USA
2	/organization/-the one-of-them-inc	Apps Games Mobile	USA
3	/organization/O-6-com	Curated Web	CHN
4	/organization/O04-technologies	Software	USA

```
In [41]: rounds2_df_treated['company_permalink'] = rounds2_df_treated['company_permalink'].str.encode('ascii', 'ignore').str.decode('ascii')
rounds2_df_treated.head()
```

```
Out[41]:
```

	company_permalink	funding_round_type	raised_amount_usd
0	/organization/-fame	venture	1000000.0
1	/organization/-quarter	venture	5000000.0
2	/organization/-quarter	seed	700000.0
3	/organization/-the one-of-them-inc	venture	3406878.0
4	/organization/O-6-com	venture	2000000.0

```
In [42]: set(rounds2_df_treated['company_permalink']) - set(companies_df_treated['permalink'])

Out[42]:
```

```
['/organization/innovatiff-reklam-ve-tant-hizmetleri-tic',
'/organization/patrofin']
```

```
In [43]: set(companies_df_treated['permalink']) - set(rounds2_df_treated['company_permalink'])

Out[43]:
```

```
['/organization/innovatiff-reklam-ve-tant-hizmetleri-tic',
'/organization/patrofin']
```

Partial Match to verify if strings are same - see #3

```
In [44]: companies_df[companies_df['permalink'].str.contains('/Organization/Patrof')]

Out[44]:
```

	permalink	category_list	country_code
42529	/Organization/Patrof	Software	TUR
42530	/Organization/Patron-Technology	Art CRM Enterprise Software Music Sports Techn...	USA
42531	/Organization/Patronpath	Software	USA
42532	/Organization/Patronus-Medical	Medical	USA

```
In [45]: companies_df_treated[companies_df_treated['permalink'].str.contains('/organization/patrof')]

Out[45]:
```

	permalink	category_list	country_code
42529	/organization/patrofin	Software	TUR
42530	/organization/patron-technology	Art CRM Enterprise Software Music Sports Techn...	USA
42531	/organization/patronpath	Software	USA
42532	/organization/patronus-medical	Medical	USA

```
In [46]: rounds2_df_treated[rounds2_df_treated['company_permalink'].str.contains('/organization/patrof')]

Out[46]:
```

	company_permalink	funding_round_type	raised_amount_usd
73633	/organization/patrofin	grant	42607.0
73634	/organization/patron-technology	angel	2500000.0
73635	/organization/patron-technology	debt_financing	600000.0
73636	/organization/patronpath	venture	700000.0
73637	/organization/patronus-medical	venture	1000000.0

```
In [47]: companies_df_treated[companies_df_treated['permalink'].str.contains('reklam')]['permalink'][66367]

Out[47]: '/organization/innovatiff-reklam-ve-tant-hizmetleri-tic'
```

```
In [48]: rounds2_df_treated[rounds2_df_treated['company_permalink'].str.contains('reklam')]['company_permalink'][114948]

Out[48]: '/organization/innovatiff-reklam-ve-tant-hizmetleri-tic'
```

We can see that the permalinks are for same company, but character 'Y' is encoded a bit differently in rounds2. We will fix these 2 entries manually

```
In [49]: rounds2_df_treated.loc[[114948], 'company_permalink'] = '/organization/innovatiff-reklam-ve-tant-hizmetleri-tic'
rounds2_df_treated[rounds2_df_treated['company_permalink'].str.contains('reklam')]['company_permalink'][114948]

Out[49]: '/organization/innovatiff-reklam-ve-tant-hizmetleri-tic'
```

```
In [50]: companies_df_treated.loc[[42529], 'permalink'] = '/organization/patrofin'
companies_df_treated[companies_df_treated['permalink'].str.contains('/organization/patrof')]

Out[50]:
```

	permalink	category_list	country_code
42529	/organization/patrofin	Software	TUR

Permalink values have been cleaned and normalized for both dataframes

```
In [51]: set(rounds2_df_treated['company_permalink']) - set(companies_df_treated['permalink'])

Out[51]: set()

In [52]: rounds2_df_treated.company_permalink.describe()

Out[52]:
```

	count	unique	top	freq	Name
company_permalink	114949	66368	/organization/solarflare	19	company_permalink, dtype: object

```
In [53]: companies_df_treated_permalink.describe()

Out[53]:
```

	count	unique	top	freq	Name
company_permalink	66368	66368	/organization/solarflare	1	company_permalink, dtype: object

We find out that both frames contain the same unique permalinks

Merging rounds2 & companies dataframes after treatment

```
In [54]: # changing name of column from company_permalink to permalink in rounds2 dataframe
rounds2_df_treated = rounds2_df_treated.rename(columns={'company_permalink': 'permalink'})
rounds2_df_treated.head()
```

```
Out[54]:
```

	permalink	funding_round_type	raised_amount_usd
0	/organization/-fame	venture	1000000.0
1	/organization/-quarter	venture	5000000.0
2	/organization/-quarter	seed	700000.0
3	/organization/-the one-of-them-inc	venture	3406878.0
4	/organization/O-6-com	venture	2000000.0

```
In [55]: master_frame = pd.merge(rounds2_df_treated, companies_df_treated, on='permalink')

In [56]: master_frame.head()
```

```
Out[56]:
```

	permalink	funding_round_type	raised_amount_usd	category_list	country_code
0	/organization/-fame	venture	1000000.0	Media	IND
1	/organization/-quarter	venture	5000000.0	Application Platforms Real Time Social Network...	USA
2	/organization/-quarter	seed	700000.0	Apps Games Mobile	USA
3	/organization/-the one-of-them-inc	venture	3406878.0	Curated Web	CHN
4	/organization/O-6-com	venture	2000000.0		

```
In [57]: master_frame.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 114949 entries, 0 to 114948
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---  ---
0 permalink 114949 non-null object
1 funding_round_type 114949 non-null object
2 raised_amount_usd 114949 non-null float64
3 category_list 66368 non-null object
4 country_code 59439 non-null object
dtypes: float64(1), object(4)
memory usage: 4.3+ MB
```

Understand the Data Set

```
unique companies in rounds2 66368
unique companies in companies 66368
unique key in companies permalink
unique key in rounds2
permalink round2 - permalink companies 0 entries - all entries in rounds2 permalink are present in company_permalink and vice versa
total observations after merge 114949
```


	permalink	raised_amount_usd	category_list	country_code	sector
0	/organization/fame	10000000.0	media	IND	Entertainment
1	/organization/90min	15000000.0	media	GBR	Entertainment
2	/organization/90min	5800000.0	media	GBR	Entertainment
3	/organization/90min	18000000.0	media	GBR	Entertainment
4	/organization/all-def-digital	5000000.0	media	USA	Entertainment

```
In [105]: top3_countries_filtered = top3_countries_sectors[(top3_countries_sectors.raised_amount_usd <= 15000000) & (top3_countries_sectors.country_code != "IND")]
```

```
In [106]: top3_countries_filtered.head()
```

	permalink	raised_amount_usd	category_list	country_code	sector
0	/organization/fame	10000000.0	media	IND	Entertainment
1	/organization/90min	15000000.0	media	GBR	Entertainment
2	/organization/90min	5800000.0	media	GBR	Entertainment
3	/organization/all-def-digital	5000000.0	media	USA	Entertainment
4	/organization/chefs-feed	5000000.0	media	USA	Entertainment
16	/organization/chefs-feed	5000000.0	media	USA	Entertainment

Understand the data for the Top 3 Countries

```
In [107]: ind_df = top3_countries_filtered[top3_countries_filtered['country_code'] == "IND"]
us_df = top3_countries_filtered[top3_countries_filtered['country_code'] == "USA"]
uk_df = top3_countries_filtered[top3_countries_filtered['country_code'] == "GBR"]
```

```
In [108]: ind_df.head()
```

	permalink	raised_amount_usd	category_list	country_code	sector
0	/organization/fame	10000000.0	media	IND	Entertainment
243	/organization/birds-eye-systems	5000000.0	apps	IND	News, Search and Messaging
250	/organization/cobbie-app	5000000.0	apps	IND	News, Search and Messaging
279	/organization/chilli-2	6000000.0	apps	IND	News, Search and Messaging
415	/organization/fmnoz	5000000.0	apps	IND	News, Search and Messaging

```
In [109]: ind_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 498 entries, 0 to 55378
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---
0 permalink 498 non-null object
1 raised_amount_usd 498 non-null float64
2 category_list 498 non-null object
3 country_code 498 non-null object
4 sector 498 non-null object
dtypes: float64(1), object(4)
memory usage: 23.3+ KB
```

```
In [110]: us_df.head()
```

	permalink	raised_amount_usd	category_list	country_code	sector
4	/organization/all-def-digital	5000000.0	media	USA	Entertainment
16	/organization/chefs-feed	5000000.0	media	USA	Entertainment
25	/organization/huffingtonpost	5000000.0	media	USA	Entertainment
26	/organization/huffingtonpost	5000000.0	media	USA	Entertainment
36	/organization/matchmine	10000000.0	media	USA	Entertainment

```
In [111]: us_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 15843 entries, 4 to 55473
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---
0 permalink 15843 non-null object
1 raised_amount_usd 15843 non-null float64
2 category_list 15843 non-null object
3 country_code 15843 non-null object
4 sector 15843 non-null object
dtypes: float64(1), object(4)
memory usage: 742.6+ KB
```

```
In [112]: uk_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 876 entries, 1 to 55466
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---
0 permalink 876 non-null object
1 raised_amount_usd 876 non-null float64
2 category_list 876 non-null object
3 country_code 876 non-null object
4 sector 876 non-null object
dtypes: float64(1), object(4)
memory usage: 41.1+ KB
```

```
In [113]: uk_df.head()
```

	permalink	raised_amount_usd	category_list	country_code	sector
1	/organization/90min	15000000.0	media	GBR	Entertainment
2	/organization/90min	5800000.0	media	GBR	Entertainment
102	/organization/common-interest-communities	10000000.0	application platforms	GBR	News, Search and Messaging
118	/organization/geopack-td-	5460000.0	application platforms	GBR	News, Search and Messaging
119	/organization/geopack-td-	5400000.0	application platforms	GBR	News, Search and Messaging

```
In [114]: uk_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 876 entries, 1 to 55466
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---
0 permalink 876 non-null object
1 raised_amount_usd 876 non-null float64
2 category_list 876 non-null object
3 country_code 876 non-null object
4 sector 876 non-null object
dtypes: float64(1), object(4)
memory usage: 41.1+ KB
```

```
In [115]: us_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 15843 entries, 4 to 55473
Data columns (total 5 columns):
# Column Non-Null Count Dtype
---
0 permalink 15843 non-null object
1 raised_amount_usd 15843 non-null float64
2 category_list 15843 non-null object
3 country_code 15843 non-null object
4 sector 15843 non-null object
dtypes: float64(1), object(4)
memory usage: 742.6+ KB
```

Calculating Total Invested Amount in each country

```
In [116]: us_df.raised_amount_usd.sum()
```

```
Out[116]: 129657400937.0
```

```
In [117]: uk_df.raised_amount_usd.sum()
```

```
Out[117]: 6676843539.0
```

```
In [118]: ind_df.raised_amount_usd.sum()
```

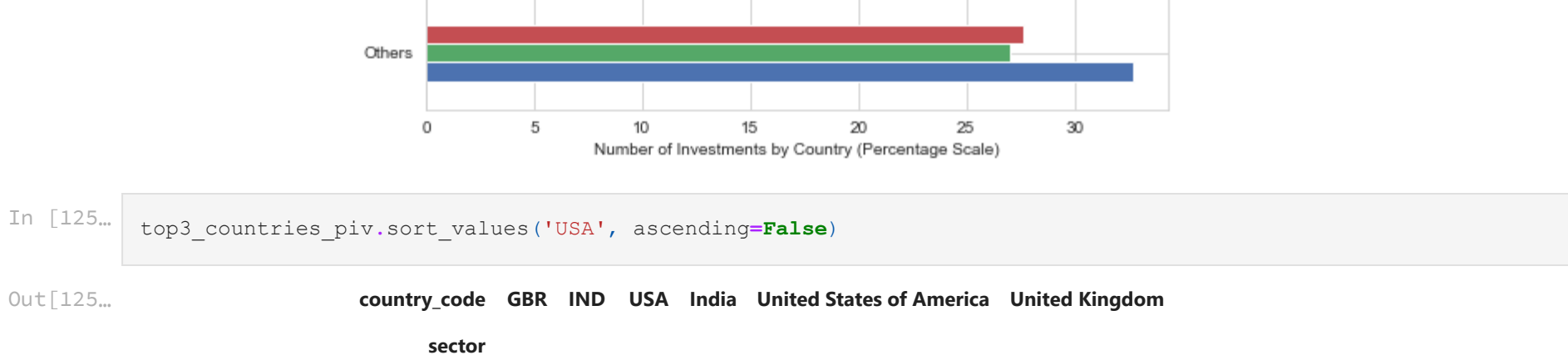
```
Out[118]: 3816543602.0
```

```
In [119]: top3_countries_piv = top3_countries_filtered.pivot_table(
    index='sector',
    values='raised_amount_usd',
    aggfunc='count',
    columns='country_code',
)
```

```
In [120]: top3_countries_piv
```

	country_code	GBR	IND	USA
	sector			
	Automotive & Sports	26	14	256
	Cleantech / Semiconductors	154	29	2745
	Entertainment	80	47	842
	Health	32	33	1123
	Manufacturing	52	30	1076
	News, Search and Messaging	97	76	2025
	Others	242	163	4285
	Social, Finance, Analytics, Advertising	193	106	3491

```
Out[120]: top3_countries_piv.plot(kind='barh')
plt.show()
```



Since the total count is very high for USA, it is difficult to understand the graph properly.

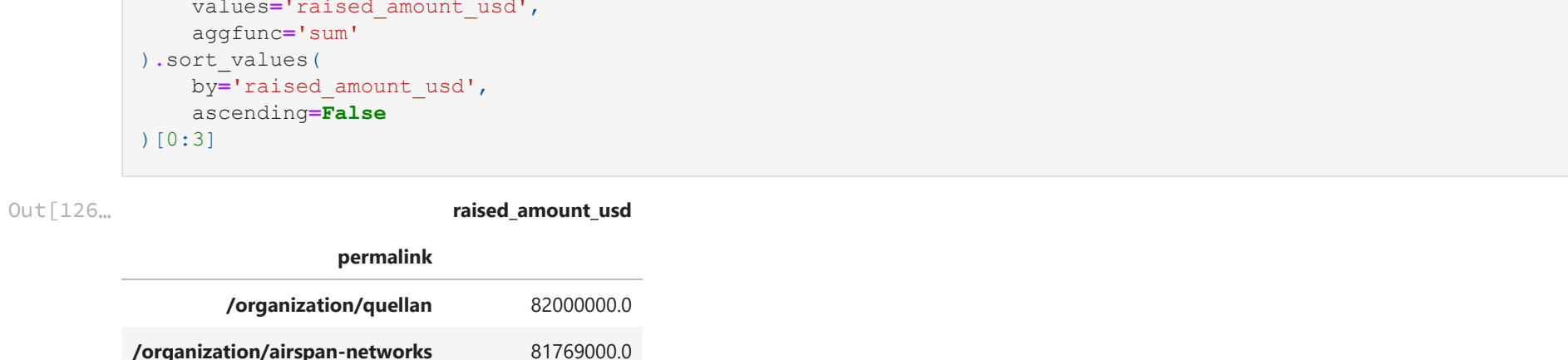
We need to analyze which sector is more active in which country, so we can generate percentage distribution of sector for each country.

```
In [122]: top3_countries_piv['India'] = round(top3_countries_piv['IND']/sum(top3_countries_piv['IND'])*100, 2)
top3_countries_piv['United States of America'] = round(top3_countries_piv['USA']/sum(top3_countries_piv['USA'])
top3_countries_piv['United Kingdom'] = round(top3_countries_piv['GBR']/sum(top3_countries_piv['GBR'])*100, 2)
```

```
In [123]: top3_sectors_in_countries_percent = top3_countries_piv.sort_values('India', ascending=False)[['India', 'United States of America', 'United Kingdom']]
```

Plotting the Sector Wise Distribution of Companies in the Top 3 English speaking countries

```
In [124]: plt.style.use('seaborn-deep')
fig, ax = plt.subplots(figsize=(8,10))
top3_sectors_in_countries_percent.plot(kind='barh', ax=ax)
ax.set_xlabel('Number of Investments by Country (Percentage Scale)')
ax.set_ylabel('')
# ax.legend(title='Country Code')
ax.set_title('Percentage Distribution of Number of Investments in Each Sector')
fig.savefig('investment_by_country.png', bbox_inches='tight')
plt.show()
```



```
In [125]: top3_countries_piv.sort_values('USA', ascending=False)
```

	country_code	GBR	IND	USA	United States of America	United Kingdom	
	sector						
	Others	242	163	4285	32.73	27.05	27.63
	Social, Finance, Analytics, Advertising	193	106	3491	21.29	22.03	22.03
	Cleantech / Semiconductors	154	29	2745	5.82	17.33	17.58
	News, Search and Messaging	97	76	2025	15.26	12.78	11.07
	Health	32	33	1123	6.63	7.09	3.65
	Manufacturing	52	30	1076	6.02	6.79	5.94
	Entertainment	80	47	842	9.44	5.31	9.13
	Automotive & Sports	26	14	256	2.81	1.62	2.97

```
In [126]: # Sorted by Total Invest - US
us_df.pivot_table(
    index='permlink',
    values='raised_amount_usd',
    aggfunc='sum',
).sort_values(
    by='raised_amount_usd',
    ascending=False
)[0:3]
```

	raised_amount_usd
/organization/quellan	82000000.0
/organization/alrspan-networks	81769000.0
/organization/biodesix	75300000.0

```
In [127]: # Sorted by Total Invest - UK
uk_df.pivot_table(
    index='permlink',
    values='raised_amount_usd',
    aggfunc='sum',
).sort_values(
    by='raised_amount_usd',
    ascending=False
)[0:3]
```

	raised_amount_usd
/organization/greenroad-technologies	52500000.0
/organization/myupermarket	43400000.0
/organization/topa	37900000.0

```
In [128]: # Sorted by Total Investment Amount - IND
ind_df.pivot_table(
    index='permlink',
    values='raised_amount_usd',
    aggfunc='sum',
).sort_values(
    by='raised_amount_usd',
    ascending=False
)[0:3]
```

	raised_amount_usd
/organization/azure-power	57200000.0
/organization/manthan-systems	50700000.0
/organization/firstcry-com	39000000.0

```
In [129]: # Sanity check
ind_df.pivot_table(
    index='permlink',
    values='raised_amount_usd',
    aggfunc='count',
).sort_values(
    by='raised_amount_usd',
    ascending=False
)[0:3]
```

	raised_amount_usd
/organization/azure-power	6
/organization/manthan-systems	4
/organization/mynta	4

```
In [130]: # sanity check
ind_df[ind_df['permlink'] == '/organization/mynta'].raised_amount_usd.sum() - ind_df[ind_df['permlink'] ==
-1000000.0
```

```
In [131]: # sanity check
us_df[us_df.permlink == '/organization/quellan'].raised_amount_usd.sum() - us_df[us_df.permlink == '/organization
17700000.0
```

Conclusion

- Spark Funds can become a Venture Fund Investor
- It can invest in USA, UK & India
- In USA, it can invest in Others, Cleantech/Semiconductors & Social, Finance, Analytics, Advertising sectors
- In UK, it can invest in Others, Cleantech/Semiconductors & Social, Finance, Analytics, Advertising sectors
- In UK, it can invest in Others, Social, Finance, Analytics, Advertising & News, Search and Messaging sectors

Note: We have chosen Others as one of the sectors since it was mentioned that this is one of the 8 sectors that can be analyzed. However, we should dig deeper and figure what prominent sub-sectors/categories can be identified after discussing the findings thoroughly with Spark Fund

References

- #1 - <https://stackoverflow.com/a/24047916>
 - #2 - <https://stackoverflow.com/a/42789818/7048915>
 - #3 - <https://stackoverflow.com/a/11531402/7048915>
 - #4 - <https://stackoverflow.com/a/56744855/7048915>
 - #5 - https://www.reddit.com/r/learnpython/comments/3cnpq/seaborn_axis_as_index_cuo79g?utm_source=share&utm_medium=web2&context=3
 - #6 - <https://discuss.analyticsvidhya.com/t/difference-between-wide-and-long-data-format/8110>
 - #7 - <https://pandas.pydata.org/docs/reference/api/pandas.melt.html>
- <https://stackabuse.com/change-tick-frequency-in-matplotlib>
 - <https://stackoverflow.com/a/33382750/7048915>
 - <https://github.com/rmwaskon/seaborn/issues/871#issuecomment-399701820>
 - https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.axes.Axes.set_title.html
 - <https://stackoverflow.com/a/7125157/7048915>
 - <https://www.geeksforgeeks.org/how-to-drop-rows-in-pandas-dataframe-by-index-labels/>
 - <https://stackoverflow.com/questions/44723377/pandas-combining-two-dataframes-horizontally>
 - <https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html>