ae3/numpy-1.19.5-cp36-cp36m-manylinux1\_x86\_64.whl (13.4MB) ■| 13.4MB 1.4MB/s eta 0:00:01 24% | 100% | 3.2M B 33.5MB/s eta 0:00:01 36% | | 4.8MB 32.1MB/s eta 0:00:01 47% | | 6.4MB 31.2MB/s eta 0:00:01 81% | | 10.9MB 32.1MB/s eta 0:00:01 tensorflow 1.3.0 requires tensorflow-tensorboard<0.2.0,>=0.1.0, which is not installed. Installing collected packages: numpy Found existing installation: numpy 1.12.1 Uninstalling numpy-1.12.1: Successfully uninstalled numpy-1.12.1 Successfully installed numpy-1.19.5 Collecting missingno Downloading https://files.pythonhosted.org/packages/87/22/cd5cf999af21c2f97486622c551ac3d07361ced8125121e907f588ff5 f24/missingno-0.5.2-py3-none-any.whl Requirement already satisfied: numpy in /opt/conda/lib/python3.6/site-packages (from missingno) (1.19.5) Requirement already satisfied: scipy in /opt/conda/lib/python3.6/site-packages (from missingno) (1.2.1) Requirement already satisfied: seaborn in /opt/conda/lib/python3.6/site-packages (from missingno) (0.8.1) Requirement already satisfied: matplotlib in /opt/conda/lib/python3.6/site-packages (from missingno) (2.1.0) Requirement already satisfied: pandas in /opt/conda/lib/python3.6/site-packages (from seaborn->missingno) (0.23.3) Requirement already satisfied: six>=1.10 in /opt/conda/lib/python3.6/site-packages (from matplotlib->missingno) (1.1 Requirement already satisfied: python-dateutil>=2.0 in /opt/conda/lib/python3.6/site-packages (from matplotlib->missi ngno) (2.6.1) Requirement already satisfied: pytz in /opt/conda/lib/python3.6/site-packages (from matplotlib->missingno) (2017.3) Requirement already satisfied: cycler>=0.10 in /opt/conda/lib/python3.6/site-packages/cycler-0.10.0-py3.6.egg (from m atplotlib->missingno) (0.10.0) Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /opt/conda/lib/python3.6/site-packages (fr om matplotlib->missingno) (2.2.0) Installing collected packages: missingno Successfully installed missingno-0.5.2 In [2]: import pandas as pd import matplotlib.pyplot as plt import os import configparser import datetime as dt from pyspark.sql.functions import isnan, when, count, col, udf, dayofmonth, dayofweek, month, year, weekofyear, avg, monotonically\_increasing\_id from pyspark.sql.types import \* import requests requests.packages.urllib3.disable\_warnings() from pyspark.sql.functions import year, month, dayofmonth, weekofyear, date\_format from pyspark.sql import SparkSession, SQLContext, GroupedData, HiveContext from pyspark.sql.functions import \* from pyspark.sql.functions import date\_add as d\_add from pyspark.sql.types import DoubleType, StringType, IntegerType, FloatType from pyspark.sql import functions as F from pyspark.sql.functions import lit from pyspark.sql import Row import datetime, time import seaborn as sns import numpy as np import tools as tools import tables\_func as ct import matplotlib.pyplot as plt %matplotlib inline **Scope of the Project:** • In this project I will gather the data from two sources. I will load this data into staging dataframes. I will clean the raw data, write it to parquet files and perform an ETL process using a Spark cluster. Then I will write the data into Fact & Dimension tables to form a star schema. The star schema can then be used by the relevant parties to perform data analytics, correlation and ad-hoc reporting in an effective and efficient manner. ### Data Descriptions: • Edu\_ministry is Sample data of saudi scholarships students records from the Ministry of Education. This data source will serve as the Fact table in the schema. This data comes from <a href="https://od.data.gov.sa/Data/ar/group/education\_and\_training?page=2">https://od.data.gov.sa/Data/ar/group/education\_and\_training?page=2</a>: البوابة الوطنية للبيانات المفتوحة! Countries of the World :This dataset contains Country names linked to region, population, area size, GDP, mortality and more. This data comes from Kaggle: <a href="https://www.kaggle.com/datasets/fernandol/countries-of-the-world">https://www.kaggle.com/datasets/fernandol/countries-of-the-world</a> spark = SparkSession\ In [3]: .builder \ .config("spark.jars.packages", "org.apache.hadoop:hadoop-aws:2.7.0") \ .getOrCreate() In [5]: #import shutil #shutil.rmtree('../Project/parquet\_tables/countries') **Spark configuration parameters** • ('spark.executor.id', 'driver'), • ('spark.app.name', 'pyspark-shell'), ('spark.driver.port', '33753'), • ('spark.jars.packages', 'org.apache.hadoop:hadoop-aws:2.7.0'), • ('spark.rdd.compress', 'True'), • ('spark.serializer.objectStreamReset', '100'), ('spark.master', 'local[\*]'), • ('spark.driver.host', 'ad2cb8c219e1'), ('spark.submit.deployMode', 'client'), • ('spark.jars', In [6]: fname = '../Project/Sources/Edu\_minisrty.csv' edu\_df = pd.read\_csv(fname) In [7]: edu\_df.head() Out[7]: المرحلة الدراسية فئة الإبتعاث القطاع الدولي السنة الدراسية العدد الجنس المجالالواسع 0 2015.0 الأعمال، الإدارة والقانون طالب مبتعث استراليا أخرى 11.0 اً خرى 1 2015.0 6.0 طالب مبتعث استراليا الأعمال، الإدارة والقانون 2 4.0 أنثى التعليم 2015.0 طالب مبتعث استراليا 3 2015.0 طالب مبتعث استراليا الصحة والرفاه 3.0 2015.0 1.0 ذكر العلوم الطبيعية, والرياضايات والإحصاء أخرى طالب مبتعث استراليا In [8]: fname = '../Project/Sources/countries of the world.csv' country\_df = pd.read\_csv(fname) In [9]: country\_df.head() Out[9]: Infant Pop. Coastline mortality GDP (\$ **Phones** Area Density Net Literacy **Arable Crops Other** Climate Birthrate Country **Region Population** (coast/area migration (sq. mi.) (per sq. (%) ratio) 1000 capita) 1000) births) ASIA (EX. 0 Afghanistan NEAR 31056997 647500 48,0 0,00 23,06 163,07 700.0 36,0 3,2 12,13 0,22 87,65 46,6 EAST) EASTERN 28748 71,2 21,09 3581655 124,6 1,26 -4,93 21,52 4500.0 4,42 74,49 15,11 EUROPE 32930091 2381740 13,8 0,04 31 6000.0 70,0 78,1 3,22 0,25 96,53 17,14 Algeria -0,39 1 **AFRICA** American OCEANIA 9,27 8000.0 22,46 Samoa WESTERN 71201 4,05 19000.0 100,0 497,2 2,22 0 97,78 Andorra 468 152,1 0,00 8,71 EUROPE **Exploration and Assessing the Data Data Cleaning:** • Drop columns containing over 90% missing values • Drop duplicate values In [10]: # checking missing values edu\_df.isnull().sum() Out[10]: 0 السنة الدراسية القطاع الدوليي فئة الإبتعاث المرحلة الدراسية المجال الواسع الجنس العدد dtype: int64 In [11]: # Drop columns with over 90% missing values clean\_edu = tools.eliminate\_missing\_data(edu\_df) Dropping missing data... Cleaning complete! In [12]: | clean\_edu = tools.drop\_duplicate\_rows(clean\_edu) Dropping duplicate rows... 0 rows dropped. In [13]: # Drop columns with over 90% missing values clean\_country = tools.eliminate\_missing\_data(country\_df) Dropping missing data... Cleaning complete! In [14]: | clean\_country = tools.drop\_duplicate\_rows(clean\_country) Dropping duplicate rows... 0 rows dropped. Now I'll translate the dataset to english language In [15]: clean\_edu = ct.translate\_en(clean\_edu) clean\_edu.head() Out[15]: Academic\_year Country Scholarship\_type Educational\_level Specialization Sex Count 2015.0 Australia Scholarship student Business, Management and Law 1 2015.0 Australia Scholarship student Other Business, Management and Law 6.0 2015.0 Australia Scholarship student Other Education 4.0 3 2015.0 Australia Scholarship student Other Health and well-being 2015.0 Australia Scholarship student Other Natural sciences, mathematics and statistics In [16]: clean\_edu['Country'].count() Out[16]: 3148 Males represent No.1 and females No.2 In [15]: clean\_edu['Sex'].value\_counts().plot(kind='pie', autopct='%.2f%%') Out[15]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7f501e43fda0> 1 52.03% 47.97% 2 In [17]: fig, ax = plt.subplots(figsize=(16,4)) clean\_edu['Scholarship\_type'].value\_counts().plot(kind='bar', fontsize =12, rot = 0) plt.title('Scholarship Types') Out[17]: Text(0.5,1,'Scholarship Types') Scholarship Types 1600 1400 1200 1000 800 600 400 200 Scholarship student Expatriate employee on his own In [18]: fig, ax = plt.subplots(figsize=(18,6)) clean\_edu['Country'].value\_counts().plot(kind='bar', fontsize =12, rot = 0) plt.title('Countries') Out[18]: Text(0.5,1,'Countries') Countries 600 500 400 300 200 100 -**European Countries Arab Countries** United States United Kingdom Australia Canada Asian Countries In [19]: fig, ax = plt.subplots(figsize=(14,6)) clean\_edu['Academic\_year'].value\_counts().plot(kind='bar', fontsize =12, rot = 0) plt.title('Year') Out[19]: Text(0.5,1,'Year') Year 1000 800 600 400 200 2018.0 2015.0 2017.0 2016.0 In [20]: fig, ax = plt.subplots(figsize=(12,12)) clean\_edu['Specialization'].value\_counts().plot(kind='pie',autopct='%.2f%%') Out[20]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7f6a26e11588> Business, Management and Law Natural sciences, mathematics and statistics Communication and IT 8.80% 8.67% Education 9.91% Health and well-being 8.58% 10.74% IN THE PROPERTY OF THE PROPERT 7.72% Arts and #Jumanities Manufacturing and processing Press and media Agriculture, forestry, fisheries and veterinary 6.58% Personal services Non-specialized programs and qualifications well-being Engineering, manufacturing and construction Behavioral Social Sciences Mathematics and statistics Biological and related sciences Social Sciences, Journalism and Media Architecture and construction Health Languages Engineering and engineering crafts Services Not Specified Businesses and Management **Defining the Data Model** 3.1 Conceptual Data Model In accordance with Kimball Dimensional Modelling Techniques, laid out in this document (<a href="http://www.kimballgroup.com/wp-content/uploads/2013/08/2013.09-">http://www.kimballgroup.com/wp-content/uploads/2013/08/2013.09-</a> <u>Kimball-Dimensional-Modeling-Techniques11.pdf</u>), the following modelling steps have been taken: Select the Business Process: The Scholarships department follows their business process of recording scholarships students. This process generates events which are captured and translated to facts in a fact table • Declare the Grain: The grain identifies exactly what is represented in a single fact table row. • In this project, the grain is declared as a single occurrence of a batch of scholarships students out of the KSA. • Identify the Dimensions: Dimension tables provide context around an event or business process. • The dimensions identified in this project are: dim country dim\_specialization • Identify the Facts: • Fact tables focus on the occurrences of a singular business process, and have a one-to-one relationship with the events described in the grain. The fact table identified in this project is: fact\_scholarships For this application, I have developed a set of Fact and Dimension tables in a Relational Database Management System to form a Star Schema. This Star Schema can be used by Data Analysts and other relevant business professionals to gain deeper insight into various scholarships figures, trends and statistics recorded historically. 3.2 Mapping Out Data Pipelines List the steps necessary to pipeline the data into the chosen data model: 1 Load the data into staging tables • 2 Create Dimension tables • 3 Create Fact table · 4 Write data into parquet files • 5 Perform data quality checks **Running Pipelines to Model the Data** 4.1 Create the data model Build the data pipelines to create the data model In [17]: | output\_path = "parquet\_tables/" In [18]: clean\_edu.head() Out[18]: Academic\_year Country Scholarship\_type Educational\_level Specialization Sex Count 0 2015.0 Australia Scholarship student 2 11.0 Other Business, Management and Law 1 2015.0 Australia Scholarship student Other Business, Management and Law 1 6.0 2015.0 Australia Scholarship student Other 4.0 Education 2 3 2015.0 Australia Scholarship student Health and well-being 2015.0 Australia Scholarship student Other Natural sciences, mathematics and statistics In [19]: # creating schema scholarships\_schema = StructType([StructField("Academic\_year", StringType(), **True**)\ ,StructField("Country", StringType(), True)\ ,StructField("Scholarship\_type", StringType(), True)\ ,StructField("Educational\_level", StringType(), True)\ ,StructField("Specialization", StringType(), True)\ ,StructField("Sex", StringType(), True)\ ,StructField("Count", FloatType(), True)]) scholarships\_spark = spark.createDataFrame(clean\_edu, schema=scholarships\_schema) scholarships\_spark.toPandas().head() Out[19]: Academic\_year Country Scholarship\_type Educational\_level Specialization Sex Count 2 11.0 2015.0 Australia Scholarship student Business, Management and Law 1 2015.0 Australia Scholarship student Other Business, Management and Law 6.0 2015.0 Australia Scholarship student Other Education 4.0 3 2015.0 Australia Scholarship student Other Health and well-being 3.0 2015.0 Australia Scholarship student Other Natural sciences, mathematics and statistics In [31]: clean\_country.head() Out[31]: Infant GDP (\$ Coastline mortality Phones Area Density Literacy **Arable Crops Other** Climate Birthrate Country **Region Population** (coast/area (sq. mi.) (per sq. migration (%) (%) (%) 1000) 1000 ratio) capita) births) ASIA (EX. **0** Afghanistan NEAR 31056997 647500 48,0 0,00 23,06 163,07 700.0 36,0 3,2 12,13 0,22 87,65 46,6 EAST) EASTERN Albania 3581655 28748 124,6 1,26 -4,93 21,52 4500.0 86,5 71,2 21,09 4,42 74,49 15,11 **EUROPE** NORTHERN Algeria 32930091 2381740 13,8 0,04 -0,39 31 6000.0 70,0 78,1 3,22 0,25 96,53 17,14 **AFRICA** American OCEANIA 199 290,4 58,29 -20,71 9,27 8000.0 97,0 259,5 10 15 75 22,46 Samoa WESTERN 0 97,78 Andorra 71201 468 152,1 0,00 4,05 19000.0 100,0 497,2 2,22 8,71 EUROPE In [21]: #"Country", "Region", "Population", "Net migration", "Coastline (coast/area ratio)", "GDP (\$ per capita) country\_schema = StructType([StructField("Country", StringType(), True)\ ,StructField("Region", StringType(), True)\ ,StructField("Population", IntegerType(), True)\ ,StructField("Area (sq. mi.)", IntegerType(), True)\ ,StructField("Pop. Density (per sq. mi.)", StringType(), True)\ ,StructField("Coastline (coast/area ratio", StringType(), True)\ ,StructField("Net migration)", StringType(), True)\ ,StructField("Infant mortality (per 1000 births)", StringType(), True)\ ,StructField("GDP (\$ per capita", FloatType(), True) ,StructField("Literacy (%)", StringType(), True)\ ,StructField("Phones (per 1000)", StringType(), True)\ ,StructField("Arable (%)", StringType(), True)\ ,StructField("Crops (%)", StringType(), True)\ ,StructField("Other (%)", StringType(), True)\ ,StructField("Climate", StringType(), True)\ ,StructField("Birthrate", StringType(), True)\ ,StructField("Deathrate", StringType(), True)\ ,StructField("Agriculture", StringType(), True)\ ,StructField("ndustry", StringType(), True)\ ,StructField("Service", StringType(), True)]) country\_spark = spark.createDataFrame(clean\_country, schema=country\_schema) country\_spark.toPandas().head() Out[21]: Infant GDP (\$ Literacy Coastline mortality Net Arable Crops Other Area Density Region Population Climate Birthrate (coast/area Country per (per (sq. mi.) (per sq. migration) (%) ratio 1000 capita 1000) mi.) births) ASIA (EX. NEAR 31056997 647500 48,0 163,07 700.0 3,2 12,13 0,22 87,65 O Afghanistan 0,00 23,06 36,0 46,6 EAST) EASTERN 71.2 21.09 15,11 Albania 3581655 28748 124,6 1,26 -4,93 21,52 4500.0 4,42 74,49 **EUROPE** NORTHERN 17,14 32930091 2381740 13,8 0,04 -0,39 31 6000.0 70,0 78,1 3,22 0,25 96,53 Algeria American OCEANIA 199 290,4 58,29 -20,71 9,27 8000.0 97,0 259,5 22,46 Samoa WESTERN 71201 468 152,1 0,00 4,05 19000.0 100,0 497,2 2,22 0 97,78 8,71 Andorra **EUROPE** In [22]: # inserting missed rows # Creating the DataFrame newRow = spark.createDataFrame(added\_row, schema=None) country\_spark= country\_spark.union(newRow) 1. Create dim\_Country table In [23]: ct.create\_country\_dim(country\_spark, output\_path) Writing table countries to parquet\_tables/countries Write complete! Out[23]: DataFrame[country\_id: int, country\_name: string, Region: string, Population: bigint, Coastline\_area\_ratio: string, Ne t\_migration: string, GDP\_\$\_per\_capital: double] In [24]: | country = spark.read.parquet("parquet\_tables/countries") country.toPandas().head() Out[24]: Region Population Coastline\_area\_ratio Net\_migration GDP\_\$\_per\_capital country\_name Afghanistan ASIA (EX. NEAR EAST) 31056997 23,06 700.0 1 EASTERN EUROPE 3581655 1,26 -4,93 4500.0 Albania NORTHERN AFRICA 32930091 0,04 -0,39 6000.0 Algeria 57794 58,29 8000.0 3 3 American Samoa **OCEANIA** -20,71 WESTERN EUROPE 6,6 19000.0 2. Create dim\_Specialization table In [25]: ct.create\_specialization\_dim(scholarships\_spark, output\_path) Writing table specializations to parquet\_tables/specializations Write complete! Out[25]: DataFrame[spec\_id: int, spec\_desc: string] In [41]: | specialization = spark.read.parquet("parquet\_tables/specializations") specialization.toPandas().head() Out[41]: spec\_id spec\_desc Services Engineering, manufacturing and construction Education 3 3 Farming 4 Others in engineering, manufacturing and const... 3. Create fact\_Scholarships table In [27]: scholarships = ct.create\_scholarships\_fact(scholarships\_spark, output\_path, spark) Writing table scholarships to parquet\_tables/scholarships Write complete! In [28]: | scholarships = spark.read.parquet("parquet\_tables/scholarships") scholarships.toPandas().head() Out[28]: scholarships\_id spec\_id country\_id Academic\_year Educational\_level Sex Count 0 1 30.0 228 2015.0 1 33 228 2015.0 2 33 2015.0 228 1.0 3 3 228 2015.0 2.0 16 228 2015.0 Fellowship 1 2.0 In [29]: scholarships.count() Out[29]: 1738 **Data Quality Checks** In [5]: tables = { "specialization": specialization, "country": country, "scholarships": scholarships for table\_name, table in tables.items(): tools.perform\_quality\_check(table, table\_name) NameError Traceback (most recent call last) <ipython-input-5-56fc45cfc6ef> in <module>() **1** tables = { "specialization": specialization, ---> 2 "country": country, 3 "scholarships": scholarships 4 **5** } NameError: name 'specialization' is not defined In [6]: countriesTable = spark.read.parquet("parquet\_tables/countries") specializationTable = spark.read.parquet("parquet\_tables/specializations") scholarshipsTable = spark.read.parquet('parquet\_tables/scholarships') In [30]: #Females scholarshipsTable.filter(scholarshipsTable.Sex == '2') \ .select('scholarships\_id', 'Sex') \ .dropDuplicates() \ .count() Out[30]: 828 In [31]: # Males scholarshipsTable.filter(scholarshipsTable.Sex == '1') \ .select('scholarships\_id', 'Sex') \ .dropDuplicates() \ .count() Out[31]: 910 In [25]: query = scholarshipsTable.select(["\*"])\ .join(countriesTable, (countriesTable.country\_id == scholarshipsTable.country\_id), how='inner')\ .join(specializationTable, (specializationTable.spec\_id == scholarshipsTable.spec\_id), how='inner')\ .select(['scholarships\_id' , specializationTable.spec\_desc.alias('Specialization Name'), countriesTabl e.country\_name.alias('Country Name'),countriesTable.Population.alias('Country Population'),\ 'Sex', 'Educational\_level','Academic\_year','count'])\ .filter(countriesTable.country\_name != ('Asian Countries'))\ .filter(countriesTable.country\_name != ('European Countries'))\ .filter(countriesTable.country\_name != ('Arab Countries')) query.toPandas().head(50) Out[25]: Specialization Name Country Name Country Population Sex Educational\_level Academic\_year count scholarships\_id 170 298444215 Business, Management and Law United States Other 2015.0 1.0 1 171 Business, Management and Law **United States** 298444215 Other 4.0 2015.0 172 298444215 2 **United States** Other 2015.0 1.0 Education 173 3 298444215 Other 6.0 Education United States 2015.0 174 Health and well-being 298444215 **United States** Other 2015.0 17.0 175 5 Arts and humanities **United States** 298444215 2 Other 2015.0 8.0 176 Arts and humanities **United States** 298444215 Other 2015.0 2.0 7 **United States** 298444215 1 Other 177 Engineering, manufacturing and construction 2015.0 1.0 178 Communication IT 298444215 United States Other 2015.0 1.0 179 9 Communication\_IT **United States** 298444215 Other 2015.0 2.0 10 180 298444215 2 Not Specified **United States** Other 2015.0 1.0 11 181 298444215 Other 5.0 Not Specified **United States** 2015.0 12 182 298444215 Health and well-being **United States** Fellowship 2015.0 98.0 13 183 298444215 141.0 Health and well-being **United States** Fellowship 2015.0 14 184 298444215 Not Specified **United States** Fellowship 2015.0 2.0 15 185 Business, Management and Law **United States** 298444215 2 Bachelor 2015.0 208.0 16 186 Business, Management and Law **United States** 298444215 Bachelor 2015.0 1287.0 17 187 Non-specialized programs and qualifications **United States** 298444215 Bachelor 2015.0 1.0 18 188 298444215 2 **United States** Bachelor 2015.0 29.0 Education 189 19 298444215 13.0 Education United States Bachelor 2015.0 20 190 298444215 Services **United States** Bachelor 2015.0 7.0 21 191 298444215 Bachelor 47.0 Services **United States** 2015.0 22 192 298444215 2015.0 305.0 Health and well-being **United States** Bachelor 23 193 **United States** 298444215 Bachelor 298.0 Health and well-being 2015.0 24 298444215 194 Social Sciences, Journalism and Media **United States** Bachelor 2015.0 52.0 25 195 Social Sciences, Journalism and Media **United States** 298444215 Bachelor 2015.0 190.0 26 196 Natural sciences, mathematics and statistics **United States** 298444215 2 Bachelor 2015.0 139.0 27 197 Natural sciences, mathematics and statistics United States 298444215 1 Bachelor 2015.0 289.0 Arts and humanities United States 298444215 2 2015.0 125.0 29 199 298444215 1 2015.0 88.0 Arts and humanities United States Bachelor 30 298444215 2015.0 129.0 200 Engineering, manufacturing and construction **United States** Bachelor 31 201 Engineering, manufacturing and construction **United States** 298444215 1 Bachelor 2015.0 1677.0 32 202 Communication IT 298444215 United States Bachelor 2015.0 83.0 33 203 Communication\_IT 298444215 Bachelor 2015.0 642.0 **United States** 34 204 298444215 2 Not Specified **United States** Bachelor 2015.0 14.0 35 205 Not Specified **United States** 298444215 Bachelor 2015.0 39.0 36 206 Business, Management and Law **United States** 298444215 Ph.d 2015.0 26.0 37 207 Business, Management and Law 298444215 1 Ph.d 2015.0 37.0 **United States** 38 208 298444215 2 Ph.d 12.0 **United States** 2015.0 Education 39 209 **United States** 298444215 1 Ph.d 21.0 Education 2015.0 40 210 298444215 2 Ph.d 19.0 Health and well-being **United States** 2015.0 41 211 298444215 Ph.d 2015.0 9.0 Health and well-being **United States** 42 212 298444215 2 Ph.d Social Sciences, Journalism and Media **United States** 2015.0 5.0 43 213 298444215 1 Ph.d 2015.0 2.0 Social Sciences, Journalism and Media **United States** 44 298444215 Ph.d 214 Natural sciences, mathematics and statistics **United States** 2015.0 20.0 215 Natural sciences, mathematics and statistics 45 **United States** 298444215 1 Ph.d 2015.0 13.0 46 216 Arts and humanities **United States** 298444215 2 Ph.d 2015.0 7.0 47 217 298444215 1 Arts and humanities **United States** Ph.d 1.0 2015.0

48

49

218 Engineering, manufacturing and construction

219 Engineering, manufacturing and construction

298444215

298444215 1

Ph.d

Ph.d

2015.0

2015.0

1.0

21.0

**United States** 

**United States** 

In [1]: ! pip install -U numpy

Collecting numpy

! pip install missingno

Downloading https://files.pythonhosted.org/packages/45/b2/6c7545bb7a38754d63048c7696804a0d947328125d81bf12beaa692c3