Programming Fundamentals for Data Engineers

Final Project Report UN Cell Phone and Internet Data Wrangling and Analysis

Table of Contents

Initial Data Review:	3
Preliminary Design/Planning:	
Program execution/flow of control overview:	
Design Process:	
The final/master dataframe was merged in 2 steps:	
Data Statistics Design:	
Masking:	
Dataframe Manipluation:	6
Charting:	6
Quick Find Guide for Code (Where to find things):	7
Citation:	7
Screenshots:	8
Execution:	8
Handling of Incorrect Inputs Screenshot	10

Initial Data Review:

Before choosing a dataset we reviewed various potential datasets provided; we ultimately landed on the UN data for cell phones and internet usage because they provided a large amount of relatively clean numeric data that we could use for data/stats work.

Preliminary Design/Planning:

Before coding we laid out the main blocks of our code as well as the user inputs we would request, as well as some high-level stats we would be providing.

•	
Main	code blocks:
	Main (also contains data/stats coding) User input Find nulls
User in	nputs:
	UN Sub-region Dataset (cell phones or internet usage) Years for specific data column operations
Data/S	Stats:
	Print summary of final merged dataframe Mean Sum
	am execution/flow of control overview: → user_input(args) → main() → find_null(args) → main() → data/stats work (done in

Design Process:

Three (3) separate original dataframes were merged into a master dataframe with a multiindex for work on the project; Excel was <u>not</u> used to modify the files prior to loading. General overview of steps employed for each individual dataframe prior to merging:

- 1. High-level review of data in Jupyter Lab first to get a better understanding of NaNs, data types (e.g. df.info()).
- Loaded data with pd.read_excel()
- 3. Dropped duplicate rows
- 4. Dropped unnecessary columns (e.g. sparse data)
- 5. Casted dataframe elements as needed (e.g. int to float)
- 6. Checked and changed column labels as necessary (e.g. "Country" was changed to "country" to facilitate dataframe merging later.

The final/master dataframe was merged in 2 steps:

First the two large dataset frames (cell phones and internet user usage data) were merged. Second the UN Codes dataframe was merged into the first merged dataframe. In both merges, Pandas pd.merge() was used. Following each merge duplicate rows and columns were checked for and removed if they existed.

The final/master dataframe (tech_df) was provided a multi-index in the following manner:

The row multi-index consisting of UN sub-regions and countries was created by using
df.set_index() to effectively move the "un sub-region" and "country" columns into a 2 level row index.
The column multi-index was prepared using a list of valid years and a list of valid data
type labels. These 2 lists where then used to construct a 2-level multi-index using
Pandas pd.MultiIndex.from_product() function. (Note: an intermediate step of
transposing the dataframe was required, since Pandas does not have a
"set_column_index()" function, and the "set_index()" function is specific to rows).

user_input(args) function has primarily 4 main parts:

- 1. The 1st prompt is preceded with an overview and notes/instructions. The initial prompt than allows the user to "quit" or enter any other key to continue. (Also of note, user inputs are case insensitive).
- 2. If the user does not quit, a sub-region **mask** is created to be used in a nested while loop try/catch block for validation. In this first try and catch block, the user is prompted again to enter a valid sub-region name or index number then this input is checked against the validation mask created (the "if" statement checks if a valid name was entered; and the "elif" checks whether a valid index number was entered). If an invalid entry was given, user has the option to get "help".
- 3. Once a valid sub-region has been entered, the user is prompted for a valid data type ('1' for cell phone or '2' internet usage data). This part of the user input is also implemented with a try-catch block, with if/elif /else statements in the try-body: acceptable entries are "1", "2", or "quit" and any other entries raise an exception to be handled giving the user another chance to enter.
- 4. The last 2 user input blocks are very similar as they both prompt/obtain a specific year input (e.g. 2000, 2001. These blocks also employ nested while loops with try-catch blocks with conditional if/elif/else statements. Valid entries include dates from 1995-2017, or "quit" other entries cause an exception and user has another opportunity to make a valid entry.

Note: Nested while loops were used to ensure that the if/elif/else statements were evaluated/used first (i.e. to prevent the else-statement being associated with the try-catch.)

Data Statistics Design:

Masking:

We used a **masking** operation to find if any missing values were in the sub region - the user then chooses and a statement is printed out based off if there was any missing data. This notice to the user was just to let them know if there were missing data.

Dataframe Manipluation:

For the statistics portion of the project, the first thing we did was show the user the dataframe filtered down to their sub-region and data-type selections. Our first statistic was calculating the **mean** of their data type selected from all data from 1995 to 2017 and presenting it to the user.

Next we **created two columns** taking the differences from 1995 to 2017 for both number of cell phones and internet users and appended them to the final/master merged dataframe. We then presented these columns to the users based off the sub region they choose.

For the global data statistics, we used the **groupby** function to take the **mean and sum** of the number of cell phones and presented this data to the user. We then calculated the total number of cell phones across all un sub-regions and presented to the user.

Next the **.describe** function was used to show some basic statistics and percentiles for the merged dataframe. From the two (2) user year inputs (e.g. 2000, 2002), we took the average of these two columns individually and presented them back to the user based on their chosen sub region and data-type.

The **Pivot Table** was next printed (note: it was created above when the dataframes were being merged). The pivot table information isn't great due to the data set not being a good application for use with a pivot table, but we put something together to create one with the aggregate function of sum.

Charting:

A **Bar Plot** was created using the sub region the user chose and presented them the internet users % for 2017 in the chart. We added labels and a grid with a legend. This plot is saved to the folder of the .py file and then showed to the user.

A **Pie Chart** was created using the sub region the user chose and presented them the number of cell phones for 2017 in the chart. We added labels and a grid with a legend. This plot is saved to the folder of the .py file and then showed to the user.

To finish off the project we saved the merged dataframe to excel using the **.to_excel** function and choose its given name.

Quick Find Guide for Code (Where to find things):

•	•
Imports:	line 22
user_input function:	line 28
find_null function:	line 173
main function:	line 183
load first dataframe:	line 206
load second dataframe:	line 220
load third dataframe:	line 230
Merge dataframes:	line 243
Pivot Table Creation:	line 263
Multi-Index Creation:	line 277
Request User Input:	line 299
Statistics:	
Mean:	line 312
Adding columns:	line 318
Groupby:	line 335
Sum with Groupby:	line 344
.Describe	line 347
Averages of user inputs:	line 350
Bar Plot1:	line 361
Pie Chart1:	line 373
Exporting:	line 385
Masking Operations:	line 76
	line 173

Citation:

Cell Phones (per 100 people), United Nations, June 2021, [Online] Available:

https://www.gapminder.org/data/

Cell Phones (total), United Nations, June 2021, [Online] Available:

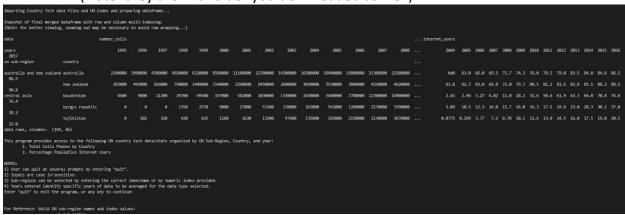
https://www.gapminder.org/data/

Individuals using the Internet (% of population), United Nations, June 2021, [Online] Available: https://www.gapminder.org/data/

^{*}Note these datasets where combined to create the given datasets*

Screenshots:

Execution: (Note: they are in a folder you downloaded as well)



For Reference: Valid UN sub-region names and index values:
un sub-region
0 australia and new zealand
1 central asia
2 eastern asia
3 estern europe 4 latin america and the caribbean
4 data merata ma the caracement
6 niconesia
7 northern africa
8 northern aserica
9 northern europe
10 polynesia
11 south-eastern asia
12 southern asía
13 southern europe
14 sub-saharan africa
15 western africa 16 western saia
Jo western assa 17 western europe
1/ RESIGNITION OF
Please enter a UN Sub-Region name or numeric index value: western europe
western europe
For Cell Phone data enter "1", for Internet Usage (%) enter "2", or "quit" to exit: 1
Selection: number of cell phones data
Victoria projecti de la compansa del compansa del compansa de la c
Enter the 1st specific year of data for avarage calc: 1995
1995
Enter 2nd specific year of data for avarage calc: 2017
chier on specific year or data for also age cast: 2017
6931
User inputs summary:
Quit? (or any key to continue): , UN Sub-Region: western europe, Data: number cells, Year 1: 1995, Year 2: 2017
There are values missing within the data, calculations were performed on the data assuming a zero value was present
Sub-Region & Data Type You Choose***********************************

	***********	**********	********	******51	b-Region &	Data Type	You Choos	·c******	******	*******	*******	********	*******	**										
ta		umber_cells																						
ars sub-region o		1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017
stern europe a		384000	599000	1160000	2298888	4250000	6120000	6540000	6740000	7270000	7990000	8670000	9280000	9910000	10000000	11400000	12200000	13000000	12600000	12200000	12000000	13500000	11100000	10900000
	belgium	235000	478000	974000	1760000	3190000	5630000	7700000	8100000	8610000	9130000	9600000	9850000	10700000	11300000	11800000	12200000	12500000	12300000	12300000			12600000	11400000
	france	1300000					29100000	37999999	38600000	41700000	44500000	48100000	51799999	55400000	58000000	57900000	57898999	59888888	62300000	63300000			67600000	69000000
	gernany	3730000	5510000	8280000	13900000	23400000	48200000	56100000	59100000	64800000	71300000	79300000	85700000	96200000	106000000	105000000	88400000	90900000	92400000	100000000	99500000	96400000	103000000	110000000
	liechtenstein	<na></na>	<na></na>	<na></na>	7500	9000	10000	11000	11400	25000	25500	27500	28800	32000	34000	35000	35500	37000	36100	38400	40700	41000	44300	46400
	luxembourg	26888	45000	67200	131000	209000	303000	409000	473000	539000	479999	510000	713000	685000	707000	720000	727000	765000	761000	788000	802000	807000	764000	794000
r.	monaco	3010	5400	7200	11500	13100	13900	14300	14900	15100	15800	17200	18300	28488	22000	23000	23400	31800	33200	35500	33700	34000	33300	33000
	netherlands	539000	1020000	1720000	3350000		10000000	12200000	12100000	13200000	14800000	15800000	17300000	19300000	28600000	20100000	19200000	19800000	19700000	19500000	19600000	28899999	20900000	20500000
	switzerland	447999	663000	1040000	1700000	3060000	4640000	5288888	5740000	6190000	6270000	6830000	7440000	8210000	8900000	9320000	9640000	10100000	10600000	11000000	11200000	11200000	11200000	11100000
			ettellana o	6 433 Det	n From 100	to 2017	Con number	colleges																
			nean o	t vii bai	a rrom 199	5 to 201/	tor number	_cerrs																
	data from 1995	to 2017 for	number_cel																					
	country																							
stern europe		8.434913e																						
	belgium	8.602478e																						
	france	4.417739e																						
	germany	6.987478e																						
	liechtenstein	2.880500e 5.311304e																						
	luxembourg	2.056565e																						
	netherlands	1,432952e																						
	switzerland	7.033478e																						
type: float64																								
			***Adding	Columns 1	o The Data	frame for	Difference	s From 19	95 to 2017															
																								internet
			number_cel												intern	et_users							number_cell	
					os 100	7 1005	1000	2000	2001	2002	2002	2004	2005	2006			000 2010	2011 201	2 2012 2	014 2015	2016 201			
ata rs ears			number_cel		196 199	7 1998	1999	2000	2001	2002	2003	2004	2005	2006			009 2010	2011 201	2 2013 2	914 2015	2016 201			: Internet Diffe
ears	coun				196 199	7 1998	1999	2000	2001	2902	2003	2004	2005				909 2010	2011 201	2 2013 2	914 2015	2016 201			
rs ears ce	coun				196 199	7 1998	1999	2000	2001	2002	2003	2004	2905				009 2010	2011 201	.2 2013 2	014 2015	2016 201			
rs ears ce n sub-region ustralia and ne	coun			95 19				2000								2988 2				914 2915		7 Cell Phon		: Internet Diffe
rs ears ce n sub-region			19	95 19										19800000		2998 2 71.7 7	4.3 76.0		0 83.5 8	4.0 84.6	86.5 86.	7 Cell Phon	e Differenc 2446000	: Internet Diffe
ears te te to sub-region ustralia and ne	new zealand aust		19	95 19 00 39906	100 458000	9 4920006	6320000									2998 2 71.7 7	4.3 76.0		0 83.5 8		86.5 86.	7 Cell Phon	e Differenc	: Internet Diffe
rs cears ce n sub-region ustralia and ne 44	new zealand aust	try ralia zealand	19 22400 3650	95 19 90 39900 90 4930	100 458000 100 56600	3 4920006 3 790006	6320000	8560000 1540000	11100000 2290000	12700000 2450000	14300000 2600000	16500000 3030000	18400000 3530000	19800000 3800000		71.7 7 72.0 7	4.3 76.0 9.7 80.5	79.5 79. 81.2 81.	9 83.5 8 6 82.8 8	4.0 84.6 5.5 88.2	86.5 86. 88.5 90.	7 Cell Phon 5	2446000 603500	Internet Differ
s ears te sub-region stralia and ne 4 2 entral asia	new zealand aust		19 22400	95 19 90 39900 90 4930	100 458000 100 56600	3 4920006 3 790006	6320000	8560000 1540000	11100000	12700000	14300000	16500000	18400000	19800000		71.7 7 72.0 7	4.3 76.0 9.7 80.5	79.5 79. 81.2 81.	9 83.5 8 6 82.8 8	4.0 84.6	86.5 86. 88.5 90.	7 Cell Phon 5	e Differenc 2446000	Internet Differ
s ars e sub-region stralia and ne 4 2 ntral asia	ew zealand aust new : kaza	try ralia zealand khstan	19 22400 3650	95 19 90 39900 90 4930	900 458000 900 56600 900 1120	9 4929996 9 799996 3 29796	6320000 1400000 49500	8560000 1540000 197000	11100000 2290000 582000	12700000 2450000 1030000	14300000 2600000 1330000	16500000 3030000 2450000	18400000 3530000 5400000	19800000 3800000 7780000		71.7 7 72.0 7 11.0 1	4.3 76.0 9.7 80.5 8.2 31.6	79.5 79. 81.2 81. 50.6 61.	0 83.5 8 6 82.8 8 9 63.3 6	4.0 84.6 5.5 88.2 6.0 70.8	86.5 86. 88.5 90. 74.6 76.	7 Cell Phon 5 8	2446000 603500 2669540	Internet Differ
s ars e sub-region stralia and ne 4 2 ntral asia 7	new zealand aust new : kaza	try ralia zealand	19 22400 3650	95 19 90 39900 90 4930	100 458000 100 56600	9 4920006 9 790006 3 29706	6320000 1400000 49500	8560000 1540000	11100000 2290000	12700000 2450000	14300000 2600000	16500000 3030000	18400000 3530000	19800000 3800000		71.7 7 72.0 7 11.0 1	4.3 76.0 9.7 80.5 8.2 31.6	79.5 79. 81.2 81. 50.6 61.	0 83.5 8 6 82.8 8 9 63.3 6	4.0 84.6 5.5 88.2	86.5 86. 88.5 90. 74.6 76.	7 Cell Phon 5 8	2446000 603500	Internet Differ
s ars e sub-region stralia and ne 4 2 ntral asia 7	new zealand aust new : kaza kyrg	try ralia zealand khstan yz republic	19 22400 3650	95 19 90 39900 90 4930 90 98	900 458999 900 56699 900 1129	9 4920006 9 790006 9 29706 9 1356	6320000 1400000 49500 2570	8560900 1540000 197000 9000	11199909 2299900 582909 27909	12790900 2450900 1030900 53100	14300000 2600000 1330000 133000	16500000 3030000 2450000 263000	18400000 3530000 5400000 542000	19800000 3800000 7780000 1260000		71.7 7 72.0 7 11.0 1 15.7 1	4.3 76.0 9.7 80.5 8.2 31.6 6.0 16.3	79.5 79. 81.2 81. 50.6 61. 17.5 19.	0 83.5 8 6 82.8 8 9 63.3 6 8 23.0 2	4.0 84.6 5.5 88.2 6.0 79.8 8.3 39.2	86.5 86. 88.5 90. 74.6 76. 37.0 38.	7 Cell Phon 5 8 4	2446000 603500 2669540 847000	: Internet Differ
s ears e sub-region estralia and ne	new zealand aust new : kaza kyrg	try ralia zealand khstan	19 22400 3650	95 19 90 39900 90 4930 90 98	900 458000 900 56600 900 1120	9 4920006 9 790006 9 29706 9 1356	6320000 1400000 49500 2570	8560000 1540000 197000	11100000 2290000 582000	12700000 2450000 1030000	14300000 2600000 1330000	16500000 3030000 2450000	18400000 3530000 5400000	19800000 3800000 7780000		71.7 7 72.0 7 11.0 1 15.7 1	4.3 76.0 9.7 80.5 8.2 31.6 6.0 16.3	79.5 79. 81.2 81. 50.6 61. 17.5 19.	0 83.5 8 6 82.8 8 9 63.3 6 8 23.0 2	4.0 84.6 5.5 88.2 6.0 70.8	86.5 86. 88.5 90. 74.6 76. 37.0 38.	7 Cell Phon 5 8 4	2446000 603500 2669540	: Internet Differ

		Total number of	cell phones a	across all UN	Sub Regions in	2017*******										
un sub-region																
australia and new zealand	33100000															
central asia	78750000															
eastern asia	1714400000															
eastern europe	488950000															
latin america and the caribbean	668581100															
melanesia	5743000															
micronesia	95200															
northern africa	241620000															
northern america	431813000															
northern europe	125641000															
polynesia	240000															
south-eastern asia	915995000															
southern asia	1642831000															
southern europe	187323300															
sub-saharan africa	736334000															
western africa	612000															
western asia	271300000															
western europe	233773400															
Name: (number_cells, 2017), dtype	: int64															
***************************************	********Total nu	umber of Cell p	thones in the w	world under Ut	N-Sub Regions i	n 2017******		***************************************								
7697102000																
100E 100E										2012	2012	2014	2015	2016	2017	Call Chana Difference
years 1995 1996	1997	1998	1999	2999	2001	2002	2003	2010	2011	2012	2013	2014	2015	2016		Cell Phone Difference
years 1995 1996 count 1.800000e+01 1.800000e+01	1997 1.8800000e+01	1998 1.800000e+01	1999 1.800000e+01	2000 1.800000e+01	2001 1.800000e+01	2002 1.800000c+01	2003 1.800000e+01	2010 1.800000e+01	2011 1.800000e+01	1.800000e+01	1.800000e+01	1.800000e+01	1.8000000=+01	1.800000e+01	1.800000e+01	1.800000e+01
years 1995 1996 count 1.800000e+01 1.800000e+01 mean 1.102660e+06 1.681392e+06	1997 1.800000e+01 2.324167e+06	1998 1.800000e+01 3.165700e+06	1999 1.800000e+01 4.477693e+06	2000 1.800000e+01 6.362844e+06	2001 1.800000e+01 8.218179e+06	2002 1.800000e+01 9.836579e+06	2003 1.800000e+01 1.174858e+07	2010 1.800000e+01 3.637401e+07	2011 1.800000e+01 4.029215e+07	1.800000e+01 4.280315e+07	1.800000e+01 4.541216e+07	1.800000e+01 4.804369e+07	1.800000e+01 4.918868e+07	1.800000e+01 5.179235e+07	1.800000e+01 5.398858e+07	1.800000e+01 5.372774e+07
years 1995 1996 count 1.800000e+01 1.800000e+01 mean 1.102660e+06 1.681392e+06 std 2.876833e+06 3.944024e+06	1997 1.800000e+01 2.324167e+06 5.153486e+06	1998 1.800000e+01 3.165700e+06 6.717112e+06	1999 1.800000e+01 4.477693e+06 8.801550e+06	2000 1.8000000e+01 6.362844e+05 1.178165e+07	2001 1.800000e+01 8.218179e+06 1.501800e+07	2002 1.800000e+01 9.836579e+06 1.796993e+07	2003 1.8000000e+01 1.174858e+07 2.126517e+07	2010 1.800000e+01 3.637401e+07 5.372767e+07	2011 1.800000e+01 4.029215e+07 6.082340e+07	1.800000e+01 4.280315e+07 6.614046e+07	1.800000e+01 4.541216e+07 7.179184e+07	1.800000e+01 4.804369e+07 7.632656e+07	1.800000c+01 4.918868e+07 7.764204e+07	1.800000e+01 5.179235e+07 8.258637e+07	1.800000e+01 5.398858e+07 8.805482e+07	1.800000e+01 5.372774e+07 8.683329e+07
years 1995 1996 count 1.800000e+01 1.8000000e+01 mean 1.102660e+06 1.681392e+06 std 2.876833e+06 3.944024e+06 min 0.000000e+00 0.000000e+00	1997 1.800000e+01 2.324167e+06 5.153486e+06 2.000000e+01	1998 1.800000e+01 3.165700e+06 6.717112e+06 3.042500e+02	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02	2000 1.8000000e+01 6.362844e+06 1.178165e+07 4.867500e+02	2001 1.800000e+01 8.218179e+06 1.501800e+07 5.960000e+02	2002 1.800000e+01 9.836579e+06 1.796993e+07 9.017500e+02	2003 1.800000e+01 1.174858e+07 2.126517e+07 2.728500e+03	2010 1.800000e+01 3.637401e+07 5.372767e+07 1.470000e+04	2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+04	1.800000e+01 4.280315e+07 6.614046e+07 1.830000e+04	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04	1.800000c+01 4.918868e+07 7.764204e+07 2.242000e+04	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04	1.890900e+01 5.398858e+07 8.805482e+07 2.380900e+04	1.800000e+01 5.372774e+07 8.683329e+07 2.360900e+04
years 1995 1996 count 1.89999e+01 1.899999e+01 mean 1.102669e+06 1.681392e+96 std 2.876833e+06 3.944024e+96 min 0.99999e+09 0.999999e+09 25% 3.114683e+03 5.856383e+03	1997 1.800000e+01 2.324167e+06 5.153486e+06 2.000000e+01 1.213133e+04	1998 1.800000e+01 3.165700e+06 6.717112e+06 3.042500e+02 2.141550e+04	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04	2000 1.8000000e+01 6.362844e+06 1.178165e+07 4.867500e+02 1.021361e+05	2001 1.8000000e+01 8.218179e+06 1.501800e+07 5.960000e+02 2.040218e+05	2002 1.800000e+01 9.836579e+06 1.796993e+07 9.017500e+02 3.289628e+05	2903 1.809090e+01 1.174858e+07 2.126517e+07 2.728500e+03 4.680576e+05	2010 1.800000e+01 3.637401e+07 5.372767e+07 1.470000e+04 8.639559e+06	2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+04 1.024594e+07	1.800000e+01 4.280315e+07 6.614046e+07 1.830000e+04 1.143364e+07	1.899999e+91 4.541216e+97 7.179184e+97 2.319999e+94 1.252127e+97	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07	1.800000c+01 4.918868e+07 7.764204e+07 2.242000c+04 1.283887e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.300000e+04 1.276813e+07	1.899000e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07
years 1995 1996 count 1.890909e-81 1.890909e-81 mean 1.192669e-96 1.681392e-96 std 2.876833e-96 3.944924e-96 min 0.090909e-90 0.090909e-90 25% 3.114693e-93 5.856383e-93 56% 6.785714e-94 1.434934e-95	1997 1.800000e+01 2.324167e+06 5.153486e+06 2.000000e+01 1.213133e+04 2.970314e+05	1998 1.800000e+01 3.165700e+06 6.717112e+06 3.042500e+02 2.141550e+04 5.434697e+05	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06	2999 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+02 1.021361e+05 1.706570e+06	2001 1.800000e+01 8.218179e+06 1.501800e+07 5.960000e+02 2.040218e+05 2.261685e+06	2892 1.809090e+01 9.836579e+06 1.796993e+07 9.017500e+02 3.289628e+05 2.790520e+06	2003 1.800000e+01 1.174858e+07 2.126517e+07 2.728500e+03 4.680576e+05 4.314390e+06	2010 1.800000e+01 3.637401e+07 5.372767e+07 1.470000e+04 8.639559e+06 1.381115e+07	2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+04 1.024594e+07 1.432473e+07	1.800000e+01 4.280315e+07 6.614046e+07 1.830000e+04 1.143364e+07 1.445760e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07	1.800000c+01 5.398858e+07 8.805482e+07 2.380000c+04 1.276813e+07 1.615000c+07	1.890999c+01 5.372774e+07 8.683329e+07 2.369990c+04 1.219548e+07 1.584093c+07
years 1955 1950 1950 1950 1950 1950 1950 1950	1997 1.800000e+01 2.324167e+06 5.153486e+06 2.00000e+01 1.213133e+04 2.970314e+05 1.762216e+06	1998 1.800000e+01 3.165700e+06 6.717112e+06 3.042500e+02 2.141550e+04 5.434697e+05 2.746639e+06	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+02 1.021361e+05 1.706570e+06 6.087775e+06	2001 1.800000e+01 8.218179e+06 1.501800e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.900450e+06	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06	2003 1.800000e+01 1.174858e+07 2.126517e+07 2.728500e+03 4.680576e+05 4.314390e+06 9.130700e+06	2010 1.800000e+01 3.637491e+07 5.37276+07 1.470000e+04 8.639559e+06 1.381115e+07 3.860517e+07	2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+04 1.024594e+07 1.432473e+07 3.760000e+07	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1955 1950 1950 1950 1950 1950 1950 1150 11	1997 1.800000e+01 2.324167e+06 5.153486e+06 2.00000e+01 1.213133e+04 2.970314e+05 1.762216e+06	1998 1.800000e+01 3.165700e+06 6.717112e+06 3.042500e+02 2.141550e+04 5.434697e+05 2.746639e+06	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+02 1.021361e+05 1.706570e+06 6.087775e+06	2001 1.800000e+01 8.218179e+06 1.501800e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.900450e+06	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06	2003 1.800000e+01 1.174858e+07 2.126517e+07 2.728500e+03 4.680576e+05 4.314390e+06 9.130700e+06	2010 1.800000e+01 3.637491e+07 5.37276+07 1.470000e+04 8.639559e+06 1.381115e+07 3.860517e+07	2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+04 1.024594e+07 1.432473e+07 3.760000e+07	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890999c+01 5.372774e+07 8.683329e+07 2.369990c+04 1.219548e+07 1.584093c+07
years 1955 1950 1950 1950 1950 1950 1950 1950	1997 1.890992+91 2.324167e+96 5.153486e+96 2.0909090+91 1.213133e+94 2.979314e+95 1.762216e+96 1.983366e+97	1998 1.800000e+01 3.165700e+06 6.717112e+06 3.042500e+02 2.141550e+04 5.434697e+05 2.746639e+06 2.485047e+07	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1996 (august 1.898080e-14) 1.898080e-14 1.8980	1997 1.890092e+91 2.324167e+06 5.153486e+06 2.080092e+91 1.213133e+04 2.970314e+05 1.762216e+06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
year's 1955 1950 1950 1950 1950 1950 1950 1950	1997 1.890092e+91 2.324167e+06 5.153486e+06 2.080092e+91 1.213133e+04 2.970314e+05 1.762216e+06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1996 (august 1.898080e-14) 1.898080e-14 1.8980	1997 1.890092e+91 2.324167e+06 5.153486e+06 2.080092e+91 1.213133e+04 2.970314e+05 1.762216e+06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1996 (court 1.8000000-e10 1.8000000-e10 1.8000000-e10 1.8000000-e10 1.8000000-e10 51 1.8000000-e10 51 1.801392-e00 51 10.81392-e00	1997 1.890092e+91 2.324167e+06 5.153486e+06 2.080092e+91 1.213133e+04 2.970314e+05 1.762216e+06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 count 1.800090-01 1.8000090-01 mean 1.1800000-02 1.8000000-01 to 2.2000000000 1.8000000-00 to 2.200000000000 to 2.000000000000000000000000000000000000	1997 1.800000e+01 2.334167e+06 5.153486e+06 2.000000e+01 1.213133e+04 2.970314e+05 1.762216e+06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1296 count 1.389389-11 1.389389-10 1 mean 1.18938-16 1.389389-10 1 2.18938-16 3.54938-16 3.54938-16 1 2.55 3.114688-16 1.56538-16 1 2.56 6.6725-16 1.45498-16 1 2.57 7.5755-16 1.24936-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-16 1.45498-16 1 2.58 1.14688-10 1 2.58 1.14688-1 2.58 1.14688-1 2.58 1.14688-1 2.58 1.14688-1 2.58 1.14688-1 2.	1997 1.3808909c+01 2.334167e+06 5.153486e+06 2.090909c+01 1.213133e+04 2.970314e+05 1.762216e+06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1.000000000000000000000000000000000000	1997 1.800090-10 2.324167e-06 5.153486e-06 5.153486e-06 1.213133e-04 2.970334e-05 1.762216e-06 1.983366e+07	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1.000000000000000000000000000000000000	1997 1.8000902+01 2.334167e+06 5.153486e+06 5.153486e+06 1.213133e+04 2.9703314e+05 1.762216e+06 1.983366e+07 4.00000000000000000000000000000000000	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+87	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
year's 1995 count 1.89999-01 1.89999-01 mon 1.107469-05 1.691972-05 to 2.79759-05 3.64092-06 55 3.70690-09 1.65197-05 55 5.310690-09 1.65197-09 57 7.709759-06 1.29567-09 Acetage of nutber_colls in 1995 at mob-region country sestern curspe bettern curspe country sestern curspe country c	1997 1.890090-10 2.324167e-96 5.153486e-96 5.153486e-96 1.213133e-94 2.970314e-95 1.762216e-96 1.983366e-97 ************************************	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+97	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1.000 1	1997 1.898989-11 2.324167e+96 5.153486e+96 5.153486e+96 1.213133e+94 1.213133e+94 1.762216e+96 1.762216e+96 d number_cells 5642808.8 5817509.8 35159009.0	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+97	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1296 count 1.0000000001 1.0000000001 count 1.000000001 1.0000000001 count 1.000000001 1.000000000000000000000000	1997 1. 8900030e-01 2. 334167e-06 5. 1153488e-06 2. 000000e-01 1. 213133e-04 1. 22316e-05 1. 762216e-06 1. 383366e-07 1. 8000000000000000000000000000000000000	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+97	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1296 count 1.389389-11 1.389389-10 1.389389-10 1.389389-10 1.389389-10 1.389389-10 1.389389-10 1.389389-10 1.389389-10 1.38938-10 1.	1997 1.898080±01 2.324167e+06 5.153465e+06 5.153465e+06 1.908030e+01 1.213133e+04 1.762216e+06 1.762216e+06 1.983366e+07 4.983366e+07 56642000.0 56665000.0 4045 410400.0	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+97	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07
years 1995 1296 count 1.0000000001 1.0000000001 count 1.000000001 1.0000000001 count 1.000000001 1.000000000000000000000000	1997 1. 8900030e-01 2. 334167e-06 5. 1153488e-06 2. 000000e-01 1. 213133e-04 1. 22316e-05 1. 762216e-06 1. 383366e-07 1. 8000000000000000000000000000000000000	1998 1.808090e+81 3.165790e+86 6.717112e+86 3.842590e+82 2.141550e+84 5.434697e+85 2.746639e+86 2.485847e+97	1999 1.800000e+01 4.477693e+06 8.801550e+06 4.107500e+02 4.960757e+04 1.010368e+06 4.245200e+06 3.097067e+07	2000 1.800000e+01 6.362844e+06 1.178165e+07 4.867500e+05 1.021361e+05 1.706570e+06 6.087775e+06 3.924477e+07	2001 1.890909e+01 8.218179e+06 1.591890e+07 5.960000e+02 2.040218e+05 2.261685e+06 6.930450e+06 4.979900e+07	2002 1.800008e+01 9.836579e+06 1.796993e+07 9.017508e+02 3.289628e+05 2.790520e+06 7.642350e+06 6.392320e+07	2993 1.899990e+91 1.174858e+97 2.126517e+97 2.728590e+93 4.690576e+95 4.314390e+96 9.130790e+96 7.812380e+97		2011 1.800000e+01 4.029215e+07 6.082340e+07 1.587500e+07 1.432473e+07 3.750000e+07 2.350880e+08	1.800000e+01 4.280315e+07 6.614846e+07 1.830000e+04 1.143364e+07 1.445760e+07 3.934317e+07	1.800000e+01 4.541216e+07 7.179184e+07 2.310000e+04 1.252127e+07 1.476361e+07 4.096575e+07	1.800000e+01 4.804369e+07 7.632656e+07 2.153333e+04 1.280546e+07 1.515250e+07 4.111658e+07	1.800000e+01 4.918868e+07 7.764204e+07 2.242000e+04 1.283887e+07 1.566056e+07 4.109058e+07	1.800000e+01 5.179235e+07 8.258637e+07 2.843333e+04 1.276089e+07 1.587591e+07 4.097583e+07	1.800000e+01 5.398858e+07 8.805482e+07 2.380000e+04 1.276813e+07 1.615000e+07 4.073875e+07	1.890090e+01 5.372774e+07 8.683329e+07 2.360900e+04 1.219548e+07 1.594093e+07 4.069642e+07

ountry 917	afghanistan			andorra		ntigua and barbuda					united states	uruguay					yemen	zambia	zimbabwe
wi/ un sub-region	23900000	3630000	45800000	80.900	13300000	184000	61900000	3490000	26700000		400000000	100000	24300000	228999	24500000	120000000	15400000	13400000	14100000
.995 australia and new zealand	0.0	0.0	0.0	0.0	0.0	0.6	0.0	9.0	2240000.0		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
central asia	0.0	0.0	0.0	0.0	0.0	0.6					0.0	0.0		0.0	0.0	0.0	0.0		0.0
eastern asia	0.0		0.0		0.0	0.6					0.0	0.0		0.0		0.0	9.0		
eastern europe	0.0		0.0		0.0	0.6					0.0	0.0		0.0		0.0	0.0		
latin america and the caribbean	9.0	9.0	0.0		0.0		405000.0				0.0	39900.0			484080.0	0.0	9.9		
916 southern europe	0.0		0.0	76100.0	0.0	0.6					0.0	0.0		0.0		0.0	0.0		0.0
sub-saharan africa	9.9	0.0	0.0	9.0	13000000.0	0.6	0.0	0.0	0.0		0.0	0.0	0.0	0.0	0.0	0.0	0.0	12000000.0	12900000.0
western africa	0.0	0.0	0.0		0.0	0.6					0.0	0.0		0.0	0.0	0.0	0.0	0.0	0.0
western asia	0.0	9.0	0.0	0.0	0.0	0.6	0.0	3430000.0	0.0		0.0	0.0	0.0	0.0	0.0	0.0	16400000.0	0.0	0.0
western europe	9.0	9.0	0.0	0.0	0.0	0.6	0.0	0.0	0.0		0.0	0.0	0.0	0.0	0.0	0.0	9.9	0.0	0.0
ustria 87.9 elgium 87.7 rance 80.5 ermany 84.4 ichtentrin 90.5 ucombourg 97.4 enaco 97.1 etherlands 93.2 eldtzerland 89.7 mane: (internet_usors, 2017), dtype:																			
ountry	************		••••••	***Plot 2 E	ata******	•••••	•••••	•••••	••••••	•••••	••••								
ustria 10900000																			
rance 69000000																			
ermany 110000000																			
iechtenstein 46400																			
iechtenstein 46400 uxembourg 794000																			
iechtenstein 46400 .uxembourg 794000 ionaco 33000																			
iechtenstein 46400 uxembourg 794000																			

Handling of Incorrect Inputs Screenshot

```
NOIEs:
1) User can quit at several prompts by entering "quit".
2) Inputs are case in-sensitive.
3) Sub-regions can be selected by entering the correct labe/name or by numeric index provided.
4) Years entered identify specific years of data to be averaged for the data type selected.
Enter "quit" to exit the program, or any key to continue:
For Reference: Valid UN sub-region names and index values:
             un sub-region
australia and new zealand
                               central asia
eastern asia
3 eastern europe
4 latin america and the caribbean
                          melanesia
micronesia
northern africa
                        northern america
northern europe
                     polynesia
south-eastern asia
                          southern asia
southern europe
                      sub-saharan africa
western africa
                           western asia
western europe
Please enter a UN Sub-Region name or numeric index value: atlantis Invalid UN Sub-Region.
Enter "help" to see the UN Sub-Region table, or any other key to continue: Please enter a UN Sub-Region name or numeric index value: western europe western europe
For Cell Phone data enter "1", for Internet Usage (%) enter "2", or "quit" to exit: 1 Selection: number of cell phones data
Enter the 1st specific year of data for avarage calc: 2 Invalid year entered. Please try again, or enter "quit" to exit...
Enter the 1st specific year of data for avarage calc: 1995 1995
Enter 2nd specific year of data for avarage calc: 2017
User inputs summary:
Quit? (or any key to continue): , UN Sub-Region: western europe, Data: number_cells, Year 1: 1995, Year 2: 2017
There are values missing within the data, calculations were performed on the data assuming a zero value was present
```