Cleaning and Storing Supply Chain Delivery Data

DSCS6020 Collect/Store/Retrieve Data

Abhiram Paruchuri 100933135

Introduction

Project Proposal:

The proposed project idea is to clean a supply chain delivery data set, obtained from The United States President's Emergency Plan for Aids Relief (pepfar). The organization's main goal is to provide better health care and testing facilities to the places that are lacking them. The data consists of all the deliveries made by the pepfar delivery network over the past few years. Cleaning and preprocessing the data can help in using the data to observe trends in the medical needs and help in effectively predicting the inventory. As supply chain is a constantly evolving process, effectively forecasting the changes and trends can help in developing a better plan to handle the needs. In this case, it is the HIV medications and test kits. In addition, storing the dataset into a database can make visualizing the data easier.

In addition, I have also retrieved the HIV/AIDS prevalence rates, deaths caused by AIDS, number of people infected from the data on the size of HIV epidemic from WHO global health observatory data repository. The idea is to collect all the available data on the epidemic and combine it into a single repository. This will make analysis easier and improve the scope of the supply chain network by providing answers to the questions like, which is the most affected country?, how many are in need of medical aid? Etc.

Brief introduction to prepfar:

The U.S. President's Emergency Plan for AIDS Relief (PEPFAR) is the U.S. Government initiative to help save the lives of those suffering from HIV/AIDS around the world. This historic commitment is the largest by any nation to combat a single disease internationally, and PEPFAR investments also help alleviate suffering from other diseases across the global health spectrum. PEPFAR is driven by a shared responsibility among donor and partner nations and others to make smart investments to save lives.

PEPFAR is the cornerstone and largest component of the U.S. President's Global Health Initiative. With a special focus on improving the health of women, newborns and children, the Global Health Initiative's goal is to save the greatest number of lives by increasing and building upon what works and, then, supporting countries as they work to improve the health of their own people.

Brief introduction to SCMS:

The SCMS aim is to provide cost-effective, reliable, secure and sustainable supply chains for developing countries, which can be a huge thing for people suffering with HIV. For more than six years, the Supply Chain Management System (SCMS) has been saving lives through stronger supply chains. In collaboration with in-country and international partners, SCMS:

- Provides quality, best-value health care products to those who need them.
- Deploys innovative solutions to help programs enhance their supply chain capacity.
- Ensures accurate supply chain information is collected, shared and used.

Cleaning the Data

Dataset Overview:

There are two data sets that are being derived. They are:

- 1. The SCMS scheduled delivery data set
- 2. Dataset created from WHO country statistics.

The SCMS scheduled delivery dataset:

First, we deal with the SCMS scheduled delivery data set, which comprises of procurement transaction data from the Supply Chain Management System (SCMS), administered by the United States Agency for International Development (USAID), i.e. data about all the deliveries made, which are related to medical needs and facilities like, drugs, tests and suspension kits.

The dataset captures deliveries of antiretroviral (ARV)drugs, rapid diagnostic tests for HIV and malaria , and anti-malaria medicines, including prices and delivery destinations. This data is particularly valuable for understanding ranges and trends in pricing as well as

First we will review the glossary of the data set and talk about the problems we have with it.

ID	FieldName	FieldDescription
1	ID	Primary key identifier of the line of data in our analytical tool
2	Project Code	Project code
3	PQ#	Price quote (PQ) number
4	PO #	Order number: Purchase order (PO) for Direct Drop deliveries, or Sales Order (SO) for from Regional Delivery Center (RDC) deliveries
5	ASN/DN#	Shipment number: Advanced Shipment Note (ASN) for Direct Drop deliveries, or Delivery Note (DN) for from RDC deliveries
6	Country	Destination country
7	Managed By	SCMS managing office: either the Program Management Office (PMO) in the U.S. or the relevant SCMS field office
8	Fulfill Via	Method through which the shipment was fulfilled: via Direct Drop from vendor or from stock available in the RDCs
9	Vendor INCO Term	The vendor INCO term (also known as International Commercial Terms) for Direct Drop deliveries
10	Shipment Mode	Method by which commodities are shipped
11	PQ First Sent to Client Date	Date the PQ is first sent to the client
12	PO Sent to Vendor Date	Date the PO is first sent to the vendor

	Scheduled Delivery	
13	Date	Current anticipated delivery date
14	Delivered to Client Date	Date of delivery to client
		Date on which delivery to client was recorded in SCMS
15	Delivery Recorded Date	information systems
16	Product Group	Product group for item, i.e. ARV, HRDT
		Identifies relevant product sub classifications, such as
		whether ARVs are pediatric or adult, whether a malaria
17	Sub Classification	product is an artemisinin-based combination therapy (ACT),
17	Sub Classification	etc.
18	Vendor	Vendor name
10	VCHUOI	Product name and formulation from Partnership for Supply
19	Item Description	Chain Management (PFSCM) Item Master
20	Molecule/Test Type	Active drug(s) or test kit type
21	Brand	Generic or branded name for the item
22	Dosage	Item dosage and unit
	Бозаве	reem dosage and ame
		Dosage form for the item (tablet, oral solution, injection,
23	Dosage Form	etc.).
	Unit of Measure (Per	
24	Pack)	Pack quantity (pills or test kits) used to compute unit price
25	Line Item Quantity	Total quantity (packs) of commodity per line item
	, ,	
26	Line Item Value	Total value of commodity per line item
		Cost per pack (i.e. month's supply of ARVs, pack of 60 test
27	Pack Price	kits)
28	Unit Price	Cost per pill (for drugs) or per test (for test kits)
		Identifies manufacturing site for the line item for direct drop
29	Manufacturing Site	and from RDC deliveries
		Designates if the line in question shows the aggregated
	<u> </u>	freight costs and weight associated with all items on the
30	First Line Designation	ASN/DN
31	Weight (Kilograms)	Weight for all lines on an ASN/DN
22	F	Freight charges associated with all lines on the respective
32	Freight Cost (USD)	ASN/DN

	Line Item Insurance	Line item cost of insurance, created by applying an annual
33	(USD)	flat rate (%) to commodity cost

Problems faced with the dataset:

As this is the data that is made available for public, there are some key limitations to the data. Multiple deliveries are consolidated into single delivery, so freight cost is inaccurate relative to the line item. The dates are not properly documented. As supply chain data vary for different transactions and the data changes overtime. So NA values are a problem. Some dates are not applicable as the goods for those order are fulfilled out of stock available at RDC (Regional Development Center). The Line Item Value/Pack Price/Unit Price values have high variability as they can be sometimes low due to donations.

Cleaning the data:

First we clean the columns with the date values. U sing the lubridate package, we convert all the values in the date columns to date classes. The columns (PO.Sent.to.Venfor.Date, Sheduled.Delivery.Date, Delivered.to .Client.Date) are converted to numeric values, so as to find the delivery time and the delay time by subtracting the PO sent date and scheduled date from the delivered date. If the delayed time value is a negative vale, then the product arrived earlier than expected.

In the program, it is done by constructing a deriveDates() function Screen shots:

Result:

PO.Sent.to.Vendor.Date	Scheduled.Delivery.Date	Delivered.to.Client.Date	Delivery.Recorded.Date
Date Not Captured	2-Jun-06	2-Jun-06	2-Jun-06
Date Not Captured	14-Nov-06	14-Nov-06	14-Nov-06
Date Not Captured	27-Aug-06	27-Aug-06	27-Aug-06
Date Not Captured	1-Sep-06	1-Sep-06	1-Sep-06
Date Not Captured	11-Aug-06	11-Aug-06	11-Aug-06
Date Not Captured	28-Sep-06	28-Sep-06	28-Sep-06
Date Not Captured	8-Jan-07	8-Jan-07	8-Jan-07
Date Not Captured	24-Nov-06	24-Nov-06	24-Nov-06
Date Not Captured	7-Dec-06	7-Dec-06	7-Dec-06
11/13/2006	30-Jan-07	30-Jan-07	30-Jan-07
12/1/2006	16-Feb-07	16-Feb-07	16-Feb-07
Date Not Captured	8-Jan-07	8-Jan-07	8-Jan-07
Date Not Captured	10-Jan-07	10-Jan-07	10-Jan-07
12/22/2006	27-Feb-07	27-Feb-07	27-Feb-07
Date Not Captured	18-Jan-07	18-Jan-07	18-Jan-07
1/10/2007	19-Mar-07	19-Mar-07	19-Mar-07
Date Not Captured	7-May-07	7-May-07	7-May-07
Date Not Captured	29-Mar-07	29-Mar-07	29-Mar-07
4/12/2007	6-Jun-07	6-Jun-07	6-Jun-07
5/13/2007	19-Jun-07	19-Jun-07	19-Jun-07
5/17/2007	19-Jun-07	19-Jun-07	19-Jun-07
7/13/2007	2-0ct-07	2-Oct-07	2-Oct-07
7/4/2007	15-Oct-07	15-0ct-07	15-Oct-07

The dates before cleaning

PO_Sent_to_Vendor_Date	Scheduled_Delivery_Date	Delivered_t_Client_Date	Year	Delivery_Time	Delay_Time
IA	2006-06-02	2006-06-02	2006	0	0
NA .	2006-11-14	2006-11-14	2006	0	Θ
NA .	2006-08-27	2006-08-27	2006	0	8
NA .	2006-09-01	2006-09-01	2006	Θ	Θ
NA .	2006-08-11	2006-08-11	2006	0	8
NA .	2006-09-28	2006-09-28	2006	Θ	Θ
IA .	2007-01-08	2007-01-08	2007	0	8
NA .	2006-11-24	2006-11-24	2006	0	Θ
IA .	2006-12-07	2006-12-07	2006	0	0
2006-11-13	2007-01-30	2007-01-30	2007	78	0
2006-12-01	2007-02-16	2007-02-16	2007	77	0
IA	2007-01-08	2007-01-08	2007	0	0
iA.	2007-01-10	2007-01-10	2007	0	0
2006-12-22	2007-02-27	2007-02-27	2007	67	8
IA	2007-01-18	2007-01-18	2007	Θ	Θ
997-91-19	2007-03-19	2007-03-19	2007	68	0
IA	2007-05-07	2007-05-07	2007	Θ	Θ
IA	2007-03-29	2007-03-29	2007	0	0
997-94-12	2007-06-06	2007-06-06	2007	55	0
2007-05-13	2007-06-19	2007-06-19	2007	37	0
1007-05-17	2007-06-19	2007-06-19	2007	33	0
1007-07-13	2007-10-02	2007-10-02	2007	81	0
997-97-94	2007-10-15	2007-10-15	2007	103	8
1997-97-94	2007-08-27	2007-08-27	2007	54	0
997-97-26	2007-08-13	2007-08-21	2007	26	8
997-97-26	2007-08-25	2007-08-25	2007	30	Θ
IA	2007-10-16	2007-10-16	2007	0	0
IA	2007-11-22	2007-11-22	2007	0	Θ
IA	2007-11-22	2007-11-22	2007	0	0
1007-10-03	2007-11-20	2007-11-20	2007	48	0
2007-08-28	2007-10-03	2007-10-03	2007	36	0
987-11-12	2008-01-29	2008-01-29	2008	78	8
2007-11-19	2008-01-21	2008-01-21	2008	63	0
2007-11-21	2008-01-21	2008-01-21	2008	61	8
2007-12-10	2008-01-31	2008-01-31	2008	52	0
IA	2008-02-05	2008-02-05	2008	0	0
2008-01-04	2008-01-21	2008-01-04	2008	Θ	-17

The cleaned dataset for dates

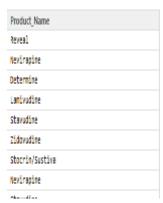
Assigning the Product name:

Now we assign a proper name to the product from the 'Molecule Test Type' column and the Brand name column. We parse through both the columns, and choose the brand name when it is present and choose the molecule name when the value is listed as generic as it means that the molecule name is generally the product name.

Result:

Molecule.Test.Type	Brand
HIV, Reveal G3 Rapid HIV-1 Antibody Test	Reveal
Nevirapine	Generic
HIV 1/2, Determine Complete HIV Kit	Determine
Lamivudine	Generic
Stavudine	Generic
Zidovudine	Generic

A .The above figure has values before merging



B. values after merging

Splitting the drug into its parts:

Several drugs contain two or more types of molecules in them. By constructing the splitElements() function, we split the Molecules and dosages into separate columns. The use of splitting the drugs into their classes is to increase the scope of analysis on the drug type and combination.

Result:

Efavirenz/Emtricitabine/Tenofovir Disoproxil Fumarate	Atripla	600/200/300mg
Quinine (as dihydrochloride)	Generic	600mg/2ml
HIV 1/2, Uni-Gold HIV Kit	Uni-Gold	N/A
Lamivudine/Zidovudine	Generic	150/300mg
Lamivudine	Generic	150mg
Lamivudine/Nevirapine/Stavudine	Generic	150/200/30mg
Lamivudinė/Zidovudine	Generic	150/300mg
Lopinavir/Ritonavir	Kaletra	80/20mg/ml
Lopinavir/Ritonavir	Generic	100/25mg

The above table consists of molecule names and drug dosages before cleaning

Lopinavir	Ritonavir	NA
Malaria Paramax-3 Kit	NA	NA
Efavirenz	Emtricitabine	Tenofovir Disoproxil Fumarate
Quinine (as dihydrochloride)	NA	NA
test	NA	NA
Lamivudine	Zidovudine	NA
Lamivudine	NA	NA
Lamivudine	Nevirapine	Stavudine
Lamivudine	Zidovudine	NA
Lopinavir	Ritonavir	NA
Lopinavir	Ritonavir	NA
test	NA	NA

Molecules split into their respective columns

263.40	200.0	50.0
158.05	300.0	0.0
68.63	30.0	0.0
13.72	0.0	0.0
H24.H9	200.0	50.0
30.07	100.0	0.0
2.12	80.0	20.0
5.52	300.0	200.0
38.75	300.0	0.0
41.16	0.0	0.0
1176.00	0.0	0.0
1002.74	0.0	0.0
757.31	600.0	0.0
215.60	0.0	0.0
1.20	100.0	0.0
76.49	150.0	300.0

Dosages split into their respective columns

Assigning the dosages to their respective units & donation status:

As many data analysis tools are sensitive to units of measure, we split the dosages in the column, which are jumbles into mg/ml, mg and g values to their respective unit columns

We also construct a new column to decide if the line item is a donation or not. The condition set is, if the pack price is less than 2.5 dollars, we designate it as a donation.

Result:

Dosage	
N/A	
10mg/ml	
N/A	
150mg	
30mg	
10mg/ml	
200mg	
200mg	
30mg	

Dosage_mg
NA
NA
NA
150
30
NA
200
200
30
200/50
200/50
NA
NA
150/300

Dosage Column

before

After

Donation_Designation
N
N
N
N
N
N
N
N
Υ
N
N
N
At .

Donation Column

Final note:

Finally the cleaned data is loaded into a data frame that can be loaded into a database. Some of the columns with redundant values and useless columns are left out while populating the new data frame as they are considered less important for data analysis or further inference.

Not considering the columns above, the other columns that are considered are, ID, project code, PO.SO, ASN, Country, Managed, Fulfill via, shipment mode, vendor, unit of measure per pack, Item Quantity, Pack price, first line designation, weight in kilograms, insurance.

The NA values are converted to zeros rather than omitting as they are very widely dispersed.

Second dataset retrieval:

The second data set is constructed using three datasets that are downloaded in program from the WHO Global Health Observatory Data Repository. The cleaning is done in the following steps:

- a. The data sets for country wise HIV prevalence, deaths and number of people infected are downloaded and opened in R in the form of a data frame.
- b. We observe that all the columns are polluted with '&It' values, which might have occurred during retrieving.
- c. We parse through each values in the three data frames and remove the '<' values from the values.
- d. Now we notice that the values for the prevalence, deaths and the number of people infected and their confidence interval values are present in the same columns.
- e. We construct a function numSplit() to separate these values into two separate columns, while removing the white spaces and the square brackets present in the value.
- f. Now we merge all the three datasets in to a single data set using merge() function. First prevalence and deaths data are merged on prevalence country column and then the merged data is merged with number of infected people data set with respect to the same prevalence country column so that we have a common binder.
- g. NA values are removed using Na.omit() function as they are not dispersed, but are present in the form of rows.
- h. All the data cleaning is done by the dataClean() function the code. The rename() function renames the data after cbind.

Screen shots:

The screen shots are provided only for the prevalence data set. It is the same for the other data sets also.

Country	X.2013	X.2009	X.2005	X.2001
Afghanistan	<0.1 [<0.1-0.1]	<0.1 [<0.1-0.1]	<0.1 [<0.1-0.1]	<0.1 [<0.1-0.1]
Albania	<0.1 [<0.1-0.1]	<0.1 [<0.1-0.1]	<0.1 [<0.1-0.1]	<0.1 [<0.1-0.1]
Algeria	0.1 [0.1-0.2]	0.1 [<0.1-0.2]	0.1 [<0.1-0.2]	0.1 [<0.1-0.2]
Angola	2.4 [1.7-3.2]	2.1 [1.4-2.9]	1.9 [1.3-2.6]	1.8 [1.2-2.4]
Argentina	No data	No data	No data	No data
Armenia	0.2 [0.1-0.3]	0.1 [0.1-0.2]	0.1 [0.1-0.2]	0.1 [<0.1-0.2]
Australia	0.2 [0.1-0.2]	0.1 [0.1-0.2]	0.1 [0.1-0.2]	0.1 [0.1-0.2]
Austria	No data	No data	No data	No data
Azerbaijan	0.2 [0.1-0.2]	0.1 [0.1-0.2]	0.1 [0.1-0.1]	<0.1 [<0.1-0.1]
Bahamas	3.2 [3.1-3.5]	3.4 [3.4-3.5]	3.4 [3.4-3.5]	3.5 [3.4-3.5]
Bangladesh	<0.1 [<0.1-0.1]	<0.1 [<0.1-<0.1]	<0.1 [<0.1-<0.1]	<0.1 [<0.1-<0.1]
Barbados	0.9 [0.7-1.2]	0.9 [0.7-1.2]	0.8 [0.7-1.1]	0.8 [0.6-1.0]
Belarus	0.5 [0.5-0.5]	0.3 [0.3-0.4]	0.2 [0.2-0.3]	0.1 [0.1-0.2]
Belgium	No data	No data	No data	No data
Belize	1.5 [1.3-1.7]	1.6 [1.5-1.7]	1.7 [1.6-1.8]	1.8 [1.5-2.2]
Benin	1.1 [1.1-1.2]	1.2 [1.1-1.3]	1.3 [1.2-1.4]	1.5 [1.4-1.7]
Bhutan	0.1 [0.1-0.4]	0.1 [0.1-0.2]	0.1 [0.1-0.1]	<0.1 [<0.1-0.1]
Bolivia (Plurinational State of)	0.2 [0.1-0.4]	0.3 [0.2-0.5]	0.4 [0.3-0.6]	0.6 [0.4-0.8]
Botswana	21.9 [20.8-23.1]	23.6 [22.5-24.9]	25.4 [24.3-26.7]	27.7 [26.5-29.0]
Brazil	0.5 [0.5-0.6]	No data	No data	No data
Bulgaria	No data	No data	No data	No data
Burkina Faso	0.9 [0.8-1.1]	1.0 [0.9-1.2]	1.3 [1.2-1.5]	2.2 [1.8-2.5]
Burundi	1.0 [0.9-1.1]	1.5 [1.3-1.6]	2.1 [1.9-2.4]	2.5 [2.2-2.7]
Câto d'Ivoino	2 7 52 4 2 01	2 ([2 2 2 0]	F 4 [4 7 F F]	C C [C O 7 2]

Prevalence dataset before cleaning

prevalence_Country	prevalence_2013	prevalence_CI_2013	prevalence_2009	prevalence_CI_2009	prevalence_2005	prevalence_CI_2005	prevalence_2001	prevalence_CI_2001
Afghanistan	0.1	0.1-0.1	0.1	0.1-0.1	0.1	0.1-0.1	0.1	0.1-0.1
Albania	0.1	0.1-0.1	0.1	0.1-0.1	0.1	0.1-0.1	0.1	0.1-0.1
Algeria	0.1	0.1-0.2	0.1	0.1-0.2	0.1	0.1-0.2	0.1	0.1-0.2
Angola	2.4	1.7-3.2	2.1	1.4-2.9	1.9	1.3-2.6	1.8	1.2-2.4
Argentina	Nodata	NA	Nodata	NA	Nodata	NA	Nodata	NA
Armenia	0.2	0.1-0.3	0.1	0.1-0.2	0.1	0.1-0.2	0.1	0.1-0.2
Australia	0.2	0.1-0.2	0.1	0.1-0.2	0.1	0.1-0.2	0.1	0.1-0.2
Austria	Nodata	NA	Nodata	NA	Nodata	NA	Nodata	NA
Azerbaijan	0.2	0.1-0.2	0.1	0.1-0.2	0.1	0.1-0.1	0.1	0.1-0.1
Bahamas	3.2	3.1-3.5	3.4	3.4-3.5	3.4	3.4-3.5	3.5	3.4-3.5
Bangladesh	0.1	0.1-0.1	0.1	0.1-0.1	0.1	0.1-0.1	0.1	0.1-0.1
Barbados	0.9	0.7-1.2	0.9	0.7-1.2	0.8	0.7-1.1	0.8	0.6-1.0
Belarus	0.5	0.5-0.5	0.3	0.3-0.4	0.2	0.2-0.3	0.1	0.1-0.2
Belgium	Nodata	NA	Nodata	NA	Nodata	NA	Nodata	NA
Belize	1.5	1.3-1.7	1.6	1.5-1.7	1.7	1.6-1.8	1.8	1.5-2.2
Benin	1.1	1.1-1.2	1.2	1.1-1.3	1.3	1.2-1.4	1.5	1.4-1.7
Bhutan	0.1	0.1-0.4	0.1	0.1-0.2	0.1	0.1-0.1	0.1	0.1-0.1
Bolivia (Plurinational State of)	0.2	0.1-0.4	0.3	0.2-0.5	0.4	0.3-0.6	0.6	0.4-0.8
Botswana	21.9	20.8-23.1	23.6	22.5-24.9	25.4	24.3-26.7	27.7	26.5-29.0
Brazil	0.5	0.5-0.6	Nodata	NA	Nodata	NA	Nodata	NA
Bulgaria	Nodata	NA	Nodata	NA	Nodata	NA	Nodata	NA
Burkina Faso	0.9	0.8-1.1	1.0	0.9-1.2	1.3	1.2-1.5	2.2	1.8-2.5
Burundi	1.0	0.9-1.1	1.5	1.3-1.6	2.1	1.9-2.4	2.5	2.2-2.7

Prevalence Dataset after cleaning

CI_2009	prevalence_2005	prevalence_CI_2005	prevalence_2001	prevalence_CI_2001	Deaths_2013	Deaths_Cl_2013	Deaths_2009	Deaths_Cl_2009	Deaths_2005	Deaths_CI_2005	Deaths_2001	Deaths_O_2001	numberinf_2013	numberinf_CI_2013	numberinf_2009	numberinf_CI_2009
	0.1	0.1-9.1	8.1	0.1-0.1	588	200-1100	588	288-1888	299	188-588	288	188-588	4588	1788-17888	3488	1589-12888
	8.1	8.1-8.1	8.1	₹.1-0.1	188	188-188	188	180-198	188	188-186	188	189-188	1999	588-1180	588	589-1888
	0.1	0.1-0.2	8.1	0.1-0.2	1488	1888-3388	1188	588-3588	1000	588-3788	1888	588-3888	25888	13888-43888	29888	9389-43888
	1.9	1.3-2.6	1.8	1,2-2.4	12990	6388-18888	9888	5388-15888	11999	7288-15888	7888	3386-12888	258888	199899-348899	298988	148686-278688
	0.1	0.1-0.2	8.1	0.1-0.2	200	288-588	288	188-588	299	188-588	188	188-288	3788	2488-5988	2998	1889-4688
	₹.1	8.1-8.2	0.1	₹.1-0.2	188	188-288	188	186-286	188	188-288	588	189-588	29800	26888-34888	25888	23866-38888
	0.1	0.1-0.1	8.1	0.1-0.1	1888	588-1888	588	588-1888	588	288-588	188	188-288	9298	6708-12808	6988	4989-9588
	3.4	3.4-3.5	3,5	3.4-3.5	1888	1888-1888	588	588-588	588	588-586	588	586-1988	7798	7388-8388	7688	7380-7980
	8.1	0.1-9.1	8.1	0.1-0.1	588	289-3888	588	289-1888	588	188-588	188	188-588	9588	4188-57888	7888	3389-38888
	8.8	8.7-1,1	8.8	8.6-1.8	188	188-188	188	180-180	188	188-186	188	189-288	1700	1398-2298	1988	1200-2000
	Đ.2	0.2-9.3	8.1	8.1-8.2	1888	1888-1288	1888	1800-1900	1888	588-1888	588	299-598	25888	24888-27888	19888	17888-21888
	1.7	1.6-1.8	1.8	1,5-2,2	280	188-288	288	200-200	288	288-286	188	189-288	3388	2988-3688	3888	2880-3380
	1.3	1.2-1,4	1.5	1.4-1.7	2798	2388-3280	3188	2789-3688	4888	4488-5488	4188	3688-4688	74888	£3000-\$0000	67888	63888-73888
	₹.1	8.1-8.1	0.1	₹.1-0.1	188	188-188	188	180-180	188	188-186	188	189-198	1900	588-2186	588	589-1888
	8.4	8.3-8.6	8.6	8.4-8.8	1298	1888-2988	1888	1389-4888	2299	1688-5388	2100	1488-5688	15888	7988-33888	19888	12000-36000
	25.4	24.3-26.7	27.7	26.5-29.8	5888	5888-6588	E200	7298-9788	14800	13000-16000	20000	19888-22888	328888	318888-348888	318886	290000-320000
	1.3	1.2-1,5	2.2	1,8-2,5	5888	4688-7388	9888	6588-9988	14888	12000-17000	19888	16889-23888	119988	100000-130000	120000	188888-138888
	2.1	1.3-2.4	2,5	1.2-2.7	4788	3988-5686	7299	6299-8599	7988	6889-9286	5788	4688-7188	E3888	76888-91888	97000	87868-118888
	5.1	4.7-5.5	6.6	5.9-7.3	29886	25888-32888	37888	33888-43888	53000	46888-68888	51888	44888-61888	378888	338889-418888	440000	418888-438888
	8.4	0.4-8.5	8.5	₹.4-8.6	188	188-188	188	180-180	188	188-286	188	189-288	1500	1300-1800	1488	1200-1700
	1.2	0.4-2.4	1.6	0.6-2.9	2298	1888-4888	3988	2000-9300	7688	1988-16888	6299	1388-18888	75888	41888-138888	82000	41888-158888
	5.2	4.9-5.5	5.3	4.9-5.7	44888	48888-45888	46888	42986-52888	45888	44888-54888	38888	34888-44888	588888	560000-650000	688888	568888-648888
	6.8	5.9-7.8	8.7	7,2-10,6	11886	9588-12888	11888	9380-13880	16888	13888-28888	15888	11888-28888	129988	110000-130000	148888	120000-150000
	3.6	3.1-4.3	3,5	3.8-4.2	15886	12008-15000	13888	11696-16888	16999	13000-15000	11888	8789-13888	219999	178888-258888	228888	190000-250000
	0.3	0.2-8.5	8.4	8.2-8.5	1888	289-1688	1388	588-2588	1788	588-3188	2100	1888-3388	38888	23888-59888	33888	20000-40000
	3.9	3.6-4.2	5.8	4.6-5.6	5488	4999-6000	£988	6299-7688	9388	3488-18888	9688	8589-11888	£9888	64888-75888	77988	71869-82888
	€.2	0.1-0.2	8.2	0.2-0.3	588	289-588	288	189-588	200	288-588	288	188-288	7688	5488-9288	7188	4889-1588

Sample of merged dataset (NA omitted)

Creating a NoSQL database

A columnar type NoSQL database (mongoDB) is used to store the two retrieved data sets and we perform a set of queries on the databases to see that all the data is loaded correctly. Two data bases are used for storage as the two data sets are of a completely different purposes but have the same aim. So the data retrieved from the two databases simultaneously can give insight to new results as one has the medical supply data and the other has the statistical facts.

Code: mongoData<-mongo("SCMS") mongoData\$insert(supplyData) mongoData1<-mongo("HIVdata") mongoData1\$insert(finalData) Queries:</pre>

```
> mongoData$count('{"Country":"vietnam"}')
[1] 2064
> yeard<-mongoData$find('{"Year":2008}')
Imported 2216 records. Simplifying into dataframe...
> dat <- mongoData$find('{"Product_Name":"Nevirapine"}', fields = '{"_id":0,"Pack_Price":1, "Country":1,"Year":1}')
Imported 1698 records. Simplifying into dataframe...
> countryd<-mongoData1$find('{"Country":"Afghanistan"}')
Imported 2 records. Simplifying into dataframe...
> mxdel<-mongoData$aggregate('[{"$group":{"_id":"$Country", "count": {"$sum":1},"max":{"$max":"$Delivery_Time"}}}]')
Imported 45 records. Simplifying into dataframe...
> view(veard)
```

Queries with head data:

```
mongoDataScount('("Country":"Vietnam")')
1] 2004
yeard--mongoDataSfind('("vear":2008)')
Imported 2216 records. Simplifying into dataFrame...
head(yeard)
ID project.code PO_SD ASN_DN COUNTRY M.
262 116-2A-701 SCMS-1070 ASN-1251 South Africa
270 108-WN-701 SCMS-14190 ASN-1192 vietnam
270 108-WN-701 SCMS-14200 ASN-1171 vietnam
274 107-RW-701 SCMS-14200 ASN-1251 SWANDA
305 123-NG-701 SCMS-13590 ASN-1258 Nigeria
343 116-ZA-701 SCMS-15950 ASN-1520 South Africa
pelivery.time pelay_time Product.proup
                                                                                                                                                                                                                                                                                                                    scheduled_pelivery_pate
2008-01-28 19:00:00
2008-01-20 19:00:00
2008-01-20 19:00:00
2008-01-30 19:00:00
2008-02-04 19:00:00
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      Delivered_t_Client_Date

2008-01-28 19:00:00

2008-01-20 19:00:00

2008-01-20 19:00:00

2008-01-30 19:00:00

2008-01-30 19:00:00

2008-01-03 19:00:00
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          2008-01-03 19:00:00
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                2008-01-20 19:00:00
                                                                                                                                                                                                                           OSDARTA AFTICA PMO - US DIFECT DEED N/A 2008-01-03 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-01-20 19:00:00 2008-0
   61
52
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            solution
Tablet
Tablet
     Line_Item_Quantity
40
1000
400
1500
650
                                                                                                                                                   7.94
17.00
8.38
0.01
0.01
     Molecule_Second Molecule_Third
                                                                         <NA>
                                                                                                                                                                     <1400
                                                                           <NA>
                                                                                                                                                                       <NA>
                                                                                                                                                                       <NA>
                                                                                                                                                                       -NAS
                                                                                                                                                                       -NAS
```

Queries being run to test the database Continuation of test Queries

```
- dat <- mongopata$find('("Product_Name":"Nevirapine")', fields = '{"_id":0,"Pack_Price":1, "country":1,"Year":1)')
Imported 1698 records. Simplifying into dataframe...
 head(dat)
       Country Year Pack_Price
Vietnam 2006 6.20
                            6.20
      Tanzania 2006
       Vietnam 2007
                            1.92
       Vietnam 2007
       vietnam 2008
                            2.90
6 South Africa 2008 4.90 
> countryd<-mongoDatalSfind('{"Country":"Afghanistan"}')
 Imported 2 records. Simplifying into dataframe...
  Country prevalence_2013 prevalence_CI_2013 prevalence_2009 prevalence_CI_2009 prevalence_2005 prevalence_CI_2005 prevalence_2001 prevalence_CI_2001 Deaths_2013 Afghanistan 0.1 0.1-0.1 0.1 0.1-0.1
> head(countryd)
                                                                                                                                                                0.1-0.1
                                            0.1-0.1
2 Afghanistan
                            0.1
                                                                   0.1
                                                                                   0.1-0.1
                                                                                                         0.1
                                                                                                                          0.1 - 0.1
                                                                                                                                                0.1
                                                                                                                                                                                   500
 Deaths_CI_2013 Deaths_2009 Deaths_CI_2009 Deaths_CI_2005 Deaths_CI_2005 Deaths_CI_2001 Deaths_CI_2001
                                                                                                                            numberInf_CI_2013 numberInf_2009 numberInf_CI_2009
        200-1100
                                     200-1000
200-1000
                                                        200
200
                                                                    100-500
100-500
                           500
                                                                                       200
                                                                                                   100-500
                                                                                                                       4500
                                                                                                                                    1700-17000
                                                                                                                                                           3400
                                                                                                                                                                        1500-12000
        200-1100
                           500
                                                                                                                                                                        1500-12000
 numberInf_2005 numberInf_CI_2005 numberInf_2001 numberInf_CI_2001
head(mxdel)
_id count max
       Angola 21 183
Nigeria 3582 341
Tanzania 1557 384
 côte d'Ivoire 1083 NA
Haiti 1965 545
```

Problems Faced:

The main problems that I faced is when I tried to combine the two data sets that I retrieved. I have tried very hard to somehow embed into each other, but then I could only complete the process till the merging of the three datasets retrieved from WHO data. I think the main problem lies in the fact that a nested loop is necessary to embed the HIV statistics in to the supply chain data by inspecting each element and then merging them. The process took a lot of time and I finally abandoned that idea to create two separate datasets

Conclusion:

- Hence we can see that the supply chain data which is available publicly can be used for conducting
 effecting analytics if proper preprocessing is performed on the data.
- We also realize the supply chain data can also be linked to various factors in a given environment.
- If properly realized, the HIV medical supplies supply chain network can be properly optimized to
 the HIV prevalence rate of a country, so that the people who need help the most will receive it.
 But this cannot be achieve without further improvement in the data cleaning and acquisition
 systems.

References:

- a. https://cran.r-project.org
- b. http://scms.pfscm.org/scms/about
- c. http://www.pepfar.gov/
- d. http://www.who.int/en/
- e. Class Notes 'Collecting, Storing and Retrieving Data'-Yatish Jain and Martin Schedlbauer