

Contents

Appendix-Power Query guide.....	2
MACROS used	2
Walkthrough	3
Produce Excel files from “.dat” files.....	3
Produce Excel files from “.pla/.plw” files.....	19
Produce Excel file from unstructured files.....	20
Appendix-Data Wrangling in Python	21
Import Packages	21
Two step process: excel files to csv & merging of csv files	21
Number of decimals changed to six	22
Drop column in for csv containing SAMPLES data	23
Create empty column for SAMPLES csv file for future FK population	23
Create empty column for SAMPLES csv file for alternative database for future incremented PK population	24
Count number of rows in csv files	24
Appendix-Zooplankton Database Creation.....	24
STAGE_RECORDS, bulk insert & RECORDS table.....	25
Creation of CRUISES table.....	28
STAGE_SAMPLES & bulk insert	28
Foreign key relationship creation between CRUISES and RECORDS tables.....	29
Fixes in STAGE_SAMPLES table, creation of FK, creation of SAMPLES table & insert of data.....	30
Creation of Species_info table	31
Changes on tables before Foreign Key relationships creation	32
Last updates on tables	32
Create remaining Foreign Key Relationships	34
Merge and Update on future data.....	34
Appendix-Alternative Database Implementation	36
STAGE_RECORDS, bulk insert & RECORDS table.....	36
STAGE_SAMPLES & bulk insert	39
Creation of CRUISES table and FK relationship with RECORDS table	40
Fixes in STAGE_SAMPLES table, creation of FK, creation of SAMPLES table & insert of data.....	41
Creation of Species_info table and make fixes in it.....	43
Insert data from Species_info table to SAMPLES table & create FK relationship with RECORDS table	43
Merge and Update on future data.....	44

Appendix-Data Mining in R	46
SCENARIO_1.....	46
SCENARIO_2.....	51
Appendix-Power BI	55

The project, “creation of a marine zooplankton database”, required the following steps:

1. data wrangling (power query, python),
2. relational database model design and implementation of it (MSSQL)
3. user-friendly queries for data mining in R programming language

Here follows the procedure of all the steps mentioned above:

Appendix-Power Query guide

MACROS used

MACROS with name ‘trX’ used:

```
Sub trX()

' TeachExcel.com

' Make the macro run faster on large data sets.
Application.ScreenUpdating = False

' Do something with the user-selected cells/range.
xRow = Selection.Rows.Count
xCol = Selection.Column

' The row that the Transposed data will be put into.
nextRow = 1

' Hard-code Column header example
'Range("C1").Value = "Column Header"

' How many rows to Transpose.
```

```

stepValue = InputBox("How many rows should be grouped together?")

' Loop through the user-selected data using a step value.
For i = 1 To xRow Step stepValue

    ' Copy the data, using the step value to determine the size of
    ' the copied range.
    Cells(i, xCol).Resize(stepValue).Copy

    'Transpose the data.
    Cells(1, xCol).Offset(nextRow, 3).PasteSpecial Paste:=xlPasteAll, Transpose:=True

    ' Increment the nextRow value so the copied data goes onto
    ' a new line.
    nextRow = nextRow + 1

Next

' Remove the "copy lines" from the Transposed data.
Application.CutCopyMode = False

' Make Excel function as expected after the macro is finished.
Application.ScreenUpdating = True

End Sub

```

Walkthrough

Abbreviations:

DA: down arrow, RA: right arrow

Produce Excel files from “.dat” files

Dat files:

Produce excel file with the results values:

Open Excel → Data tab → Get Data → From File → From Folder → Choose DANA0697 → choose Dat → OK → OK → on the pop-up Data preview window click Transform Data (Figure 1) → We are in the Power Query Editor(Figure 2) → right click on Content column → Chose Remove Other Columns → click Combine files (Figure 3) → on Combine files preview screen choose “Tab” as delimiter (Figure 4)→ click OK → click the arrow next to column1 → choose “text Filters (Figure 5)→ Does not contain “,” (Figure 6)→ Now logging data is gone (Figure 7)→ At Filter tab , under “Text filter” enter “136” and untick it (Figure 8)→ Close&Load from Home tab tool bar → Data is transferred to Excel worksheet → Right click on column → choose Table → Convert to range (Figure 9)→ OK → Delete first row → CTRL+SHIFT+DA to choose all working cells of first column → Alt+F8 to open Macros tab → Run Macro → in the pop-up window at the question “How many rows” fill “136” (Figure 10)→ Click OK

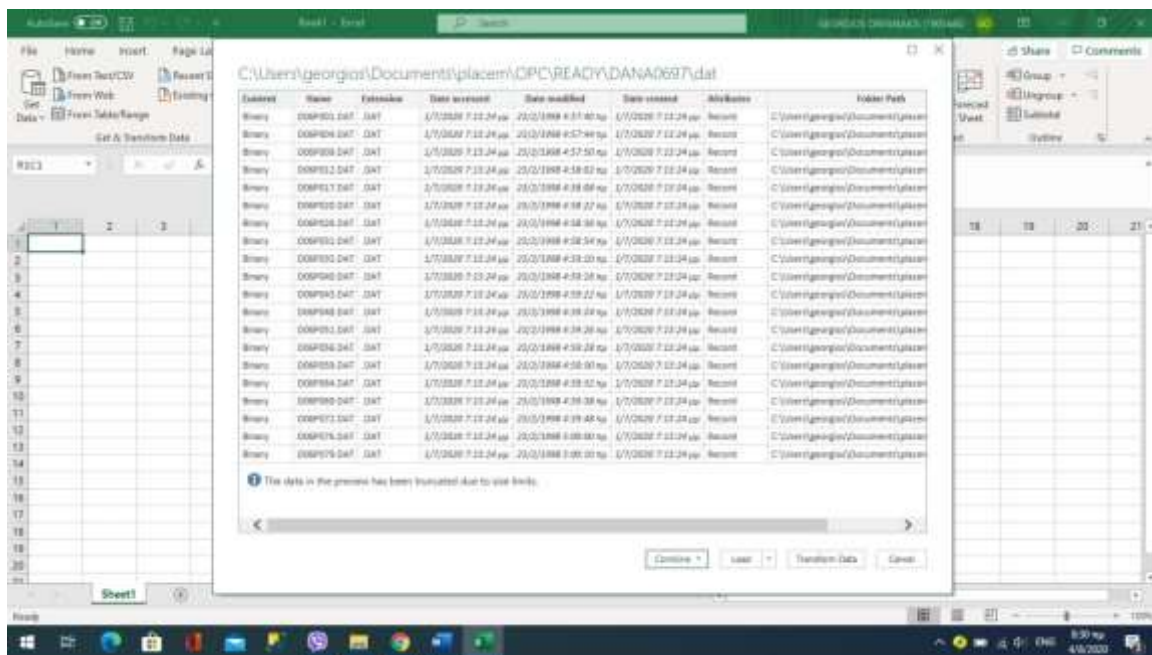


Figure 1

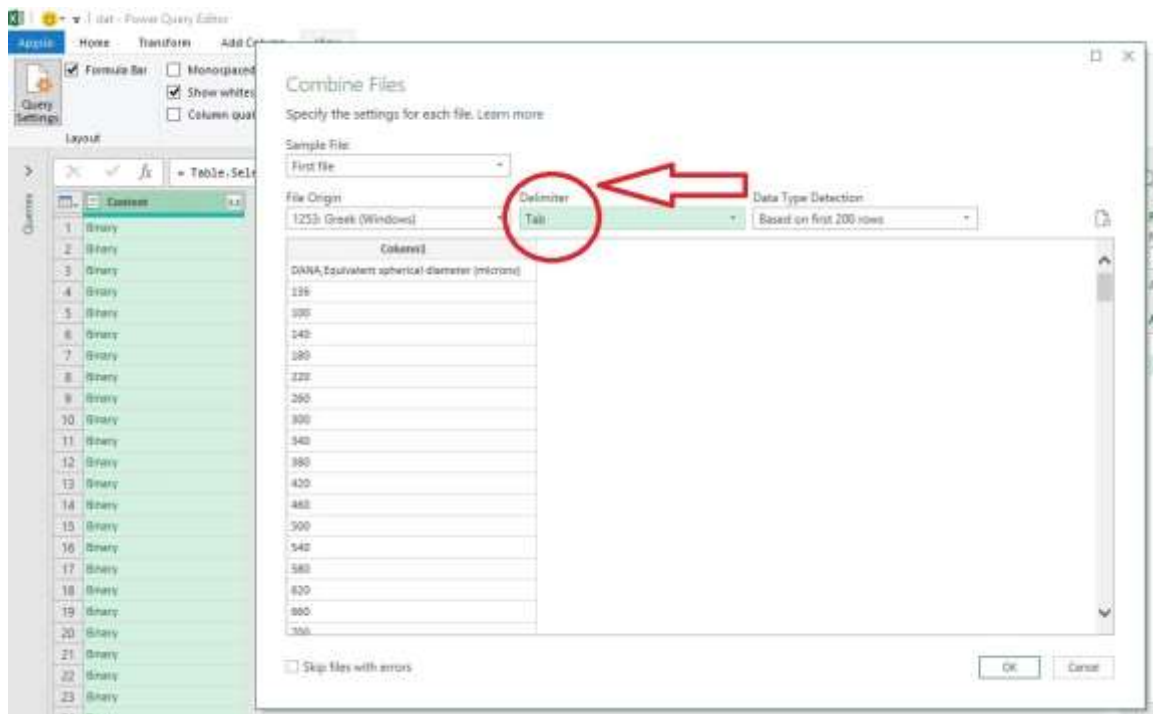


Figure 4

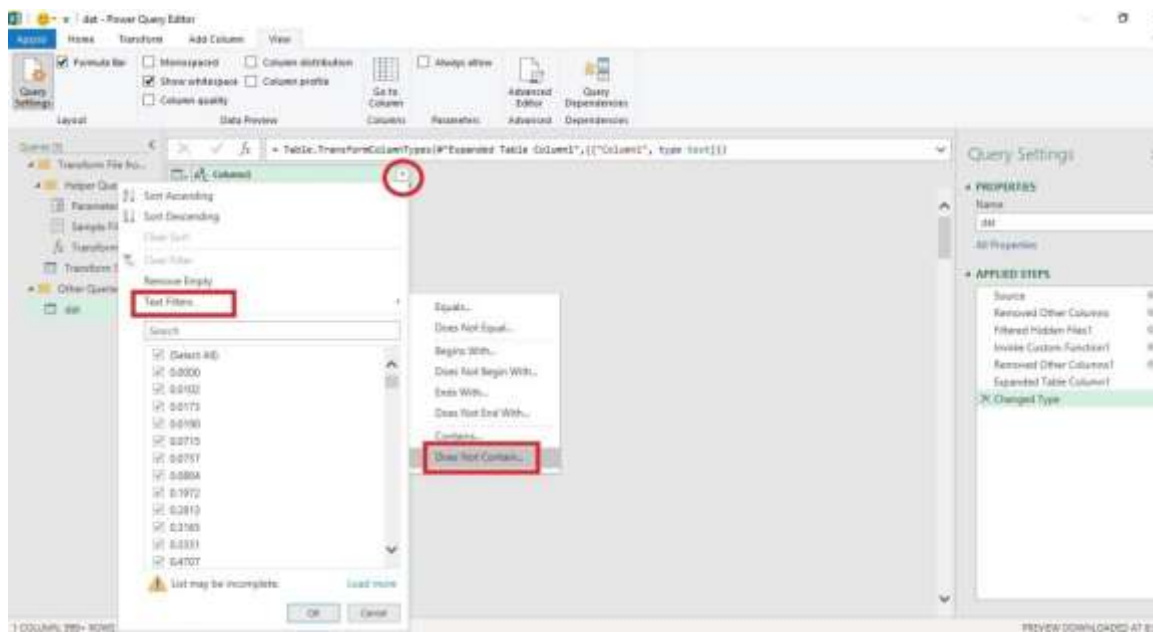


Figure 5

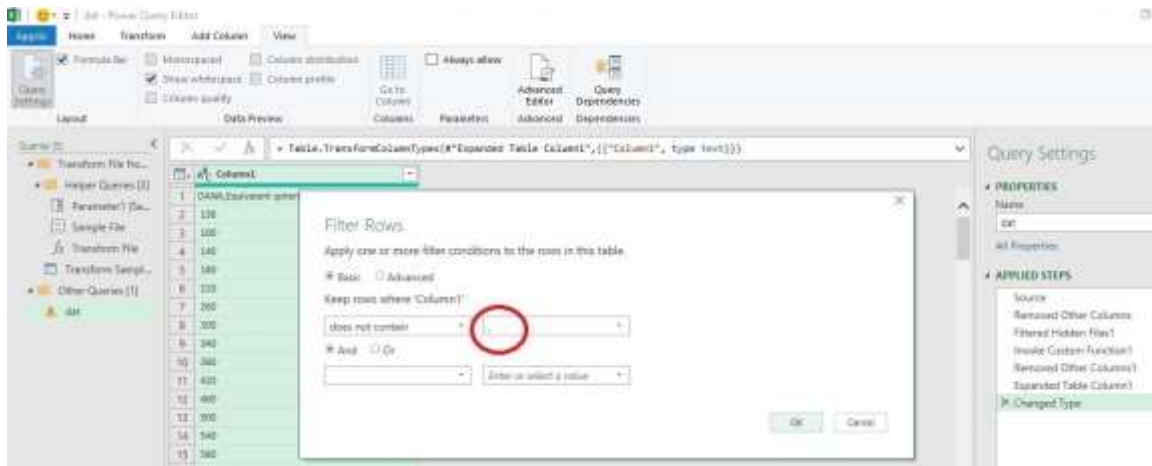


Figure 6

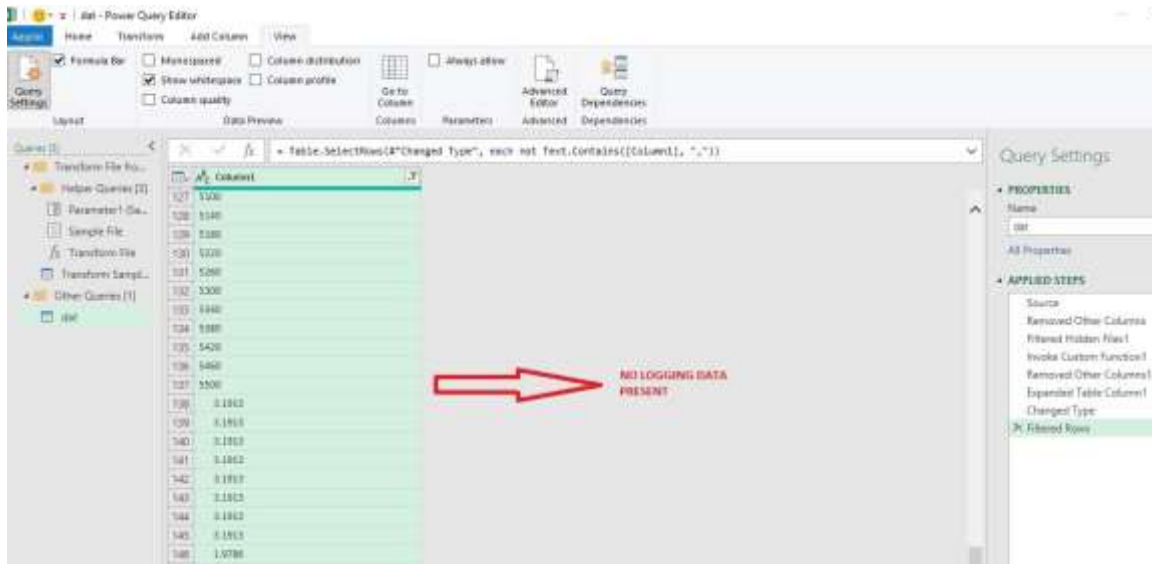


Figure 7

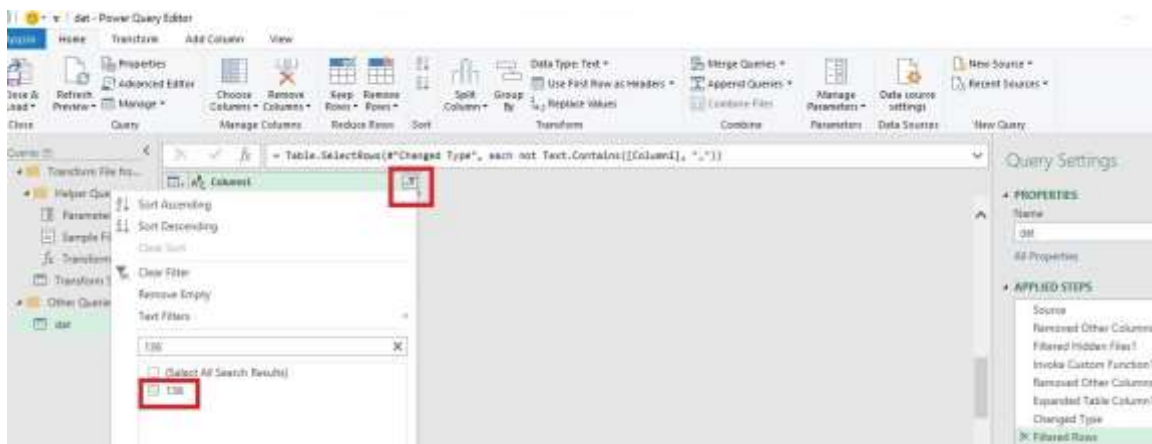


Figure 8

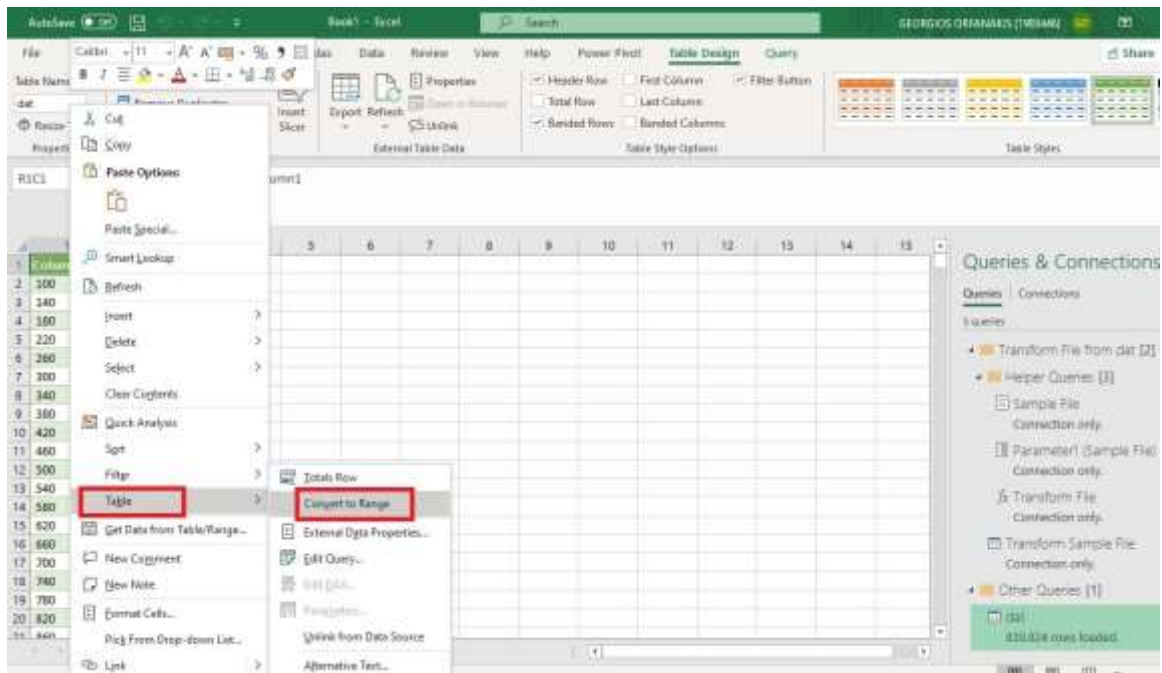


Figure 9

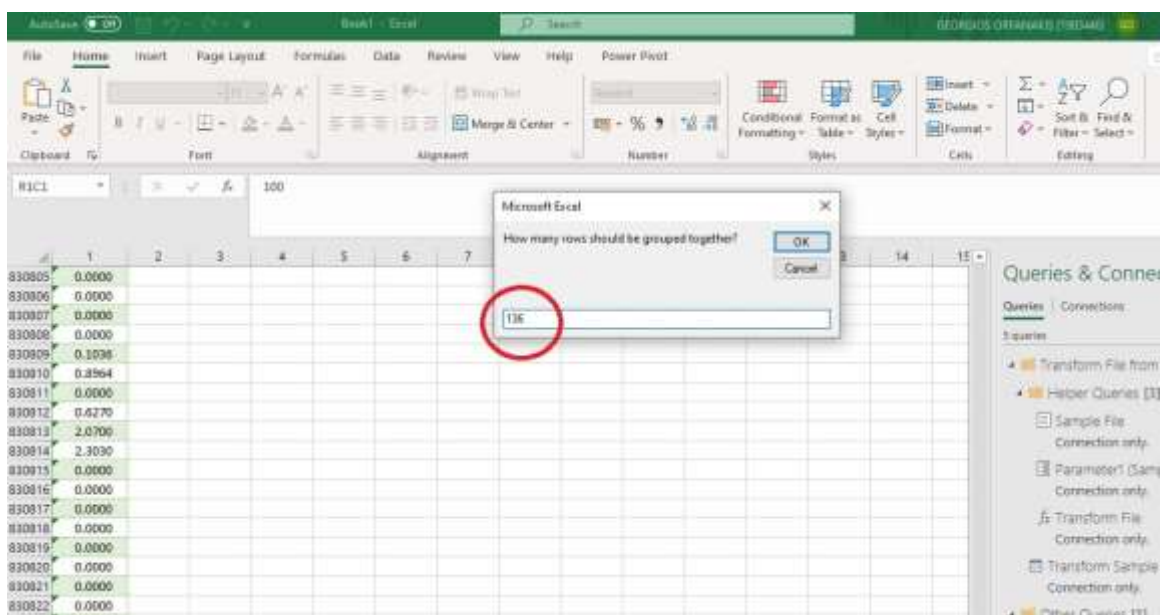


Figure 10

The initial 830824 rows were transformed into 6110 rows. The categories row though is repeated (Figure 11) so it needs to be removed.

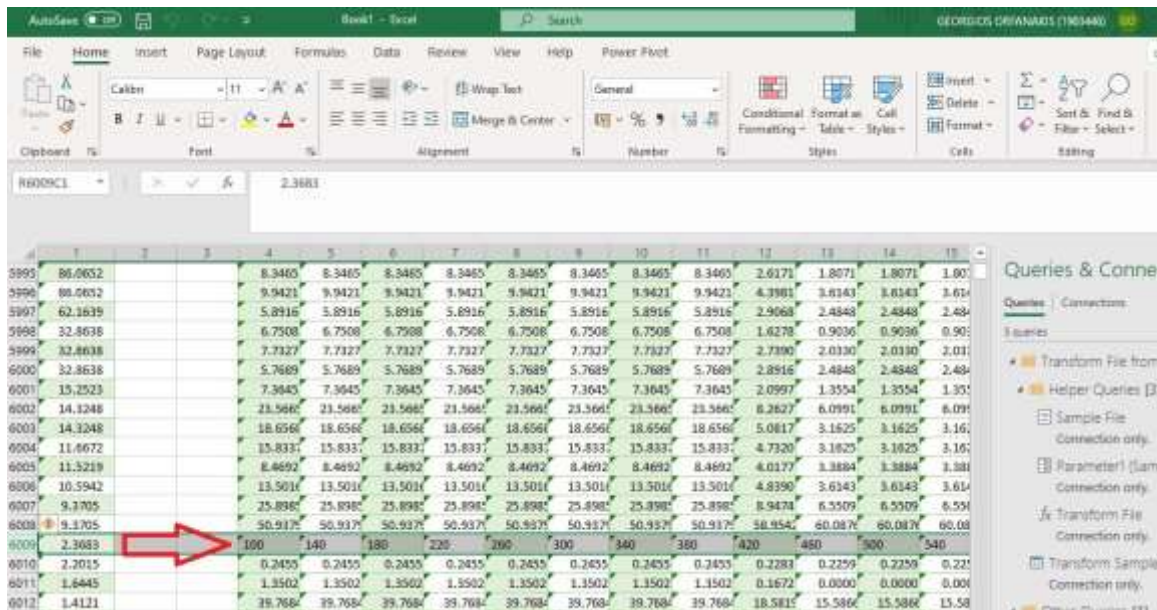


Figure 11

Right click on top of first column → Delete → OK → Click on top of first left working cell with value "100" → CTRL+SHIFT+RA then CTRL+SHIFT+DA → All working cells are chosen → CTRL+F → in "Find what" enter 1140, which is one of the categories → click "Find all" –Make sure no value other than the category is in → Click CTRL+A → hold CTRL and click on the first row in the pop-up window and on values that are not categories to exclude them (Figure 12) → Close window → On Home tool bar click Delete → Delete Sheet Rows (Figure 13). All in between rows that have categories values are omitted.

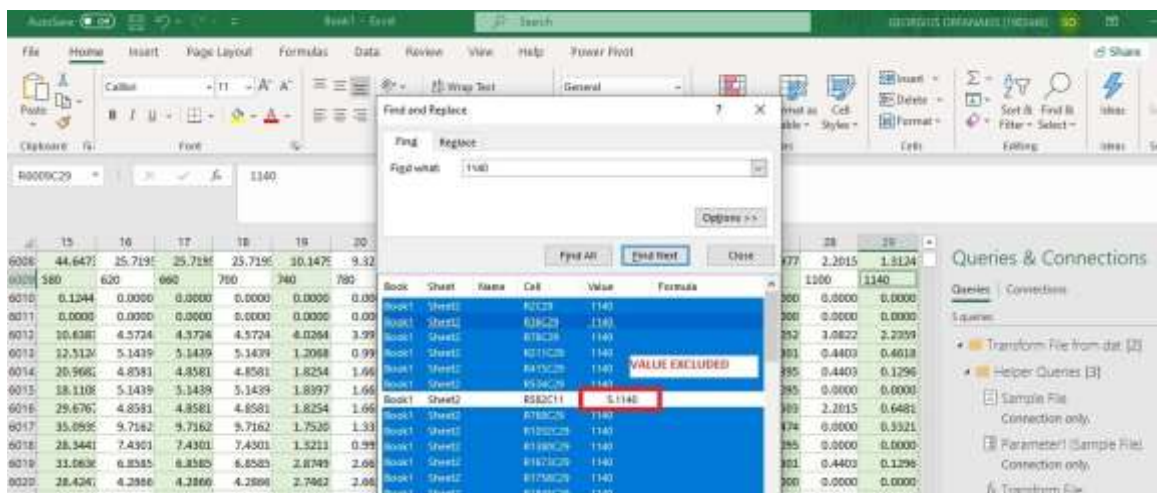


Figure 12

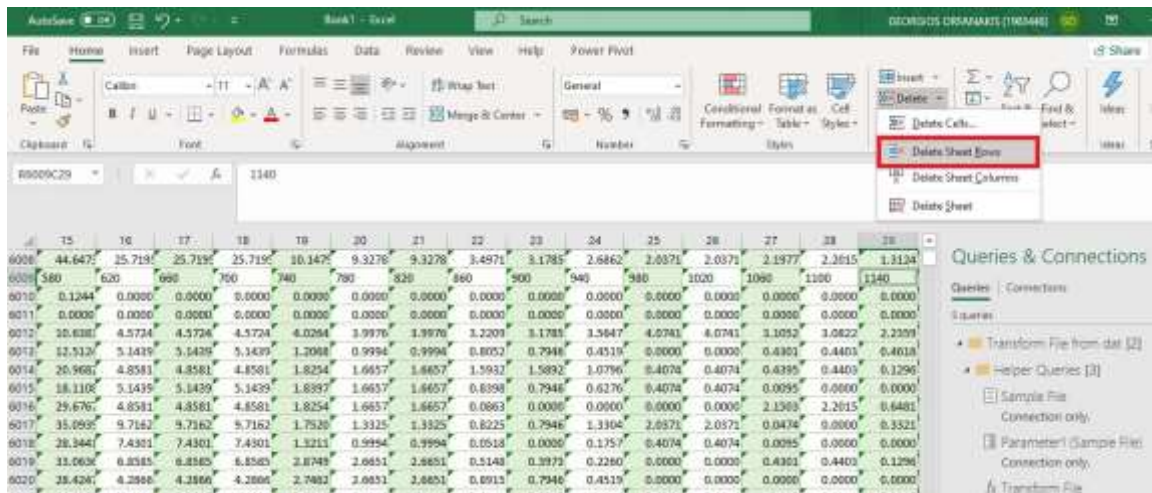


Figure 13

Choose another sheet → at the Data tab click Data → and select the same file following the steps as above. → click the arrow next to column1 → choose “text Filters” (Figure 14) → “Contains” and input “,” and “Does not contain” and input “/” (Figure 15) → The same with “Does not contain” and enter “D” (Figure 16) → Close & Load for the column containing the haul number and volume (Figure 17) → Right click on column1 → Table → Convert to range → OK → CTRL+SHIFT+DA to choose all working cells of first column → On Data tab click “Text to columns” → click the “Delimited” option → Next → choose “comma” option → Next → Choose “Text option” for both columns → Finish → click on first cell of first column → CTRL+SHIFT+DA to choose all working cells of first column → copy paste it as first column on the previous sheet where all the values are (Figure 18) → Enter column name as “HAUL_NO”. (Figure 19)

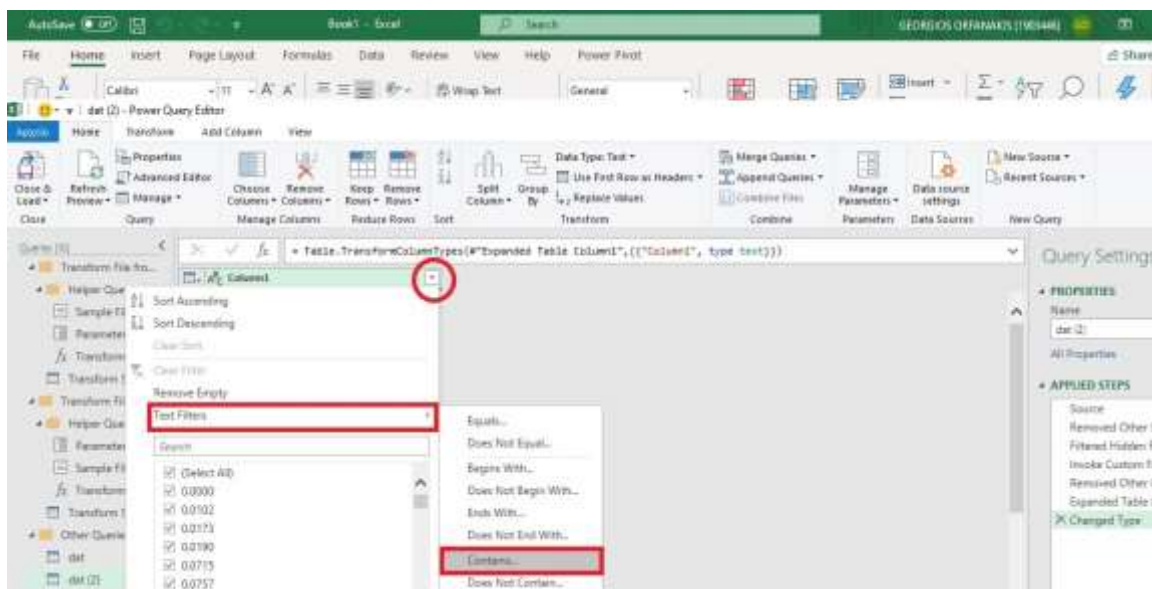


Figure 14

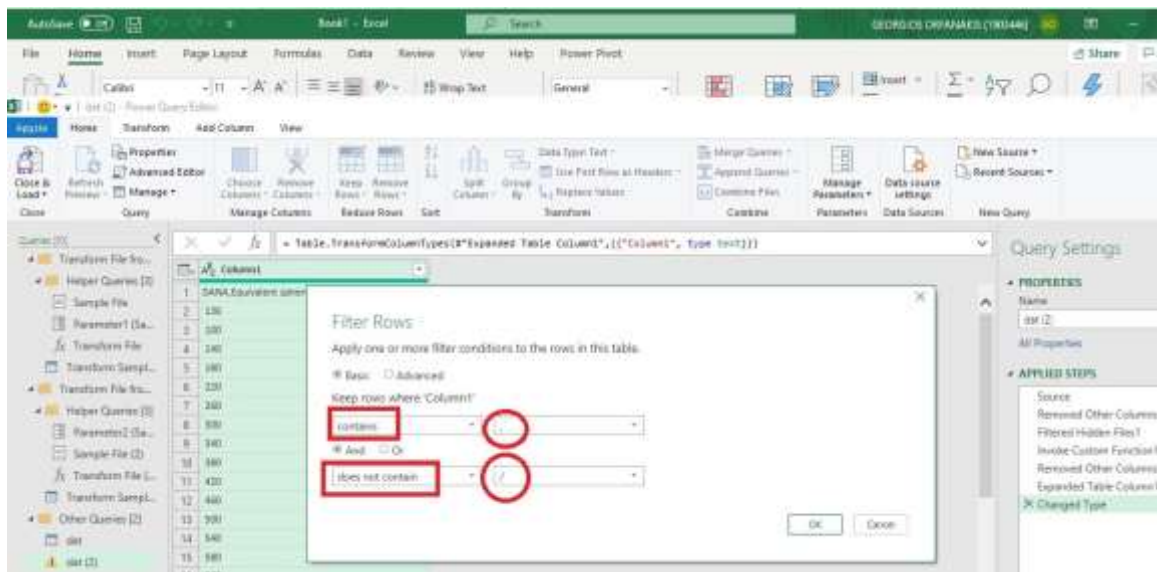


Figure 15

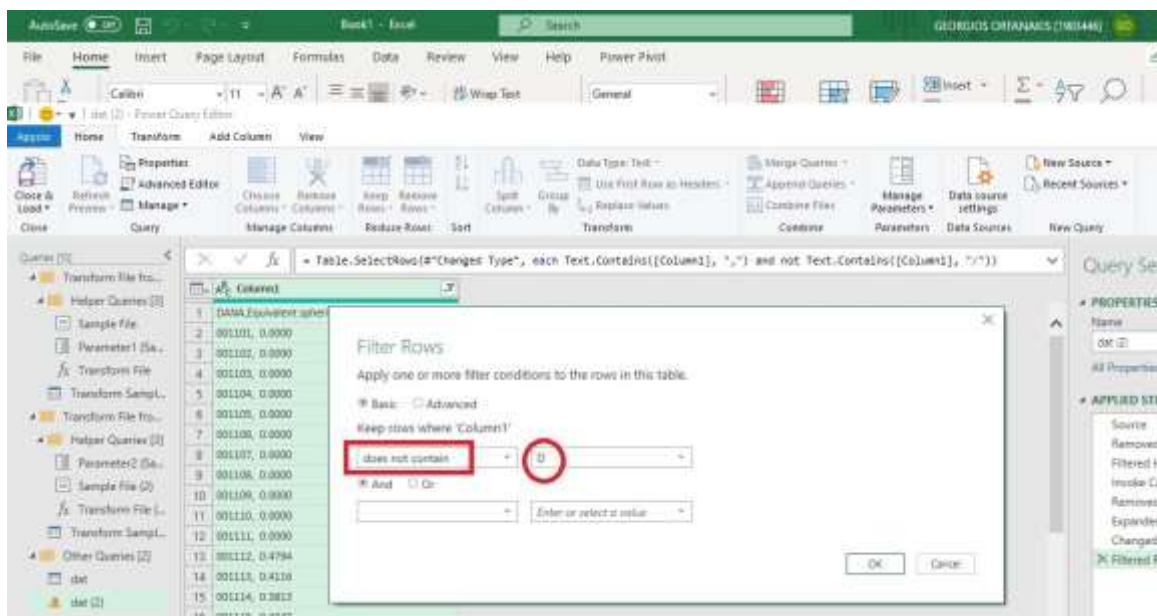


Figure 16

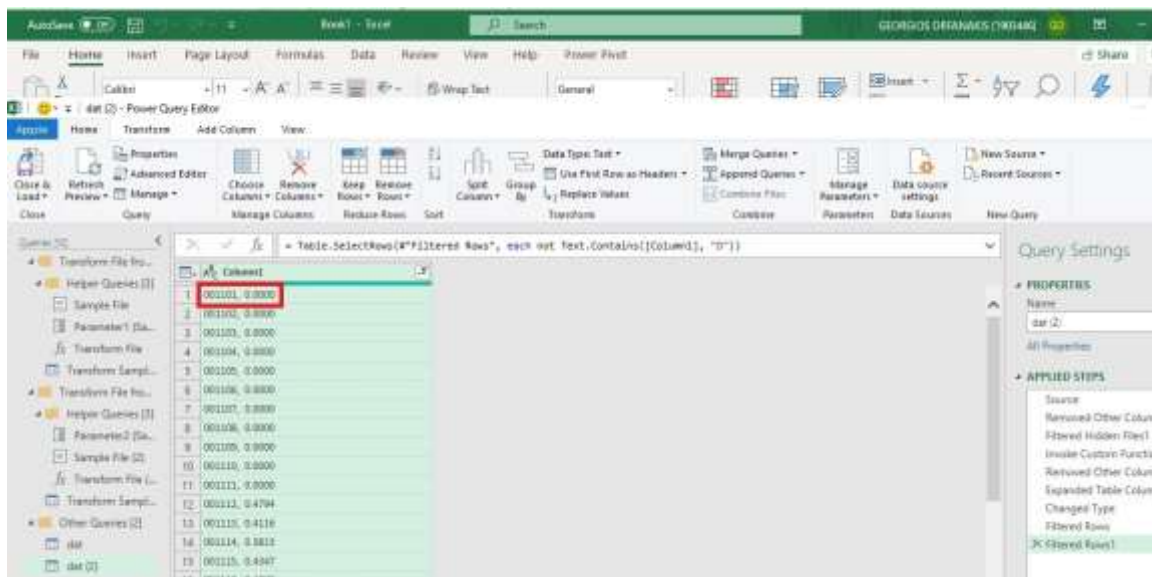


Figure 17

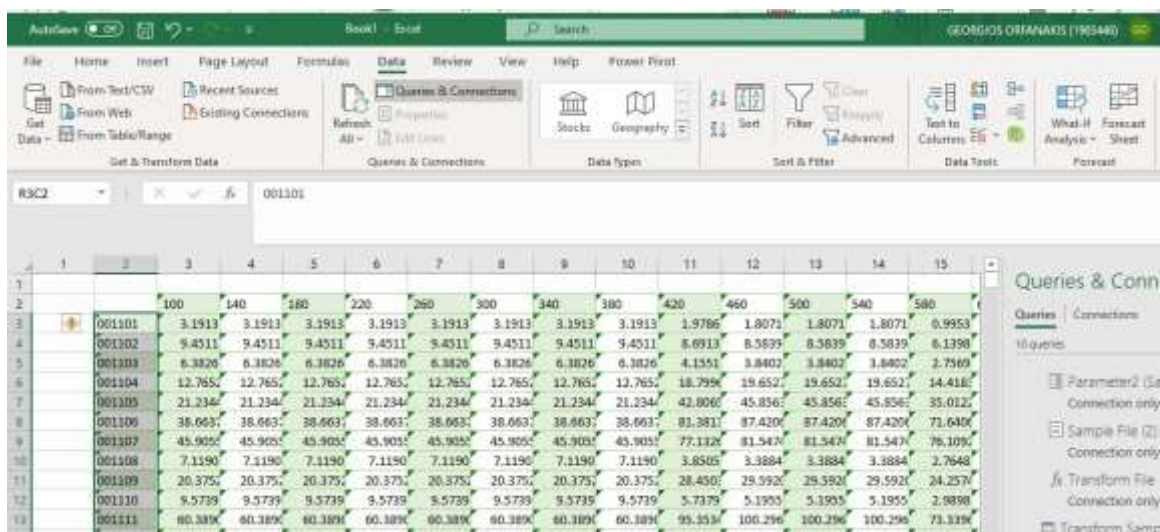


Figure 18

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
3	HAUL_NO	140	230	300	380	450	520	590	660	730	800	870	940	1010	1080	1150	1220	1290	1360	1430
4	000101	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913	8.1913
5	000102	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511	9.4511
6	000103	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826	6.3826
7	000104	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765	12.765
8	000105	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234	21.234
9	000106	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661	38.661
10	000107	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905	45.905
11	000108	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190	7.1190
12	000109	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175	20.175
13	000110	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735	9.5735
14	000111	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389	60.389
15	000112	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720	96.720

Figure 19

A way to test we are right in this process is at this step to go at the very last working cell of the first column and check if it ends where the values end. (Figure 20) The first column contains the haul numbers so there can be no value without a haul number.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
6034	240185	28.230	28.230	28.230	28.230	28.230	28.230	28.230	28.230	28.230	22.892	22.137	22.137	22.137	15.401	
6035	240186	29.212	29.212	29.212	29.212	29.212	29.212	29.212	29.212	29.212	22.618	21.685	21.685	21.685	15.409	
6036	240187	31.790	31.790	31.790	31.790	31.790	31.790	31.790	31.790	31.790	22.145	20.782	20.782	20.782	15.168	
6037	240188	29.089	29.089	29.089	29.089	29.089	29.089	29.089	29.089	29.089	22.404	21.459	21.459	21.459	14.643	
6038	240189	33.140	33.140	33.140	33.140	33.140	33.140	33.140	33.140	33.140	22.510	21.006	21.006	21.006	15.549	
6039	240190	39.154	39.154	39.154	39.154	39.154	39.154	39.154	39.154	39.154	22.860	20.556	20.556	20.556	13.888	
6040	240191	52.656	52.656	52.656	52.656	52.656	52.656	52.656	52.656	52.656	34.625	32.078	32.078	32.078	20.747	
6041	240192	47.501	47.501	47.501	47.501	47.501	47.501	47.501	47.501	47.501	34.383	32.528	32.528	32.528	20.225	
6042	240193	57.811	57.811	57.811	57.811	57.811	57.811	57.811	57.811	57.811	34.868	31.625	31.625	31.625	20.888	
6043	240194	54.865	54.865	54.865	54.865	54.865	54.865	54.865	54.865	54.865	28.764	25.074	25.074	25.074	16.120	
6044	240195	96.229	96.229	96.229	96.229	96.229	96.229	96.229	96.229	96.229	55.262	49.470	49.470	49.470	31.462	
6045	240196	101.753	101.753	101.753	101.753	101.753	101.753	101.753	101.753	101.753	70.790	66.412	66.412	66.412	40.812	
6046	240197	68.490	68.490	68.490	68.490	68.490	68.490	68.490	68.490	68.490	37.972	33.651	33.651	33.651	20.976	
6047	240198	34.981	34.981	34.981	34.981	34.981	34.981	34.981	34.981	34.981	20.363	18.297	18.297	18.297	10.847	
6048	240199	61.248	61.248	61.248	61.248	61.248	61.248	61.248	61.248	61.248	30.148	25.751	25.751	25.751	16.621	
6049	240200	75.854	75.854	75.854	75.854	75.854	75.854	75.854	75.854	75.854	34.787	34.787	34.787	34.787	25.708	
6050	240201	46.887	46.887	46.887	46.887	46.887	46.887	46.887	46.887	46.887	58.951	58.951	58.951	58.951	40.173	

Figure 20

Now, CTRL+SHIFT+RA then CTRL+SHIFT+DA → All working cells are chosen → CTRL+T to turn the working cells into a table with headers → At Data tab click From Table/Range to enter the Power Query → CTRL+A → From Home tab choose Data Type → Text (Figure 21) → Replace current → Right click on HAUL_NO column → Unpivot other columns (Figure 22) → Rename the columns as “CLASS” and “VALUE”(Figure 23) → Close&Load.

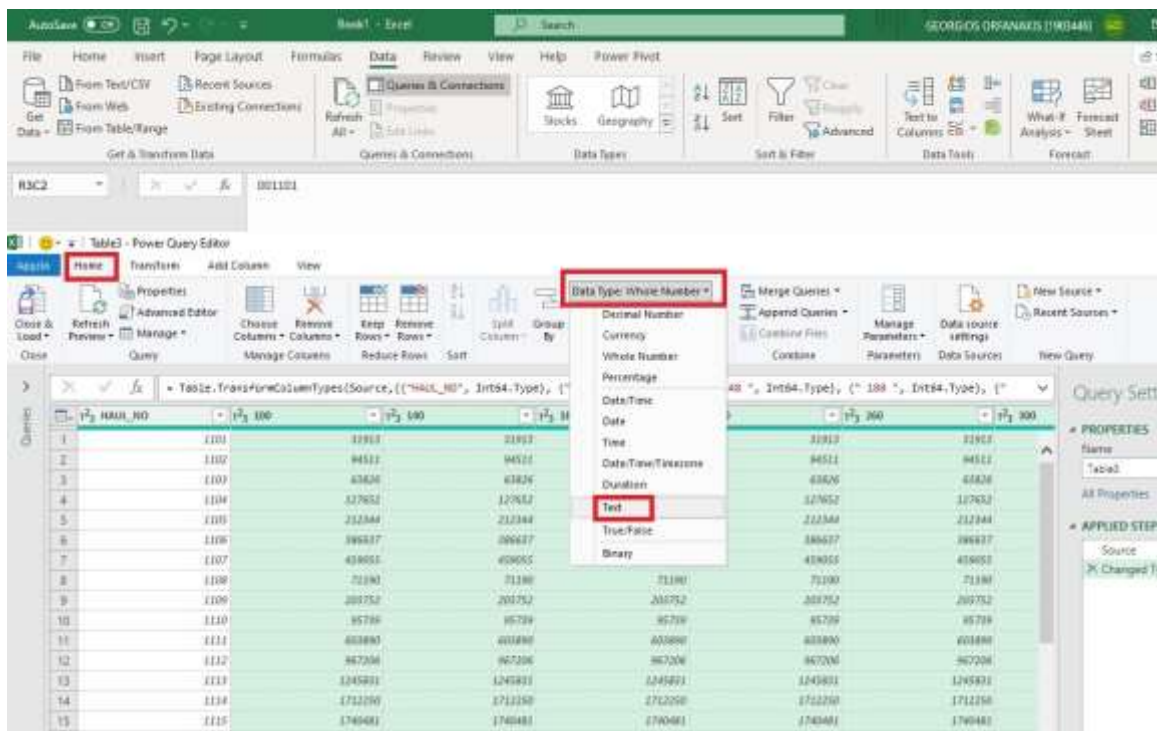


Figure 21

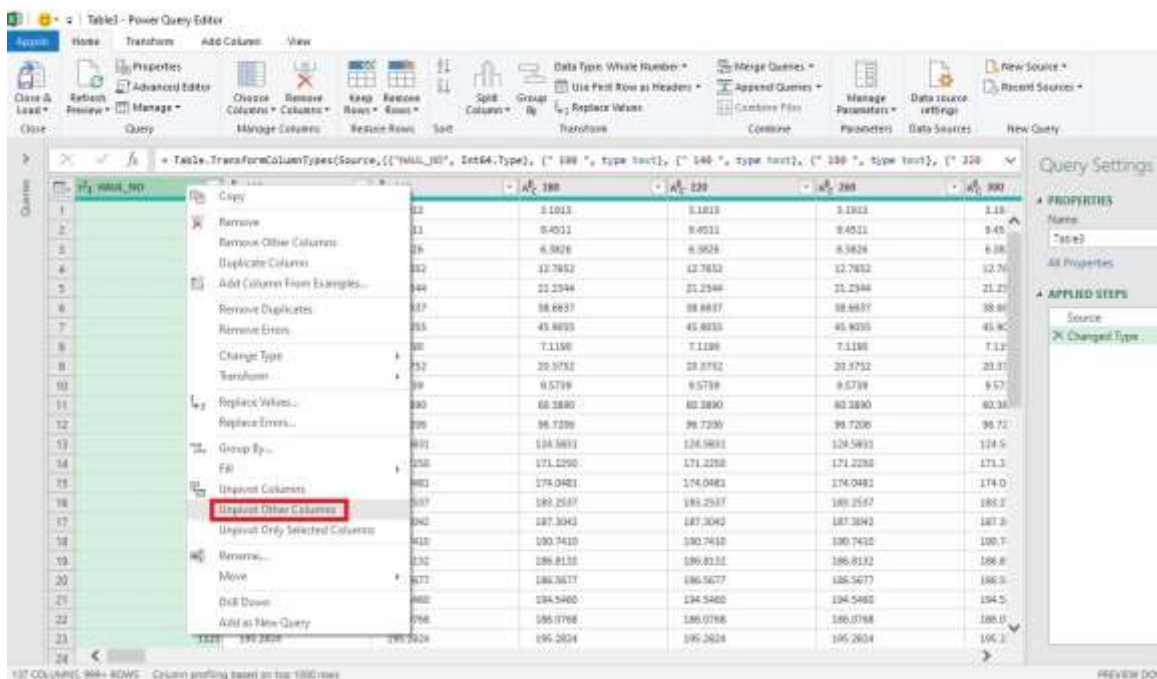


Figure 22

The screenshot shows the Power Query Editor interface. The main area displays a table with three columns: HAUL_NO, CLASS, and VALUE. The CLASS column is highlighted in red, and the VALUE column is highlighted in green. The table contains 20 rows of data. The right sidebar shows the 'Query Settings' pane with 'PROPERTIES' and 'APPLIED STEPS' sections. The 'APPLIED STEPS' section shows a list of steps: Source, Changed Type, Unpivoted Columns, and Returned Columns.

HAUL_NO	CLASS	VALUE
1	1101 100	0.1913
2	1101 540	0.1913
3	1101 580	0.1913
4	1101 230	0.1913
5	1101 260	0.1913
6	1101 300	0.1913
7	1101 540	0.1913
8	1101 580	0.1913
9	1101 420	0.0788
10	1101 440	0.0071
11	1101 500	0.0072
12	1101 540	0.0071
13	1101 580	0.0063
14	1101 620	0.0000
15	1101 660	0.0000
16	1101 700	0.0000
17	1101 740	0.3188
18	1101 780	0.3333
19	1101 820	0.3333
20	1101 860	0.0173

Figure 23

Once the data is loaded in the excel worksheet we insert the columns "TYPE" and "CRUISE_CODE" → fill "TYPE" with "CONCENTRATION" and "CRUISE_CODE" with "D0697" → copy paste "CRUISE_CODE" column as value → Delete all other sheets, but keep HAUL_NO and Volume columns on another excel workbook for later use → Save file with appropriate name eg results_opc_D0697.

Produce excel file with the logging values:

On the new workbook where the HAUL_NO and depth columns where stored on a new worksheet → At Data tab choose "From file" → From folder → OK → Choose DANA0697 → choose Dat → OK → OK → on the pop-up Data preview window click Transform Data → We are in the Power Query Editor → right click on Content column → Chose Remove Other Columns → click Combine files → on Combine files preview screen choose "Tab" as delimiter → click OK → click the arrow next to column1 → choose "text Filters → "Contains" input "/" (Figure 24) → Close&Load

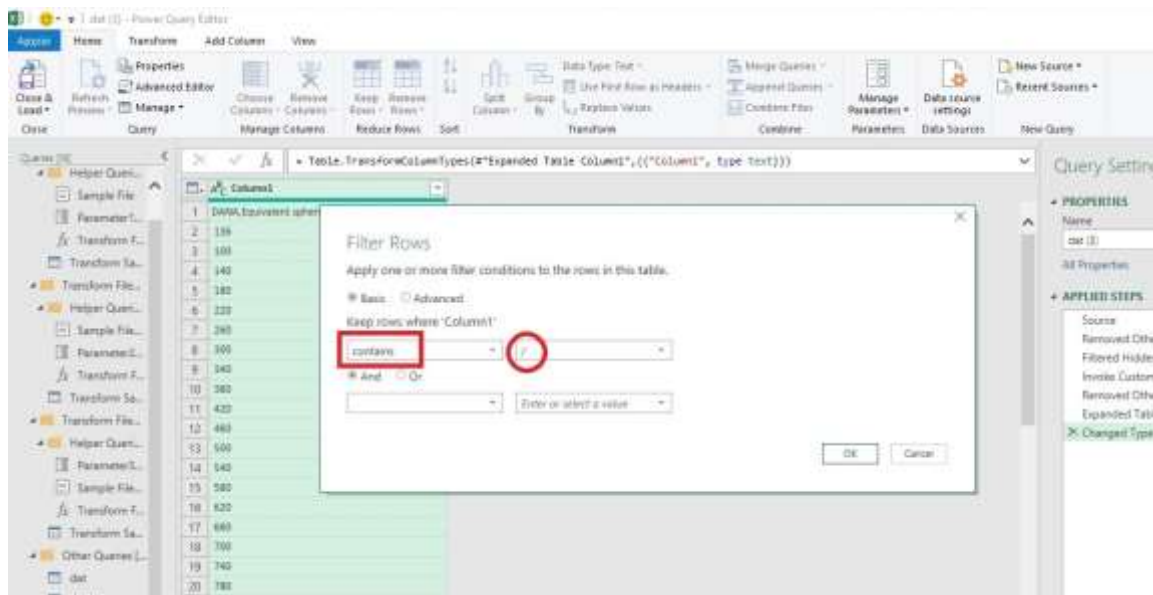


Figure 24

In the Excel sheet where the data is loaded CTRL+SHIFT+DA → All working cells of first column are chosen → At Data tab choose “Text to Data” → Delimited → Next → choose “comma” → Next → choose first column as Date in “DMY” format → Finish.

Now we copy paste at the top row the column names that we have created previously. After this the transformation of each consecutive three rows into 1 begins with Start, Mid and End columns being inserted the values as shown in Figure 25 with subsequent colors. This happens for all contents of Latitude, Longitude, Echo and Sampler_depth.

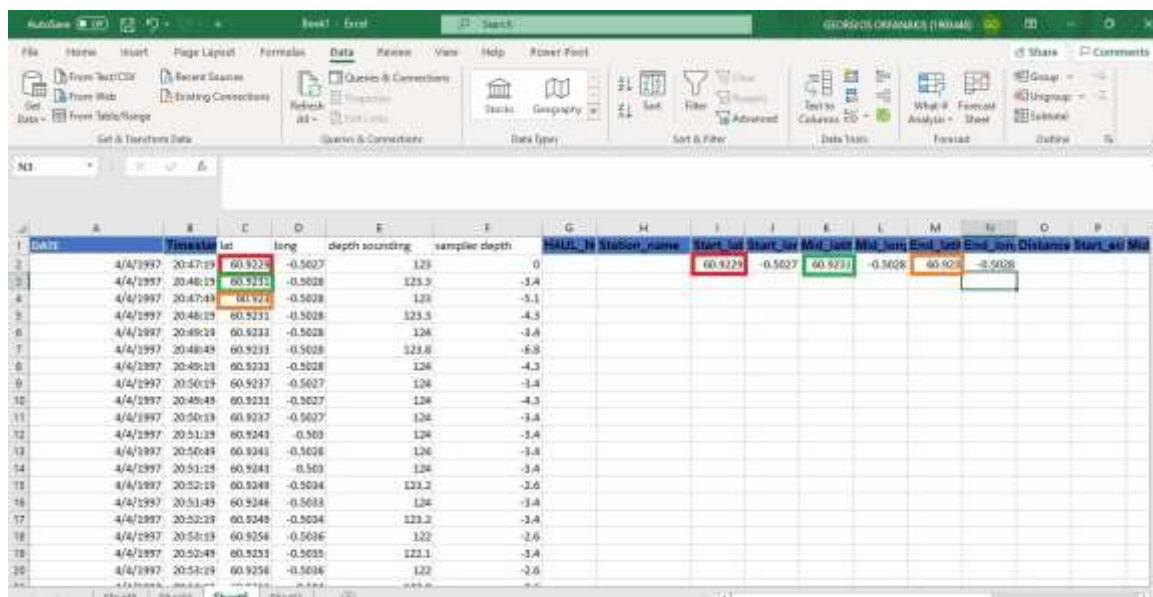


Figure 25

For Duration we calculate the difference between 3rd row and 1st row (Figure 26). We choose the above computed columns → Double click on the down left corner of the last chosen working cell and let excel iterate the calculation → Fill CRUISE_CODE column with D0697 → Timestamp column is chosen with CTRL+SHIFT+DA → right click → format cells → Time → choose “37:50:55” format → Same goes for Duration column → All computed columns are copy-pasted as values → We delete column lat-long → In the column next to

CRUISE_CODE we add on the second and third row the "@" value (Figure 27). We choose the first three rows of the column and double click an let excel iterate the values (Figure 28)→ We choose the first cell of this last column → CTRL+SHIFT+DA → the entire working column is chosen → CTRL+F → enter "@" → CTRL+A (Figure 29)→ at Home tab choose Delete → Delete sheet rows → Every second and third row is gone → the HAUL_NO from the other worksheet is added at the subsequent HAUL_NO column, same goes for depth volume_dec_cubic_m column → The sampler column is filled with "PC" as value → Valid_Y_N column is filled with "Y" as value → All other sheets are deleted → Excel file is saved with appropriate name e.g. records_opc_D0697

DATE	Time	lat	long	depth sounding	sampler depth	HAUL_NO	Station_name	Start_lon	Start_lat	Mid_lon	Mid_lat	End_lon	End_lat	Distance	Start_lon	Mid
4/4/1997	20:47:23	60.9229	-0.5027	123	0			60.9229	-0.5027	60.9229	-0.5028	60.923	-0.5026	123		
4/4/1997	20:48:19	60.9231	-0.5028	123.3	-3.4											
4/4/1997	20:47:49	60.923	-0.5028	123	-5.1											
4/4/1997	20:48:19	60.9231	-0.5028	123.3	-4.3											
4/4/1997	20:48:19	60.9231	-0.5028	124	-3.4											
4/4/1997	20:48:49	60.9231	-0.5028	123.8	-6.8											
4/4/1997	20:48:19	60.9231	-0.5028	124	-4.3											
4/4/1997	20:50:29	60.9237	-0.5027	124	-3.4											
4/4/1997	20:49:49	60.9231	-0.5027	124	-4.3											
4/4/1997	20:50:39	60.9237	-0.5027	124	-3.4											
4/4/1997	20:51:19	60.9243	-0.503	124	-3.4											
4/4/1997	20:50:49	60.9241	-0.5028	124	-3.4											
4/4/1997	20:51:19	60.9243	-0.503	124	-3.4											
4/4/1997	20:52:29	60.9249	-0.5034	123.2	-2.6											
4/4/1997	20:51:49	60.9246	-0.5033	124	-3.4											
4/4/1997	20:52:19	60.9249	-0.5034	123.2	-3.4											
4/4/1997	20:53:19	60.9256	-0.5036	123	-2.6											

Figure 26

End_lon	Distance	Start_lon	Mid_lon	End_lon	SAMPLE Mesh_m	Volume_dec	Duration	Start_lon	Mid_lon	End_lon	Flow_rate	Start_lon	End_lon	REV3	Valid_Y	CRUISE_CODE
-0.5028	123	123.3	123	123	0:00:30	0	-3.4	-5.1							D0697	D0697
-0.5028	123.3	123	123.3	123.3	0:00:00	-3.4	-5.1	-4.3							D0697	D0697
-0.5028	123	123.3	124	124	0:01:30	-5.1	-4.3	-3.4							D0697	D0697
-0.5028	123.3	124	123.8	123.8	0:00:30	-4.3	-3.4	-6.8							D0697	D0697
-0.5028	124	123.8	124	124	0:00:00	-3.4	-6.8	-4.3							D0697	D0697
-0.5027	123.8	124	124	124	0:01:30	-6.8	-4.3	-3.4							D0697	D0697
-0.5027	124	124	124	124	0:00:30	-4.3	-3.4	-4.3							D0697	D0697
-0.503	124	124	124	124	0:01:30	-4.3	-3.4	-3.4							D0697	D0697
-0.5028	124	124	124	124	0:00:30	-3.4	-3.4	-3.4							D0697	D0697
-0.503	124	124	124	124	0:00:00	-3.4	-3.4	-3.4							D0697	D0697
-0.5034	124	124	123.2	123.2	0:01:30	-3.4	-3.4	-2.6							D0697	D0697
-0.5033	124	123.2	124	124	0:00:30	-3.4	-3.4	-2.6							D0697	D0697
-0.5034	123.2	124	123.2	123.2	0:00:30	-2.6	-3.4	-3.4							D0697	D0697
-0.5034	124	123.2	122	122	0:01:30	-3.4	-3.4	-2.6							D0697	D0697
-0.5033	123.2	122	122.1	122.1	0:00:30	-3.4	-2.6	-3.4							D0697	D0697
-0.5036	122	122.1	122	122	0:00:00	-2.6	-3.4	-2.6							D0697	D0697

Figure 27

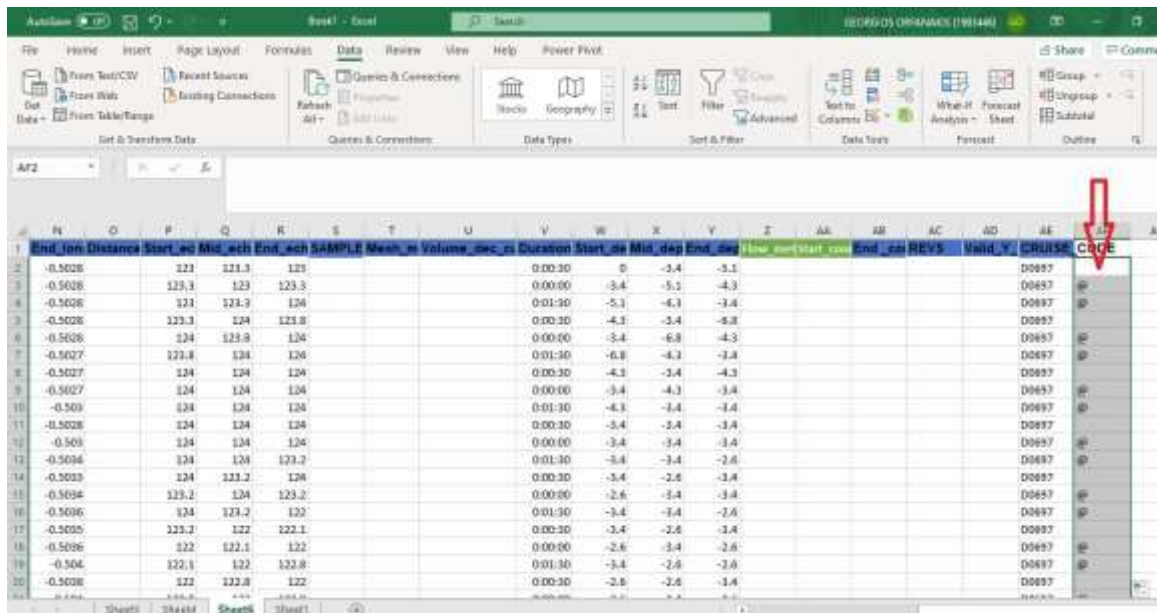


Figure 28

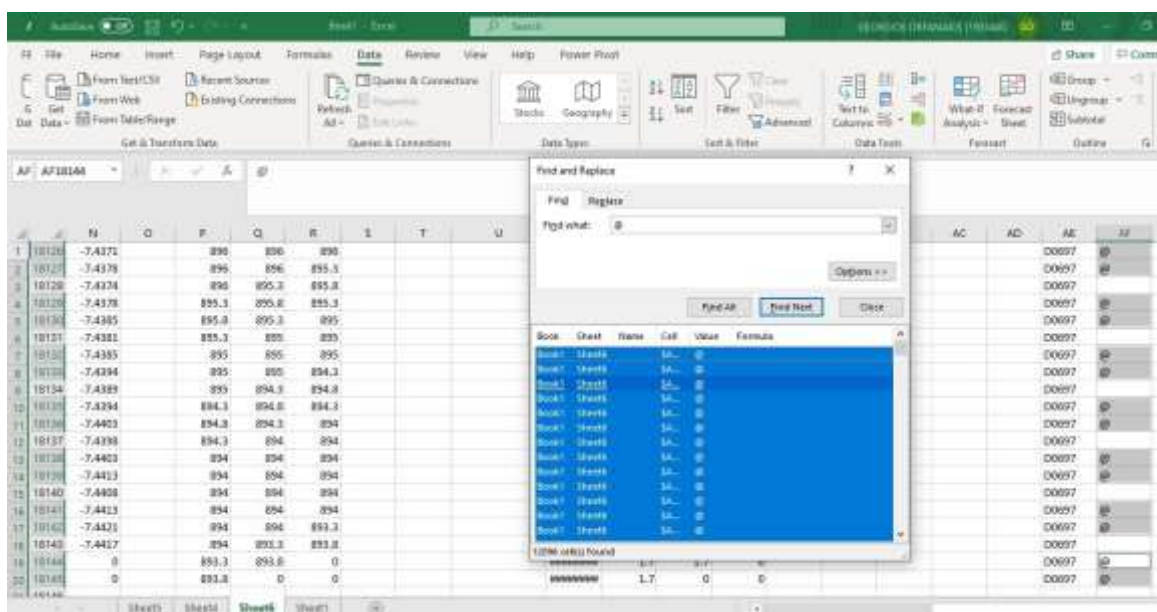


Figure 29

A way to check that the deletion has been achieved is by checking whether the HAUL_NO column's values end at the same row as the rest pre-existing values (Figure 30)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
001	18/4/1997	10:50:37	59.3877	-7.4277	906	548.5	548.5	59.3877	-7.4277	59.387	-7.4263	59.3874	-7.4279	906		
004	18/4/1997	10:52:37	59.387	-7.4283	906	548.5	548.5	59.387	-7.4283	59.3862	-7.4292	59.3867	-7.4287	906		
005	18/4/1997	10:53:37	59.3962	-7.4292	905.3	322.9	322.9	59.3962	-7.4293	59.3955	-7.4302	59.3959	-7.4297	905.3		
006	18/4/1997	10:54:37	59.3855	-7.4302	904.3	296.5	296.5	59.3855	-7.4302	59.385	-7.431	59.3852	-7.4306	904.3		
007	18/4/1997	10:55:37	59.383	-7.431	903.3	270.1	270.1	59.383	-7.431	59.3847	-7.4314	59.3848	-7.4312	903.3		
008	18/4/1997	10:56:37	59.3847	-7.4314	902.6	246.2	246.2	59.3847	-7.4314	59.3843	-7.4316	59.3846	-7.4315	902.6		
009	18/4/1997	10:57:37	59.3843	-7.4316	900.3	223.2	223.2	59.3843	-7.4316	59.3838	-7.4321	59.3841	-7.4318	900.3		
040	18/4/1997	10:58:37	59.3838	-7.4321	898.3	198.5	198.5	59.3838	-7.4321	59.3833	-7.4322	59.3836	-7.4325	898.3		
041	18/4/1997	10:59:37	59.3833	-7.4322	898.3	172.5	172.5	59.3833	-7.4322	59.3827	-7.4345	59.3831	-7.4339	898.3		
042	18/4/1997	10:10:37	59.3827	-7.4345	897.3	147.4	147.4	59.3827	-7.4345	59.382	-7.436	59.3824	-7.4352	897.3		
043	18/4/1997	10:11:37	59.382	-7.436	896.3	115	115	59.382	-7.436	59.3818	-7.4371	59.3817	-7.4366	896.3		
044	18/4/1997	10:12:37	59.3814	-7.4372	890	103.1	103.1	59.3814	-7.4372	59.3807	-7.4378	59.3811	-7.4374	890		
045	18/4/1997	10:13:37	59.3807	-7.4378	895.3	92	92	59.3807	-7.4378	59.3801	-7.4385	59.3804	-7.4381	895.3		
046	18/4/1997	10:14:37	59.3801	-7.4385	895	78.4	78.4	59.3801	-7.4385	59.3795	-7.4394	59.3799	-7.4389	895		
047	18/4/1997	10:15:37	59.3795	-7.4394	894.3	52.8	52.8	59.3795	-7.4394	59.3789	-7.4403	59.3792	-7.4398	894.3		
048	18/4/1997	10:16:37	59.3789	-7.4403	894	29.5	29.5	59.3789	-7.4403	59.3783	-7.4413	59.3786	-7.4408	894		
049	18/4/1997	10:17:37	59.3783	-7.4413	894	7.7	7.7	59.3783	-7.4413	59.3776	-7.4421	59.378	-7.4417	894		

Figure 30

For all the other collections (ARIES, MICD, OCEAN) records and sample data are stored in different files. Records are stored in pla/plw files and sample data is in excel, csv or other text format file from which it is derived. This constitutes the unstructured part of the data and it is not possible to depict each case in this report so two distinct and most common cases are shown here for sample data extraction.

Produce Excel files from “.pla/.plw” files

Pla/plw files:

On an Excel worksheet follow phase_1 and all data is imported in Power Query. Click on Filters and choose “Does not contain” with “/”. Filter out 0,0 by unticking it. Hit Close&Load. The data is loaded as one column in the Excel worksheet. CTRL+SHIFT+DA all data is chosen. On the Data tab we choose Text to Columns. On the pop-up window we choose Delimited and hit Next. Then choose “comma”. The Next. Then for the 1st column we choose Date as DMY. Hit Finish. Then we Copy Paste the columns Names on the first row. CTRL+SHIFT+RA then CTRL+SHIFT+DA to choose all working cells. CTRL+T to turn them into a table with headers. Then we need to transform the Timestamp column. For this, 4 extra columns are created after the Timestamp column, which has been formatted into a text column and its title has been deleted (Figure 31). For the first new column the function TEXT(B2;”0000”) is used to make sure all numbers have 4 digits (Figure 31). If a number is missing a digit, then 0 is added in the beginning. The result of this column is copy-pasted as Value on the next column. The third new column is using the function Timestamp (Figure 31)* to turn the previous column cells into actual time. Specifically, the function takes the text 0653 and brings it in the form “06:53”, then converts it to actual time. The columns selected, then follows right click, format cells, time, and choose format 37:30:55 (Figure 32). The cells of this column are pasted as copy-pasted as values onto the next column which gets the name Timestamp. The initial column and all the ones created except for the last one are deleted. We then fill the CRUISE_CODE column with S1801. Then we save the Excel workbook with just the working sheet using an appropriate name e.g. pla_AR_records.

The screenshot shows the Excel interface with the formula bar containing the formula: `=TIMEVALUE(CONCATENATE(LEFT(C2,2),",",RIGHT(C2,2)))`. The formula bar is highlighted with a red box. Below the formula bar, a table is visible with the following columns: DATE, Column1, Column3, Column2, TimeStamp, HAUL_NO, Station_name, and Start_latitude. The table contains 12 rows of data, with the first row having a date of 8/12/2001 and a time of 8:49:00.

DATE	Column1	Column3	Column2	TimeStamp	HAUL_NO	Station_name	Start_latitude
8/12/2001	849	0900		8:49:00	491001	FMRL	90.1
8/12/2001	904	0904		9:04:00	491002	FMRL	90.1
8/12/2001	905	0905		9:05:00	491003	FMRL	90.1
8/12/2001	906	0906		9:06:00	491004	FMRL	90.1
8/12/2001	908	0908		9:08:00	491005	FMRL	90.1
8/12/2001	913	0913		9:13:00	491006	FMRL	90.1
8/12/2001	914	0914		9:14:00	491007	FMRL	90.1
8/12/2001	916	0916		9:16:00	491008	FMRL	90.1
8/12/2001	1157	1157		11:57:00	494001	FMV3	90.3
8/12/2001	1407	1407		14:07:00	494002	FMV3	90.3
8/12/2001	1409	1409		14:09:00	494003	FMV3	90.3
8/12/2001	1412	1412		14:12:00	494004	FMV3	90.3

Figure 31

The screenshot shows the Excel interface with the 'Format Cells' dialog box open, displaying the 'Number' tab. The dialog box has a 'Category' list on the left and a 'Type' list on the right. The 'Type' list is currently set to 'Number'. The background shows a table with columns: DATE, TIMEZONE, and a table of numerical data. The table has 12 rows of data, with the first row having a date of 8/4/1997 and a time of 20:47:19.

DATE	TIMEZONE	
8/4/1997	20:47:19	
8/4/1997	20:48:18	
8/4/1997	20:49:15	
8/4/1997	20:50:13	
8/4/1997	20:51:13	
8/4/1997	20:52:13	
8/4/1997	20:53:13	
8/4/1997	20:54:12	
8/4/1997	20:55:13	
8/4/1997	20:56:13	
8/4/1997	20:57:13	
8/4/1997	20:58:13	

Figure 32

Produce Excel file from unstructured files

The results from an ARIES dataset can be rather unstructured, an example of which we see here (Figure 33). All data Chrons and Concentrations are copy pasted as values on a new workbook. All cells of the Chrons column are selected. At the Home tab we choose Fill&Select, Go To Special, then select blanks option. At the Home tab we choose Delete, Delete Sheet Rows. All rows that were initially empty are deleted. The empty columns are deleted manually. All working cells are selected, CTRL+T to turn the working cells into a table with headers. All columns are selected, then from Data tab we choose from Table/Range to enter the data into Power Query. All column except for the first are selected and from the Data Type option we choose Text. We then choose the 1st column. Right click and Unpivot other columns. Rename 1st column as HAUL_NO, Attribute column into CLASS and Value column into VALUE. Close&Load. In the Excel worksheet we choose the VALUE column and from Home tab we choose Find&Select, Replace, in Find What we enter “,” and in “Replace with” we enter “.”. Then hit Replace All. Then we select the entire column, right click, format cells, Number and fill decimal places with 6. From the Data tab we choose Text to columns and in the pop-up window we hit Finish. Then we add 2 columns, one is named TYPE and filled with “CONCENTRATION” and one named CRUISE_CODE and filled with S1801. The cruise code is copy pasted as value on the same column. The entire column is selected, right click, format cells,

Text option is chosen. From the Data tab we choose Text to columns and in the pop-up window we hit Finish. Then we save the workbook keeping only the working sheet with the appropriate name e.g. aries_results_S1801.

File	Home	Insert	Page Layout	Formulas	Data	Review	View	Help	PowerPoint																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																
------	------	--------	-------------	----------	------	--------	------	------	------------	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

Figure 33

Appendix-Data Wrangling in Python

Import Packages

```
# PACKAGES IMPORT AND CHECK OF DIRECTORY USED

import os

import glob

import pandas as pd

from pathlib import Path

import csv

os.getcwd() #check we are at the right directory, where files are kept
```

Two step process: excel files to csv & merging of csv files

Before the code below runs, we need to copy paste all excel files of a certain collection e.g. ARIES into the folder in Jupyter Notebook that has been created for this procedure. The directory of the files could have been used directly, but due to reasons of data safety this path was chosen to be followed. This makes sure that the excel files created via the data wrangling process are kept intact.

TWO STEP EXCEL FILES TO CSVS & MERGE OF CSVS TO ONE CSV FOR BULK INSERT

STEP 1

massively transform xlsx/xls to csv

for file in glob.glob("*.xlsx"):

```
df = pd.read_excel(file)

#base=os.path.basename(file)

p = Path(file).stem

df.to_csv(p+'_csv'+'.csv', index=False, header= True)

os.remove(file)
```

STEP 2

combine all csv files into 1 csv file

extension = 'csv'

all_filenames = [i for i in glob.glob('*.{}'.format(extension))]

#combine all files in the list

combined_csv = pd.concat([pd.read_csv(f, header = 0) for f in all_filenames])

#export to csv

combined_csv.to_csv("allrec.csv", index=False, encoding='utf-8-sig')

#print merged csv

print(combined_csv.head(1000))

Number of decimals changed to six

The number of decimals in columns containing needs to be down to 6 according to Marine Scotland request.

CHANGE NUMBER OF DECIMALS ON FILES

#read in initial file

dataset = pd.read_csv('allrec.csv')

#returns input if failure occurs (e.g. if not a float)

def try_cutoff(x):

try:


```

    return round(float(x), 6)

except Exception:

    return x

#loop over each column and apply try_cutoff to each element in each column

#IF VALUE COLUMN ONLY REMOVE '#' AND PUT A '#' IN SECOND hlist lines

hlist = ["VALUE"]

hlist = ['Start_latitude', 'Start_longitude', 'Mid_latitude', 'Mid_longitude', 'End_latitude',
'End_longitude','Start_echo', 'Mid_echo', 'End_echo', 'Volume_dec_cubic_m','Start_depth', 'Mid_depth',
'End_depth']

for field in hlist:

    dataset[field] = dataset[field].map(try_cutoff)

#write new dataset result to CSV file

dataset.to_csv("allrec_file.csv", index = False)

```

The code above runs for a ARIES, OPC, OCEAN and MIKD collections separately. Then the records csv file for each collection is put into the folder and the **##STEP2** and **#CHANGE NUMBER OF DECIMALS ON FILES** is run again creating a csv file that contains records of all collections. The same procedure is followed for sample data too.

Drop column in for csv containing SAMPLES data

At first the student was not sure if the volume column would be kept in the Samples data of the OPC collection and if the volume would need to be imported as column other collections too, so the OPC samples csv file created that had all OPC Sample data in it had this column too. After discussions with the client it was decided that the column would be dropped. For this to come true the following code is run for the OPC sample data csv file.

```

# DROP COLUMN IN FINAL CSV

f=pd.read_csv("samples_filename.csv ")

keep_col = ['HAUL_NO', 'CLASS', 'TYPE', 'VALUE', 'SAMPLER', 'CRUISE_CODE']

new_f = f[keep_col]

new_f.to_csv("newFile.csv", index=False)

```

Create empty column for SAMPLES csv file for future FK population

The csv file containing the sample data of all collections needed the adding of an empty column at the end, where the foreign key would be later on created. So, the following code ran.

```
# CREATE EMPTY COLUMN AT THE END OF FILE TO IMPORT TO STAGE_SAMPLES TABLE
```

```
##THE EMPTY COLUMN MATCHES THE FK_REC_ID COLUMN OF THE TABLE
```

```
import pandas as pd  
f = pd.read_csv("samples_filename.csv")  
f["FK_REC_ID"] = ""  
f.to_csv("samples_filename.csv ", index=False)
```

Create empty column for SAMPLES csv file for alternative database for future incremented PK population

For the case of the alternative database one more step needs to be added at this very point, which is adding another column at the beginning of the dataset as shown in the code below:

```
# INSERT EMPTY COLUMN AT THE BEGINNING OF THE csv FILE
```

```
g = pd.read_csv("samples_filename.csv ")  
g.insert(0, 'PK_staging_samples', "", allow_duplicates=True)  
g.head()  
g.to_csv("newStaginSamplesFile.csv", index=False)  
g.head()
```

Count number of rows in csv files

Counting the number of rows is important to make sure that all rows are uploaded in the database by comparing the number from python scripts to the number from SQL SERVER counting of rows.

```
# COUNT NUMBER OF ROWS IN CSV FILE-TESTING
```

```
with open('allrec_file.csv') as f:  
    z = sum(1 for line in f)  
    print(z)
```

Appendix-Zooplankton Database Creation

After the steps above two csv files have been created: one for the RECORDS table and one for The SAMPLES table. We open the SSMS18 → right click on databases → create new database → accept all default parameters and enter “Zooplankton” as the name of the database → Right click on the database → choose New Query. In the query tab opened the following code is inserted:

STAGE_RECORDS, bulk insert & RECORDS table

/* Entire database creation*/

---Create STAGE_RECORDS table

```
CREATE TABLE [dbo].[STAGE_RECORDS](
    [DATE] [date] NOT NULL,
    [Timestamp] [time](0) NOT NULL,
    [HAUL_NO] [int] NOT NULL,
    [Station_name] [nvarchar](50) NULL,
    [Start_latitude] [decimal](18, 6) NULL,
    [Start_longitude] [decimal](18, 6) NULL,
    [Mid_latitude] [decimal](18, 6) NULL,
    [Mid_longitude] [decimal](18, 6) NULL,
    [End_latitude] [decimal](18, 6) NULL,
    [End_longitude] [decimal](18, 6) NULL,
    [Distance] [nvarchar](50) NULL,
    [Start_echo] [decimal](18, 6) NULL,
    [Mid_echo] [decimal](18, 6) NULL,
    [End_echo] [decimal](18, 6) NULL,
    [SAMPLER] [nvarchar](50) NOT NULL,
    [Mesh_micro_m] [nvarchar](50) NULL,
    [Volume_dec_cubic_m] [decimal](18, 6) NULL,
    [Duration] [time](0) NULL,
    [Start_depth] [decimal](18, 6) NULL,
    [Mid_depth] [decimal](18, 6) NULL,
    [End_depth] [decimal](18, 6) NULL,
    [Flow_meter_type] [nvarchar](50) NULL,
    [Start_count] [nvarchar](50) NULL,
    [End_count] [nvarchar](50) NULL,
    [REVS] [nvarchar](50) NULL,
    [Valid_Y_N] [nvarchar](50) NULL,
    [CRUISE_CODE] [nvarchar](50) NOT NULL
) ON [PRIMARY]
```

```

GO

---Bulk insertion for STAGE_RECORDS table

BULK INSERT [dbo].[STAGE_RECORDS]

FROM

'C:\Users\georgios\Documents\placem\RESULTS_CSV_COMBINED\work_files\allrec_COMBfile_THE_OFFICIAL_WORKING.csv'

WITH (FIRSTROW = 2,

      FIELDTERMINATOR = ',',

      ROWTERMINATOR='\\n' ,

      MAXERRORS=4000000,

      BATCHSIZE=40000);

---Count rows in STAGE_RECORDS table

SELECT COUNT(*) FROM [dbo].[STAGE_RECORDS];

---Check for duplicates in STAGE_RECORDS table and then delete them

/*First we run the query with the SELECT* end statement. Then comment that statement out and uncomment the DELETE FROM statement and re-rerun the query with it.*/

WITH CTE ([HAUL_NO],

          [SAMPLER],

          [CRUISE_CODE],

          duplicatecount)

AS (SELECT [HAUL_NO],

          [SAMPLER],

          [CRUISE_CODE],

          ROW_NUMBER() OVER(PARTITION BY [HAUL_NO],

                               [SAMPLER],

                               [CRUISE_CODE]

                               ORDER BY [HAUL_NO]) AS DuplicateCount

FROM [dbo].[STAGE_RECORDS])

SELECT * FROM CTE WHERE DuplicateCount <> 1; /* Some duplicated have been inserted on purpose for this query to show it is working*/

---DELETE FROM CTE WHERE DuplicateCount <> 1;

---(746 rows affected)

```

---COUNT ROWS IN STAGE_RECORDS TABLE TO CHECK THE CHANGE

```
SELECT COUNT(*) FROM [dbo].[STAGE_RECORDS];
```

It is important to notice that duplicates do exist. This is not due to a problem of the data wrangling process. Some initial Excel files created had the contained data duplicated in order to the test scripts above to run and verify that they are working. This way, when Marine Scotland will use the database for future data entry, they know that they can trust this mechanism to find duplicate records.

---Create RECORDS table

```
CREATE TABLE [dbo].[RECORDS](
    [RECORDS_ID] [int] IDENTITY(1,1) NOT NULL PRIMARY KEY,
    [DATE] [date] NOT NULL,
    [Timestamp] [time](0) NOT NULL,
    [HAUL_NO] [int] NOT NULL,
    [Station_name] [nvarchar](50) NULL,
    [Start_latitude] [decimal](18, 6) NULL,
    [Start_longitude] [decimal](18, 6) NULL,
    [Mid_latitude] [decimal](18, 6) NULL,
    [Mid_longitude] [decimal](18, 6) NULL,
    [End_latitude] [decimal](18, 6) NULL,
    [End_longitude] [decimal](18, 6) NULL,
    [Distance] [nvarchar](50) NULL,
    [Start_echo] [decimal](18, 6) NULL,
    [Mid_echo] [decimal](18, 6) NULL,
    [End_echo] [decimal](18, 6) NULL,
    [SAMPLER] [nvarchar](50) NOT NULL,
    [Mesh_micro_m] [nvarchar](50) NULL,
    [Volume_dec_cubic_m] [decimal](18, 6) NULL,
    [Duration] [time](0) NULL,
    [Start_depth] [decimal](18, 6) NULL,
    [Mid_depth] [decimal](18, 6) NULL,
    [End_depth] [decimal](18, 6) NULL,
    [Flow_meter_type] [nvarchar](50) NULL,
```

```

[Start_count] [nvarchar](50) NULL,
[End_count] [nvarchar](50) NULL,
[REVS] [nvarchar](50) NULL,
[Valid_Y_N] [nvarchar](50) NULL,
[CRUISE_CODE] [nvarchar](50) NOT NULL
) ON [PRIMARY]
GO

```

---Insert data from STAGE_RECORDS table to RECORDS table

```
INSERT INTO [dbo].[RECORDS] SELECT * FROM [dbo].[STAGE_RECORDS];
```

---Count rows in STAGE_RECORDS table

```
SELECT COUNT(*) FROM [dbo].[RECORDS];
```

Creation of CRUISES table

---Create new table CRUISES using task import flat file (Figure 34)

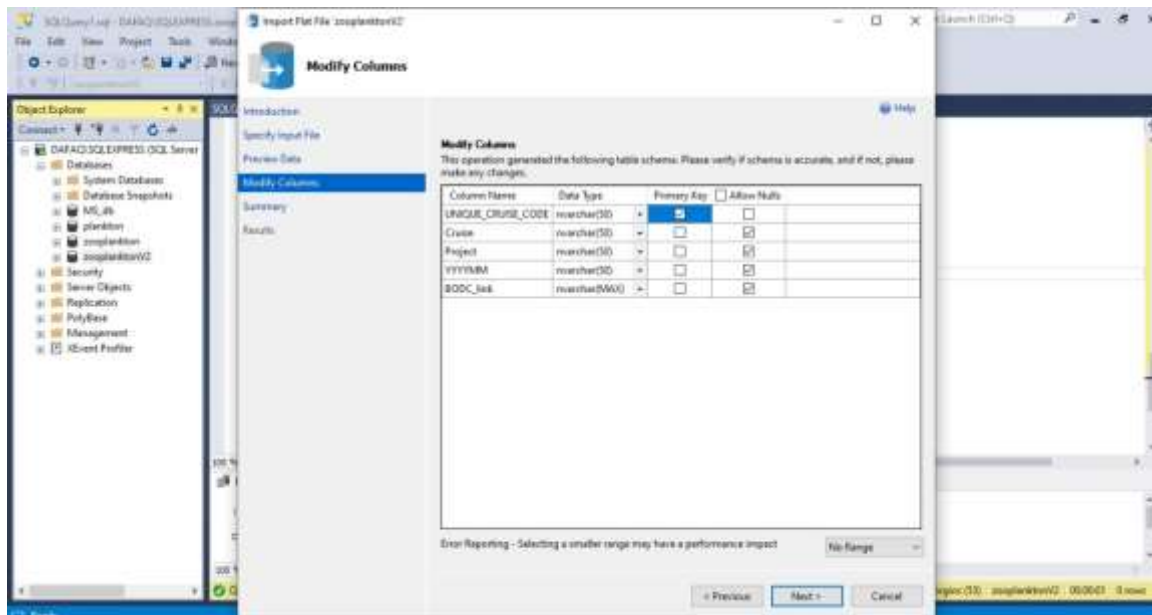


Figure 34 Cruises table

STAGE_SAMPLES & bulk insert

---Create STAGE_SAMPLES table

```

CREATE TABLE [dbo].[STAGE_SAMPLES](
    [HAUL_NO] [int] NOT NULL,
    [CLASS] [nvarchar](100) NULL,
    [TYPE] [nvarchar](50) NULL,
    [VALUE] [decimal](18, 6) NULL,
    [SAMPLER] [nvarchar](50) NOT NULL,
    [CRUISE_CODE] [nvarchar](50) NOT NULL,
    [FK_REC_ID] [int] NULL
) ON [PRIMARY]

GO

---Bulk insertion for STAGE_SAMPLES table

BULK INSERT [dbo].[STAGE_SAMPLES]

FROM

'C:\Users\georgios\Documents\placem\RESULTS_CSV_COMBINED\work_files\allSAMPLE_22_07_OFFICIAL.csv'

WITH (FIRSTROW = 2,

    FIELDTERMINATOR = ',',

    ROWTERMINATOR='\\n' ,

    MAXERRORS=4000000,

    BATCHSIZE=40000);

```

Foreign key relationship creation between CRUISES and RECORDS tables

```

---after populating CRUISES and RECORDS table create FK between them

---FK relationship between RECORDS and CRUISES tables

ALTER TABLE [dbo].[RECORDS] WITH CHECK ADD CONSTRAINT [FK_RECORDS_CRUISES] FOREIGN
KEY([CRUISE_CODE])

REFERENCES [dbo].[CRUISES] ([UNIQUE_CRUISE_CODE])

GO

ALTER TABLE [dbo].[RECORDS] CHECK CONSTRAINT [FK_RECORDS_CRUISES]

GO

```

---Problem arises/ we check problem


```

select Distinct rec.[CRUISE_CODE] from [dbo].[RECORDS] as rec
WHERE rec.[CRUISE_CODE] NOT IN
(SELECT ct.[UNIQUE_CRUISE_CODE] from [dbo].[CRUISES] as ct);
SELECT DISTINCT [CRUISE_CODE] from [dbo].[RECORDS] where [CRUISE_CODE] like 'LIN%';

```

---Make appropriate changes in RECORDS table, then run FK creation command again

```

UPDATE
    [dbo].[RECORDS]
SET
    [CRUISE_CODE] = REPLACE([CRUISE_CODE], 'LINNHE9', 'LINNHE91')
WHERE
    [CRUISE_CODE] LIKE 'LINNHE9';

```

---rerun FK creation between RECORDS nad CRYUISES table

```

ALTER TABLE [dbo].[RECORDS] WITH CHECK ADD CONSTRAINT [FK_RECORDS_CRUISES] FOREIGN
KEY([CRUISE_CODE])
REFERENCES [dbo].[CRUISES] ([UNIQUE_CRUISE_CODE])
GO

ALTER TABLE [dbo].[RECORDS] CHECK CONSTRAINT [FK_RECORDS_CRUISES]
GO

```

Fixes in STAGE_SAMPLES table, creation of FK, creation of SAMPLES table & insert of data

---Now check STAGE_SAMPLES FOR LINNHE91

```

SELECT DISTINCT [CRUISE_CODE] from [dbo].[STAGE_SAMPLES] where [CRUISE_CODE] like 'LIN%';

```

---Make appropriate change

```

UPDATE
    [dbo].[STAGE_SAMPLES]
SET
    [CRUISE_CODE] = REPLACE([CRUISE_CODE], 'LINHE91', 'LINNHE91')
WHERE
    [CRUISE_CODE] LIKE 'LINHE91';

```

---Create FK in STAGE_SAMPLES table based on constraints

```
UPDATE [dbo].[STAGE_SAMPLES]
SET [dbo].[STAGE_SAMPLES].[FK_REC_ID] = sr.[RECORDS_ID]
FROM [dbo].[RECORDS] AS sr, [dbo].[STAGE_SAMPLES] AS da
WHERE sr.HAUL_NO = da.HAUL_NO
AND sr.CRUISE_CODE = da.CRUISE_CODE
AND sr.[SAMPLER] = da.SAMPLER;
```

---Create SAMPLES Table

```
CREATE TABLE [dbo].[SAMPLES](
    [SAMPLE_ID] [int] IDENTITY(1,1) NOT NULL PRIMARY KEY,
    [HAUL_NO] [int] NOT NULL,
    [CLASS] [nvarchar](100) NOT NULL,
    [TYPE] [nvarchar](50) NULL,
    [VALUE] [decimal](18, 6) NULL,
    [SAMPLER] [nvarchar](50) NOT NULL,
    [CRUISE_CODE] [nvarchar](50) NULL,
    [FK_REC_ID] [int] NULL
) ON [PRIMARY]
GO
```

---Insert data from STAGE_SAMPLES table to SAMPLES table

```
INSERT INTO [dbo].[SAMPLES] SELECT * FROM [dbo].[STAGE_SAMPLES];
```

Creation of Species_info table

---Import Species_info table with flat file import wizard (Figure 35)

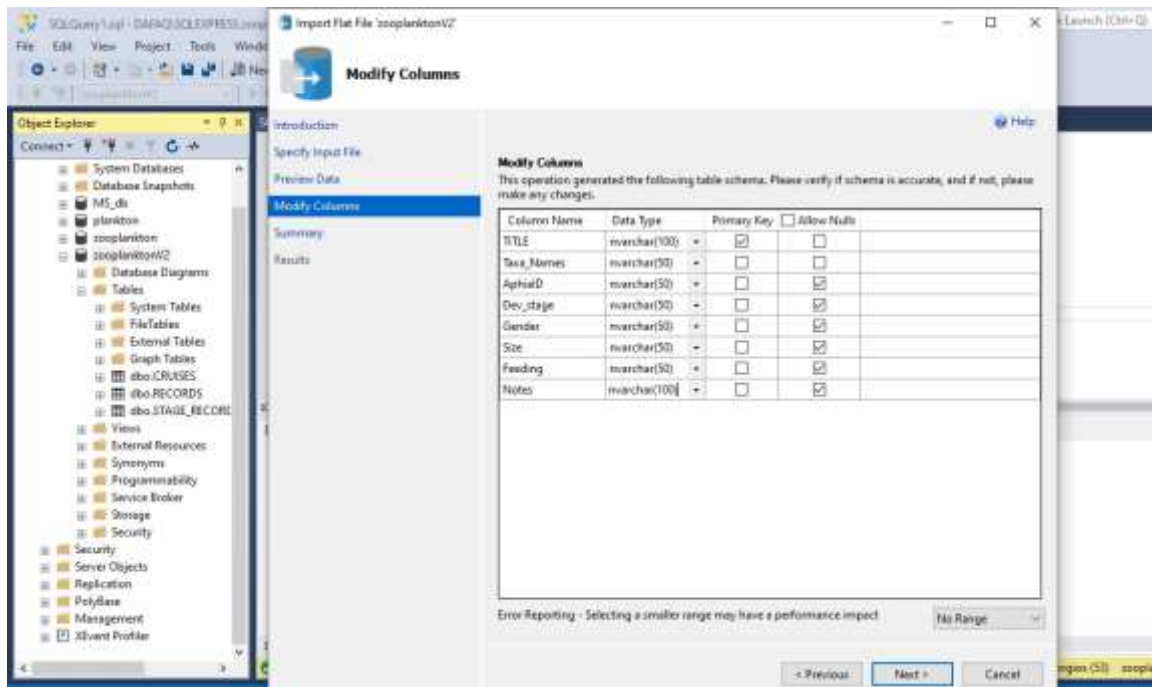


Figure 35 Species_info table

Changes on tables before Foreign Key relationships creation

---Perform left rimming on class from Samples and TITLE&Taxa_Names from Species_info table

```
UPDATE [dbo].[SAMPLES] SET CLASS = LTRIM(CLASS);
UPDATE [dbo].[Species_info] SET [Taxa_Names] = LTRIM([Taxa_Names]);
UPDATE [dbo].[Species_info] SET [TITLE] = LTRIM([TITLE]);
```

---Check not matching elements between tables for problem in FK creation

```
SELECT Distinct [CLASS] AS dCLASS FROM [dbo].[SAMPLES]
WHERE [CLASS] NOT IN
(SELECT [TITLE] from [dbo].[Species_info]) ORDER BY dCLASS;

SELECT [TITLE] FROM [dbo].[Species_info];
```

---Check for wrongs

```
select distinct CLASS from [dbo].[SAMPLES] WHERE [CLASS] LIKE '%All%';
```

Last updates on tables

Now that not matching elements have been found with the script above some changes need to be done like the one following. The “gooseberry” was identified from the initial excel file due to its name uniqueness and it was intentionally left in a wrong format to test the script above and create a script that would easily implement name changes, like the one below. Such script is important, because names of species may change occasionally due to name updates scientifically decided or due to mistakes that were not noticed in the initial files entered in the database.

---Update/replace gooseberry name

```
UPDATE
    [dbo].[Species_info]
SET
    [TITLE] = REPLACE([TITLE], '"gooseberry"', 'gooseberry')
WHERE
    [TITLE] = '"gooseberry"';

SELECT [TITLE] FROM [dbo].[Species_info];
```

---Check not matching elements between tables for problem in FK creation

```
select Distinct sa.[CRUISE_CODE] from [dbo].[SAMPLES] as sa
WHERE sa.[CRUISE_CODE] NOT IN
(SELECT ct.[UNIQUE_CRUISE_CODE] from [dbo].[CRUISES] as ct);
```

---Check and make other changes in SAMPLES table

```
SELECT DISTINCT [TYPE] AS dTYPE FROM [dbo].[SAMPLES] WHERE [TYPE] LIKE 'CO%';

UPDATE
    [dbo].[SAMPLES]
SET
    [TYPE] = REPLACE([TYPE], 'COUNTS', 'COUNT')
WHERE
    [TYPE] = 'COUNTS';

SELECT DISTINCT [TYPE] AS dTYPE FROM [dbo].[SAMPLES] WHERE [TYPE] LIKE 'CO%';
```

Create remaining Foreign Key Relationships

---Create FK relationship between Records and Samples table

```
ALTER TABLE [dbo].[SAMPLES] WITH CHECK ADD CONSTRAINT [FK_SAMPLES_RECORDS] FOREIGN  
KEY([FK_REC_ID])  
  
REFERENCES [dbo].[RECORDS] ([RECORDS_ID])  
  
GO  
  
ALTER TABLE [dbo].[SAMPLES] CHECK CONSTRAINT [FK_SAMPLES_RECORDS]  
  
GO
```

---FK relationship between SAMPLES and CRUISES tables

```
ALTER TABLE [dbo].[SAMPLES] WITH CHECK ADD CONSTRAINT [FK_SAMPLES_CRUISES] FOREIGN  
KEY([CRUISE_CODE])  
  
REFERENCES [dbo].[CRUISES] ([UNIQUE_CRUISE_CODE])  
  
GO  
  
ALTER TABLE [dbo].[SAMPLES] CHECK CONSTRAINT [FK_SAMPLES_CRUISES]  
  
GO
```

---Add FOREIGN KEY relationship with NO CHECK between SAMPLES and Species_info table

```
ALTER TABLE [dbo].[SAMPLES]  
  
WITH NOCHECK ADD CONSTRAINT [FK_SAMPLES_SpInfo]  
FOREIGN KEY (CLASS) REFERENCES [dbo].[Species_info]([TITLE])  
  
GO
```

Merge and Update on future data

As mentioned before new data may come into play at a later point in time and the client should be able to update the working RECORDS and SAMPLES tables. So once data has been uploaded in the staging tables using the bulk insertion mechanism already described, the following query would make sure that all new data is uploaded in the RECORDS table

---how to update existing and insert new

```
MERGE INTO  
  
[dbo].[RECORDS] as rec  
  
USING  
  
[dbo].[STAGE_RECORDS] as sr
```

```

ON

rec.[HAUL_NO] =sr.[HAUL_NO]

and rec.[SAMPLER] =sr.[SAMPLER] and rec.[CRUISE_CODE] =sr.[CRUISE_CODE]

WHEN MATCHED THEN

UPDATE SET

rec.[DATE] =sr.[DATE] , rec.[Timestamp] =sr.[Timestamp] , rec.[HAUL_NO] =sr.[HAUL_NO] ,
rec.[Station_name] =sr.[Station_name] ,

rec.[Start_latitude] =sr.[Start_latitude] , rec.[Start_longitude] =sr.[Start_longitude] ,
rec.[Mid_latitude] =sr.[Mid_latitude] ,

rec.[Mid_longitude] =sr.[Mid_longitude] , rec.[End_latitude] =sr.[End_latitude] ,
rec.[End_longitude] =sr.[End_longitude] ,

rec.[Distance] =sr.[Distance] , rec.[Start_echo] =sr.[Start_echo] , rec.[Mid_echo] =sr.[Mid_echo] ,
rec.[End_echo] =sr.[End_echo] ,

rec.[SAMPLER] =sr.[SAMPLER] , rec.[Mesh_micro_m] =sr.[Mesh_micro_m] ,
rec.[Volume_dec_cubic_m] =sr.[Volume_dec_cubic_m] ,

rec.[Duration] =sr.[Duration] , rec.[Start_depth] =sr.[Start_depth] , rec.[Mid_depth]
=sr.[Mid_depth] , rec.[End_depth] =sr.[End_depth] ,

rec.[Flow_meter_type] =sr.[Flow_meter_type] , rec.[Start_count] =sr.[Start_count] ,
rec.[End_count] =sr.[End_count] ,

rec.[REVS] =sr.[REVS] , rec.[Valid_Y_N] =sr.[Valid_Y_N] , rec.[CRUISE_CODE] =sr.[CRUISE_CODE]

WHEN NOT MATCHED THEN

INSERT

([DATE] , [Timestamp] , [HAUL_NO] , [Station_name] , [Start_latitude] , [Start_longitude] ,
[Mid_latitude] ,

[Mid_longitude] , [End_latitude] , [End_longitude] , [Distance] , [Start_echo] , [Mid_echo] ,
[End_echo] ,

[SAMPLER] , [Mesh_micro_m] , [Volume_dec_cubic_m] , [Duration] , [Start_depth]
, [Mid_depth] , [End_depth] ,

[Flow_meter_type] , [Start_count] , [End_count] , [REVS] , [Valid_Y_N]
, [CRUISE_CODE]

)

VALUES

(sr.[DATE] ,sr.[Timestamp] ,sr.[HAUL_NO] ,sr.[Station_name] , sr.[Start_latitude] , sr.[Start_longitude] ,
sr.[Mid_latitude] , sr.[Mid_longitude] ,

sr.[End_latitude] , sr.[End_longitude] , sr.[Distance] , sr.[Start_echo] , sr.[Mid_echo] , sr.[End_echo] ,
sr.[SAMPLER] , sr.[Mesh_micro_m] ,

sr.[Volume_dec_cubic_m] , sr.[Duration] , sr.[Start_depth] , sr.[Mid_depth] , sr.[End_depth] ,
sr.[Flow_meter_type] , sr.[Start_count] , sr.[End_count] ,

```

```
sr.[REVS] , sr.[Valid_Y_N] , sr.[CRUISE_CODE]);
```

and the following query would make sure that all new data is uploaded in the SAMPLES table

```
MERGE INTO
  [dbo].[SAMPLES] AS sa
USING
  [dbo].[STAGE_SAMPLES] AS ss
ON
  sa.[HAUL_NO] =ss.[HAUL_NO]
    and sa.[SAMPLER] =ss.[SAMPLER] and sa.[CRUISE_CODE] =ss.[CRUISE_CODE]
WHEN MATCHED THEN
  UPDATE SET
    sa.[HAUL_NO] =ss.[HAUL_NO] ,   sa.[CLASS] =ss.[CLASS] , sa.[TYPE] =ss.[TYPE] ,      sa.[VALUE]
=ss.[VALUE] ,
    sa.[SAMPLER] =ss.[SAMPLER] ,   sa.[CRUISE_CODE] =ss.[CRUISE_CODE] ,      sa.[FK_REC_ID]
=ss.[FK_REC_ID]

WHEN NOT MATCHED THEN
  INSERT
    ([HAUL_NO] ,   [CLASS], [TYPE] , [VALUE],[SAMPLER],      [CRUISE_CODE], [FK_REC_ID]
)
  VALUES
    (ss.[HAUL_NO] ,ss.[CLASS] ,ss.[TYPE] ,ss.[VALUE] , ss.[SAMPLER] , ss.[CRUISE_CODE] , ss.[FK_REC_ID]);
```

Appendix-Alternative Database Implementation

The process in the SSMS and code used to create it is very similar to the originally created database with the exceptions documented here with blue:

STAGE_RECORDS, bulk insert & RECORDS table

```
/*final zooplankton_db_alternative creation*/
CREATE TABLE [dbo].[STAGE_RECORDS](
    [DATE] [date] NOT NULL,
    [Timestamp] [time](0) NOT NULL,
    [HAUL_NO] [int] NOT NULL,
    [Station_name] [nvarchar](50) NULL,
    [Start_latitude] [decimal](18, 6) NULL,
    [Start_longitude] [decimal](18, 6) NULL,
    [Mid_latitude] [decimal](18, 6) NULL,
```



```

[Mid_longitude] [decimal](18, 6) NULL,
[End_latitude] [decimal](18, 6) NULL,
[End_longitude] [decimal](18, 6) NULL,
[Distance] [nvarchar](50) NULL,
[Start_echo] [decimal](18, 6) NULL,
[Mid_echo] [decimal](18, 6) NULL,
[End_echo] [decimal](18, 6) NULL,
[SAMPLER] [nvarchar](50) NOT NULL,
[Mesh_micro_m] [nvarchar](50) NULL,
[Volume_dec_cubic_m] [decimal](18, 6) NULL,
[Duration] [time](0) NULL,
[Start_depth] [decimal](18, 6) NULL,
[Mid_depth] [decimal](18, 6) NULL,
[End_depth] [decimal](18, 6) NULL,
[Flow_meter_type] [nvarchar](50) NULL,
[Start_count] [nvarchar](50) NULL,
[End_count] [nvarchar](50) NULL,
[REVS] [nvarchar](50) NULL,
[Valid_Y_N] [nvarchar](50) NULL,
[CRUISE_CODE] [nvarchar](50) NOT NULL
) ON [PRIMARY]
GO

---Bulk insertion for STAGE_RECORDS table

BULK INSERT [dbo].[STAGE_RECORDS]

FROM

'C:\Users\georgios\Documents\placem\RESULTS_CSV_COMBINED\work_files\allrec_COMBfile_THE_OFFICIAL_WORKING.csv'

WITH (FIRSTROW = 2,

    FIELDTERMINATOR = ',',

    ROWTERMINATOR='\\n' ,

    MAXERRORS=4000000,

    BATCHSIZE=40000);

```

---COUNT ROWS IN STAGE_RECORDS TABLE

```
SELECT COUNT(*) FROM [dbo].[STAGE_RECORDS];
```

---Check for duplicates in STAGE_RECORDS table and then delete them

```
WITH CTE ([HAUL_NO],
          [SAMPLER],
          [CRUISE_CODE],
          duplicatecount)
AS (SELECT [HAUL_NO],
          [SAMPLER],
          [CRUISE_CODE],
          ROW_NUMBER() OVER(PARTITION BY [HAUL_NO],
                               [SAMPLER],
                               [CRUISE_CODE]
                              ORDER BY [HAUL_NO]) AS DuplicateCount
   FROM [dbo].[STAGE_RECORDS])

SELECT * FROM CTE WHERE DuplicateCount <> 1; /* Some duplicated have been inserted on purpose for this
query to show it is working*/

---DELETE FROM CTE WHERE DuplicateCount <> 1;

---(746 rows affected)
```

---COUNT ROWS IN STAGE_RECORDS TABLE TO CHECK THE CHANGE

```
SELECT COUNT(*) FROM [dbo].[STAGE_RECORDS];
```

---Create RECORDS table

```
CREATE TABLE [dbo].[RECORDS](
    [RECORDS_ID] [int] IDENTITY(1,1) NOT NULL PRIMARY KEY,
    [DATE] [date] NOT NULL,
    [Timestamp] [time](0) NOT NULL,
    [HAUL_NO] [int] NOT NULL,
    [Station_name] [nvarchar](50) NULL,
    [Start_latitude] [decimal](18, 6) NULL,
    [Start_longitude] [decimal](18, 6) NULL,
    [Mid_latitude] [decimal](18, 6) NULL,
```

```

[Mid_longitude] [decimal](18, 6) NULL,
[End_latitude] [decimal](18, 6) NULL,
[End_longitude] [decimal](18, 6) NULL,
[Distance] [nvarchar](50) NULL,
[Start_echo] [decimal](18, 6) NULL,
[Mid_echo] [decimal](18, 6) NULL,
[End_echo] [decimal](18, 6) NULL,
[SAMPLER] [nvarchar](50) NOT NULL,
[Mesh_micro_m] [nvarchar](50) NULL,
[Volume_dec_cubic_m] [decimal](18, 6) NULL,
[Duration] [time](0) NULL,
[Start_depth] [decimal](18, 6) NULL,
[Mid_depth] [decimal](18, 6) NULL,
[End_depth] [decimal](18, 6) NULL,
[Flow_meter_type] [nvarchar](50) NULL,
[Start_count] [nvarchar](50) NULL,
[End_count] [nvarchar](50) NULL,
[REVS] [nvarchar](50) NULL,
[Valid_Y_N] [nvarchar](50) NULL,
[CRUISE_CODE] [nvarchar](50) NOT NULL
) ON [PRIMARY]
GO

```

---Insert data from STAGE_RECORDS table to RECORDS table

```
INSERT INTO [dbo].[RECORDS] SELECT * FROM [dbo].[STAGE_RECORDS];
```

---COUNT ROWS IN STAGE_RECORDS TABLE

```
SELECT COUNT(*) FROM [dbo].[RECORDS];
```

STAGE_SAMPLES & bulk insert

---Create STAGE_SAMPLES table

```
CREATE TABLE [dbo].[STAGE_SAMPLES](
    [PK_staging_samples] [int] IDENTITY(1,1) NOT NULL PRIMARY KEY,
```

```

[HAUL_NO] [int] NOT NULL,
[CLASS] [nvarchar](100) NULL,
[TYPE] [nvarchar](50) NULL,
[VALUE] [decimal](18, 6) NULL,
[SAMPLER] [nvarchar](50) NOT NULL,
[CRUISE_CODE] [nvarchar](50) NOT NULL,
[FK_REC_ID] [int] NULL
) ON [PRIMARY]

```

GO

---Bulk insertion for STAGE_SAMPLES table

```
BULK INSERT [dbo].[STAGE_SAMPLES]
```

```
FROM
```

```
'C:\Users\georgios\Documents\placem\RESULTS_CSV_COMBINED\work_files\for_alternativedb\newStaginSamplesFile.csv'
```

```
WITH (FIRSTROW = 2,
```

```
    FIELDTERMINATOR = ',',
```

```
    ROWTERMINATOR='\\n' ,
```

```
    MAXERRORS=4000000,
```

```
    BATCHSIZE=40000);
```

Creation of CRUISES table and FK relationship with RECORDS table

---Create new table CRUISES using task import flat file (Figure 34)

---after populating CRUISES and RECORDS table create FK between them

---FK relationship between RECORDS and CRUISES tables

```
ALTER TABLE [dbo].[RECORDS] WITH CHECK ADD CONSTRAINT [FK_RECORDS_CRUISES] FOREIGN
KEY([CRUISE_CODE])
```

```
REFERENCES [dbo].[CRUISES] ([UNIQUE_CRUISE_CODE])
```

GO

```
ALTER TABLE [dbo].[RECORDS] CHECK CONSTRAINT [FK_RECORDS_CRUISES]
```

GO

---Problem arises/ we check problem

```
SELECT DISTINCT rec.[CRUISE_CODE] from [dbo].[RECORDS] as rec
WHERE rec.[CRUISE_CODE] NOT IN
(SELECT ct.[UNIQUE_CRUISE_CODE] FROM [dbo].[CRUISES] AS ct);
SELECT DISTINCT [CRUISE_CODE] FROM [dbo].[RECORDS] WHERE [CRUISE_CODE] like 'LIN%';
```

---Make appropriate changes in RECORDS table, then run FK creation command again

```
UPDATE
    [dbo].[RECORDS]
SET
    [CRUISE_CODE] = REPLACE([CRUISE_CODE], 'LINNHE9', 'LINNHE91')
WHERE
    [CRUISE_CODE] LIKE 'LINNHE9';
```

---rerun FK creation between RECORDS and CRUISES table

```
ALTER TABLE [dbo].[RECORDS] WITH CHECK ADD CONSTRAINT [FK_RECORDS_CRUISES] FOREIGN
KEY([CRUISE_CODE])
REFERENCES [dbo].[CRUISES] ([UNIQUE_CRUISE_CODE])
GO

ALTER TABLE [dbo].[RECORDS] CHECK CONSTRAINT [FK_RECORDS_CRUISES]
GO
```

Fixes in STAGE_SAMPLES table, creation of FK, creation of SAMPLES table & insert of data

---Now check STAGE_SAMPLES FOR LINNHE91

```
SELECT DISTINCT [CRUISE_CODE] from [dbo].[STAGE_SAMPLES] where [CRUISE_CODE] like 'LIN%';
```

---Make appropriate change

```
UPDATE
    [dbo].[STAGE_SAMPLES]
SET
    [CRUISE_CODE] = REPLACE([CRUISE_CODE], 'LINHE91', 'LINNHE91')
```

WHERE

[CRUISE_CODE] LIKE 'LINHE91';

---populate fk_rec column based on RECORDS tables constraints

UPDATE [dbo].[STAGE_SAMPLES]

SET [dbo].[STAGE_SAMPLES].[FK_REC_ID] = sr.[RECORDS_ID]

FROM [dbo].[RECORDS] AS sr, [dbo].[STAGE_SAMPLES] AS da

WHERE sr.HAUL_NO = da.HAUL_NO

AND sr.CRUISE_CODE = da.CRUISE_CODE

AND sr.[SAMPLER] = da.SAMPLER;

---Create SAMPLES Table

CREATE TABLE [dbo].[SAMPLES](

[SAMPLE_ID] [int] NOT NULL PRIMARY KEY,

[CLASS] [nvarchar](100) NULL,

[TYPE] [nvarchar](50) NULL,

[VALUE] [decimal](18, 6) NULL,

[Taxa_Names] [nvarchar](50) NULL,

[AphialID] [nvarchar](50) NULL,

[Dev_stage] [nvarchar](50) NULL,

[Gender] [nvarchar](50) NULL,

[Size] [nvarchar](50) NULL,

[Feeding] [nvarchar](50) NULL,

[Notes] [nvarchar](100) NULL,

[FK_REC_ID] [int] NULL

) ON [PRIMARY]

GO

---Insert data from STAGE_SAMPLES table SAMPLES table

INSERT INTO [dbo].[SAMPLES]([SAMPLE_ID], [CLASS], [TYPE], [VALUE], [FK_REC_ID])

SELECT [PK_staging_samples], [CLASS], [TYPE], [VALUE], [FK_REC_ID]

FROM [dbo].[STAGE_SAMPLES];

---Make appropriate changes to SAMPLES table

```
UPDATE [dbo].[SAMPLES] SET CLASS = LTRIM(CLASS);
```

```
UPDATE
```

```
    [dbo].[SAMPLES]
```

```
SET
```

```
    [TYPE] = REPLACE([TYPE], 'COUNTS', 'COUNT')
```

```
WHERE
```

```
    [TYPE] = 'COUNTS';
```

Creation of Species_info table and make fixes in it

---Import species_info table with flat file import wizard (Figure 35)

---Perform left rimming on class from Samples and TITLE&Taxa_Names from Species_info table

```
UPDATE [dbo].[Species_info] SET TITLE = LTRIM(TITLE);
```

---Update/replace gooseberry name in Species_info table

```
UPDATE
```

```
    [dbo].[Species_info]
```

```
SET
```

```
    [TITLE] = REPLACE([TITLE], "'gooseberry'", 'gooseberry')
```

```
WHERE
```

```
    [TITLE] = "'gooseberry'";
```

```
SELECT [TITLE] FROM [dbo].[Species_info];
```

Insert data from Species_info table to SAMPLES table & create FK relationship with RECORDS table

---UPDATE SAMPLES table from Species_info table

```
UPDATE [dbo].[SAMPLES]
```

```
    SET [SAMPLES].[Taxa_Names] = sp.[Taxa_Names],
```

```
    [SAMPLES].[AphiaID] = sp.[AphiaID],
```

```
    [SAMPLES].[Dev_stage] = sp.[Dev_stage],
```

```
    [SAMPLES].[Gender] = sp.[Gender],
```



```

[SAMPLES].[Size] = sp.[Size],
[SAMPLES].[Feeding] = sp.[Feeding],
[SAMPLES].[Notes] = sp.[Notes]
FROM [dbo].[SAMPLES] AS sa INNER JOIN [dbo].[Species_info] as sp ON sa.[CLASS] = sp.[TITLE]

```

---Create FK relationship between Records and Samples table

```

ALTER TABLE [dbo].[SAMPLES] WITH CHECK ADD CONSTRAINT [FK_SAMPLES_RECORDS] FOREIGN
KEY([FK_REC_ID])
REFERENCES [dbo].[RECORDS] ([RECORDS_ID])
GO

```

Merge and Update on future data

---UPDATE SAMPLES TABLE

---STEP1

```

MERGE INTO
    [dbo].[SAMPLES] AS sa
USING
    [dbo].[STAGE_SAMPLES] AS ss
ON
    sa.[SAMPLE_ID] = ss.[PK_staging_samples]
WHEN MATCHED THEN
    UPDATE SET
        sa.[HAUL_NO] =ss.[HAUL_NO] , sa.[CLASS] =ss.[CLASS] , sa.[TYPE] =ss.[TYPE] , sa.[VALUE] =ss.[VALUE] ,
        sa.[FK_REC_ID] =ss.[FK_REC_ID]

WHEN NOT MATCHED THEN
    INSERT
        ([SAMPLE_ID], [CLASS], [TYPE] , [VALUE], [FK_REC_ID]
    )
VALUES
    (ss.[PK_staging_samples] ,ss.[CLASS] ,ss.[TYPE] ,ss.[VALUE] , , ss.[FK_REC_ID]);

```

---STEP 2

```

MERGE INTO

```

```

[dbo].[SAMPLES] AS sa
USING
[dbo].[Species_info] AS sp
ON
sa.[CLASS] = sp.[TITLE]
WHEN MATCHED THEN
UPDATE SET
sa.[Taxa_Names] =sp.[Taxa_Names] ,    sa.[AphialD]  =sp.[AphialD] , sa.[Dev_stage]  =sp.[Dev_stage] ,
sa.[Gender] =sp.[Gender],
sa.[Size] =sp.[Size] , sa.[Feeding] =sp.[Feeding], sa.[Notes] =sp.[Notes]

WHEN NOT MATCHED THEN
INSERT
([Taxa_Names],[AphialD],      [Dev_stage], [Gender], [Size], [Feeding], [Notes]
)
VALUES
(sp.[Taxa_Names],sp.[AphialD],  sp.[Dev_stage], sp.[Gender], sp.[Size], sp.[Feeding], sp.[Notes]);

```

Appendix-Data Mining in R

SCENARIO_1

First the function that will handle the scripting that follows needs to be created.

function to run scripts

```
g <- function(sql){  
  sqlQuery(dbhandle,sql)  
}
```

The researcher first needs to verify existing records that are available in both tables according to cruise codes
Figure 36

SCENARIO_1

SAMPLE TABLE ON GENERIC SEARCH BASED ON CRUISE CODE

```
sql <- 'SELECT DISTINCT rec.CRUISE_CODE AS records_Cruise, sa.CRUISE_CODE as samples_Cruise FROM [dbo].[RECORDS] AS rec  
FULL JOIN [dbo].[SAMPLES] as sa ON sa.FK_REC_ID = rec.RECORDS_ID;'
```

```
df <- g(sql)
```

```
df
```

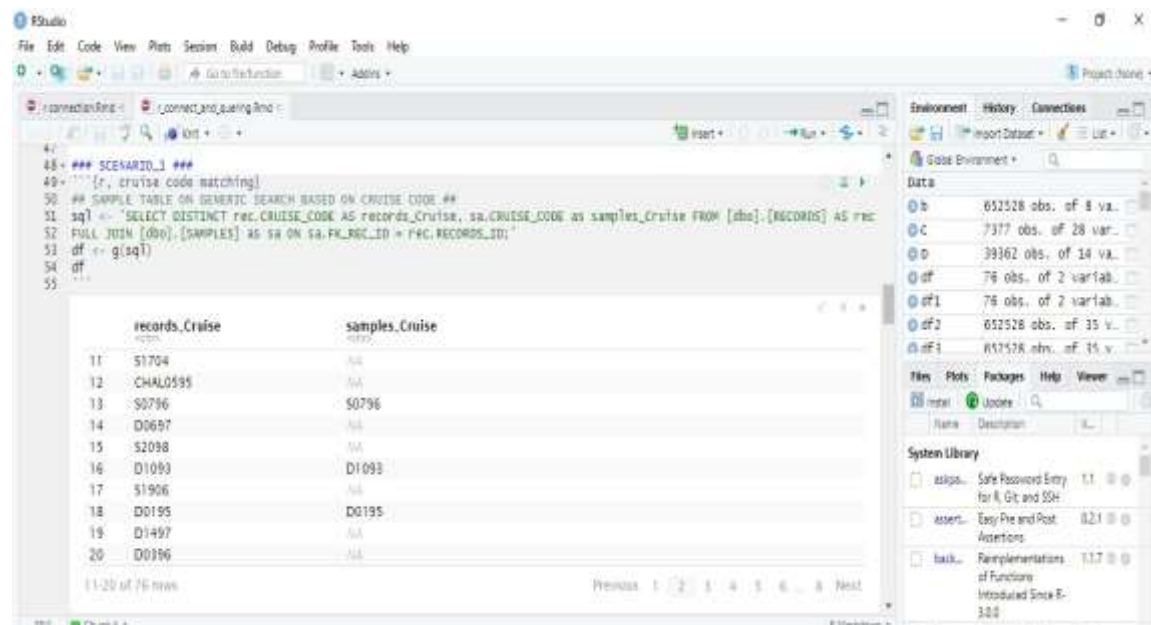


Figure 36 Available cruises on both tables

The researcher is then asked to provide in the dialogue box the cruise code that they are interested in. Here 'S1704'. Figure 37

search on cruise code

```

cruise <- readline(prompt = "ENTER CRUISE CODE IN ' ': ")

message1 <- paste("SELECT * FROM [dbo].[SAMPLES] WHERE CRUISE_CODE = ", cruise, ";")

b <- g(message1)

b

```

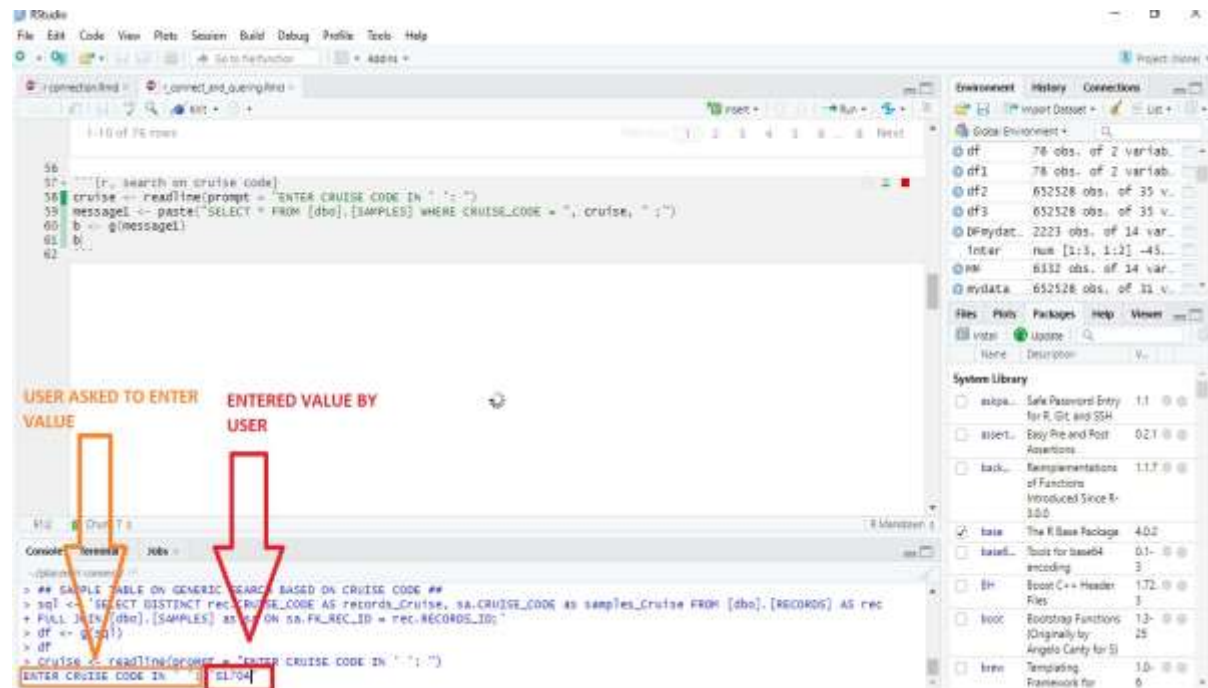


Figure 37 User asked to provide cruise code

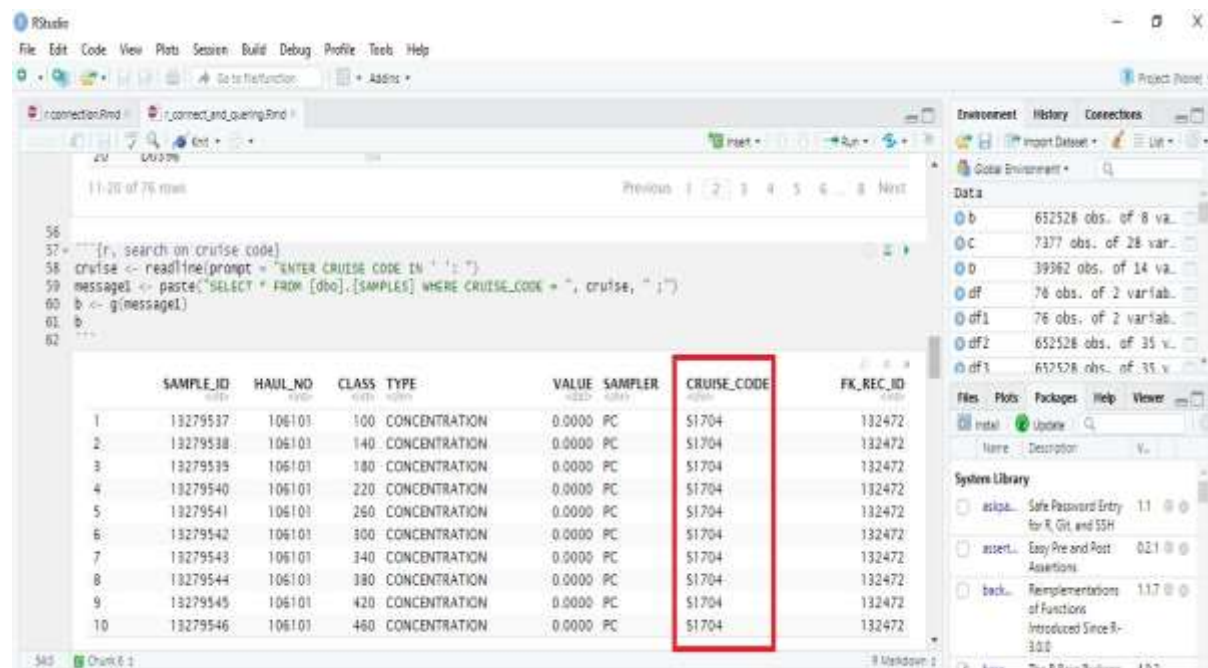


Figure 38 Retrieved data of interest

Researcher sees something interesting on retrieves data (Figure 38) decides to work on it. This implies combination with logging records. This translates into joining RECORDS, SAMPLES and Species_info tables. (Figure 39) The cruise code previously entered is automatically imported here.

combine tables of interest

```
message2 <- paste("SELECT rec.*, sa.[CLASS], sa.[TYPE], sa.[VALUE], sp.[Taxa_Names], sp.[AphiaID], ct.[Project],
ct.[BODC_link]

FROM [dbo].[CRUISES] AS ct JOIN [dbo].[RECORDS] AS rec ON ct.[UNIQUE_CRUISE_CODE] = rec.[CRUISE_CODE]

JOIN [dbo].[SAMPLES] AS sa ON rec.[RECORDS_ID] = sa.[FK_REC_ID]

FULL JOIN [dbo].[Species_info] AS sp ON sp.[TITLE] = sa.[CLASS] WHERE rec.CRUISE_CODE = ", cruise, ";")

df2 <- g(message2)

df2
```

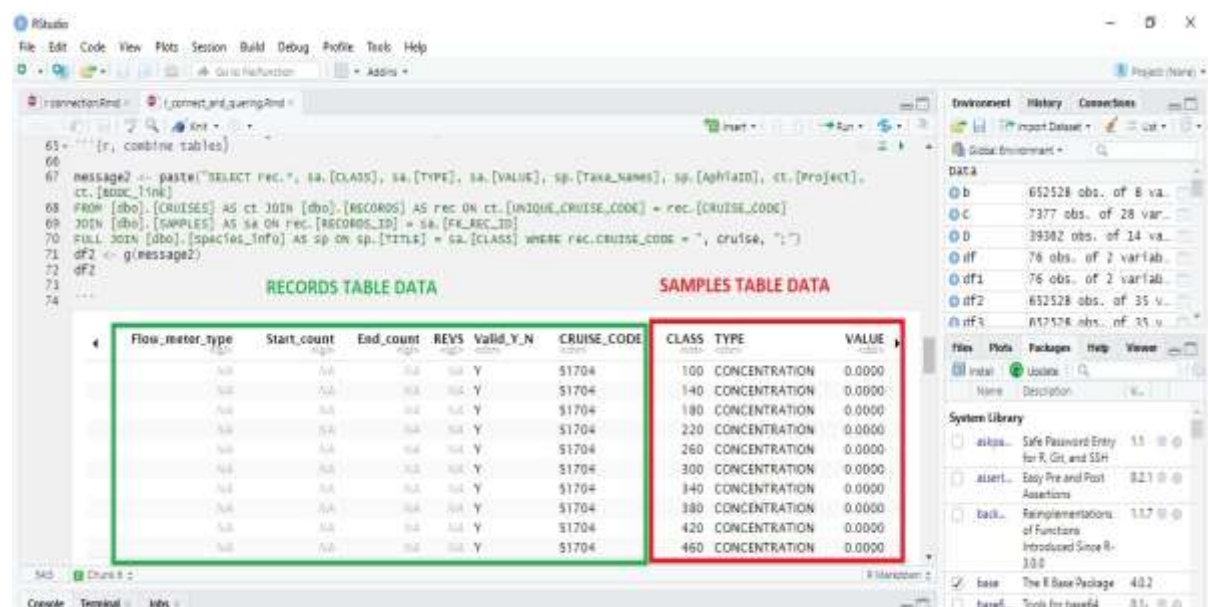


Figure 39 Combined tables

The researcher decides that adequate information exists for further research. The retrieved information above is kept intact and copied as a data frame in R for further data mining. At this point connection with the database is not used which will speed up the search and data mining process. In the code box below excluded columns are for the researcher to be chosen as well as the filtering out of 0s.(Figure 40)

the retrieved data is now and R data frame

```
df3 <- data.frame(df2)
```

exclude columns of no use

```
mydata = select(df3, -c(RECORDS_ID, Station_name, Taxa_Names, AphiaID))
```

```
mydata
```

filter out 0s

```
mydata2 <- filter(mydata, VALUE != 0.0)
```

mydata2

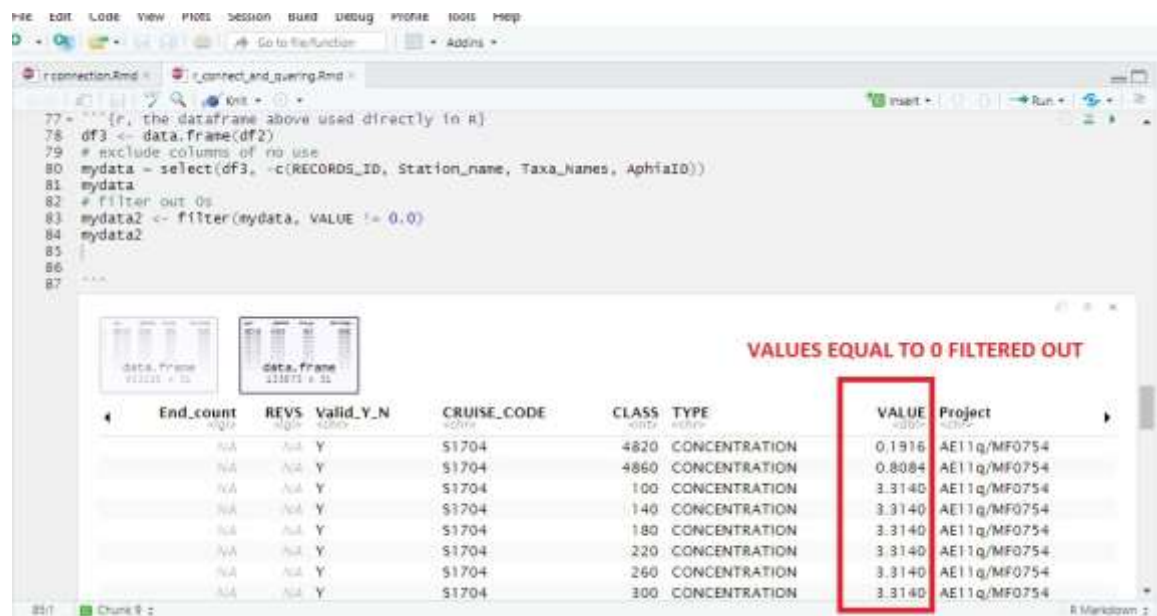


Figure 40 Filtered out data

The researcher needs to narrow down results according to some perimeter. As mentioned by Marine Scotland the echo attribute is of high importance. The script below is run and the researcher is asked to provide numbers for the echo attributes (Figure 41). As can be seen the operators ">" and "<" are fixed and this is correct given the nature of the data stored in them. The result comes up as a data frame (Figure 42)

work on columns based on interest

```
st_echo <- readline(prompt = "Enter start_echo quantity: ")
en_echo <- readline(prompt = "Enter end_echo quantity: ")
mydata3 <- filter(mydata2, Start_echo > as.numeric(st_echo))
mydata3 <- filter(mydata3, End_echo < as.numeric(en_echo))
mydata3
```



```

points(max(mydata3$End_longitude), max(mydata3$End_latitude, na.rm = TRUE), col = "blue", cex = 1.5)
points(mydata3$Start_longitude, mydata3$Start_latitude, col = "yellow", cex = .8)
points(mydata3$End_longitude, mydata3$End_latitude, col = "yellow", cex = .8)
# Add a legend
legend("topright", legend=c("Start lon/lat", "End long/lat", "In between points"),
      col=c("red", "blue", "yellow"), lty=1:2, cex=0.8)

```

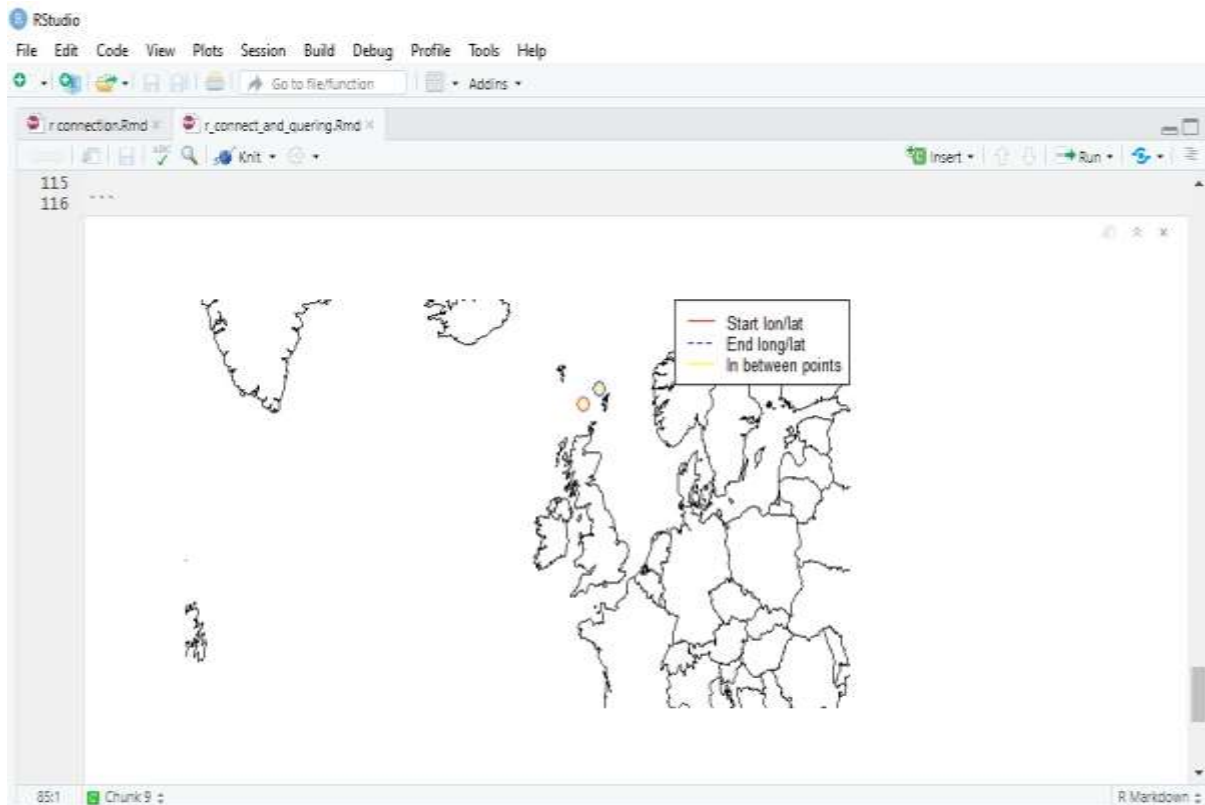


Figure 43 Values depicted on map

According to the interest of the researcher a step may be skipped or rerun if they wish.

SCENARIO_2

First the function that will handle the scripting that follows needs to be created.

function to run scripts

```

g <- function(sql){
  sqlQuery(dbhandle,sql)
}

```

The researcher first needs to verify existing records that are available in both tables according to taxa names (Figure 44)

```
## SAMPLE TABLE ON GENERIC SEARCH BASED ON CRUISE CODE ##

sql <- 'SELECT DISTINCT sp.[Taxa_Names], sa.Cruise_code AS S_CRUISES, rec.Cruise_code AS REC_CRUISES FROM
[dbo].[Species_info] AS sp
JOIN [dbo].[SAMPLES] as sa ON sa.CLASS = sp.[TITLE]
full JOIN [dbo].[RECORDS] AS rec ON sa.FK_REC_ID = rec.RECORDS_ID WHERE sp.[Taxa_Names] IS NOT NULL
ORDER BY sp.[Taxa_Names];'

df <- g(sql)

df
```

40 The user sees specific Taxa names and whether their data exists in the SAMPLES and RECORDS tables.
 41 (r, cruise code matching)
 42 ## SAMPLE TABLE ON GENERIC SEARCH BASED ON CRUISE CODE ##
 43 sql <- 'SELECT DISTINCT sp.[Taxa_Names], sa.Cruise_code AS S_CRUISES, rec.Cruise_code AS REC_CRUISES FROM
 44 [dbo].[Species_info] AS sp
 45 JOIN [dbo].[SAMPLES] as sa ON sa.CLASS = sp.[TITLE]
 46 full JOIN [dbo].[RECORDS] AS rec ON sa.FK_REC_ID = rec.RECORDS_ID WHERE sp.[Taxa_Names] IS NOT NULL
 47 ORDER BY sp.[Taxa_Names];'
 48 df <- g(sql)
 49 df

	Taxa_Names	S_CRUISES	REC_CRUISES
141	Crustacea	DISCO267	DISCO267
142	Crustacea	DISCO262	DISCO262
143	Crustacea	DISCO258	NA
144	Ctenocalanus vanus	D0697	D0697
145	Ctenophora	DISCO258	DISCO258
146	Ctenophora	DISCO264	DISCO264
147	Ctenophora	DISCO262	DISCO262
148	Ctenophora	DISCO267	DISCO267
149	Cumacea	D0495	D0495
150	Cyclopoida	DISCO258	NA

141-150 of 437 rows

Figure 44 Checking existing records for taxa names

The user decides to search for "Ctenophora" (Figure 46)

```
taxa_name <- readline(prompt = "ENTER TAXA NAME IN ' ': ")

message2 <- paste("SELECT rec.*, sa.[CLASS], sa.[TYPE], sa.[VALUE], sp.[Taxa_Names], sp.[AphiaID], ct.[Project],
ct.[BODC_link]
FROM [dbo].[CRUISES] AS ct JOIN [dbo].[RECORDS] AS rec ON ct.[UNIQUE_CRUISE_CODE] = rec.[CRUISE_CODE]
JOIN [dbo].[SAMPLES] AS sa ON rec.[RECORDS_ID] = sa.[FK_REC_ID]
FULL JOIN [dbo].[Species_info] AS sp ON sp.[TITLE] = sa.[CLASS] WHERE sp.[Taxa_Names] = ", taxa_name, ";")

df2 <- g(message2)

df2
```



Figure 45 User asked to enter value for attribute 'Taxa_name'

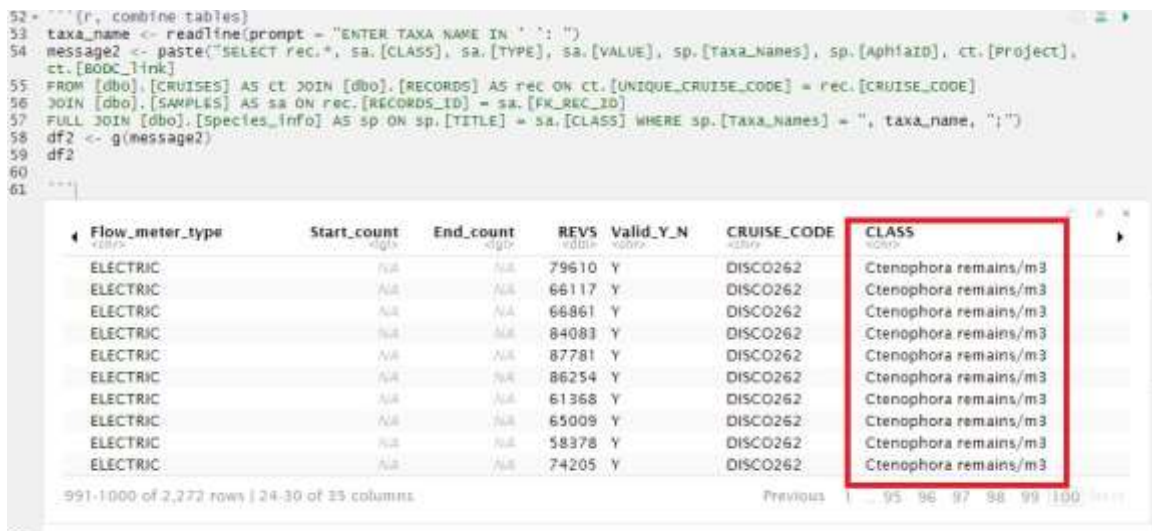


Figure 46 Result of search for 'Ctenophora'

The researcher decides that adequate information exists for further research. The retrieved information above is kept intact and copied as a data frame in R for further data mining (Figure 47).

```

df3 <- data.frame(df2)

# exclude columns of no use

mydata = select(df3, -c(RECORDS_ID, Start_count, End_count))

mydata

# filter out 0s

mydata2 <- filter(mydata, VALUE != 0.0)

mydata2

# filter according to Date

```

```
mydata3 <- tibble::as_tibble(mydata2)
mydata3 <- filter(mydata3, DATE >= '2002')
mydata3
```

```

70 mydata2 <- filter(mydata, VALUE != 0.0)
71 mydata2
72 # filter according to date
73 mydata3 <- tibble::as_tibble(mydata2)
74 mydata3 <- filter(mydata3, DATE >= '2002')
75 mydata3
76
77

```

DATE	Timestamp	HAUL_NO	Station_name	Start_latitude	Start_longitude	Mid_latitude
2002-05-08	02:22:00	438011	D5	60.92612	-36.84950	60.92477
2002-05-08	02:59:00	438018	D5	60.90391	-36.87566	60.90246
2002-05-09	02:27:00	451011	D3	61.31113	-38.24468	61.30958
2002-05-09	02:58:00	451016	D3	61.29323	-38.24807	61.29085
2002-08-23	20:36:00	1040024	D6	57.13510	-34.33061	57.13394

5 rows | 1-7 of 32 columns

Figure 47 Filter out unwanted columns, 0s and filter for date

The researcher decides to see on the map the data of interest (Figure 48).

```
newmap <- getMap(resolution = "low")
plot(newmap, xlim = c(-35, 10), ylim = c(45, 65), asp = 2)
points(min(mydata3$Start_longitude, na.rm = TRUE), min(mydata3$Start_latitude, na.rm = TRUE), col = "red",
cex = 1.5)
points(max(mydata3$End_longitude), max(mydata3$End_latitude, na.rm = TRUE), col = "blue", cex = 1.5)
points(mydata3$Start_longitude, mydata3$Start_latitude, col = "brown", cex = .8)
points(mydata3$End_longitude, mydata3$End_latitude, col = "brown", cex = .8)
# Add a legend
legend("topright", legend=c("Start lon/lat", "End long/lat", "In between points"),
col=c("red", "blue", "brown"), lty=1:2, cex=0.8)
```

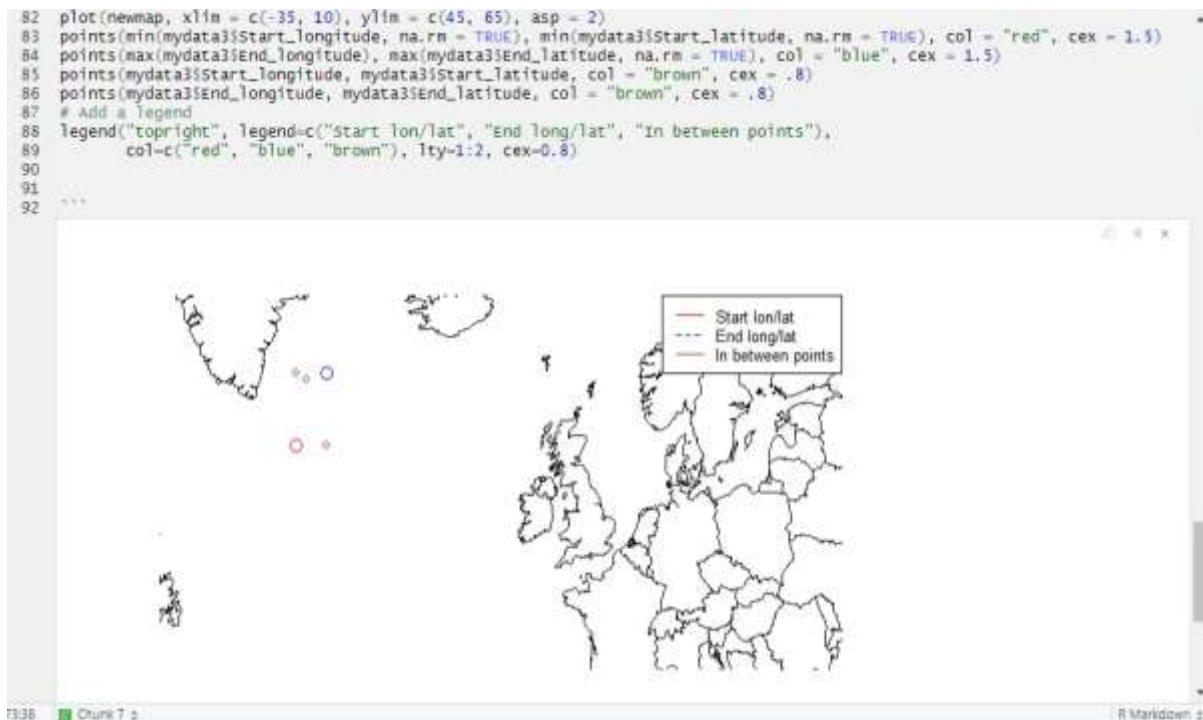


Figure 48 Look data of interest on map

Appendix-Power BI

First connectivity between the database and Power BI need to be established. We open Power BI from desktop. At Home tab we choose Get Data → SQL server (Figure 49) → in the pop-up window we enter the server name

and database name (Figure 50)→ At the new pop-up window we tick the tables we want and click Load (Figure 51)→ once the tables are successfully loaded we see them all on the right of our panel (Figure 52). From this point, one proceeds with the filtering and visualizations they prefer.



Figure 49 Import from SQL server

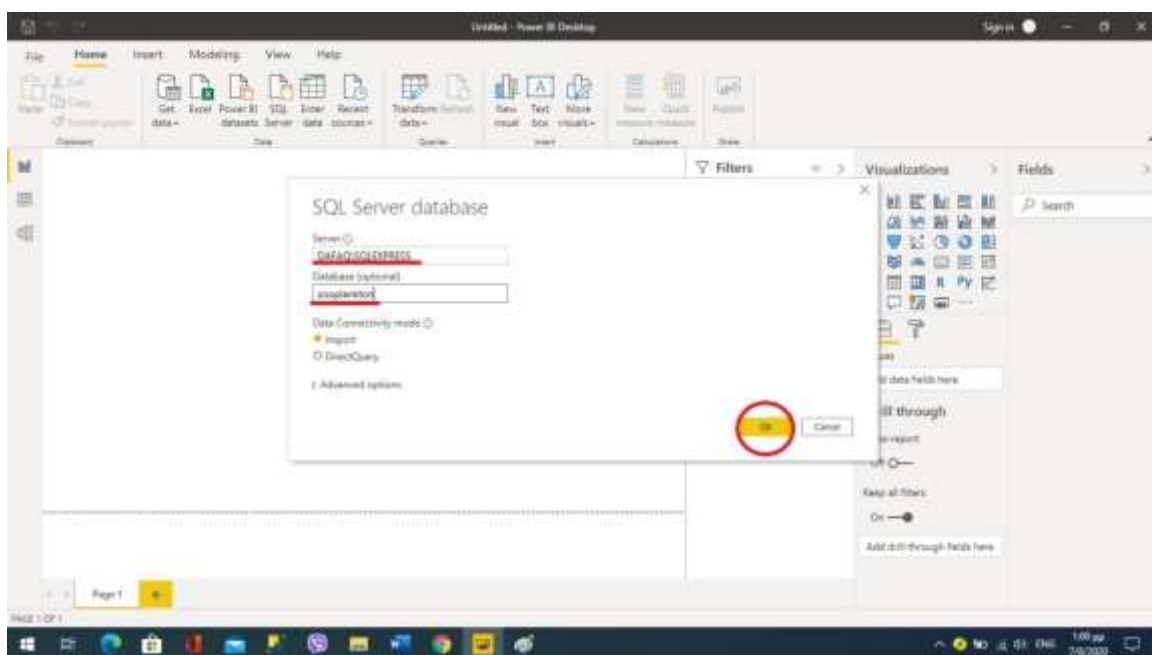


Figure 50 Insert credentials

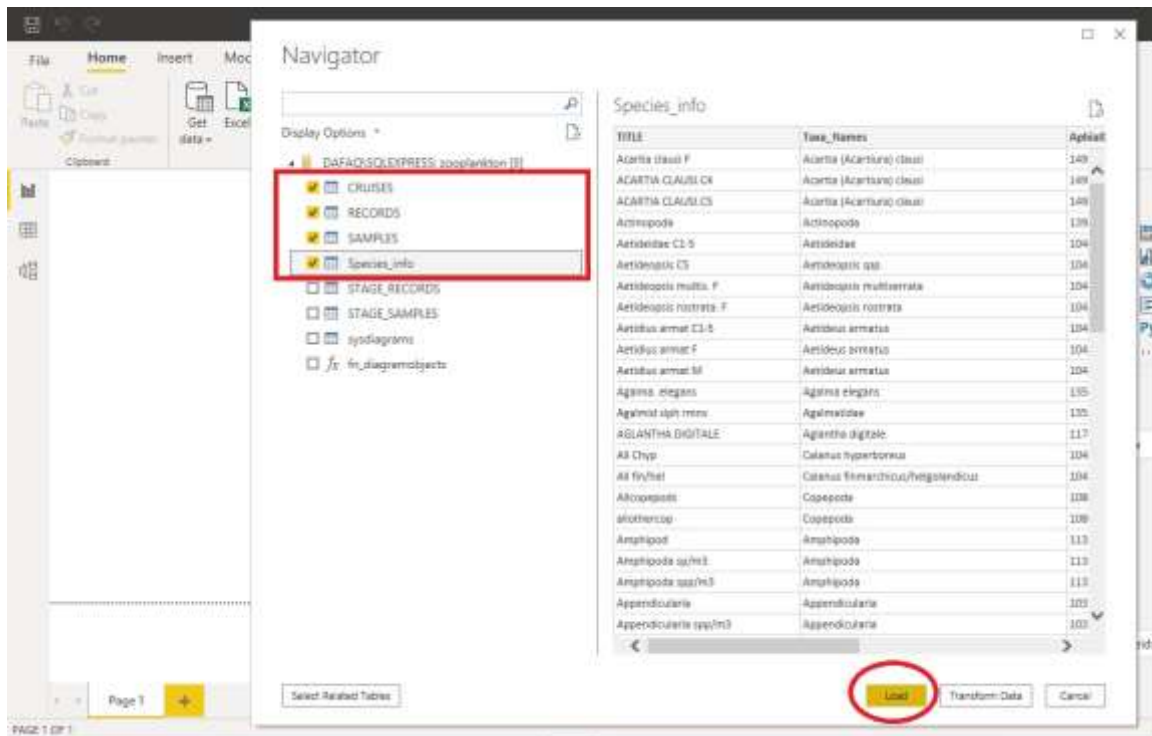


Figure 51 Load tables

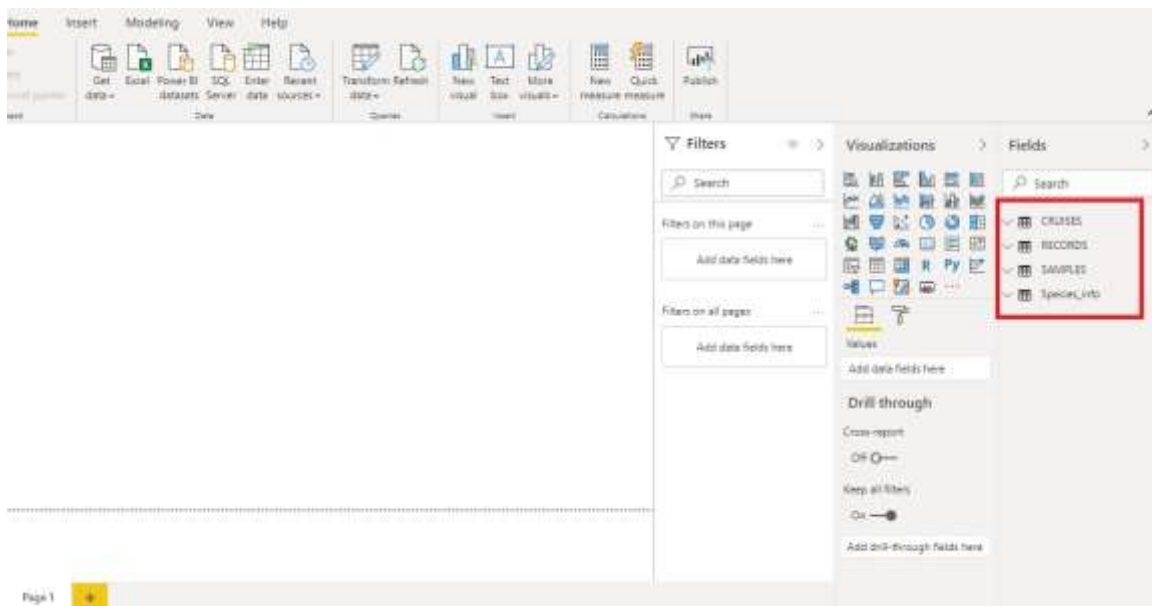


Figure 52 Successful import of tables

When we try to combine the data of two tables and for this example put stigmas on a an map we get the following response from Power BI (Figure 53). We click on “Fix this”.

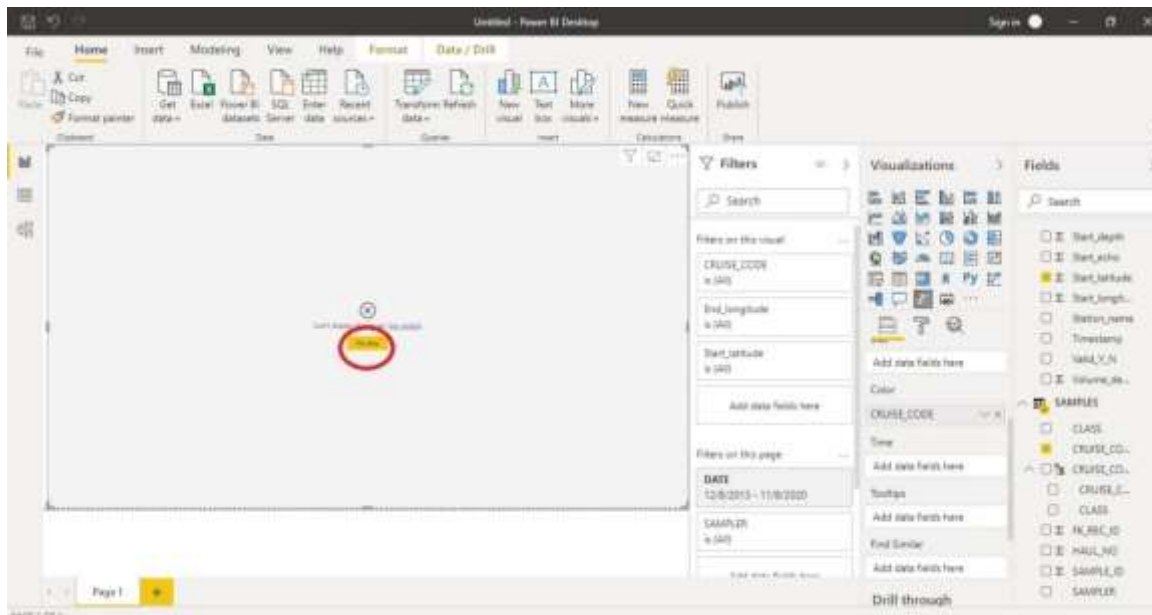


Figure 53 Map tables

We then click on “Autodetect” (Figure 54). Yet Power BI is incapable of auto-detecting relationships so we choose “Manage relationships” (Figure 55).

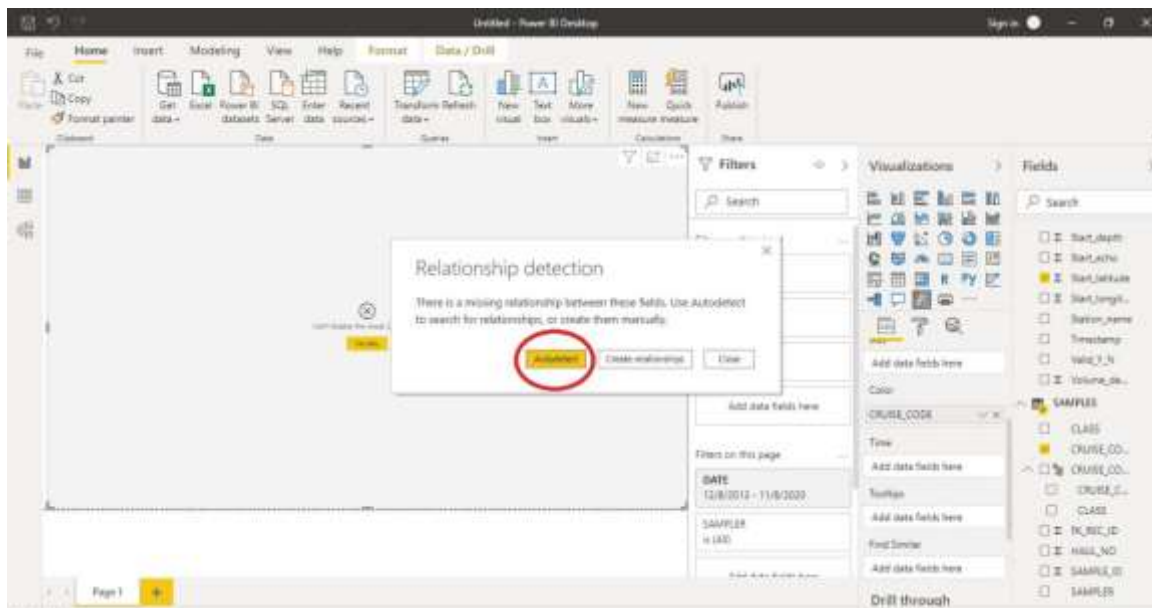


Figure 54 Relationship auto-detection

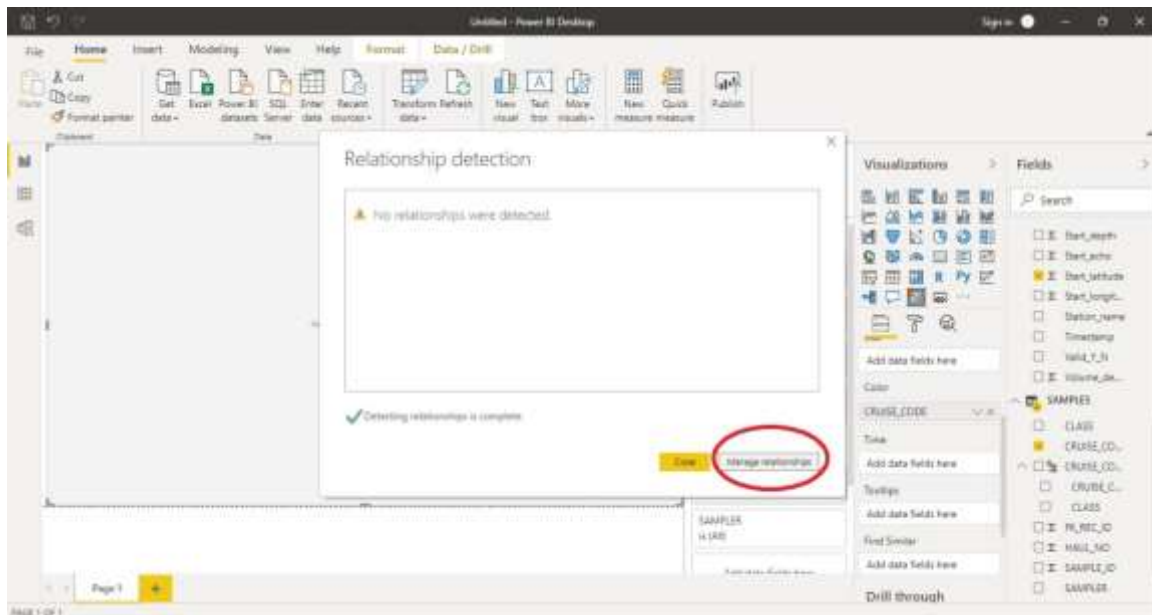


Figure 55 Manage relationship option

In the “Manage relationships” pop-up window we choose the relationship that is ticked, which is the one we have created in SSMS for this database, the RECORDS_ID to FK_REC_ID relationship (Figure 56).

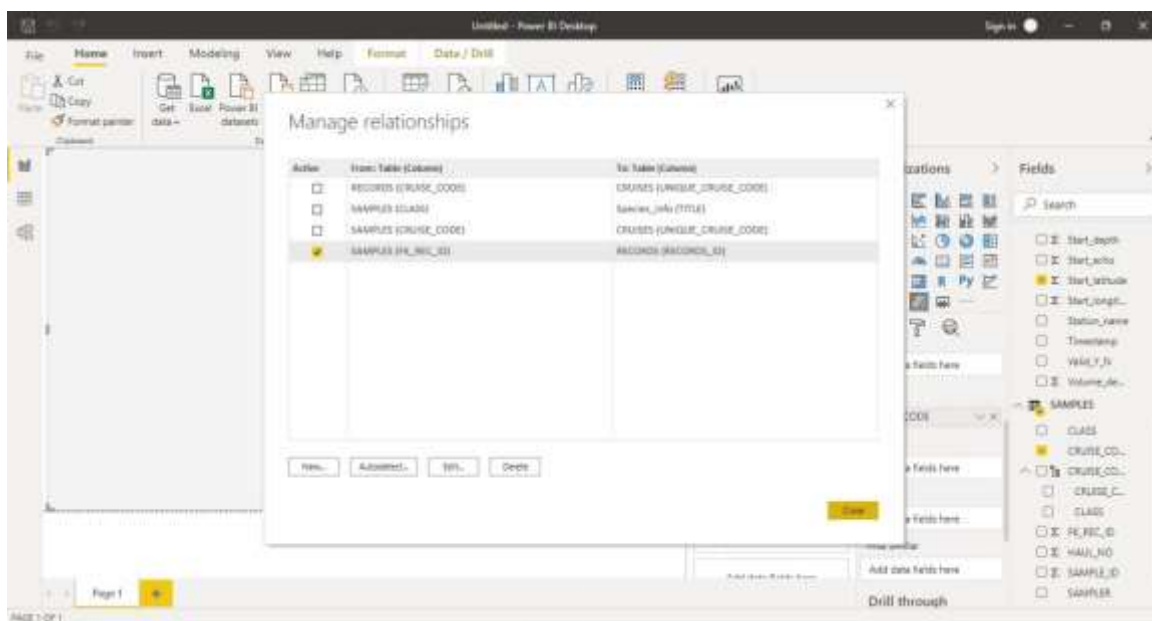


Figure 56 Choosing the correct relationship

Then we filter according to Sampler and Date and depict Cruise_code on the map (Figure 57).

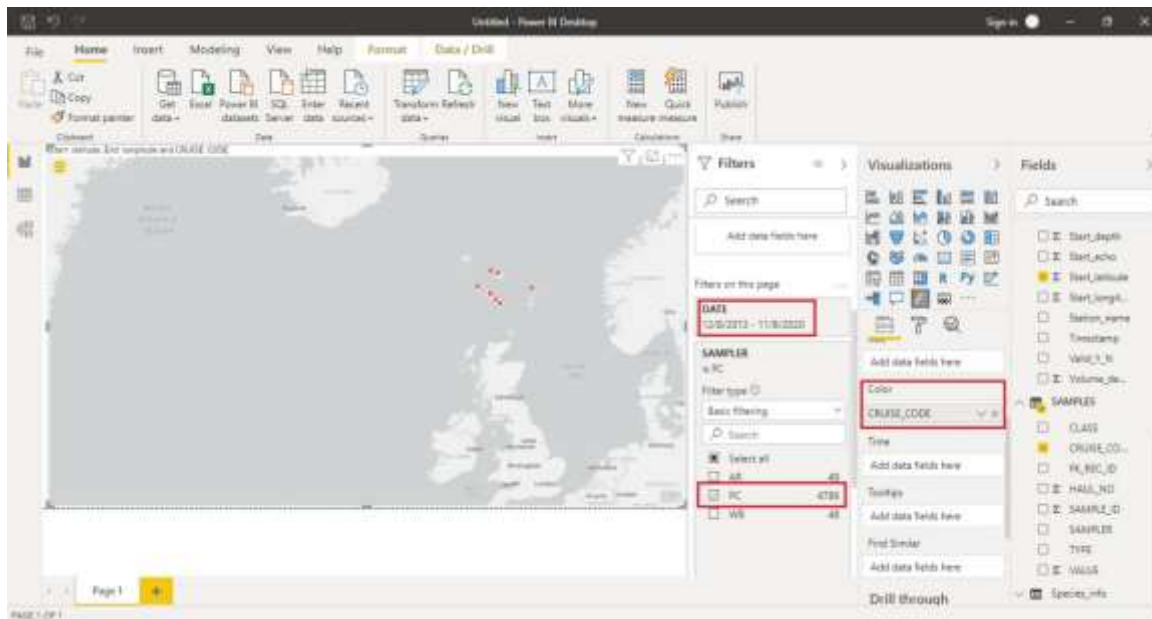


Figure 57 Multiple filtering

This is an example of how Power BI can easily provide initial information that can lead to insights for further data mining and analysis in RStudio.