# Documentation for Cleaning of Dataset Using Excel- Assignment No. 1
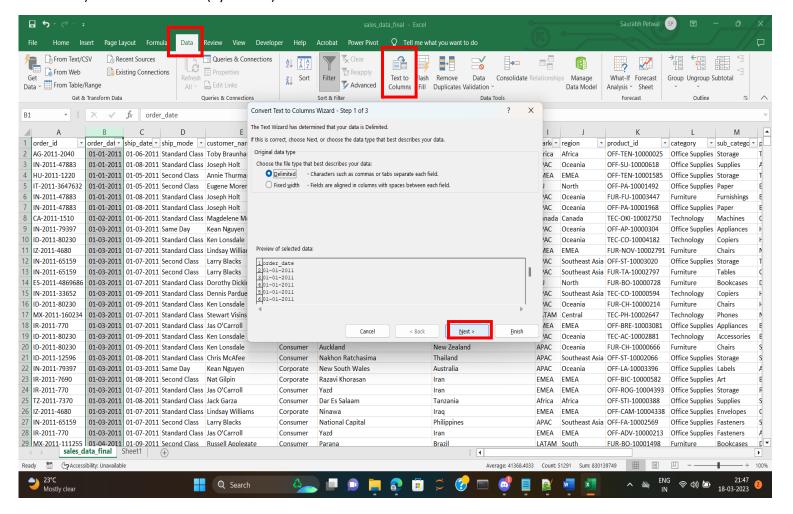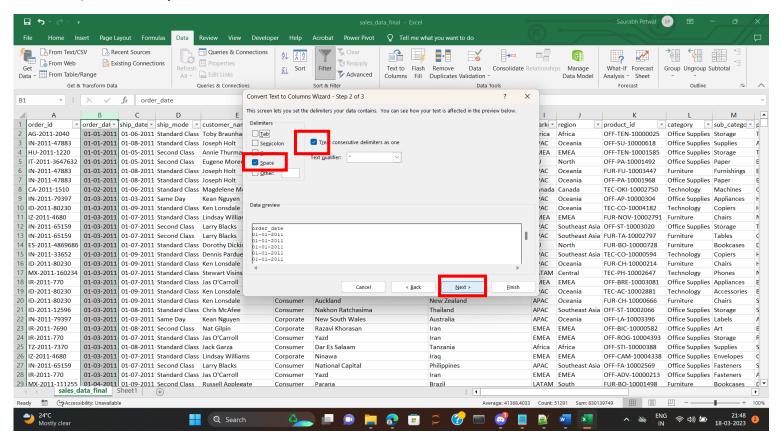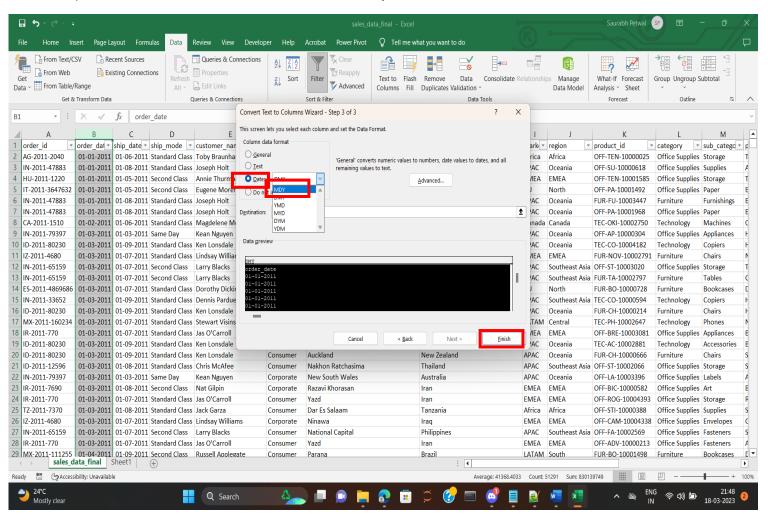
1. "order_date" and "ship_Date" columns in dataset are not in proper format. Some dates are in general format while rest are in date format. So changed the "order_date" and "ship_date" column format from "General" format to "date" format via the procedure shown below:

1.a) Select the **order_date** column (whole column)

1.b) Go to **Data** tab and then select **Text to columns**

1.c) Select • Delimited (by default) and click on **Next** button
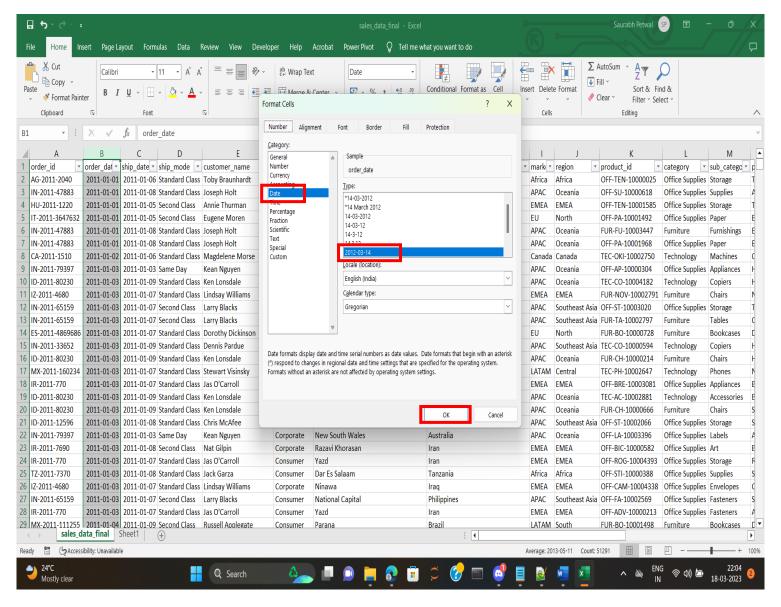
**1.d)** Check on √ **Space** and click on **Next** button

1.e) Check on √ **Date** and select **DMY** format from drop down and then click on Finish button
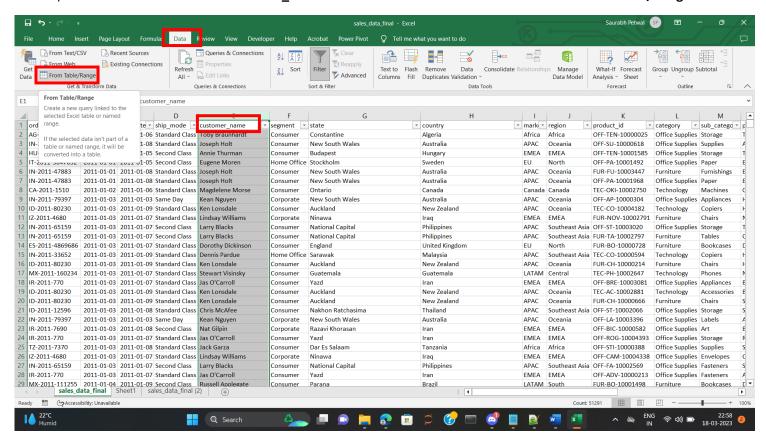
1.f) Now, Go to Format type (**ctrl +1** : shortcut key) and select date format as YYYY-MM-DD, as snow flake takes date in the same format
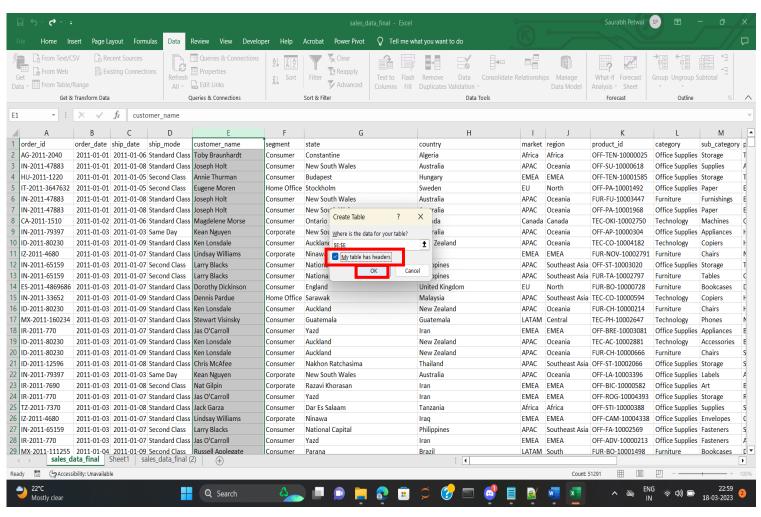


**Note: This is the effective way to turn the format to data type using excel. Similarly follow same steps for column "ship_date".**
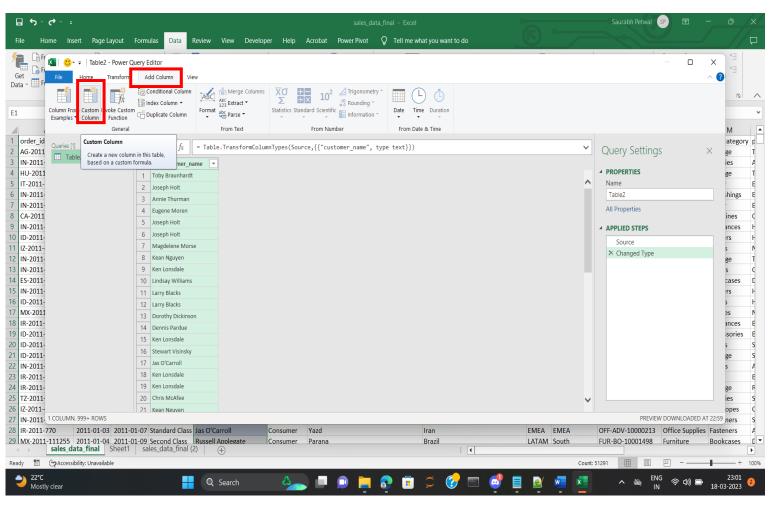
2. Using Custom filter checking data in columns "customer_name", "product_name", "state" as these columns have special characters so removed the special characters and UTF 8 characters step wise step. Using below procedure we can clean the data in alpha numeric datatype but since it also contain UTF8 characters and some latin and german characters, so to keep in mind to not to loose data we will manually remove and replace them. We will create a Flag for rows having special characters and can remove them manually using replace function ( **ctrl +f** >> then click Replace tab).

   2.a) Select the whole column **"customer_name"** and then click on **"Data"** tab. Now select From **Table/Range**

2.b) Check on popup **My table has headers** and then click on **OK**

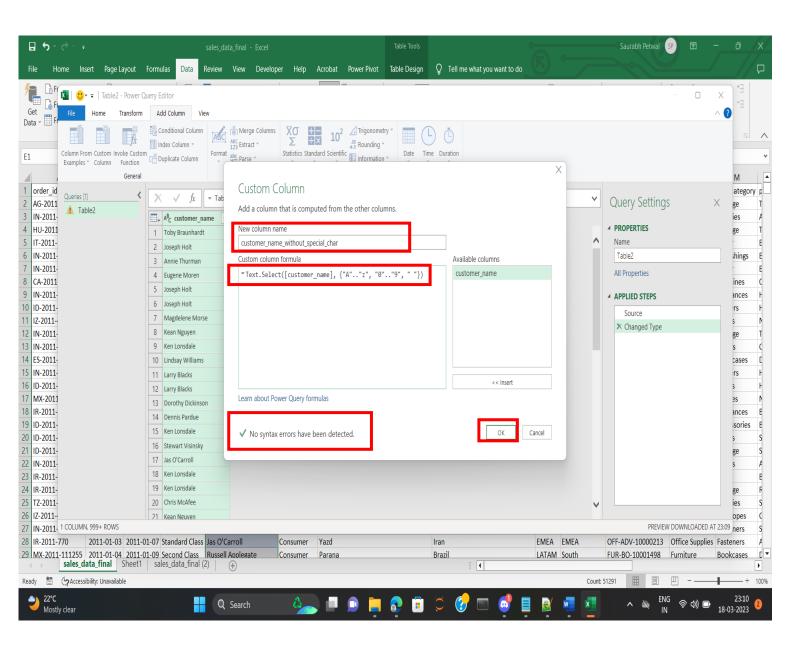2.c) Click on **"Add column"** tab then click on **"Custom Column"**

2.d) Now type a new column name, then write a syntax to clean the existing column and the values will be saved to new_column, syntax used to convert the data to alphanumeric datatype is as:
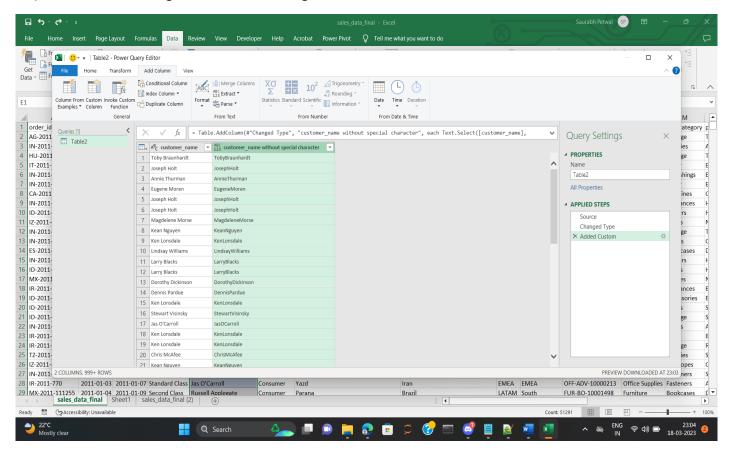
> **= Text.Select([customer_name], {"A".."z", "0".."9", " "})**

Here,

- "customer_name" is the existing column in the table to which we are cleaning and creating a new column to save values in it, it is named below as "customer_name_with_special_char"
- "A".."z" in the syntax is used for all values of "A-z and a-z " since ASCII char are in sequence A,B,C.....Y,Z,a,b...y,z (capital then small case) so it gives range for capital as well as small case alphabets
- "0".."9" in the syntax used for all numeric values from 0 to 9
- " " (white space) is used for white space between two words as we have data with space separated as first and last name of the customer in customer_name column. If we don't use this " " in our sytntax then forst and last name will not having space between them in the new column generated by cleaning of data i.e "customer_name_with_special_char"

2.e) New column will be generated from existing column with cleaned data and can be used further.



2.f) Now, the new generated column can be used in place of original column in original dataset table.
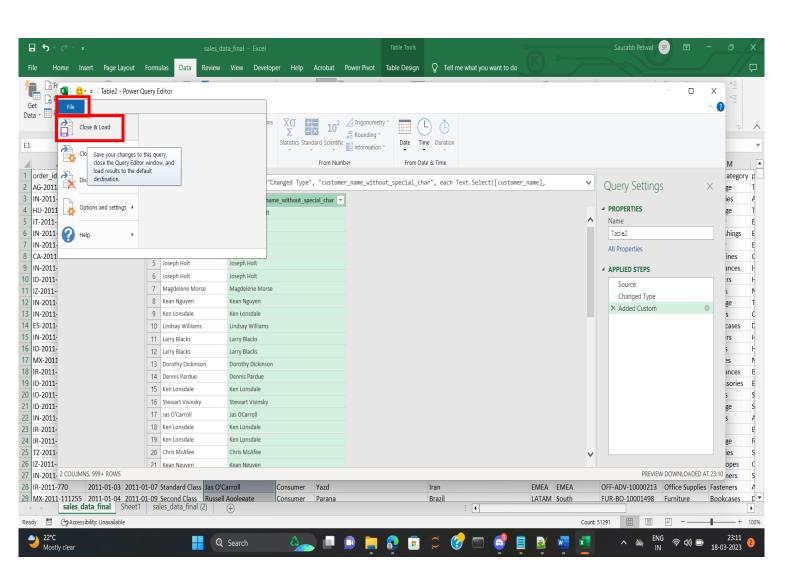
Steps for copying the data:

a. Click on file and select close and load.
b. A separate sheet will be created on the existing workbook and from there the new cleaned column can be used in place of original uncleaned column. But in our case there are some german, latin etc characters were there so we haven't replaced new column with old existing column "customer_name" in order to save the information from getting lost, but used the new column for identifying the rows having special characters and german characters from existing column by creating a "Flag" column and comparing both existing and new column contents using "if" formula in excel as shown below:

=if (cell of 'customer_name_with_special_character' = cell of 'customer_name', 1, 0)

Then copying the formula to all rows by dragging the bottom corner of cell where the above formula is applied.
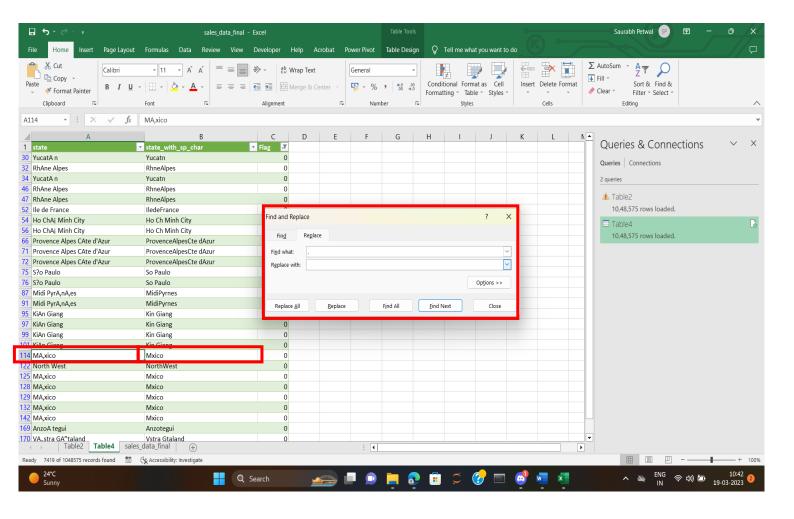
**Note: The above procedure is applied to three columns to clean the data and are as follows:**

1. **state**
2. **customer_name**
3. **Product_name**

2.g) As from the below example we can clearly see that by using the above procedure to clean the data there might be chance of loosing information from data (characters similar to alphabets except special characters – **"MÂ,xico"** >> **"Mxico"**. So from help of flag column where the values are 0 (zero) we can check the particular columns and can refill the alphabets in the data where such Latin, German or etc words was removed.

Also UTF 8 characters were in form of space was present in the data so it was removed by replacing with "[null]" manually one by one using filter and replace function.



------------------------------------------------------------------**END**------------------------------------------------------------------