# Exercise- Data Cleansing Exercise

**Student Name:** Animesh Johri                    **Student Id:** 2021292753

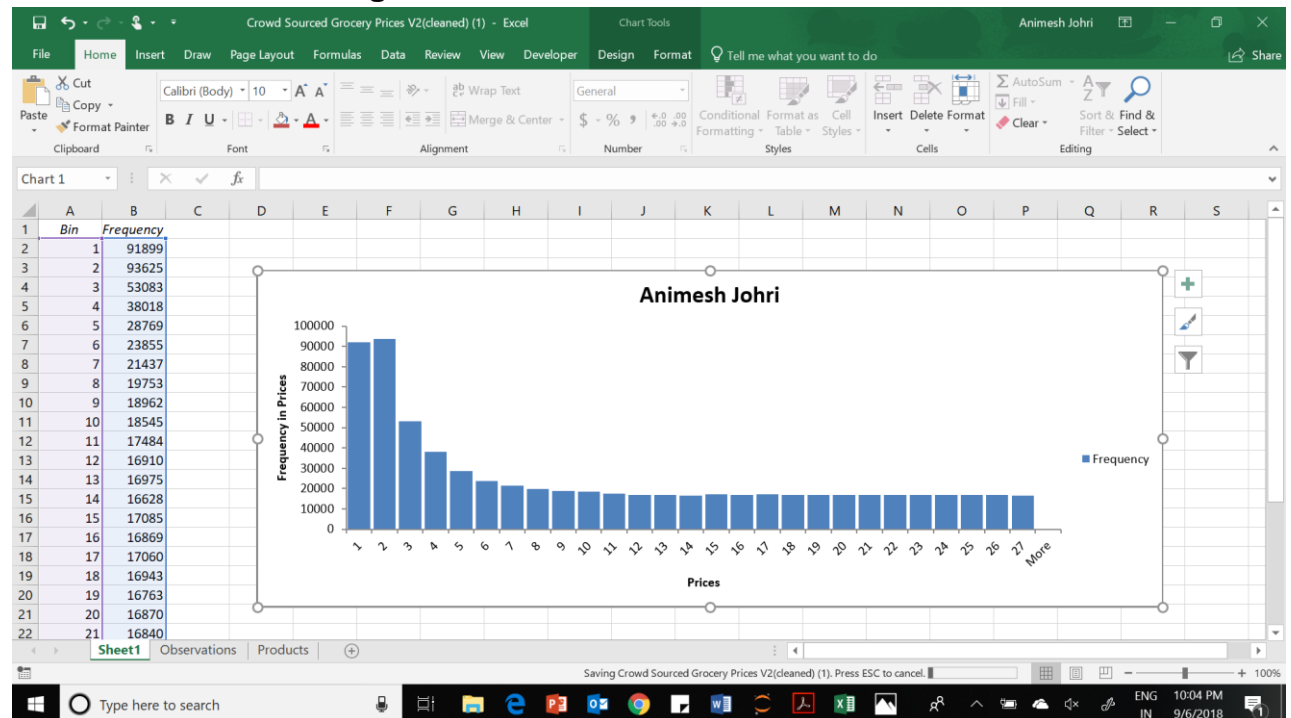**Date:** 09/07/2018

**a)** The value in cell D439899 is 35

**Screenshot:**



**b)** **The methodology chosen for imputation:** I made a histogram based on the price column, the histogram was skewed, therefore I chose to impute the median of the price column in the blank cells. To impute -> I selected the whole price column using "ctrl + shift + down arrow" -> I clicked "ctrl + H" to find and replace -> I found all the blank cells within the price column and replaced it with the median value of 8.91 (To find median, I selected all the values in price column and used the formula "=median" over the price column)

**Screenshot of the histogram:**



**c) A screenshot of the table with the imputed values:**

## d) A screenshot of the table after initial cleaning:



| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Country | City | Obs Date | Product Code | Quantity | Price | | BINS | | | | | | | |
| 2 | Kenya | Eldoret | 2009-01-01 | 52 | 1 kg | 11.92 | | 1 | | | Number of | 0 | | | |
| 3 | Pakistan | Attock | 2009-01-01 | 44 | 1 kg | 9.36 | | 2 | | | | | | | |
| 4 | India | Nashik | 2009-01-01 | 47 | 1 kg | 6.47 | | 3 | | | | | Here the missing blanks are | | |
| 5 | Kenya | Eldoret | 2009-01-01 | 52 | 1 kg | 5.91 | | 4 | | | Percentage | 3.122384 | more than 2%(3.122384 ) | | |
| 6 | India | Nashik | 2009-01-01 | 62 | 100 g | 1.80 | | 5 | | 0.376005 | | | for the price, therefore we | | |
| 7 | Pakistan | Attock | 2009-01-02 | 46 | 1 kg | 16.97 | | 6 | | | | | need to see the histogram | | |
| 8 | Kenya | Eldoret | 2009-01-02 | 52 | 1 kg | 16.50 | | 7 | | | | | of the data. If histogram is | | |
| 9 | India | Nashik | 2009-01-02 | 47 | 1 kg | 10.43 | | 8 | | | | | skewed then median is | | |
| 10 | Kenya | Eldoret | 2009-01-02 | 52 | 1 kg | 9.77 | | 9 | | | | | imputed if histogram is | | |
| 11 | Pakistan | Attock | 2009-01-02 | 47 | 1 kg | 1.61 | | 10 | | | | | normal distribution then | | |
| 12 | Pakistan | Rawalpindi | 2009-01-03 | 52 | 1 kg | 17.80 | | 11 | | | | | mean is imputed in the | | |
| 13 | Pakistan | Rawalpindi | 2009-01-03 | 44 | 1 kg | 8.77 | | 12 | | | Median | 8.91 | missing blanks. | | |
| 14 | Pakistan | Attock | 2009-01-04 | 47 | 1 kg | 26.22 | | 13 | | | | | | | |
| 15 | Pakistan | Attock | 2009-01-04 | 46 | 1 kg | 20.07 | | 14 | | | | | | | |
| 16 | Pakistan | Islamabad | 2009-01-04 | 45 | 1 kg | 19.83 | | 15 | | | | | | | |
| 17 | Pakistan | Lahore | 2009-01-04 | 44 | 1 kg | 19.35 | | 16 | | | | | | | |
| 18 | Pakistan | Islamabad | 2009-01-04 | 47 | 1 kg | 15.86 | | 17 | | | | | | | |
| 19 | Pakistan | Islamabad | 2009-01-04 | 46 | 1 kg | 9.42 | | 18 | | | | | | | |
| 20 | Pakistan | Karachi | 2009-01-04 | 46 | 1 kg | 8.53 | | 19 | | | | | | | |
| 21 | India | Jabalpur | 2009-01-04 | 49 | 500 g | 0.99 | | 20 | | | | | | | |
| 22 | India | Indore | 2009-01-05 | 45 | 1 kg | 26.13 | | 21 | | | | | | | |

## e) The total number of rows (including header) present in the data sheet are 961733

## f) Countries with spell check error:

1. Brazil (wrong spelling: Brazel)
2. Indonesia (wrong spelling: Indonseia)
3. Philippines (wrong spelling: Phillippines).

**The countries with duplicate entries are:**
1. India
2. Bangladesh
3. Kenya
4. Pakistan

### g. Screenshot showing the changes made:



### h. Product Name is "Carrots" and the Quantity of Product Code 44 in Islamabad is "1kg".
### Screenshot: