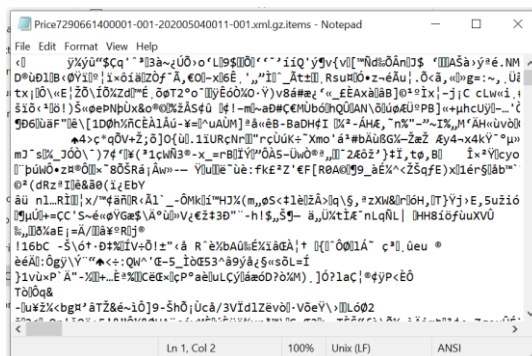


Ex2 – Submission practical part

סעיף א': איתור בעיות בקבצי הנתונים.

1. **קבצים מושחתים** (Not approachable):
 - תיאור הבעיה: קבצים מסוימים מוגדרים מושחתים ועל כן אין אפשרות לעבוד איתם.
 - דוגמות: בפרט הקבצים בתיקיית "freshmarket", ובסה"כ כ-20 קבצים שאין אפשרות לפתוח אותם (UNZIP).
2. **קבצים לא קריאים** (Bad Encoding):
 - תיאור הבעיה: ישנם קבצים שניתן לפתוח אותם (UNZIP), אבל הקידוד שלהן לא בוצע כהלכה והם אינם ניתנים לעיבוד.
 - דוגמות: צפינו ב-1415 מקרים בתיקייה hashook\2020-05-04, כאשר לשמות הקבצים נוספו .item.

דוגמא לקובץ שכזה –



3. **פער בקבצי נתונים** (Missing Data files):
 - תיאור הבעיה: בהתאם למבנה העקרוני הנדרש של הנתונים, בחברות רבות אין את כל הקבצים שאמורים להיות, דבר שיפגע ביכולת ההשוואה בין רשתות ובניתוח הנתונים בנוגע לכל רשת, בגלל היעדר בסיס להשוואה.
 - דוגמות:
 - א. קיום קובץ "חנויות" בלבד (ללא פירוט מחירים ומבצעים כנדרש), ברשתות: "קשת", "אורש עד", "רמי לוי", טיב טעם ועל כן לא ניתן לעשות השוואת מחירים ברשתות אלו
 - ב. קיום רק פלט קבצי "מבצעים" (ללא פירוט מחירים ורשימת חנויות כנדרש) ברשת "מגה".

4. בעיית נתוני זמן (Temporality):

- תיאור הבעיה: טווחי תאריכים בשדות עצמם (מועד עדכון מחיר המוצרים) לא חופפים, ולא מתואמים עם נתוני תיארוך נוספים (בשם הקובץ). הדבר מונע בחינה ברורה בהתאם לתאריך, או לכל הפחות ידרוש הסדרה של השימוש בתאריכים לפי אחד מהשדות המדוברים ובמודל אחיד (מועד הפקת הדו"ח – בשם הקובץ, מועד העדכון של הפרטים). זאת, כדי לאפשר השוואת מחירים של מוצרים בין חנויות שונות (באותה תקופה).
- דוגמות: לשני מוצרים שונים תחת אותו הקובץ יש שנת עדכון מחיר שונה.

| Price7290661400001-001-202005040407-001.xml | Price7290661400001-001-202005040407-001.xml |
|-------------------------------------------------------|--------------------------------------------------------|
| 1 <?xml version="1.0" encoding="utf-8"?> | 28 <Product> |
| 2 <Prices> | 29 <PriceUpdateDate>2020/01/14 11:01</PriceUpdateDate> |
| 3 <ChainID>7290661400001</ChainID> | 30 <ItemCode>42209</ItemCode> |
| 4 <SubChainID>001</SubChainID> | 31 <ItemType>0</ItemType> |
| 5 <StoreID>001</StoreID> | 32 <ItemName>בשר טחון סר' חנה חובר</ItemName> |
| 6 <SikoretNo>000</SikoretNo> | 33 <ManufactureName>טמ' יורסה בטמ'</ManufactureName> |
| 7 <Product> | 34 <ManufactureCountry /> |
| 8 <Product> | 35 <ManufactureItemDescription /> |
| 9 <PriceUpdateDate>2016/11/14 07:18</PriceUpdateDate> | 36 <UnitQty>'p'</UnitQty> |
| 10 <ItemCode>403</ItemCode> | 37 <Quantity>1</Quantity> |
| 11 <ItemType>0</ItemType> | 38 <UnitMeasure /> |
| 12 <ItemName>דג מ'לה טחון סר'</ItemName> | 39 <Bisweighted>1</Bisweighted> |
| 13 <ManufactureName>טמ' יורסה בטמ'</ManufactureName> | 40 <QtyInPackage>1</QtyInPackage> |
| 14 <ManufactureCountry /> | 41 <ItemPrice>59.90</ItemPrice> |
| 15 <ManufactureItemDescription /> | 42 <UnitOfMeasurePrice>59.90</UnitOfMeasurePrice> |
| 16 <UnitQty>'p'</UnitQty> | 43 <AllowDiscount>1</AllowDiscount> |
| 17 <Quantity>1</Quantity> | 44 <ItemStatus>1</ItemStatus> |
| 18 <UnitMeasure /> | 45 <LastUpdateDate /> |
| 19 <Bisweighted>1</Bisweighted> | 46 <LastUpdateTime /> |
| 20 <QtyInPackage>1</QtyInPackage> | 47 </Product> |
| 21 <ItemPrice>99.90</ItemPrice> | 48 <Product> |
| 22 <UnitOfMeasurePrice>99.90</UnitOfMeasurePrice> | 49 <PriceUpdateDate>2019/06/25 07:27</PriceUpdateDate> |
| 23 <AllowDiscount>1</AllowDiscount> | 50 <ItemCode>7290000230029</ItemCode> |
| 24 <ItemStatus>1</ItemStatus> | 51 <ItemType>1</ItemType> |
| 25 <LastUpdateDate /> | 52 <ItemName>ק'1 טעונים דרמ'</ItemName> |
| 26 <LastUpdateTime /> | 53 <ManufactureName>בשר יורסה בטמ'</ManufactureName> |

5. ריבוי ערכים חסרים בקובץ (Missing Data):

- תיאור הבעיה: על פי המבנה הנדרש של הקבצים, בכל רשומה אמורים להיות נתונים מסוימים. היעדרו בולטת של נתונים בקבצים מסוימים פוגעת משמעותית ביכולת הניתוח שלו.
- דוגמות: בקבצים עם הסיומת FULL נמצאו ערכים חסרים. בריצה על 20 קבצים עם הסיומת FULL נתקלנו ב-9 קבצים עם ערכים חסרים, כלומר 45% מהקבצים שנדגמו לא היו בפורמט המצופה. בבדיקה ידנית שעשינו (פתיחת הקבצים באקסל) אכן קיימות עמודות או ערכים ריקים, כפי שהקוד התריע (כלומר, לא מדובר בבעיית פרסור). להלן פלט הריצה:

```
C:\Users\ASUS\PycharmProjects\HW2\converted\PriceFull7290696200003-086-202005040327-001.csv
there are missing values in this file True
```

```
C:\Users\ASUS\PycharmProjects\HW2\converted\PriceFull7290696200003-088-202005040300-001.csv
there are missing values in this file True
```

```
C:\Users\ASUS\PycharmProjects\HW2\converted\PriceFull7290696200003-089-202005040331-001.csv
there are missing values in this file True
```

6. סוג ערך לא תקין (Datatype invalid):

- תיאור הבעיה: על פי המבנה הנדרש של הקבצים, הערך של עמודת UnitQty בקבצי המחירים חייבים להיות אותיות בלבד. פער במאפיין זה פוגע הן בתהליך עיבוד המידע וכן מעיד על שימוש לא תקין בעמודה זו (על חשבון נתונים שאמורים להימצא בעמודת Quantity).
- דוגמות: פלט של קבצים שנמצאה בהם טעות:

```
in file: PriceFull7290696200003-095-202005040324-001.csv
found: non alphabetic chars in UnitQty
['1"ק', '1']
in file: PriceFull7290696200003-095-202005040324-001.csv
found: non alphabetic chars in UnitQty
['1"ק', '1']
in file: PriceFull7290696200003-095-202005040324-001.csv
found: non alphabetic chars in UnitQty
['1"ק', '1']
```

7. יחס בין משקלים למחירים (Uncorrelated corresponding data fields):

- תיאור הבעיה: ישנן שתי עמודות שמצופה בהן קוהרנטיות גבוהה: ItemPrice, UnitOfMeasurePrice (עד כדי סטייה הנובעת ממגבלת תמחור לפי אגורות, כלומר מיחידת המידה הקטנה ביותר). בשני אופנים:

1. אם יחידות המדידה (UnitMeasure) הן ביטוי זהה עד כדי פרמטר של יחידת המחיר הכללי (UnitQty), כמו במקרה לדוגמה של ק"ג-100 גרם, הקשר בין ItemPrice, UnitOfMeasurePrice אמור לקיים:

$$\text{ItemPrice} = \text{UnitOfMeasurePrice} * \text{UnitMeasure} * \text{quantity}$$

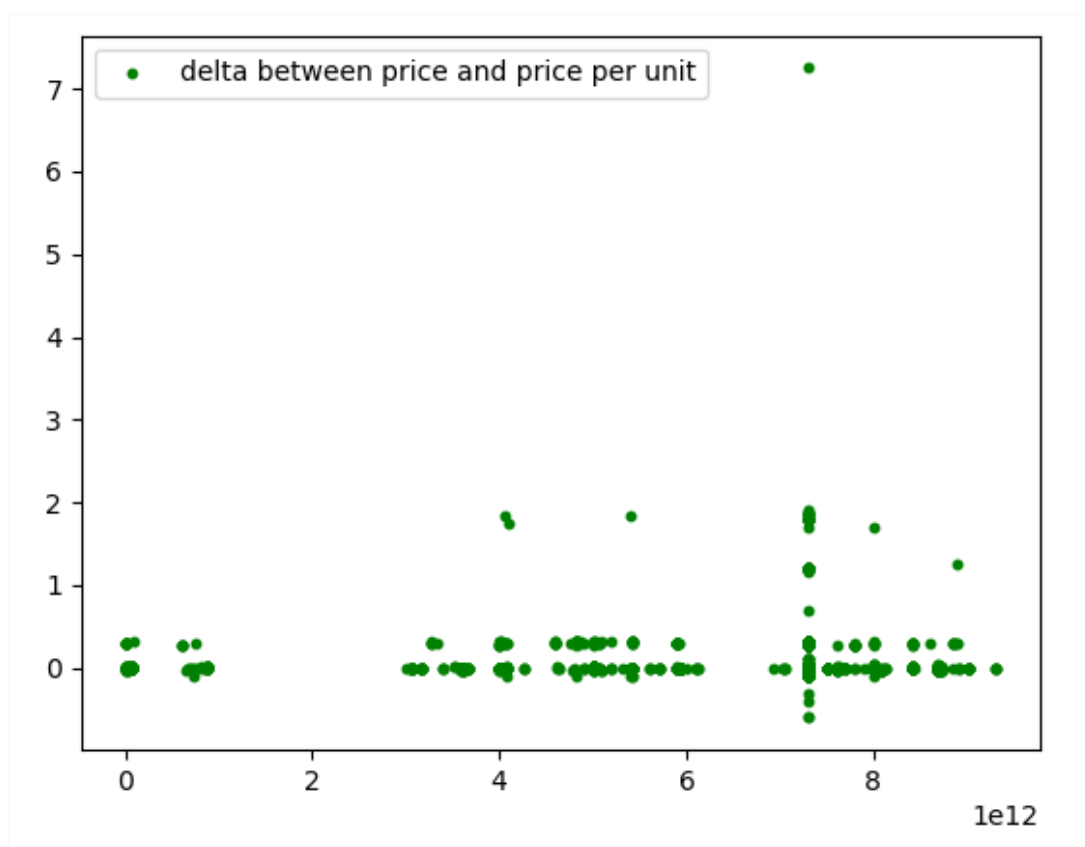
2. אם יחידות המדידה שונות בסקאלה שלהן (למשל במקרה של חישוב לפי יחידות שונות ליטר-ק"ג), הקשר בין ItemPrice, UnitOfMeasurePrice אמור לקיים:

$$\text{ItemPrice} = \text{UnitOfMeasurePrice} * \text{quantity}$$

בהתאם לאמור, חריגות מקשר בין נתונים, שאמור להיות שיקוף טכני (ערך נגזר מערך בין אם ידנית או במסגרת הדו"ח), מעידות ששישנם פערים במידע. יתכן שלאור טעויות הזנה – ובכך פגיעה באיכות הנתונים; או לחלופין, פער מחירים במציאות – ובכך, תדרש הסדרה של ניתוח המחירים לפי אחד הערכים באופן עקבי, כדי לנתחם בצורה אפקטיבית (מעבר לכך שהנתון משקף פגיעה בלקוחות).

- דוגמות: להלן ניתוח בסיסי של נתונים החורגים מקשר זה וגרף פיזור של הפערים בין הערכים הנצפים.

| price coha | hitOfMea | mPrice | ylnPack | sWeight | hitMeasu | quantity | hitQty | mCode |
|------------|----------|--------|---------|---------|----------|----------|--------|-----------------|
| | 49.9 | 49.9 | 1 | 1 | ק"ג | 1 | ק"ג | 9391 |
| invalid | 59.9 | 59.9 | 1 | 1 | גרם | 100 | ק"ג | 23489 |
| | 49.9 | 49.9 | 1 | 1 | ק"ג | 1 | ק"ג | 47997 |
| invalid | 5 | 50 | 1 | 1 | גרם | 100 | ק"ג | 58078 |
| invalid | 2.2 | 22 | 1 | 1 | גרם | 100 | ק"ג | 58238 |
| invalid | 6.9 | 69 | 1 | 1 | גרם | 100 | ק"ג | 318103 |
| invalid | 7.9 | 79 | 1 | 1 | גרם | 100 | ק"ג | 318165 |
| invalid | 7.9 | 79 | 1 | 1 | גרם | 100 | ק"ג | 318288 |
| invalid | 8.5 | 85 | 1 | 1 | גרם | 100 | ק"ג | 318349 |
| invalid | 7.9 | 79 | 1 | 1 | גרם | 100 | ק"ג | 318363 |
| invalid | 6.9 | 69 | 1 | 1 | גרם | 100 | ק"ג | 318387 |
| invalid | 7.5 | 75 | 1 | 1 | גרם | 100 | ק"ג | 318851 |
| invalid | 7.2 | 72 | 1 | 1 | גרם | 100 | ק"ג | 318868 |
| invalid | 4.9 | 49 | 1 | 1 | גרם | 100 | ק"ג | 318899 |
| invalid | 4.9 | 49 | 1 | 1 | גרם | 100 | ק"ג | 319049 |
| invalid | 7.2 | 72 | 1 | 1 | גרם | 100 | ק"ג | 7.29E+12 |
| invalid | 4.9 | 49 | 1 | 1 | גרם | 100 | ק"ג | 319759 |
| invalid | 8 | 80 | 1 | 1 | גרם | 100 | ק"ג | 7.29E+12 |
| valid | | 0.43 | 6.5 | 1 | 0 | מ"ל 100 | 15 | יח' 7.29002E+12 |
| invalid | | 0.82 | 12 | 1 | 0 | מ"ל 100 | 13.2 | יח' 7.29002E+12 |
| invalid | | 0.82 | 12 | 1 | 0 | מ"ל 100 | 13.2 | יח' 7.29002E+12 |
| invalid | | 0.82 | 12 | 1 | 0 | מ"ל 100 | 13.2 | יח' 7.29002E+12 |
| valid | | 4.6 | 6.9 | 1 | 0 | ליטר 1 | 1.5 | יח' 7.29002E+12 |
| valid | | 4.6 | 6.9 | 1 | 0 | ליטר 1 | 1.5 | יח' 7.29002E+12 |
| valid | | 4.6 | 6.9 | 1 | 0 | ליטר 1 | 1.5 | יח' 7.29002E+12 |
| invalid | | 0.82 | 12 | 1 | 0 | מ"ל 100 | 13.2 | יח' 7.29002E+12 |
| invalid | | 1.27 | 13.9 | 1 | 0 | מ"ל 100 | 10 | יח' 7.29002E+12 |
| invalid | | 1.27 | 13.9 | 1 | 0 | מ"ל 100 | 10 | יח' 7.29002E+12 |
| invalid | | 1.27 | 13.9 | 1 | 0 | מ"ל 100 | 10 | יח' 7.29002E+12 |
| valid | | 426.67 | 6.4 | 1 | 0 | מ"ל 100 | 0.015 | יח' 7.29002E+12 |
| valid | | 7.98 | 39.9 | 1 | 0 | מ"ל 100 | 5 | יח' 7.29002E+12 |



הסבר על הגרף: כל ערך שגדול מאפס מעיד על פער נצפה בין המחיר לבין הערך המחושב של המחירים כפונקציה של המחיר ליחידה והכמות הנמכרת

סעיף ב': בדיקת מציאות

בדקנו את התפלגות המחירים של כ-30% מקבצי PricesFull (10 קבצים מתיקיית victory ו-10 נוספים מתיקיית hashook). בחרנו לדגום את הקבצים האלו מכיוון שהם מאפשרים דגימה של יחסית הרבה מאוד מחירים. בבדיקה שלנו בחנו שני מאפיינים:

- ראשית, האם המחירים עומדים בבדיקת ולידיות בסיסית: ערכם גדול מאפס.
- שנית, האם המחירים מתפלגים כמצופה: ציפייה לטווחי מחירים של סופר, לפיהם ישנם הרבה מוצרים זולים ומעט מוצרים יקרים, עד כדי אלפים בודדים. ובהתאם, התפלגות המחירים הצפויה הינה לפי Power-law.

כדי לבחון את הנתונים בנוחות, לאחר שאכן התקבלה התפלגות התאמנו את גודל הbins והכנסנו ערכי קצה ל-collection bin כדי שהפריסה שתתקבל תהיה נוחה ככל הניתן.

לסיכום ניתן לראות שאכן הערכים עמודים בשני הקריטריונים והמחירים בכללם נראים כערכים תקינים.

