

Mattia Danese  
CS 119 - Big Data  
Professor Singh  
Quiz 2: The Opioid Files

Problem 1

The column names of the opioid dataset are:

- REPORTER\_DEA\_NO
- REPORTER\_BUS\_ACT
- REPORTER\_NAME
- REPORTER\_ADDL\_CO\_INFO
- REPORTER\_ADDRESS1
- REPORTER\_ADDRESS2
- REPORTER\_CITY
- REPORTER\_STATE
- REPORTER\_ZIP
- REPORTER\_COUNTY
- BUYER\_DEA\_NO
- BUYER\_BUS\_ACT
- BUYER\_NAME
- BUYER\_ADDL\_CO\_INFO
- BUYER\_ADDRESS1
- BUYER\_ADDRESS2
- BUYER\_CITY
- BUYER\_STATE
- BUYER\_ZIP
- BUYER\_COUNTY
- TRANSACTION\_CODE
- DRUG\_CODE
- NDC\_NO
- DRUG\_NAME
- QUANTITY
- UNIT
- ACTION\_INDICATOR
- ORDER\_FORM\_NO
- CORRECTION\_NO STRENGTH
- TRANSACTION\_DATE
- CALC\_BASE\_WT\_IN\_GM
- DOSAGE\_UNIT
- TRANSACTION\_ID

- Product\_Name
- Ingredient\_Name
- Measure
- MME\_Conversion\_Factor
- Combined\_Labeler\_Name
- Revised\_Company\_Name
- Reporter\_family
- dos\_str

The commands I ran:

- `zcat /comp/119/arcos_all_washpost.tsv.gz | head -1`

## Problem 2

The opioid dataset has a total of 178,598,026 rows.

The commands I ran:

- `zcat /comp/119/arcos_all_washpost.tsv.gz | wc`
- The above command returns 178,598,027 total lines, where each line represents one row, though 1 must be subtracted from this output as one of the lines is the header row and this is not an actual row in the dataset.

### Problem 3

Total number of rows for 2006: 21860398

Total number of rows for 2007: 23574939

Total number of rows for 2008: 24146453

Total number of rows for 2009: 25360920

Total number of rows for 2010: 26575386

Total number of rows for 2011: 28647123

Total number of rows for 2012: 28432806

The commands I ran:

- `zcat /comp/119/arcos_all_washpost.tsv.gz | shuf -n 5000 > temp.txt`
- I then made a python script that reads through all the lines of `temp.txt`, records how many rows correspond to each year, and then multiplies these proportions by the total number of rows in the opioid dataset

```
f = open("temp.txt")
lines = f.read().split('\n')

count = {
    '2006' : 0, '2007' : 0, '2008' : 0, '2009' : 0, '2010' : 0,
    '2011' : 0, '2012' : 0,
}

for line in lines:
    fields = line.split('\t')

    date = fields[30]
    year = date[-4:]

    count[year] += 1
print(count)

TOTAL_ROWS = 178598026
for year in count.keys():
    print("The estimated number of rows for {} is: {}".format(
        year,
        round((count[year] / 5000) * TOTAL_ROWS)))
```