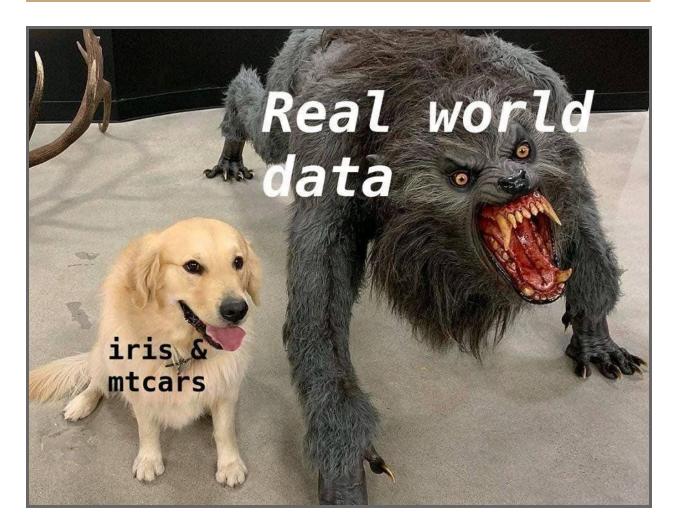# WHAT TO DO WITH DIRTY DATA



## Introduction

What is data? A better question would be "what **isn't** data?" Data is everything. Quite literally anything can be data. However, in order for it to be data we can use in a paper like this one, we have to take what is often noisy, non-numeric information and carefully clean it (often transforming things like words or photos into tokens and pixels) into something that ends up looking very different from the "data" we started with.

We usually think of "data" as a nice clean spreadsheet. This, however, is a lie. A fantasy. Something you must shoo from your mind immediately or else you will enter into a void of

pain and existentialism eventually asking yourself "what even is a number and why do I care?" But, before we leave this dream world entirely, let's consult the **Dream Data.** This dream data is clearly labeled, filled with mostly numeric, normalized data. There are no missing fields, no strange characters, no instances of nesting and definitely no emojis. This dream data is definitely not missing large swaths of information. It's not an entire endless matrix of 0s with the occasional 1s. And it is most definitely not a hexadecimal representation of a photo of handwriting from a serial killer. However, to reiterate, all of that is data. So, how can we, as scientists, turn all of this flotsam and jetsam into something our computers can begin to "understand?"

# Analysis & Models

## ABOUT THE DATA

MoviesRAW.csv is a "dirty" data file. If the "ideal data file" is something like a two column spreadsheet, this is a non-ideal data file. It has many columns for each review and a number of non-alpha characters cluttering each review string. Finally, the label for each review is tacked on to the end, wherever that end may be.

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | text | reviewclass | | | | | | | | | |
| 2 | 'plot : two te | drink and th | but his girlfr | and has nigh | but presents | since i gene | mess with y | but there ar | and these fo | but execute | its main pro |
| 3 | 'the happy ba | virus still fe | like a movie | we don\'t kr | and | of course | we don\'t kr | it\'s just \" | let\'s chase | even from th | well |
| 4 | 'it is movies | the mod squ | things go w | the ads mak | cool music | claire dane\ | car chases | stuff blowin | and the like | does it not ? | it quickly be |
| 5 | ' \" quest for | fully-animat | but the mou | if flawed | 20th century | \" but disne | \" with its li | had her bea | it\'s no cont | as \" quest f | the early-te |
| 6 | 'synopsis : a | a fledgling r | he takes pic | both theatri | no big name | he\'s rejecte | ex-wife | or ex-husba | stalked is ju | though that | for instance |
| 7 | 'capsule : in | especially as | that the nev | 0 ( -4 to +4 ) | supposedly a | it is a very b | but it is mor | carpenter w | the thing | and prince o | in fact |
| 8 | 'so ask yours | sleazy unde | how | bubbling jus | there\'s a sc | those who a | supposed \" | joel schuma | \" \" a time | \" \" batma | \" \" the clie |
| 9 | 'that\'s exac | mr . hugh gr | a huge dork | not me ) tha | it\'s the fact | we\'re talkir | on the other | since when | the obstetri | tells grant\' | \" referring |
| 10 | 'call it a road | boozed-out | a sentiment | unable to ex | focusing on | tomas ( skar | openly hosti | kaisa ( lena | gossip ) . \n | and wouldn' | scotland by |

FIGURE 1: The original csv file

## CLEANING THE DATA

There are many different ways up this mountain but for this exercise, we will demonstrate two different approaches. The first way, Ami's way, takes the data and keeps it in text form while cleaning, ultimately exporting a new, cleaner text file for a fresh import. The second way, Kendra's way, takes the data and immediately turns it into a pandas dataframe. There are pros and cons to each way.

## AMI'S WAY

Overview:

1. Read in the dirty file
2. Prep new clean files
3. Clean the data
   a. For each row in the data, clean the row
   b. For each word in the row, clean the word
4. Export clean data to new clean files
5. Re-import the cleaned data
6. Turn cleaned data into a pandas df

*Code to clean each row*

```python
def display_rows(file_data):
    for row in file_data:
        row = row.lstrip()
        row = row.rstrip()
        row = row.strip()
        raw_row = "\n\nROW:" + row + "\n"
        outfile.write(raw_row)
        row_list = row.split(" ")
        new_list = []
        for word in row_list:
            to_put_in_outfile = "The next word BEFORE is: "+ word +"\n"
            outfile.write(to_put_in_outfile)
            word = clean_word(word)
            if word:
                new_list.append(word)
        label = ''.join(char for char in new_list[-1] if char.isalpha())
        new_list.pop()
        just_text = ' '.join(new_list)
        to_write = label + ',' + just_text + '\n'
        cleanfile.write(to_write)
```

*Code to clean each word*

```python
def clean_word(word):
    word=word.lower()
    word=word.lstrip()
    word=word.lstrip("\\n")
    word=word.strip("\n")
    word=word.replace(",","")
    word=word.replace(" ","")
    word=word.replace("_","")
    word=re.sub('\+', ' ',word)
```

```python
word=re.sub('.*\+\n', '',word)
word=re.sub('zz+', ' ',word)
word=word.replace("\t","")
word=word.replace(".","")
word=word.strip()
word = word.replace("\\'","")
if word not in ["", "\\", '"', "'", "*", ":", ";"]:
    if len(word) >= 3:
        if not re.search(r'\d', word): ##remove digits
            return word
```



FIGURE 2: An example of the outfile generated by Ami's way

## KENDRA'S WAY

Overview:

1. Read in the dirty file

2. Turn it into a pandas data frame

3. Merge all the rows to get the review text

4. Remove the last characters to get the "labels"

5. Clean the review text

*Code to merge all the rows together*

```python
import pandas as pd
dirtyFile = pd.read_csv('moviereviewRAW.csv')
df = pd.DataFrame()
```

```
df['all'] = dirtyFile[dirtyFile.columns[0:]].apply(
    lambda x: ','.join(x.dropna().astype(str)),
    axis=1)
```

*Code to get the label*

```
df['label'] = df.apply(lambda x: x['all'][-3], axis=1)
```

In the same way there are many ways up the "how to process files" mountain, there are many ways to clean text. As we discovered in HW1, there are some things we want to keep, some things we want to discard. We rarely want to keep strange characters and unnecessary white spaces.

*Code to clean excess characters*

```
def clean_rogue_characters(string):
    exclude = ['\\','"\'"]
    string = '.'.join(string.split('\\n'))
    string = ''.join(ch for ch in string if ch not in exclude)
    return string

df['all'] = df['all'].apply( lambda x: clean_rogue_characters(x) )
```

# Results

To get 'results' for this quick-and-dirty assignment, the researchers compared the 'dirty data' to past data sets in their 'sentiment analysis' pipeline to answer the question -- does this newly cleaned data behave very similarly, slightly similarly or not at all similarly to a cleaner dataset from the wild?

## Text Blob

|  | Kendra's Data | Ami's Data | Cornell Data | Dirty Data | Joker Data |
|---|---|---|---|---|---|
| **CORRECT NEG** | 5 | 1 | 229 | 227 | 64 |
| **CORRECT POS** | 0 | 4 | 971 | 972 | 114 |

## VADER

|  | Kendra's Data | Ami's Data | Cornell Data | Dirty Data | Joker Data |
|---|---|---|---|---|---|
| **CORRECT NEG** | 2 | 3 | 445 | 454 | 64 |
| **CORRECT POS** | 5 | 3 | 828 | 824 | 114 |

## NLTK

|  | Kendra's Data | Ami's Data | Cornell Data | Dirty Data | Joker Data |
|---|---|---|---|---|---|
| **CORRECT NEG** | -- | -- | 89% | 86% | 81% |
| **CORRECT POS** | -- | -- | 74% | 70% | 35% |
| **ACCURACY** | -- | -- | 81% | 77% | 58% |

Looking at the "Dirty Data" compared to the Cornell data, it's clear to see that once cleaned, the dirty data performed almost eerily similarly to the clean data. This is without additional NLP cleaning, simply a baseline analysis. For the purposes of this exploration, the dirty data was cleaned enough to perform as well as its cleaned counterpart.

# Conclusion

Dirty data is everywhere. Having a pipeline, (or multiple pipelines!) with which to quickly clean, format and export data is essential to being an effective data scientist. The researchers suggest having multiple different (and mutable!) pipelines for different tasks -- is the dirty data coming from the web? Is it littered with HTML? Is the dirty data coming from a poorly formatted csv? There will never be a one-size-fits-all cleaner, however, there will be ways to quickly format the data for easier more in-depth "post cleaning."