

DATA 601-W6

Data Loading and Wrangling

Data Wrangling

- We will learn ways and techniques to getting data into Python
- We already see one of these ways: Working with .csv files

Types of Data: Structured Data

- Conforms a data model
- Fields of data will be stored
- How that data will be stored: data type (numeric, alphabetic, date, Boolean)
- Any restrictions on the data input (number of characters, etc.).

Types of Data: Semi-Structured

- There is a structure in the data in terms of tags, rows, columns, hierarchies, fields.
- But not necessarily there are unified restrictions and pre-defined types for fields and tags.

Types of Data: Unstructured Data

- No pre-defined data model
- Not organized in fields, tags, attributes, etc.

Examples:

- social media posts, images, texts.

File Formats (How we record data)

What is a file format?

A file format is a standard way that information is encoded for storage in a computer file. It specifies how bits are used to encode information in a digital storage medium. File formats may be either proprietary or free and may be either unpublished or open.

Why do we need formatting the files?

- Share (conventions)
- Store - Preserve (efficient ways - speed, different needs, loose of data)
- Access (speed)

Why to use so many different formats?

- Speed
- Different nature of the data
- Each has different advantages and disadvantages
- Each has different use cases

File Formats: CSV Files

- Uses .csv extension.
- Data is organized in columns and each column is separated with a column.
- Very common format in data science.
- Looks very similar to spreadsheet files however note that you cannot encode formulas and format.
- Because there is no formatting and formulas it is easier to work with this format for bigger file sizes.
- If we use tab instead of , then the file format is called tsv (tab separated values)

Libraries to work with CSV files

- Pandas
 - Faster to load for larger datasets
 - Mainly for analyzing, applying functions to data
- CSV module in Python
 - For reading and writing csv files only.
 - Built in Python, you don't have to install.
 - [Faster with smaller datasets \(<1K rows\)](#)

File Formats: JSON Files

- **JavaScript Object Notation**
- Based on JavaScript but it is completely language independent
- Consists of name-value pairs
- It is a very common format especially in data interchange of web applications.
- JSON is **not** a programming language

Libraries to work with JSON files

- Python has a built in library called json to read and parse json files.
- We can also use pandas library to read json files.

Other File Formats...(but we will not be covering today)

- XML
- HTML
- Audio (mp4, wav, etc.)
- Image (png, jpeg, jpg, etc.)
- Other (pickle, hdf5, parquet, etc.,.)

Character Encoding

- Do we know bytes? Bits?
- Character encoding is basically mapping each character to bytes (Simple huh?)
- As anything this is not also that straightforward!!

Character Encoding: ASCII

- [ASCII Character Set](#)

bits	character
01000001	A
01000010	B
01000011	C
01000100	D
01000101	E
01000110	F

Then What Happened?

- The need for way bigger character sets needed!
- Can you see why?
- [Let's take a look at the wikipedia again!](#)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
	Latin Extended-B															
0180	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
0190	Ę	ę	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
01A0	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć	Ń	ń
01B0	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć	Ń	ń
01C0	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
01D0	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
01E0	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
01F0	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
0200	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
0210	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
0220	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
0230	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
0240	Ł	ł	Œ	œ	Š	š	Ž	ž	Ć	ć	Ń	ń	Ą	ą	Ć	ć
	IPA Extensions															
0250	ɐ	ɑ	ɒ	ɓ	ɔ	ɔ̃	ɔ̄	ə	ə̃	ə̄	ɛ	ɛ̃	ɛ̄	ɛ̅	ɛ̆	ɛ̇
0260	ɔ̈	ɔ̉	ɔ̊	ɔ̋	ɔ̌	ɔ̍	ɔ̎	ɔ̏	ɔ̐	ɔ̑	ɔ̒	ɔ̓	ɔ̔	ɔ̕	ɔ̖	ɔ̗
0270	ɔ̘	ɔ̙	ɔ̚	ɔ̛	ɔ̜	ɔ̝	ɔ̞	ɔ̟	ɔ̠	ɔ̡	ɔ̢	ɔ̣	ɔ̤	ɔ̥	ɔ̦	ɔ̧
0280	ɔ̨	ɔ̩	ɔ̪	ɔ̫	ɔ̬	ɔ̭	ɔ̮	ɔ̯	ɔ̰	ɔ̱	ɔ̲	ɔ̳	ɔ̴	ɔ̵	ɔ̶	ɔ̷
0290	ɔ̸	ɔ̹	ɔ̺	ɔ̻	ɔ̼	ɔ̽	ɔ̾	ɔ̿	ɔ̻	ɔ̼	ɔ̽	ɔ̾	ɔ̿	ɔ̻	ɔ̼	ɔ̽
02A0	ɔ̿	ɔ̻	ɔ̼	ɔ̽	ɔ̾	ɔ̿	ɔ̻	ɔ̼	ɔ̽	ɔ̾	ɔ̿	ɔ̻	ɔ̼	ɔ̽	ɔ̾	ɔ̿

Why Should we care?

```
pd.read_csv('/content/PoliceShootingsUS.csv')
```

```
-----  
UnicodeDecodeError                                Traceback (most recent call last)  
pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._convert_tokens()  
  
pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._convert_with_dtype()  
  
pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._string_convert()  
  
pandas/_libs/parsers.pyx in pandas._libs.parsers._string_box_utf8()
```

```
UnicodeDecodeError: 'utf-8' codec can't decode byte 0x96 in position 2: invalid start byte
```

During handling of the above exception, another exception occurred:

```
UnicodeDecodeError                                Traceback (most recent call last)  
<ipython-input-180-04c68a038be7> in <module>()  
----> 1 pd.read_csv('/content/PoliceShootingsUS.csv')
```


Solution?

- Chardet library
- We will discuss in the second part of the lecture...