
Week 3 Slide Deck

Wrangling, Filtering; Formats

Jack Bandy 2026

Data Formats

How Do We Store This Data?



Spreadsheet??

- You've probably seen this
- Grid with cells, rows, columns (and sheets/tabs)
- Excel, Google Sheets, Numbers, and LibreOffice Calc
- Convenient for manual entry, review, sharing
- But don't keep your data in a spreadsheet...

	A	B	C	D	E
1	name	company	street	city	phone
2	Tyler Durden	Paper Street Soap Co.	537 Paper Street	Bradford	(288) 555-0153

CSV

- Comma-separated values
- Plain-text table: one row per record, one delimiter between fields
- Often used for spreadsheets, exports, and simple datasets
- Easy to inspect
- types and hierarchy are usually implicit

```
1 name,company,product,street,city,postalCode,phone
2 Tyler Durden,Paper Street Soap Co.,All Natural Handmade,537 Paper Stre
```

TSV

- Tab-separated values
- Plain-text table like CSV, but fields are separated by tabs
- Still flat: hierarchy and data types need outside context

```
1 name      company product street  city      postalCode phone
2 Tyler Durden  Paper Street Soap Co. All Natural Handmade 537 Pa
```

JSON

- JavaScript Object Notation
- Text format for data interchange
- Built from objects, arrays, strings, numbers, booleans, and null
- Common for web APIs and configuration files

```
1 {
2   "name": "Tyler Durden",
3   "company": "Paper Street Soap Co.",
4   "product": "All Natural Handmade",
5   "address": {
6     "street": "537 Paper Street",
7     "city": "Bradford",
8     "postalCode": "19808"
9   },
10  "phone": "(288) 555-0153"
11 }
```

XML

- Extensible Markup Language
- Nested tags represent elements and attributes
- Verbose, but widely used by older document systems

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <person>
3   <name>Tyler Durden</name>
4   <company>Paper Street Soap Co.</company>
5   <product>All Natural Handmade</product>
6   <address>
7     <street>537 Paper Street</street>
8     <city>Bradford</city>
9     <postalCode>19808</postalCode>
10  </address>
11  <phone>(288) 555-0153</phone>
12 </person>
```

YAML

- YAML Ain't Markup Language
- Human-readable format based on indentation
- Supports mappings, lists, scalars, and comments
- Common for configuration files and data pipelines

```
1 name: Tyler Durden
2 company: Paper Street Soap Co.
3 product: All Natural Handmade
4 address:
5   street: 537 Paper Street
6   city: Bradford
7   postalCode: 19808
8 phone: (288) 555-0153
```

Parquet

- Binary columnar storage format
- Efficient for large datasets
- Stores schema and data types with the data
- Common in Spark, DuckDB, Polars, and “data lakes”

```
1 message business_card {
2   required binary name (STRING);
3   required binary company (STRING);
4   required binary product (STRING);
5   required binary street (STRING);
6   required binary city (STRING);
7   required binary postalCode (STRING);
8   required binary phone (STRING);
9 }
10
11 row 1:
12 Tyler Durden | Paper Street Soap Co. | All Natural Handmade |
13 537 Paper Street | Bradford | 19808 | (288) 555-0153
```

Format Summary

Format	Best fit	Short example
CSV	Flat tables, spreadsheet exports, simple datasets	<code>name,company,phone</code> <code>Tyler Durden,Paper Street Soap Co.,(288) 555-0153</code>
TSV	Flat text tables where commas may appear in fields	<code>name company phone</code>
JSON	APIs, nested records, web data	<code>{"name":"Tyler Durden","city":"Bradford"}</code>
XML	Document-like data with tags and attributes	<code><name>Tyler Durden</name></code>
YAML	Human-edited configuration and pipeline settings	<code>name: Tyler Durden</code> <code>city: Bradford</code>
Parquet	Typed, compressed analytics data	<code>name: STRING</code> <code>city: STRING</code>

Sources

1. GitHub source: <https://github.com/jackbandy/data-science-fun/blob/main/docs/slides/week3.md>.
2. Slide deck built with [Quarto](#) revealjs.
3. Format examples adapted from Wikipedia and project documentation: [JSON](#), [XML](#), [Comma-separated values](#), [YAML](#), [Tab-separated values](#), [Spreadsheet](#), and [Apache Parquet](#).
4. Tyler Durden business card image: Wikimedia Commons remake by Michaelpreid, modified, [CC BY-SA 4.0](#), https://commons.wikimedia.org/wiki/File:Tyler_Durden_Business_Card.png.
5. Title font is Big Shoulders; Body font is [Libre Franklin](#).