

Data exploration

Collecting data

Differentiate between data formats and structures

- ✓ **Video:** Discover data formats
5 min
- ✓ **Reading:** Data formats in practice
10 min
- ✓ **Practice Quiz:** Self-Reflection: Unstructured data
2 questions
- ✓ **Video:** Understanding structured data
1 min
- ✓ **Reading:** The structure of data
10 min
- ✓ **Ungraded Plugin:** Differentiating Data Types
30 min
- 📖 **Reading:** Data modeling levels and techniques

The structure of data

Data is everywhere and it can be stored in lots of ways. Two general categories of data are:

- **Structured data:** Organized in a certain format, such as rows and columns.
- **Unstructured data:** Not organized in any easy-to-identify way.

For example, when you rate your favorite restaurant online, you're creating structured data. But when you use Google Earth to check out a satellite image of a restaurant location, you're using unstructured data.

Here's a refresher on the characteristics of structured and unstructured data:

Structured data



- Defined data types
- Most often quantitative data
- Easy to organize
- Easy to search
- Easy to analyze
- Stored in relational databases & data warehouses
- Contained in rows and columns
- Examples: Excel, Google Sheets, SQL, customer data, phone records, transaction history

Unstructured data



- Varied data types
- Most often qualitative data
- Difficult to search
- Provides more freedom for analysis
- Stored in data lakes, data warehouses, and NoSQL databases
- Can't be put in rows and columns
- Examples: Text messages, social media comments, phone call transcripts, various log files, images, audio, video

Structured data

As we described earlier, **structured data** is organized in a certain format. This makes it easier to store and query for business needs. If the data is exported, the structure goes along with the data.

Unstructured data

Unstructured data can't be organized in any easily identifiable manner. And there is much more unstructured than structured data in the world. Video and audio files, text files, social media content, satellite imagery, presentations, PDF files, open-ended survey responses, and websites all qualify as types of unstructured data.

The fairness issue

The lack of structure makes unstructured data difficult to search, manage, and analyze. But recent advancements in artificial intelligence and machine learning algorithms are beginning to change that. Now, the new challenge facing data scientists is making sure these tools are inclusive and unbiased. Otherwise, certain elements of a dataset will be more heavily weighted and/or represented than others. And as you're learning, an unfair dataset does not accurately represent the population, causing skewed outcomes, low accuracy levels, and unreliable analysis.

✓ Completed

Go to next item

👍 Like 🗨 Dislike 📄 Report an issue

