

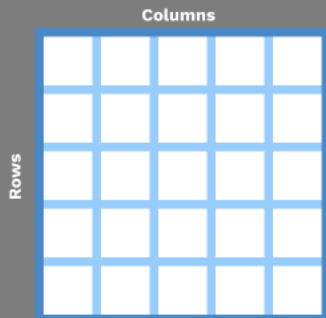
Structured data

- **Only ~20% of the data in an organisation is structured data.**
- **However constitutes ~80% of the data science use cases.**

Common sources/types of structured data



Data is organized in rows and columns



Unstructured data

~80% - 90% of the data in an organization is unstructured

Common Types of unstructured data

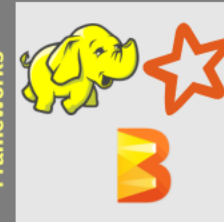


Best and most economical storage



Pre-processing and processing

Big Data Processing Frameworks



- Scaling images
- Text from audio
- Video annotation
- Tokenization for text

NoSQL Databases



- Storing metadata and location of thumbnails and image in data lake

Common data science use cases

- Computer vision
- Speech to text
- Natural language processing (NLP)

Semi structured data

Common sources of semi-structured data



Web activity logs



Sensors



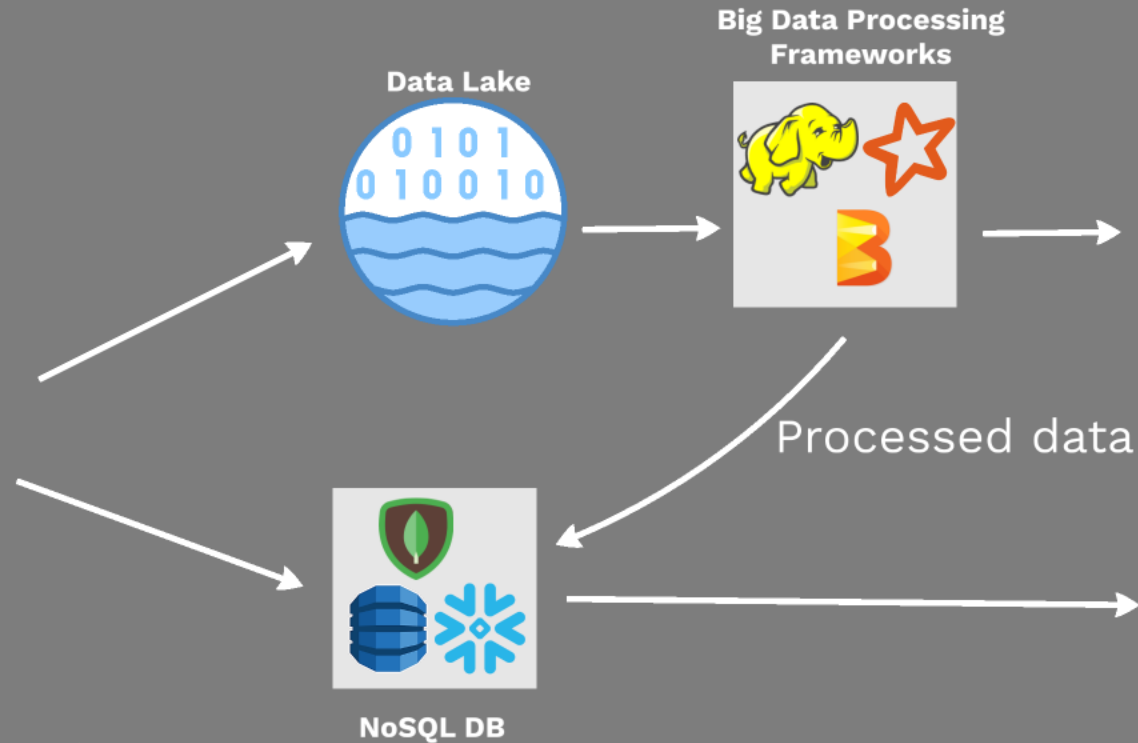
Geo loc data



Big data file formats

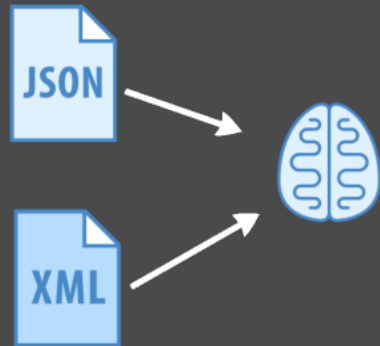
Data Formats

- XML - eXtensible markup language
- JSON - JavaScript Object Notation (JSON)
- Big data file formats like ORC, Parquet and Avro.

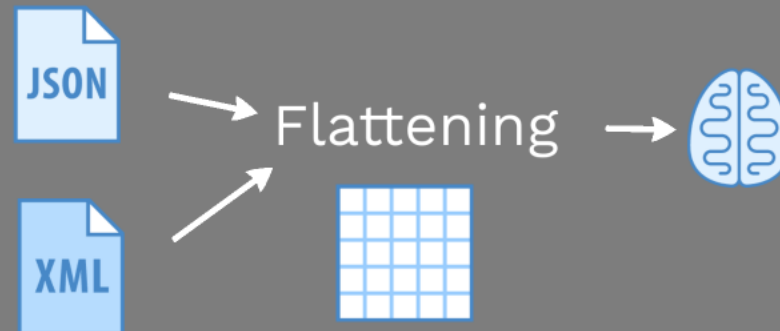


- Semi-structured data can be loaded either into a data lake or a NoSQL database.
- The choice depends on
 - data volume
 - amount of processing required
 - existing data engineering skill set
- In some cases data is processed by big data frameworks and the NoSQL database is used as a serving layer.

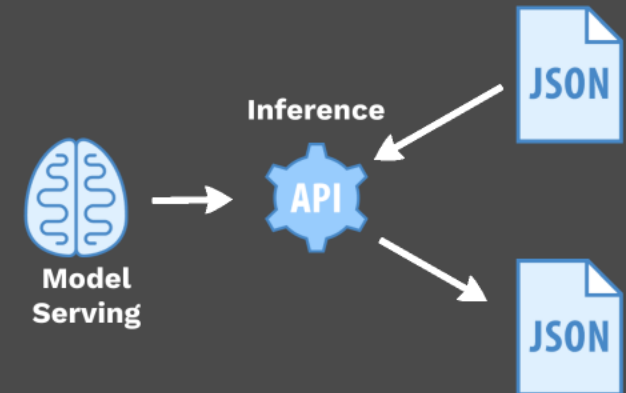
Machine learning with semi-structured data?



Semi-structured data is never used to train models directly.



Semi-structured data is first flattened and then used for model training.



Model serving API's generally accept semi-structured data and also output data in the same format