

COURSE03: Prepare data for Exploration

Instructors Hallie Analytical lead
Healthcare

Page No.
Date 07.01.23

Week1 Data types and structures:→

W4 Data Exploration

➤ Introduction:→

• Understanding the different types of data and data structures.

- what type of data is right for the question you're answering
- practical skills about how to extract, use, organise and protect your data

Story: she analyzes Medicare enrollment data over time and make connections to how people research Medicare plans on google.

This is helpful for medical professionals & patients

Prepare steps:→

- How data is generated
- Different formats, types, and structures of data
- Analyze data for bias and credibility
(cause not all data fits each need)
- what "clean data" means
- Databases
- Extract your own data using spreadsheets & SQL
- Basics of data organization
- The process of protecting your data

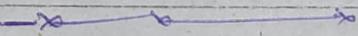
You still with me?:) → Yes ma'am!

* Fascinating data insights!

Healthcare & data together is a newer concept. Using ML, AI, Big data help healthcare industry.

"It's really fascinating that we can take all of these data sets and synthesize them and allow us to deliver some cool insights and trends to our hospital systems."

The creativity scale of analysis grows with exp, you can look at data one way now and a week later you might see a different trend altogether!



2/4 Collecting data :-

Every digital photo online is one piece of data. Every photo itself holds even more data from number of pixel to color contained in them.

US Census Bureau use forms to collect data about the country's population.

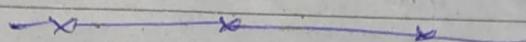
This data is used for numerous reasons like funding for schools, hospitals & fire departments.

~~* As a DA you'll have access to all kinds of data and lots of it. Knowing how its generated can add more context to it and knowing how to collect it, can make DA process more efficient.~~

No matter what kind of data you use it must be checked for accuracy & Trustworthiness

How data is collected.

- > Interviews
- > Observations
- > forms
- > Questionnaires
- > Surveys
- > Cookies : (When tracking people's online activity and interests (most effective))



* What Data to collect ? →

* First-party data: Data collected by individual or group using their own resources.



Preferred method, as you know where the data came from.

* Second-party data: Data collected by a group directly from its audience and then sold



* You buy from someone. Still reliable.

* Third party data: →

Data collected from outside sources who did not collect it directly

Might not be reliable BUT still useful

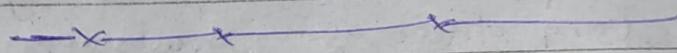
You must check it for accuracy, bias & credibility.

populations All possible values in ? Page 10 Certain data
sample: part of poplⁿ that is representative of poplⁿ

* Data collection Considerations →

- > How the data will be collected
- > Choose data sources (first, 2ⁿ, 3rd-Party)
- > Decide what data to use
- > How much data to collect (poplⁿ & sample)
- > Select the right data type.
- > Determine the time frame.

ex. if you need something immediately ; You'll need to use historical data



08-01-23

Week 1

3/4

Differentiate between date formats and structures.

* Discover data formats →

* > Qualitative Data: Name, category, description

* > Quantitative Data: Numeric data

stars, point → Discrete: counted & limit values

runtime → Continuous: measured & any num value

half not allowed

110.0356 minutes just time

* Nominal Data: qualitative data categorized without set order.



Doesn't have sequence

What?

→ (Yes, No, Not sure) Doesn't specify likelihood, etc.

Did you like? → 4, 5, 3, 4, 1, 5
 ↓

* Ordinal Data: Qualitative data with set order/scale

* Internal data: →

Data that lives within a company's own system.

↓
more reliable & easier to collect

* External data: →

Data that lives and is generated outside of an organisation.

↓
valuable when you need as much data as possible.

* Structured data: →

Data organized in a certain format such as rows & columns.

↓
Spreadsheets & relational databases

* Unstructured data: →

Data that is not organised in any easily identifiable manner. (might have internal structure but doesn't fit neatly in rows & columns).
 ↓

Audio, Video as there's no clear way to identify or organize their content

Understanding structured data

Most of the data that is being generated is unstructured: → Audio, video, emails, photos, social media.

Structured data works nicely in a data model

Data model →

A model that is used for organizing data elements and how they relate to one another.

Data elements

Pieces of information, such as people's names, acc no's and addresses.

Data model help keep the data consistent and provide a map of how data is organised.

Structured data is also useful for databases, this way data analysts can enter, query & analyze the data whenever they need to.

Makes visualization easy!

as structured data can be directly applied to charts, maps, graphs, dashboards & most other visual representations of data.

Sources: spreadsheets & databases

* Data modelling :

- Process of creating diagrams that visually represent how data is organised and structured.
- These visual representations are called as Data models

* 3 most common types :-

- ① Conceptual (Business concepts)
- ② Logical (Data entities)
- ③ Physical (Physical labels)

3 schema Architecture

Common methods :-

Entity Relationship Diagram (ERD).

visual way to understand relationship

Unified Modelling Language (UML)

detailed diagrams that describe structure

more about it:

- * Know the type of your data:

Data type:

A specific kind of data attribute that tells what kind of value the data is

Data types in spreadsheet →

> Number

> Text or String : sequence of char & punctuation

> Boolean : only 2 possibilities

Some programs have different names or include extra types.

But these cover any data you'll find in spreadsheets

Text can include numbers but they aren't used for calculation

(12)

→ → →

What do these have in common?

- > A music player
- > A calendar agenda
- > Email inbox

Rows = Records

All are arranged in tables

Columns → Fields

Rows and columns are usually reserved for spreadsheets

Record & Columns are universal.

* Wide & long data →

Wide data:

Data in which every data subject has a single row with multiple columns to hold the values of various attributes of the objects.

long data Subject : year : 2010
2011

WIDE DATA

so	country	Y2010	Y2011	Y2012
A	pa akistan palestine	x	y	z
B	co ontry			
C				

long data →

Data in each row is one time point per subject, so each subject will have data in multiple rows.

Country	Year	Population	long Data
A	2010	x	
B	2011	y	
C	2012	z	
A	2010	p	
B	2012	q	
C	2010	r	
A	2012	s	
B	2020	t	
C	2011	u	

* Long data is a great format for storing and organizing data when there are multiple variables at each time for each subject.



Nice & Compact

* But it depends on project which format you should be using.

Learned: + Identifying bias in data and how to embrace credibility, integrity and ethics.

* Data Transformation:-

Data organization: Better organized data is easy to use.

Data compatibility: Different apps & system can then use the same data.

Data migration: Data with matching formats can be moved from one system to another.

Data merging: Data with the same organization can be merged together.

Data enhancement: Data can be displayed with more detailed fields.

Data comparison: Apple-to-Apples comparison or data can then be made.

Week 2 | Ensuring data integrity

Page No. 10
Date 01/23

- * > Analyze data for bias and credibility
Important because even the most sound data can be skewed or misinterpreted.

Good data vs Bad data

Data ethics, privacy and access

We need to ask question like:

- ① Who owns this data?
- ② How much control do we have over the privacy of data?
- ③ Can we use & reuse data however we want to?

Bias: A preference in favor of or against a person, group of people or thing.

Once we know & accept that we have bias, we can start to recognize our own patterns of thinking & learn to manage it.

Data Bias: A type of error that systematically skews results in a certain direction.

* Bias can happen if sample group lacks inclusivity.

* Sampling Bias :

when the sample isn't representative of the population as a whole.

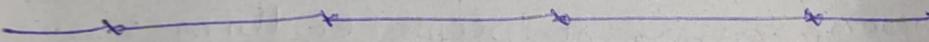


we can avoid this by choosing the sample randomly, so that all parts of the population have an equal chance of being included

* Unbiased sampling :

when a sample is representative of the population being measured.

visualizing the data will reflect if it's biased or unbiased.



* More Types of Bias

* Observer Bias →

Also called as Experimenter/ Researcher Bias

The tendency for different people to observe things differently.

① Ex. scientists see different things under microscope

② medical professional see different pp readings. (Growth/Ht)

4 Types of Data Bias

- (1) Sampling Bias
- (2) Observer Bias
- (3) Interpretation Bias
- (4) Confirmation Bias

Page No.

Date

* Interpretation Bias:

The tendency to always interpret ambiguous situations in a positive or negative way

→ Two people seeing and hearing the exact same thing and interpreting it in a variety of different ways (due to different background & experiences).

* Confirmation Bias

The tendency to search for or interpret information in a way that confirms pre-existing beliefs.

"people see what they want to see"

week 2

2/4

measure the credibility of data sets →

ROCCC

Reliable

Original

Comprehensive (contains all info)

Current (relevant)

Cited

When choosing a datasource, think these 3 things:

(1) Who created the dataset

(2) Is it part of credible organization?

(3) When was the data last updated?

Good data location! Vetted public datasets, Academic papers, financial data, Governmental Agency data

Bad data sources don't ROCCC

R → Inaccurate, incomplete or biased

could be sampling bias data
may lead to misleading insights

D → From 2nd / 3rd party

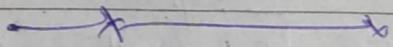
so you need to be extra careful while
understanding your data

C → ~~too~~ missing important information needed
not comprehensive to answer the question
may contain human error

C → not current, out of date & irrelevant

C → not cited, not vetted then it's a no go

"Every good solution is found by avoiding bad
data"



Introduction to data ethics:-

Ethics:

Well-founded standards of right and wrong that prescribe what humans ought to do, usually in terms of rights, obligations, benefits to society, fairness, or specific virtues.

Data Ethics:

Well-founded standards of right and wrong that dictate how data is collected, shared and used.

GDPR

General Data Protection Regulation of the European Union

Data protection legislation

Six aspects of data ethics:-

> Ownership

Individuals own the raw data they provide and they have primary control over its usage, how it's processed, and how it's shared

> Transaction Transparency:

All data-processing activities and algorithms should be completely explainable and understood by the individual who provides their data.

Let's people judge whether outcome is fair and unbiased & allows to raise concerns.

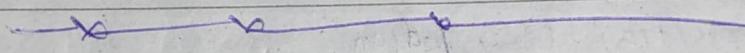
> Consent :

An individuals right to know explicit details about how and why their data will be used before agreeing to provide it.

Nowadays, it just looks like terms & conditions with just a checkbox.

> currency :

Individuals should be aware of financial transactions resulting from the use of their personal data and the scale of these transactions.



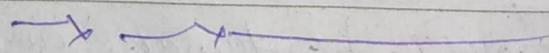
Alex : Research Scientist with the Ethical AI team

Optional!

"How do we actually improve the lives of people by using data?" → [Don't lose this view]

"Data Are people"

Take responsibility of those people



six ETHICS Are → Ownership
Transaction Transparency
Consent
Currency
Privacy
Openness

Privacy:

Preserving a data subject's information and activity any time a data transaction occurs.

Also called information privacy / data protection.

* It's about access, use & collection of data

Legal rights to one's data →

- > Protection from unauthorized access to our private data
- > Freedom from inappropriate use of our data
- > The right to inspect, update, or correct our data
- > Ability to give consent to use our data
- > Legal right to access the data

* Openness →

Free access, usage, and sharing of data

Andrew: senior developer Advocate, ethical AI research group

"One consequence of not using this technology responsibly is the possibility of amplifying or enforcing unfair biases."

Week 2
9/4

Open data :

Free access, usage and sharing of data

But we should still be transparent, respect privacy and make sure we have consent for data that's owned by other

In simple terms, if data meets these conditions then only we may access, use & share it.

* Availability & access

Open data must be available as a whole, preferably downloadable over internet in convenient & modifiable form, (ex. data.gov)

* reuse & redistribution:-

It should allow reuse & re-distribution also the ability to use it with other data sets.

* Universal participation:-

Everyone should be eligible & there shouldn't be discrimination

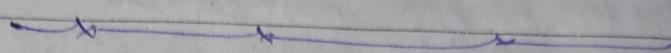
Ex. in Health sector, openness allows us to access & combine diverse data to detect diseases earlier & earlier.

Data Interoperability \rightarrow

The ability of data systems and services to openly connect and share data.



Requires lot of cooperation



Week
3
VS

12/01/23

Working with datasets.

All about datasets \rightarrow

Database \rightarrow

A collection of data stored in a computer system.

meta \rightarrow Referencing back to itself / Being completely self aware.
Metadata : Data about data

A character knows its inside a book

Documentary about making documentary
Analyzing how you analyze data

Once you have data in a spreadsheet the possibilities are endless.

Relational Database \rightarrow

A database that contains a series of related tables that can be connected via their relationship.

for Rel/PN: one or more fields must exist inside both tables.

Primary Key ↗

A identifier that references a column in which each value is unique.



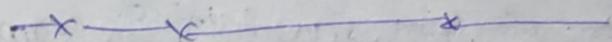
Unique, not null ^{or} blank

Foreign key ↗

A foreign key is a field within a table that's a primary key in another table.

Simple,

↗ table can have only 1 Primary key but many foreign keys.



Normalization ↗

Process of organizing RDB. It is applied to eliminate data redundancy, increase data integrity and reduce complexity in a database.

- ↗ Some table's don't require PK.
- ↗ A pk can also be constructed using multiple columns of a table its called composite key.



meta data Example:

① Photo

② mail

"Metadata is used in database management to help data analysts interpret the contents of the data within the database."

3 Types of metadata:-

(1) Descriptive metadata:

metadata that describes a piece of data and can be used to identify it at a later point in time.

Ex. descriptive metadata of book: Title, author and Unique International Standard Book Number ISBN

(2) Structural metadata:-

metadata that indicates how a piece of data is organized and whether it is part of one, or more than one, data collection.

Ex. Contents page of book, digital document → printed one

This type of metadata also keeps track of relationship between 2 things.

(3) Administrative metadata:-

metadata that indicates the technical sources of a digital asset.

Ex. photo information

* Using metadata as an analyst →

Putting data into context is probably most valuable thing that metadata does.

Benefits of using metadata →

single
version
of truth

→ metadata creates a single source of truth by keeping things consistent and uniform.

Reliable > metadata also makes data more reliable by making sure it's accurate, precise, relevant and timely.

~~metadata repository~~ Data Analyst use Metadata repository to make sure their data is consistent & reliable →

⇒ A database specifically created to store metadata

Metadata repository can be stored in a physical location or can be virtual (cloud)

→ metadata repository makes it easier to and faster to bring together multiple sources for data analysis.

- Describes state & location of metadata
- describes structure of tables inside
- describes how data flows through repository
- keeps track of who accesses metadata & when

* Metadata Management ↗

Metadata is stored in a single, central location and gives the company standardized information about all of its data.

Data Governance ↗

A process to ensure the formal management of a company's data assets.

↳ Gives better control of data & helps company manage issues related to data security, privacy, integrity and usability & internal-external data flows.

* Metadata Analyst *

Megan: measurement lead

"Metadata is kind of a shorthand or a CliffNotes version of a much more complex set of information."

Data Lake

>Data sources that may overlap & have something in common