# Using databases for social scientists

Damian Trilling

d.c.trilling@uva.nl
@damian0604
www.damiantrilling.net

Afdeling Communicatiewetenschap
Universiteit van Amsterdam

Computational Social Science Amsterdam
30 November 2018

## Today

Günther, Elisabeth; Trilling, Damian; van de Velde, Bob: But how do we store it? Data architecture in the social-scientific research process. In: Stuetzer, C.M. (Hrsg.); Welker, M. (Hrsg.); Egger, M. (Hrsg.): Computational social science in the age of Big Data. Concepts, methodologies, tools, and applications. Cologne: Herbert von Halem, 2018, pp. 161–187

ADD LINK TO PDF

When and why databases?

## The "traditional" approach

### Example: Analysis of a couple of thousands articles/speeches/reports/...

- Store as seperate `.txt` files

If metadata (beyond what can be inferred from filename and location) are important

- Store as tabular dataset (`.csv` or proprietary format)
- (possibly: store as seperate `.json` files)

## The "traditional" approach

### Example: Analysis of a couple of thousands articles/speeches/reports/. . .

- Store as seperate .txt files

If metadata (beyond what can be inferred from filename and location) are important

- Store as tabular dataset (.csv or proprietary format)
- (possibly: store as seperate .json files)

## The "traditional" approach

### Example: Analysis of a couple of thousands articles/speeches/reports/...

- Store as seperate `.txt` files

If metadata (beyond what can be inferred from filename and location) are important

- Store as tabular dataset (`.csv` or proprietary format)
- (possibly: store as seperate `.json` files)

## The "traditional" approach

> ### Example: Analysis of a couple of thousands articles/speeches/reports/...
>
> - Store as seperate `.txt` files
>
> If metadata (beyond what can be inferred from filename and location) are important
>
> - Store as tabular dataset (`.csv` or proprietary format)
> - (possibly: store as seperate `.json` files)

## The "traditional" approach

### pro

- easy to understand
- no dependencies
- works on all platforms, also in the future → good fallback/backup option

### con

- inefficient
- requires loading whole dataset into memory *or* reading all files from disk to query/aggregate/etc.
- requires *you* to deal with file I/O

...

# The "traditional" approach

## pro

- easy to understand
- no dependencies
- works on all platforms, also in the future → good fallback/backup option

## con

- inefficient
- requires loading whole dataset into memory *or* reading all files from disk to query/aggregate/etc.
- requires *you* to deal with file I/O

...

When and why databases?   Data architecture   Database types   MongoDB and Elastic Search   Practical example
○○○●○○                     ○○○                 ○○○○○○○○○        ○○○○○                        ○○
Because you have to

Databases: because you have to

- Your data is too big to fit in RAM (and you need to do some querying)
- Because using files would be prohibitively inefficient
- Because you need to scale horizontally

When and why databases?   Data architecture   Database types   MongoDB and Elastic Search   Practical example
○○○●○○                     ○○○              ○○○○○○○○○        ○○○○○                         ○○
Because you have to

## Databases: because you have to

- Your data is too big to fit in RAM (and you need to do some querying)
- Because using files would be prohibitively inefficient
- Because you need to scale horizontally

When and why databases?    Data architecture    Database types    MongoDB and Elastic Search    Practical example
○○○●○○                     ○○○                 ○○○○○○○○○         ○○○○○                         ○○
Because you have to

Databases: because you have to

- Your data is too big to fit in RAM (and you need to do some querying)
- Because using files would be prohibitively inefficient
- Because you need to scale horizontally

When and why databases?   Data architecture   Database types   MongoDB and Elastic Search   Practical example
○○○○●○                     ○○○                000000000         ○○○○○                      ○○
Because you want to

## Databases: because you want to

- You do not want to take care of I/O yourself

- Because it allows you to do better searches and queries

- Because you can easily aggregate your data

- Because you can easily join/merge data

- Because you want to be able to easily modify/update records
  without rewriting this whole CSV table

When and why databases?   Data architecture   Database types   MongoDB and Elastic Search   Practical example
○○○○●○                    ○○○            ○○○○○○○○○     ○○○○○                       ○○
Because you want to

## Databases: because you want to

- You do not want to take care of I/O yourself

- Because it allows you to do better searches and queries

- Because you can easily aggregate your data

- Because you can easily join/merge data

- Because you want to be able to easily modify/update records without rewriting this whole CSV table

When and why databases?    Data architecture    Database types    MongoDB and Elastic Search    Practical example
○○○○●○                     ○○○              ○○○○○○○○○         ○○○○○                          ○○
Because you want to

## Databases: because you want to

- You do not want to take care of I/O yourself
- Because it allows you to do better searches and queries
- Because you can easily aggregate your data
- Because you can easily join/merge data
- Because you want to be able to easily modify/update records without rewriting this whole CSV table

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| oooo●o | ooo | ooooooooo | ooooo | oo |

Because you want to

## Databases: because you want to

- You do not want to take care of I/O yourself
- Because it allows you to do better searches and queries
- Because you can easily aggregate your data
- Because you can easily join/merge data
- Because you want to be able to easily modify/update records without rewriting this whole CSV table

When and why databases?  Data architecture  Database types  MongoDB and Elastic Search  Practical example
○○○○●○                    ○○○          ○○○○○○○○○       ○○○○○                      ○○
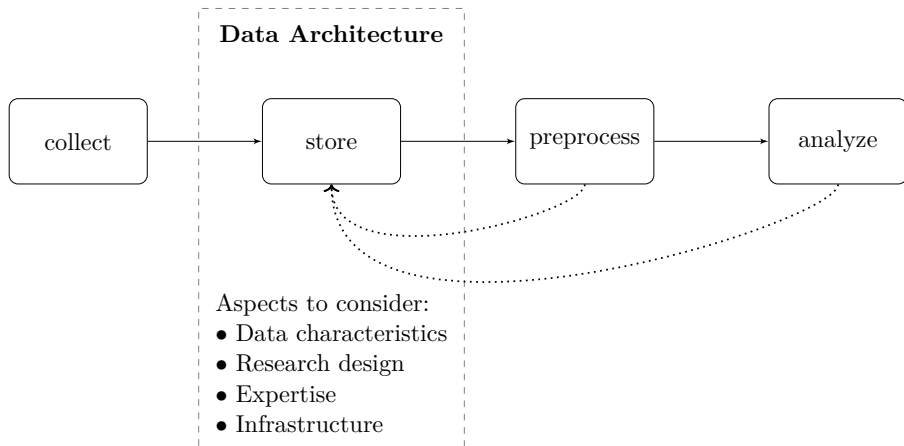Because you want to

Databases: because you want to

- You do not want to take care of I/O yourself
- Because it allows you to do better searches and queries
- Because you can easily aggregate your data
- Because you can easily join/merge data
- Because you want to be able to easily modify/update records without rewriting this whole CSV table

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| ○○○○○● | ○○○ | ○○○○○○○○○ | ○○○○○ | ○○ |

Considerations

I *Do you need pre-processing / cleaning?*

*What's your input?* II

a) Technical properties (data types)
b) Source properties (static vs. dynamic)
c) Corpus properties (quantity, messiness homogeneity)

Data

Research design

Data Architecture

Expertise

Infra-structure

III *What's your goal?*

a) Project setup (set up an archive vs. specific research interest)
b) Project priorities (availability, consistency, performance)
c) Scientific standards (accessibility, sustainability)

*What are your skills?* IV

a) Previous experiences
b) External support (collaboration vs. contracting)
c) Future demands

V *What else do you use?*

a) Software compatibility (technical interface)
b) Hardware scalability (vertical vs. horizontal)

Data architecture

## Data architecture

- File formats
- Linkage
- Internal structure

**Data Architecture**

collect → store → preprocess → analyze

Aspects to consider:
• Data characteristics
• Research design
• Expertise
• Infrastructure

Database types

Database types

When and why databases?    Data architecture    **Database types**    MongoDB and Elastic Search    Practical example

Relational databases

## Relational databases

- Tabular data structure

- Relational: tables are linked by keys

- Well-defined data types per column

- Good at joining and aggregating

When and why databases?   Data architecture   **Database types**   MongoDB and Elastic Search   Practical example
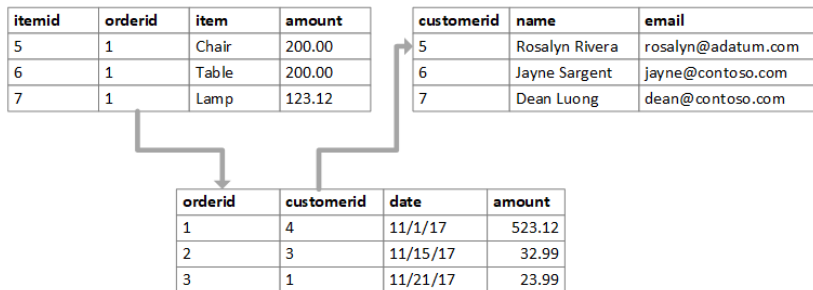
Relational databases

## Relational databases

- Tabular data structure
- Relational: tables are linked by keys
- Well-defined data types per column
- Good at joining and aggregating

## Relational databases

- Tabular data structure
- Relational: tables are linked by keys
- Well-defined data types per column
- Good at joining and aggregating

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| ○○○○○○ | ○○○ | ○○●○○○○○○ | ○○○○○ | ○○ |

Relational databases

## Relational databases

- Tabular data structure
- Relational: tables are linked by keys
- Well-defined data types per column
- Good at joining and aggregating

When and why databases?
○○○○○○

Data architecture
○○○

Database types
○○○●○○○○○

MongoDB and Elastic Search
○○○○○

Practical example
○○

Relational databases

| itemid | orderid | item  | amount |
|--------|---------|-------|--------|
| 5      | 1       | Chair | 200.00 |
| 6      | 1       | Table | 200.00 |
| 7      | 1       | Lamp  | 123.12 |

| customerid | name           | email              |
|------------|----------------|--------------------|
| 5          | Rosalyn Rivera | rosalyn@adatum.com |
| 6          | Jayne Sargent  | jayne@contoso.com  |
| 7          | Dean Luong     | dean@contoso.com   |

| orderid | customerid | date     | amount |
|---------|------------|----------|--------|
| 1       | 4          | 11/1/17  | 523.12 |
| 2       | 3          | 11/15/17 | 32.99  |
| 3       | 1          | 11/21/17 | 23.99  |

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| 000000 | 000 | 000000000 | 00000 | 00 |

Relational databases

## Relational databases and text

- You can store text

- In fact, often used as backend for for instance tweets and sometimes news articles

- BUT (1): Not optimized for searching in text

- BUT (2): Not good for messy data

- BUT (3): Hard to add extra columns or change specifications of existing ones

# Relational databases and text

- You can store text
- In fact, often used as backend for for instance tweets and sometimes news articles
- BUT (1): Not optimized for searching in text
- BUT (2): Not good for messy data
- BUT (3): Hard to add extra columns or change specifications of existing ones

# Relational databases and text

- You can store text
- In fact, often used as backend for for instance tweets and sometimes news articles
- BUT (1): Not optimized for searching in text
- BUT (2): Not good for messy data
- BUT (3): Hard to add extra columns or change specifications of existing ones

| When and why databases? | Data architecture | **Database types** | MongoDB and Elastic Search | Practical example |
| 000000 | 000 | 000●00000 | 00000 | 00 |

Relational databases

# Relational databases and text

- You can store text
- In fact, often used as backend for for instance tweets and sometimes news articles
- BUT (1): Not optimized for searching in text
- BUT (2): Not good for messy data
- BUT (3): Hard to add extra columns or change specifications of existing ones

| When and why databases? | Data architecture | **Database types** | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| 000000 | 000 | 0000●0000 | 00000 | 00 |

Relational databases

# Relational databases and text

- You can store text
- In fact, often used as backend for for instance tweets and sometimes news articles
- BUT (1): Not optimized for searching in text
- BUT (2): Not good for messy data
- BUT (3): Hard to add extra columns or change specifications of existing ones

# NoSQL databases

- optimized for messy data

- can be schema-free $\Rightarrow$ we do not have to enforce a specific format at insertation time

- new entries do not necessarily have to follow the specifications of old ones

- Allow you to just throw in an arbitrary JSON object

- CAP-theorem: we trade consistency for availability and performance

# NoSQL databases

- optimized for messy data
- can be schema-free $\Rightarrow$ we do not have to enforce a specific format at insertion time
- new entries do not necessarily have to follow the specifications of old ones
- Allow you to just throw in an arbitrary JSON object
- CAP-theorem: we trade consistency for availability and performance

When and why databases?      Data architecture      **Database types**      MongoDB and Elastic Search      Practical example
000000                 000              000000●000            00000                00

NoSQL databases

# NoSQL databases

- optimized for messy data
- can be schema-free ⇒ we do not have to enforce a specific format at insertion time
- new entries do not necessarily have to follow the specifications of old ones
- Allow you to just throw in an arbitrary JSON object
- CAP-theorem: we trade consistency for availability and performance

When and why databases?    Data architecture    **Database types**    MongoDB and Elastic Search    Practical example
000000        000       00000●000      00000        00

NoSQL databases

# NoSQL databases

- optimized for messy data
- can be schema-free ⇒ we do not have to enforce a specific format at insertion time
- new entries do not necessarily have to follow the specifications of old ones
- Allow you to just throw in an arbitrary JSON object
- CAP-theorem: we trade consistency for availability and performance

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| 000000 | 000 | 0000000000 | 00000 | 00 |

NoSQL databases

# NoSQL databases

- optimized for messy data
- can be schema-free $\Rightarrow$ we do not have to enforce a specific format at insertion time
- new entries do not necessarily have to follow the specifications of old ones
- Allow you to just throw in an arbitrary JSON object
- CAP-theorem: we trade consistency for availability and performance

When and why databases?   Data architecture   **Database types**   MongoDB and Elastic Search   Practical example
oooooo                    ooo                ooooooo●oo          ooooo                        oo
NoSQL databases

We can store retrieved web pages (1) together with some first roughly parsed extracted data (2), and do some cleaning and enrichment (3,4) later.

| 1. Retrieved | 2. Structured | 3. Cleaned | 4. Enriched |
|---|---|---|---|
| `<div>`<br>`<span id='author'`<br>`class='big'>`<br>`Author: John`<br>`Doe </span>`<br>`<span id='viewed'>`<br>`seen 42 times`<br>`</span>`<br>`</div>` | `{`<br>`"author":`<br>`" Author: John Doe ",`<br>`"viewed":`<br>`" seen 42 times "`<br>`}` | `{`<br>`"author":`<br>`" john doe ",`<br>`"viewed":`<br>`42 ,`<br>`}` | `{`<br>`"author":`<br>`" john doe ",`<br>`"author_gender":`<br>`" M ",`<br>`"viewed":`<br>`42 ,`<br>`}` |

# NoSQL databases and text

- Often optimized for indexing text

- Internal preprocessing ("analysis") under the hood: e.g., search based on stemmed text

- Nested data: e.g., comments within articles

- Texts can have additional keys that others don't (e.g., online news have urls, offline news have page numbers; some may have)

- You may have even different language versions of the same text (and all having analyzed accordingly)

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| oooooo | ooo | ooooooooeo | ooooo | oo |

NoSQL databases

# NoSQL databases and text

- Often optimized for indexing text
- Internal preprocessing ("analysis") under the hood: e.g., search based on stemmed text
- Nested data: e.g., comments within articles
- Texts can have additional keys that others don't (e.g., online news have urls, offline news have page numbers; some may have)
- You may have even different language versions of the same text (and all having analyzed accordingly)

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| 000000 | 000 | 000000000●0 | 00000 | 00 |

NoSQL databases

# NoSQL databases and text

- Often optimized for indexing text
- Internal preprocessing ("analysis") under the hood: e.g., search based on stemmed text
- Nested data: e.g., comments within articles
- Texts can have additional keys that others don't (e.g., online news have urls, offline news have page numbers; some may have)
- You may have even different language versions of the same text (and all having analyzed accordingly)

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
| 000000 | 000 | 000000000●0 | 00000 | 00 |

NoSQL databases

# NoSQL databases and text

- Often optimized for indexing text
- Internal preprocessing ("analysis") under the hood: e.g., search based on stemmed text
- Nested data: e.g., comments within articles
- Texts can have additional keys that others don't (e.g., online news have urls, offline news have page numbers; some may have)
- You may have even different language versions of the same text (and all having analyzed accordingly)

| When and why databases? | Data architecture | Database types | MongoDB and Elastic Search | Practical example |
|---|---|---|---|---|
| ○○○○○○ | ○○○ | ○○○○○○○●○ | ○○○○○ | ○○ |

NoSQL databases

# NoSQL databases and text

- Often optimized for indexing text
- Internal preprocessing ("analysis") under the hood: e.g., search based on stemmed text
- Nested data: e.g., comments within articles
- Texts can have additional keys that others don't (e.g., online news have urls, offline news have page numbers; some may have)
- You may have even different language versions of the same text (and all having analyzed accordingly)

When and why databases?    Data architecture    **Database types**    MongoDB and Elastic Search    Practical example
○○○○○○      ○○○      ○○○○○○○○●      ○○○○○      ○○

NoSQL databases

# SQL vs NoSQL in one slide

(for storing (textual) social-scientific data)

## SQL

- strucutre known in advance

## NoSQL

- full-text search may be relevant

MongoDB and Elastic Search

## Use case

- Scrape and store articles ($\approx$ 20M)
- We first used Mongo and later switched to ES (because of better performance for full-text search)

Use case

- Scrape and store articles ($\approx$ 20M)
- We first used Mongo and later switched to ES (because of better performance for full-text search)

# Interacting with ES via http requests

```
packer-ubuntu-16:~$ curl http://localhost:9200/inca6/_count?pretty
{
  "count" : 19054530,
  "_shards" : {
    "total" : 5,
    "successful" : 5,
    "skipped" : 0,
    "failed" : 0
  }
}
```

# Interacting with ES via Kibana

# Interacting with ES via Python

```
In [1]: from elasticsearch import Elasticsearch

In [2]: client = Elasticsearch()

In [3]: client.count('inca6')
Out[3]:
{'_shards': {'failed': 0, 'skipped': 0, 'successful': 5, 'total': 5},
 'count': 19054530}
```

Practical example (Jupyter Notebook)

Questions?

d.c.trilling@uva.nl
@damian0604
www.damiantrilling.net