

Applied Data Science - Notes

Dom Hutchinson

March 5, 2021

Contents

1	Introduction	2
2	Data Ingress & Pre-Processing	3
2.1	Data Structures	3
2.2	Data Formats	3
2.3	Web-Scraping & APIs	5
3	Privacy	6
3.1	IRL Examples	6
3.2	Security & Statistical Database	7
3.3	k -Anonymity	9
3.4	Other Approaches	10
4	Storage & Management	11
4.1	Storage Technologies	11
4.2	Databases	12
4.2.1	Graph Databases	15
5	Transformation & Integration	17
5.1	Data Wrangling	17
5.2	Data Fusion	19
6	Exploration & Visualisation	21
6.1	Data Transformation	23
6.2	Density Estimation	24
6.3	Multi-Variate Analysis	24
7	Data Science in Production	27
7.1	Failure is the Rule	27
7.2	Engineering Guarantees for Data Science	28
8	Sharing & Privacy	30

1 Introduction

Remark 1.1 - *Types of Data*

Data comes in many forms including, but not limited to, the following

- Dense & Sparse data.
- Structured/Relational Data.
- Numerical; Categorical; Ordinal; or Boolean.
- Text (Emails, Tweets, Articles).
- Records (User-Level Data, Timestamped Event Data, Log Files).
- Geo-Based Location Data.
- Data-Time Data.
- Network Data.
- Sensor Data.
- Images and Video.
- Audio and Music.

Remark 1.2 - *Big & Small Data*

Whether a dataset is big or small depends on the computational-resources available, and thus will vary over time. Here are some ways to evaluate this

	Big Data	Small Data
<i>Data Condition</i>	Always unstructured, not read for analysis, many relational database tables that need to be merged.	Ready for analysis, flat file, no need to merge tables.
<i>Location</i>	Cloud, offshore, SQL server etc.	Database, local PC.
<i>Data Size</i>	Over 50k variables, over 50k individuals, random samples, unstructured	File that is in a spreadsheet, that can be viewed on a few sheets of paper.
<i>Data Purpose</i>	No intended purpose.	Intended purpose for data collection.

Remark 1.3 - *What is Data Science?*

Data-Driven Science. An interdisciplinary field about scientific processes and systems to extract knowledge or insights from data in various forms.

Data science incorporates fields from: Mathematics, Computer Science; &, Domain Expertise.

Remark 1.4 - *Motivating Applications*

Data science is motivated by its applications. Here are some examples of such applications

- *Amazon* use recommender systems to suggest products to customers.
- *Energy Companies* use data science to try and predict future usage of customers, so that resources can be applied efficiently.

- *Agriculture* use sensors in fields to collect data in order to monitor crops and predict weather.
- *Healthcare* use sensors in homes to monitor the health of people over long periods of time (especially when the person cannot go to the hospital).

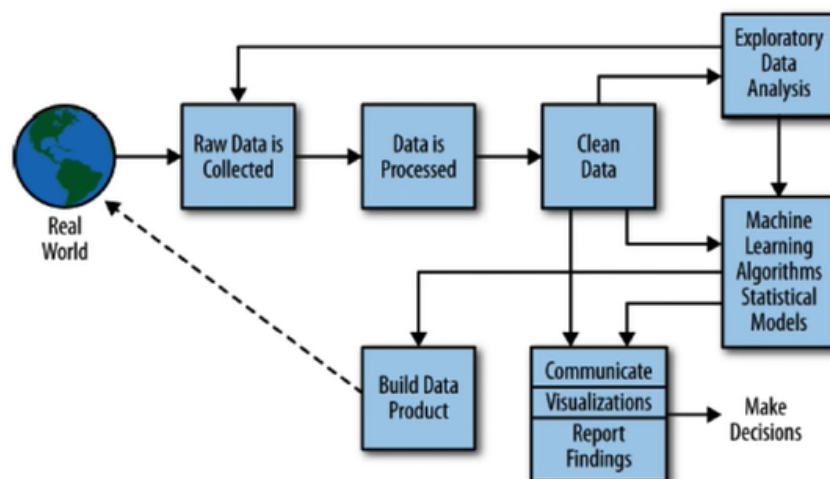


Figure 1: The pipeline for approaching problems in data science.

Proposition 1.1 - Data Science Pipeline

See Figure 1 for a pipeline for approaching data science problems.

2 Data Ingress & Pre-Processing

2.1 Data Structures

Proposition 2.1 - Native Python Data Structures

Here are some data structures which are native to python and are popular in data science

- `list` - List of elements of varying types.
- `set` - List of unique elements of varying types.
- `dict` - Key-Value pairings.

Proposition 2.2 - Non-Native Python Data Structures

Here are some data structures which are not native to python but are popular in data science.

- `np.array`.
- `pandas.DataFrame`.

2.2 Data Formats

Definition 2.1 - Object Persistence

Object Persistence is the process of ensuring that the objects which are created are kept through multiple sessions. This comes in two stages

- i). *Serialisation* - Translating data structures or objects from memory into a format which can be stored.
- ii). *Deserialisation* - The inverse. Translating data structures which have been stored, into memory.

Remark 2.1 - *Bespoke Serialisation & Deserialisation*

Bespoke serialisation and deserialisation methods can be crafted manually. (e.g. Instantiating an output file; Writing each element of a list to a different line in the file; Closing the file.)

However, there are limitations to bespoke methods:

- Methods are specific to each use case (not standardised).
- Methods may not be robust.
- Methods require testing against many test cases.
- Object metadata is not encoded.

These limitations are rarely a problem in very controlled environments.

Definition 2.2 - *Comma-Separated Values (CSV)*

Comma-Separated Values (CSV) files are well suited to tabular data (inc. matrices). Each line of a *CSV File* stores one row of the table, and each element in a row is separated by a comma.

CSV Files are generally very readable, as well as time- and space-efficient for tabular data.

Remark 2.2 - *Using CSV Files*

As *CSV Files* are very popular, there are many library methods which can interact with them. Including reading & writing to and from memory (e.g. `pandas.read_csv`).

Definition 2.3 - *JavaScript Object Notation (JSON)*

JavaScript Object Notation^[1] (JSON) is a standardised syntax for storing & exchanging data, used for serialisation. JSON uses text files which are human-readable and `dict`-like (i.e. have value-key pairs). JSON is designed to be simple so is a very robust language & suits many purposes.^[2]

Limitations of JSON are that a specific conversion process may be required to convert a JSON file into objects in memory, and JSON files can become large due to key repetition.

Definition 2.4 - *Hierarchical Data Format (HDF5)*

Hierarchical Data Format (HDF5) is a standardised format for serialisation. HDF5 is a binary format and tries to mimic file system-like access. HDF5 files have the following three components

- i). *Datasets* - Array-like collections of data. Thousands of datasets can be stored in a single file, and can be categorised and tagged however you want.
- ii). *Groups* - Folder-like structures which groups datasets & other groups.
- iii). *Metadata* - Information which pertains to all the datasets. (e.g. author, edit data & version).

^[1]JSON is not just compatible with JavaScript and is common for many languages & APIs

^[2]<https://jsonlint.com/> is a useful site for JSON validation.

HDF5 files are ideal for large numerical data sets, and can easily be manipulated by `numpy`. HDF5 files support a variety of transparent storage features (inc. compression, error-detection and chunked I/O), which `numpy.array` do not.^[3]

Proposition 2.3 - *HDF5 files vs Traditional File Systems*

There are a few key differences between *HDF5 Files* and *Traditional File Systems*.

- i). *HDF5 Files* are portable, as the entire structure is contained in the file independent of the underlying file system.^[4]
- ii). Datasets in *HDF5 Files* are all homogeneous hyper-rectangular numerical arrays, whereas files in traditional file system can be anything.
- iii). Metadata can be added to groups in *HDF5 Files*. This is not possible in traditional file systems.

Remark 2.3 - *Other Standardised Serialisation Method*

XML, *Protocol Buffers* & YAML are other popular standardised serialisation methods.

2.3 Web-Scraping & APIs

Remark 2.4 - *Terms of Use*

Web Scraping should be done within the website's terms of use.

Definition 2.5 - *Web Scraping*

Web Scraping is the practice of collecting data from websites. This can take many forms, but typically involves taking a raw webpage and parsing the desired data.

Proposition 2.4 - *Approaching Web Scraping*

To perform *Web Scraping* successfully, you need a good idea of how a webpage is structure. Typically webpages are based around a `html` file, which are well structured^[5]. Identifying combinations of tags, classes & ids in the `html` file can help locate the desired data.

It is harder, sometimes impossible, to navigate poorly designed websites as they are less structured and inconsistent.

Remark 2.5 - *Tools for Web Scraping*

Web Scraping technologies need to be tolerant to several artefacts of real-world data (known as "*wrangling*") as-well-as errors in the website.

Some popular tools for *Web Scraping* are

- BeautifulSoup - A python library for parsing XML & HTML.
- scrapy - A python library. Generally faster than BeautifulSoup.
- Selenium - A web-browser plugin, generally used to test web services.

^[3]<http://docs.h5py.org/en/latest/quick.html> provides a quick-start guide to using HDF5 files in python.

^[4]However, it does depend on the HDF5 library.

^[5]Webpages are often interpreted to have a tree structure, with each tag being a node

Definition 2.6 - Web APIs

Web APIs greatly simplify *Web Scraping* by provide a portal for explicit data acquisition, and are generally less prone to the issues which arise when *Web Scraping*.^[6]

- The code running *Web APIs* is optimised for data requesting & retrieval. It does not waste time on visualisation or aesthetics. This means the bandwidth required for an *API* request is much lower than for a similiar *Web Scraping* process (As images etc. do not need to be loaded).
- *Web API* querying is robust, reliable, well maintained and documented with a static schema (HTML-based *Web Scraping* is not).
- *Web APIs* use standardised *Serialisation Tools* (e.g. JSON).
- *Web APIs* have already extracted and organised the desired data, however this does mean the user can only access what the operator will allow. This is much better than HTML-based *Web Scraping* where you rely upon fickle naming conventions of tags.

Definition 2.7 - RESTful APIs

Representational State Transfer APIs (REST/RESTful) are a popular form of *Web API* and generally require an *API Key* to access data.

Request to *RESTful APIs* generally involves constructing a URL which contains your keys and the parameters of your query.

Remark 2.6 - Regular Expression (RegEx)

Regular Expression (ReGex) queries are useful for extracting data, either while web scraping or from API requests.^[7] Python has the `re` library for *RegEx* queries, two popular methods from this library are

- i). `re.match` - Attempts to match a *RegEx* pattern to the whole string.
- ii). `re.search` - Searches for the first occurrence of a *RegEx* pattern in a string.

3 Privacy

3.1 IRL Examples

Remark 3.1 - AOL Data Incident, 2006

In 2006, *AOL Research* released to the public, for the purpose of research, a text file containing 20mn search terms from over 650k users, over a 3-month period.

The users were anonymised in the in the data, but personally identifiable information was present in many of the search terms.

AOL Research retracted the document shortly after publishing.

Remark 3.2 - Netflix Prize, 2007

In 2007, *Netflix* ran a competition where \$1mn was offered to anyone who was able to produce a system which outperformed their existing recommender system by at least 10%.

The training data offered by *Netflix* contained `{user,movie,date,grade}` and customer ids were anonymised.

^[6]See <https://github.com/public-apis/public-apis> for a categorised list of public APIs.

^[7]<http://pythex.org/> is a website which can perform *RegEx*.

It was later found that some users had their *Netflix* and *IMDb* accounts linked (so reviews were posted on both). This meant the training data could be partially de-anonymised. This led to lawsuits and *Netflix* cancelling future competitions.

Remark 3.3 - Pay Attention

The moral of these examples is that you have to be really careful and diligent when anonymising data as there may be non-obvious ways for people to de-anonymise it. (Often through little fault of your own).

3.2 Security & Statistical Database

Definition 3.1 - Statistical Database

A *Statistical Database* is a database for statistics. They can come in several forms

- *Tabular Data* - Tables with counts or magnitudes.
- *Queryable Databases* - On-line databases which accept statistical queries (e.g. min, max, sum etc.)
- *Microdata* - Files where each record contains information about an individual.

Remark 3.4 - Privacy Trade-Off

Statistical Databases have to trade-off privacy and functionality. Generally, more privacy means less functionality.

Definition 3.2 - Identifiers

Identifiers are attributes that unambiguously identify someone (e.g. passport number, NI number, name etc.)

Definition 3.3 - Quasi-Identifiers

Quasi-Identifiers are attributes which identify someone, but there is some ambiguity. Often, having multiple *Quasi-Identifiers* are enough to confidently identify someone. (e.g. address, age, gender).

Definition 3.4 - Confidential Attributes

Confidential Attributes are attributes which contain sensitive information about an individual (e.g. salary, religion, medical conditions etc.)

Definition 3.5 - Non-Confidential Attributes

Non-Confidential Attributes are attributes which do not contain sensitive information about an individual.

Remark 3.5 - Anonymisation, Identifiers & Quasi-Identifiers

When anonymising a data set:

- *Identifiers* need to be suppressed from the data sets (i.e. removed).
- *Quasi-Identifiers* hold disclosure risk as they can often not be suppress due their high analytical value, but can often be linked with other non-anonymised datasets to de-anonymise your dataset.

Definition 3.6 - Attribute Disclosure

Attribute Disclosure is when the value of a confidential attribute is made available, and in doing so the individual is easier to identifier (compared to the value not being disclosed).

Definition 3.7 - Identity Disclosure

Identity Disclosure is when an entire record in an anonymised data set can be linked with that individual's identity.

Definition 3.8 - Membership Disclosure

Membership Disclosure is whether or not data about an individual is contained in a dataset.

Definition 3.9 - External Attack on a Table

An *External Attack* on a table occurs when another dataset is released which contains data which comprises more table, or can be combined with your table to comprise others.

Definition 3.10 - Internal Attack on a Table

An *Internal Attack* on a table occurs when those within the table (i.e. with access to it and know which data represents them) are able to deduce which data belongs to other people by a process of elimination.

Only really possible when the table has few respondents.

Definition 3.11 - Dominance Attack on a Table

An *Internal Attack* on a table occurs a (or a few) respondents dominate the contribution in a particular cell of a magnitude table, these *Dominant Respondents* are then able to upper-bound the contributions from the rest due it being easier to differentiate other users from them.

Proposition 3.1 - Combatting Table Attacks

There are two main approaches to combatting *Table Attacks*

- i). *Non-Perturbative* techniques which do not modify the values in the cells, but they may suppress or recode them (e.g. cell suppression, recoding of categorical attributes).
- ii). *Perturbative* techniques do modify the values in the cells (e.g. controlled rounding).

Proposition 3.2 - Combatting Database Attacks

There are three main approaches to combatting attacks on databases, all focused around adjusting/controlling queries

- i). *Query Input Perturbation* - Perturbation is applied to records on which queries are computed.
- ii). *Query Output Perturbation* - Perturbation is applied to the results of a query.
- iii). *Query Restriction* - The database refuses certain queries.
- iv). *Camouflage* - Returns answers which have been made non-exact in some deterministic fashion.

Remark 3.6 - Differential Attack^[8]

^[8]Not proper name.

Suppose you know that someone is going to be changing doctors soon (and thus will be removed from the doctors database), if you copied the database before and after this person moves you could likely deduce details about the individual by comparing the two versions.

How do we prevent this?

Definition 3.12 - Differential Privacy

Consider a learner implements a summary statistic A (e.g. sum of all elements); an adversary proposes two datasets S and S' which differ by only one row; and a test set Q . The learner wants the summary statistic to have the same value if only one row has changed.

Summary statistic A is ε -Differentially Private iff

$$\left| \ln \left(\frac{\mathbb{P}(A(S) \in Q)}{\mathbb{P}(A(S') \in Q)} \right) \right| \leq \varepsilon$$

Note that *Differential Privacy* is a condition on the release mechanism of the data A , not on the data S, S' itself.

Remark 3.7 - Differential Privacy in Practice

In practice we avoid differential privacy issues by adding noise to our responses (either to all responses, or incorrectly responding to some queries).

3.3 k -Anonymity

Definition 3.13 - k -Anonymity

A dataset satisfies k -Anonymity if each combination of values of *Quasi-Identifier Attributes* in the table are shared by at least k records.

Measuring k -Anonymity requires a brute force approach which compares every pair of records in the table.

Proposition 3.3 - Achieving k -Anonymity

There are two main approaches to achieving k -Anonymity:

- *Suppression* - Replace the values of attributes with an asterisk.
- *Generalisation* - Values of attributes are replaced with a broader category. (e.g. age-ranges rather than exact age)

Proposition 3.4 - k -Anonymity Attacks

k -Anonymity does not protect against attribute disclosure attacks in general if the values of *Confidential Attributes* are very similar in the group of k records who share *Quasi-Identifier* values.

Here are two specific attacks that can happen

- *Homogeneity Attack* - If all records in the set of k identical records have the same *Sensitive Value*, then the value can likely be predicted (and would comprise all k).
- *Background Knowledge Attack* - If associations exist between one or more *Quasi-Identifier Attributes* and a *Sensitive Attribute* then the set of possible values of the *Sensitive Attribute* is reduced.

Definition 3.14 - l -Diversity

l -Diversity is an extension of k -Anonymity.

l -Diversity requires that the values of all *Confidential Attributes* within any group of k -anonymous records contain at least l clearly distinct values.

Proposition 3.5 - Limitations of l -Diversity

Here are two limitations of l -Diversity

- Not every value shows equal sensitivity. (A rare positive indicator for a disease may provide more information than a common negative indicator).
- It ensures diversity of sensitive values in each group, but it does not recognise that values may be semantically similar.

Definition 3.15 - t -Closeness

t -Closeness requires that the distribution of the confidential attribute within a group of k records is similar to the distribution of the confidential attribute in the entire data set (at most distance t between the two distributions).

Remark 3.8 - Tools for Anonymisation

Here are three popular tools for anonymisation

- i). *ARX* - Applies k -anonymity, l -diversity and t -closeness in Java.
- ii). *Argus* - Creates safe microdata files.
- iii). *SdcMicro* - Provides disclosure control methods for anonymisation and risk-estimation.

3.4 Other Approaches

Definition 3.16 - Privacy-Preserving Data Mining, PPDM

Privacy-Preserving Data Mining (PPDM) seeks privacy for the data's owner when several owners wish to co-operate without giving away their data to each other

Definition 3.17 - Private Information Retrieval

Private Information Retrieval (PIR) seeks user privacy to allow the user of a database to retrieve information without the database knowing which item was retrieved.

4 Storage & Management

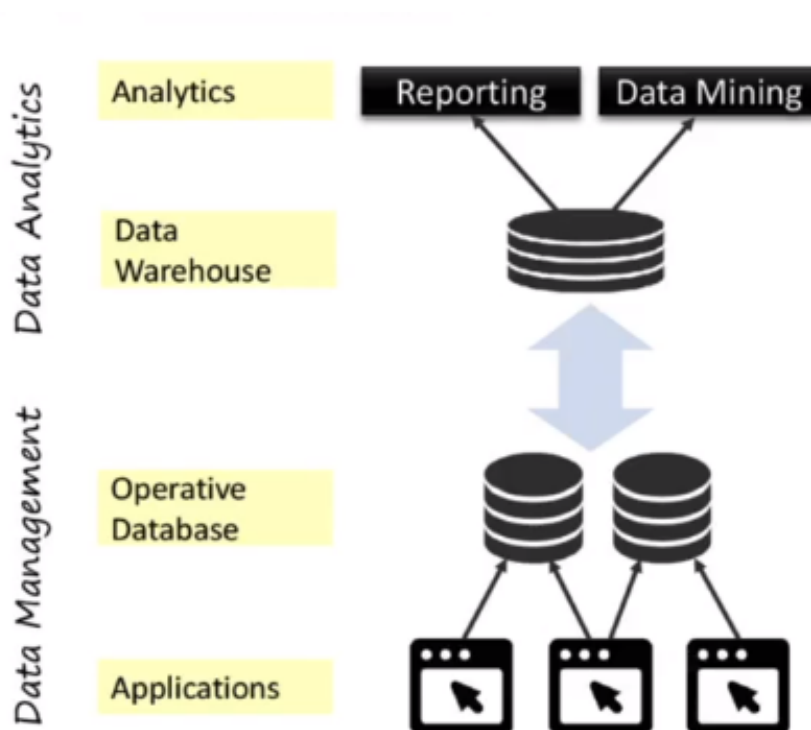


Figure 2: Typical data architecture.

Proposition 4.1 - Typical Data Architecture

Organisations tend to have two main data stores

- i). A *Data Warehouse* which contains copies of historic data. Data analysis is performed on this data as it is not critical to the operation of the organisation. *Data Warehouses* usually allow data analysts read-only access.
- ii). Several *Operative Databases* which the organisation's core products use and rely on. Data analysis should not be performed on these databases directly as it could interfere with customers (or worse, break the company's products).

See Figure 2 for an illustration.

Definition 4.1 - Datalake

TODO

4.1 Storage Technologies

Definition 4.2 - BLOB Storage

A *Binary Large Object* (BLOB) storage is typically used for large, unstructured files (e.g. images & backup files).

Due to being unstructured, *BLOB* files are not quick to query and are best used just as a record of the data.

BLOB files are generally accessed using a REST-API or an SDK, they are not mountable as a disk and drivers need to be build into applications in order for the app to access the *BLOB Storage*. *Spark* is a tool used by many datalakes for *BLOB* storage.

Data storage providers typically offer two tiers of *BLOB* storage

- i). *Hot Access* - Data is accessed frequently.
- ii). *Cold Access* - Data is not accessed very often. (Cheaper, but less availability).

Definition 4.3 - Disk Storage

Disk Storage^[9] are an alternative form of storage, but are essentially a facade for *BLOB-backed storage* and tend to be up to 1TB in size.

Disk Storage can be mounted to a virtual machine. Once mounted and formatted, the file is accessible like any other local file, but with greater throughput than standard-file-storage (60/sec per blob)^[10]

Disk Storage is ideal when you want to move applications which you were using locally into cloud storage, and you want the data to be decoupled from the virtual machine^[11]. But, only one person can access a *Disk Storage* at a time.

Definition 4.4 - File Storage

File Storage^[12] uses bytes across a network and a using a protocol (e.g SMB), this creates the illusion of the data being a local files. This means several instances of an application can access the data at the same time.

File Storage can be upto 5TB each and has a maximum throughput of 60MB/sec across a share^[13]. Throughput is better in Local-Area networks.

File Storage is backed by *BLOB Storage* and thus is just as durable.

Capabilities such as concurrency and file locking can be introduced into *File Storage*, availability depends on the protocol being used.

4.2 Databases

Remark 4.1 - Database Classification

There are two features typically used to classify databases

- i). The model they use to store data.
- ii). Their “Consistency/Availability” trade-off.

^[9]AKA *Block Storage*.

^[10]Throughput is lower than *BLOB Storage*.

^[11]As, otherwise, if the machine was to die then you would loose your data too!

^[12]AKA *Network Attached Storage* (NAS).

^[13]Slower than BLOB & Disk

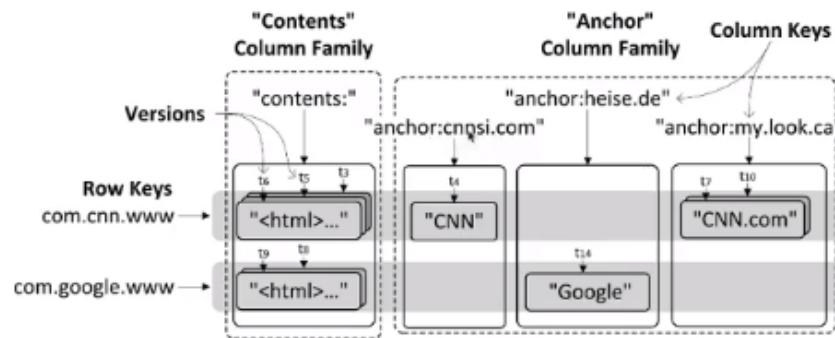


Figure 3: Example of data may be structured in a *Column Store*.

Definition 4.5 - Relational Database

TODO

Definition 4.6 - Database Models

Databases typically use one of the following five models to store & organise data.

- i). *Key-Value Stores* - A dictionary (keys map to values in a deterministic fashion). These databases only support get & put operations^[14], and all data is stored as a *BLOB*.
Key-Value Stores are *schemaless* meaning all assumptions about the structure of the data is implicitly encoded into the application log (and not explicitly defined by a data-definition language).
Key-Value Stores are easy to partition, query, have low latency and high throughput.
- ii). *Wide-Column Stores* - A distributed, multi-level, sorted map. This means that *Wide Column Stores* have several layers of “keys”: the first layer, called “Row Keys”, identify rows which contain a given *key-value pair*; The second layer is known as “column keys”.^[15] (Columns can be grouped into “Column Families”. Values in the same column family and same row are stored on the same disk).
Wide-Column Stores are space efficient as null values take zero space, and store data in a lexicographic order wrt their keys. Good compression rate allows for efficient querying of row subsets. Versions of each cell can be stored.
- iii). *Document Stores* - An implementation of *Key-Value Stores* but each value (known as a “document”) has to be in a semi-structured format (e.g. JSON). An entire document can be fetched using its unique key. *Document Stores* allow for more complex queries than *Key-Value Stores*, such as: retrieval of parts of documents; aggregation; query-by-example; and, full-text search.
- iv). *Graph Databases* - Store data sources in nodes and uses edges to define relationships between data sources. This is really good for complex relationships, but not very scalable as *Graph Partitioning* is NP-Hard.
- v). *Search Engines*^[16] - Essentially an inverted index, with some meta-data and query optimisation. Often will link from a value to the databases which contain that value.

^[14]There is no support for operations beyond simple CRUD (Create-Read-Update-Delete)

^[15]See Figure 3.

^[16]Not strictly a database

Remark 4.2 - Choosing Database Model

When choosing which model of database to use, we are generally trading between the models ability to handle “size” and “complexity”. Generally *Key-Value Stores* are best for the simplest datasets, then *Column Stores*, then *Document Stores* and finally, *Graph Databases* are typically best for handling complex datasets (but do not scale well).

A *Relational Database* is the default model to use for most projects, as they perform well for both size & complexity up to reasonably large projects.

Theorem 4.1 - CAP Theorem

The *CAP Theorem*^[17] states that a distributed system or database can guarantee at most two of the following three properties at any one time

- *Consistency* (C) - All users have the same view on the data at all times, with minimal *Latency*.
- *Availability* (A) - All users can always read and write to the database.
- *Partition-Tolerant* (P) - The system works well across physical network partitions (even if a machine fails or the cluster is split).

Proposition 4.2 - Database Consistency-Availability Trade-Off

Databases fulfil one of the following

- *Available & Partition Tolerant* (AP). (Never a relational database).
- *Consistent & Partition Tolerant* (CP). (Never a relational database).
- *Not Partition Tolerant* (CA).

Definition 4.7 - Sharding

Sharding specifies how data should be partitioning data across several nodes. This is an important technique in data storage as data is often stored in distributed networks.

Here are some popular approaches to *Sharding*

- *Hashing* - Apply a hash function to the data to determine which partition to place it in. Data-pieces which the same hash will be stored on the same machine. *Hashing* creates balanced partitions, but has no data locality.
- *Range* - Data-pieces are group according to whether their value falls into a given range. This allows for easier range search and sorting, but is not necessarily balanced.
- *Entity-Group* - Explicitly define how to co-locate data. This allows for easy access (equivalent to a single node), but you cannot change the partition after the fact.

When implementing these approaches we want the partitions to be as balanced as possible, for maximal efficiency.

Definition 4.8 - Replication

Replication is the practice of replicating data so that it is safe should a machine fail. The difficulty in *Replication* is *Synchronisation*.

Here are the two approaches for when to update data

^[17]Proved by Gilbert and Lynch in 2002.

- *Asynchronous* (lazy) - Data is only read on request, this means writing is fast but some replicas may not have been updated (and become stale).
- *Synchronous* (eager) - Data is consistent across users, but read-write is slower as the database need to co-ordinate data between users.

Here are the two approaches for where to store data

- *Primary Copy* - A single node accepts write requests and all other nodes just replicate this node. This is quick and easy to keep consistent, but creates a single point of failure.
- *Update Anywhere* - Any node can accept write requests, and will then propagate these requests to the rest of the network. This approach is fast to write, but requires co-ordination. This approach is more resilient to hardware failure

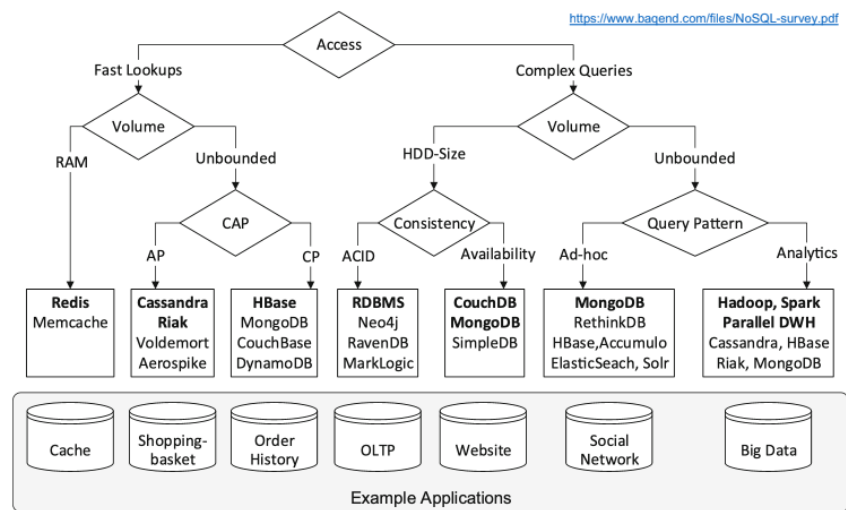


Figure 4: Flow chart to determine which data storage technology to use, given your requirements.

4.2.1 Graph Databases

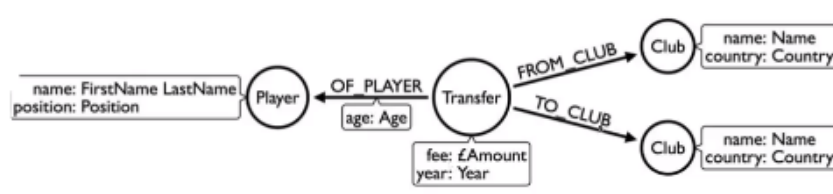


Figure 5: Example of a graph model for a database storing details about football transfers.

Definition 4.9 - Graph Model

A graph model can be used to represent data dependencies in a database

- *Nodes* are used to represent datafields or tables. Nodes are labelled with the fields which are unique to them.
- *Edges* are used to show which nodes depend on each other, often with directionality. Edges are labelled with the fields which are shared (ie connect) between two nodes.

See Figure 4.5 for an example.

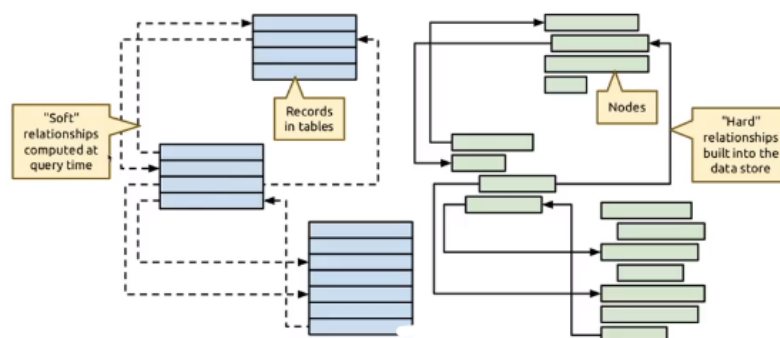


Figure 6: Comparison of the data-storage architectures for Relation Databases (Left) compared to Graph Databases (Right).

Proposition 4.3 - Relational vs Graph Database

Here are some comparisons between *Relational* and *Graph Databases*:

- In a *Relational Database* relationships are “soft” and only computed at query time.
- In a *Graph Database* relationships are “hard” as the relationships themselves are stored on the disk.

See Figure 4.6.

Proposition 4.4 - Neo4J

Neo4J is a piece of *Graph Database* management software.

Here is an example of a *Neo4J* query^[18]

```
MATCH (player:Player) WHERE player.name = 'Cristiano Ronaldo'
RETURN player;
```

In this query `player` is a variable we create, `Player` is the type of node that `player` is and `name` is a property of this variable.

We can before logical queries too. The following query returns people who have been a coach and a player.

```
MATCH (n) WHERE n:Player AND n:Coach RETURN n;
```

Proposition 4.5 - Traversals

Traversals are a feature which are unique to *Graph Databases*. *Traversals* define what relationship between several nodes to consider, and we can make queries on all these nodes^[19].

The query below fetches all transfers away from Man Utd^[20]

^[18]The equivalent SQL query is `SELECT * FROM players AS player WHERE player.name = 'Cristiano Ronaldo'`;

^[19]Similar to *Joins* in relational databases.

^[20]The equivalent SQL query is not as clean `SELECT from.name FROM transfers AS t JOIN clubs AS clubFrom ON t.from_club_id = clubFrom.id WHERE from.name="Manchester United"`;


```
MATCH (from:Club) <- [:FROM_CLUB]-(transfer:Transfer)
WHERE from.name="Manchester United"
```

We can extend this query to be all transfer from Man Utd to Real Madrid^[21]

```
MATCH (from:Club) <-
[:FROM_CLUB]-(transfer:Transfer)-[:TO_CLUB]->(to:Club)
WHERE from.name="Manchester United"
```

This query nicely shows how `<-[]-` and `-[]->` define the direction of data flow. We can extend this to find any path, including cycles.

Proposition 4.6 - Length of Relationship

Neo4j allows for variable length relationships. With `-[:__*n..]->` notation, where meaning a chain of at least `n` relationships.

The following statement defines a relationship where a player leaves a club and then after some number of transfers returns to the club.

```
(c:Club)<-[:f:FROM_CLUB]-(sale:Transfer)
-[:NEXT*1..]->(rebuy:Transfer)-[:t:TO_CLUB]->(c)
```

5 Transformation & Integration

5.1 Data Wrangling

Definition 5.1 - Data Wrangling

Data Wrangling is the process of getting to know your data and getting it into a format which you are comfortable with. This should get you interested in the project.

A big part of *Data Wrangling* is detecting inaccurate or malformed records and then correcting or removing them. (Data cleaning, validation, debugging, transforming, formatting, pre-processing).

Data Wrangling is an iterative process, you will often return to after “evaluating your model” in order to tweak the data some more.

Remark 5.1 - Data Used

Some *Data Wrangling* tasks need to be performed on the whole data set (e.g. checking & validation), but others can be performed with a random, representative subset (e.g. exploring)^[22].

Remark 5.2 - Strategies for Data Wrangling

Here are some things to consider when *Data Wrangling*

- *Type Screening* - Ensuring data is the type you expect and want.
- *Range Check* - Check data is within legal ranges or categories.
- *Illegal Values* - Check data takes legal values.

^[21]The SQL for this would require two joins

^[22]You may want to seek out some outliers, especially when your project is interest in outlier cases.

- *Multi-Column Validation* - Data in different columns is consistent.
- *Check Unique/Distinct Values* - Check that unique values are reasonable. (Not typos).
- *Remove Duplicates* - Remove duplicate records.

Remark 5.3 - Tools

There are several useful tools for *Data Wrangling*, including

- Data Wrangler <http://vis.stanford.edu/wrangler/>
- Pandas.
- Regular Expressions.

Remark 5.4 - NoSQL Databases

NoSQL Databases (e.g. MongoDB, Neo4j) are not as strongly typed as SQL databases.

Remark 5.5 - Missing Data

Missing data is common and can occur for many reasons

- Data may be sensitive and people choose to withhold.
- Data may be deemed irrelevant and skipped.
- Data may be censored.
- Communication link may fail (battery dies).
- Data collection may be aborted (boredom).
- Data may be corrupted.
- Data may not be known by the respondent.

Remark 5.6 - Informative Missing Data

Sometimes the fact that data is missing can be very informative. It can indicate that the technique/equipment used to take a reading was changed, potentially due to a change in circumstances (A patient went into/out of intensive case).

Imputing values in this case would lose this information.

Remark 5.7 - Dealing with Missing Data

Here are some approaches which can be taken to deal with missing data

- Cross reference against validated external information.
- Impute values from nearest neighbours.
- Impute value from mean/median/mode.
- Impute from a model-based calculation.
- Perform multiple imputations.
- Use a “Missing Data” indicator.

5.2 Data Fusion

Definition 5.2 - *Data Fusion*

Data Fusion is the process of taking data from different places it and representing it in a single place^[23]. This should create value by increasing the amount of available data.

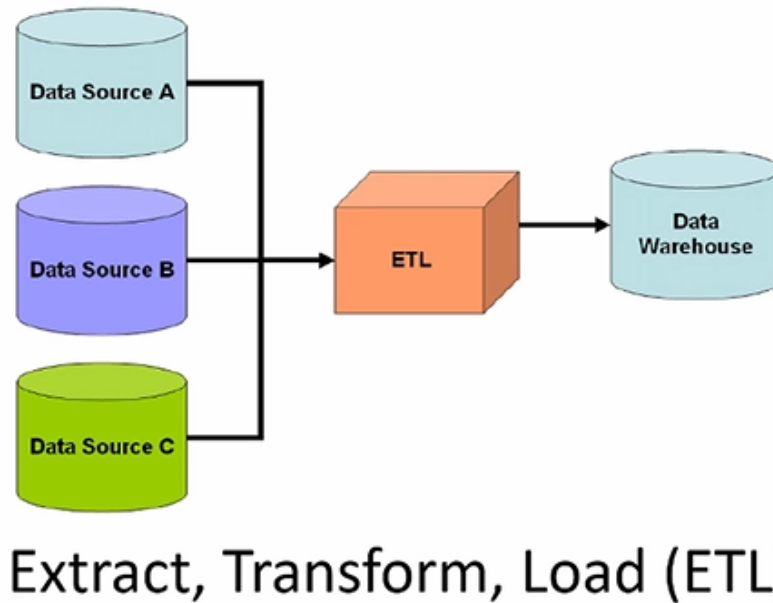


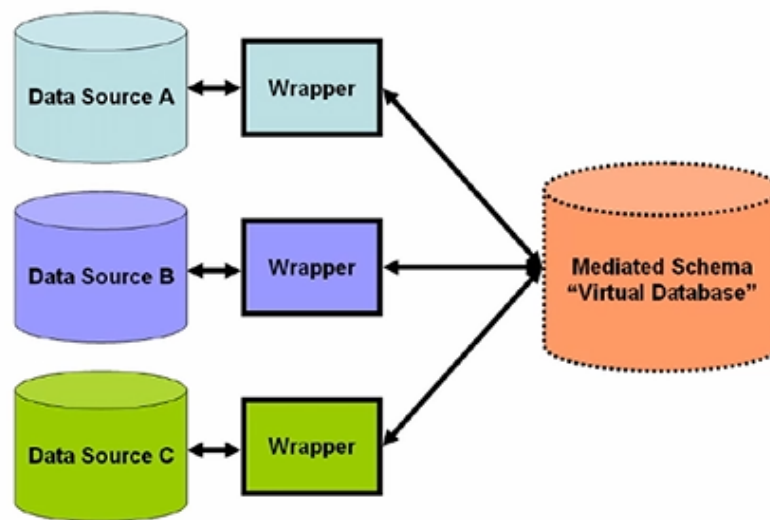
Figure 7: Example of an Extract Transform Load approach to *Data Fusion*.

Proposition 5.1 - *Extract. Transform, Load (ETL)*

Extract. Transform, Load (ETL) is an approach to *Data Fusion* where all the data sources are centralised.

See Figure 7 from an example.

^[23]e.g. Having a robot with multiple sensors, and using the accumulated data to determine the robots actions.



Data-Integration Solution

Figure 8: Example of an Data-Integration approach to *Data Fusion*.

Proposition 5.2 - *Data-Integration Solution*

Data-Integration Solution is an approach to *Data Fusion* where “wrappers” are created for each data source, allowing each to be read using the same schema.

See Figure 8 from an example.

Definition 5.3 - *Data Heterogeneity*

Data Heterogeneity is the property that datasets can vary. Data can be different in many ways, including

- Temporally.
- Spatially.
- Encoding used.
- Collection method.
- Units of measurement.
- License.
- Whether missing values have been imputed.

Remark 5.8 - *Other Issues for Data Fusion*

Here are some issues which may affect the robustness of some *Data Fusion* methods

- Bad Data - Even if a small part of the data is corrupted then the all the data is compromised.
- Lack of Storage - *Data Fusion* methods which centralise all the data will need large disk capacity.
- Original Data Sources - If you are interacting with third-party data sources there is a greater risk it will change/be taken down.^[24]

^[24]Documentation for data sources will often help alleviate this.

Definition 5.4 - Object Relational Mapping

Object Relational Mapping^[25] is an approach to aggregating data from multiple data sources into “virtual databases” by creating an interface between the user and the data sources. This makes incompatible data sources “appear” compatible. This interface is designed to be flexible and easy to use.

This should speed-up development and allow for migration to be handled automatically. But, gives the user less control (often making complex queries very difficult) and has slower execution speed than directly interacting with the data sources.

Object Relational Mappings are often optimised for certain common queries.

6 Exploration & Visualisation

Remark 6.1 - Common Questions about Data

Here are three common early questions we may have about some data

- i). What is the average of the sample?
- ii). What is the spread of the sample?
- iii). What is the effect of outliers?

Definition 6.1 - Inferential Statistics

Inferential Statistics^[26] summarise a sample and are based on probability theory.

Definition 6.2 - Descriptive Statistics

Descriptive Statistics describe or summarise a dataset. *Descriptive Statistics* are useful for assessing data quality, developing models and improving our general understanding about a dataset.

Definition 6.3 - Variables

There are several classes of variable

- *Qualitative Variables* - Categorical or discrete data. (inc. nominal, ordinal)
- *Quantitative Variables* - Continuous data. (inc. intervals & ratios)

Definition 6.4 - Univariate Measures

Univariate Measures are measures which describe only one characteristic of a sample/dataset. There are several classes of *Univariate Measures*

- *Central Tendency* (e.g. mean, median, mode).
- *Variability* (e.g. variance, quartiles, max, min, skew). See **Definition 6.1**
- *Graphical formats* (e.g. histograms).

Remark 6.2 - Deciding which “Central Tendency” measure to use

^[25]Examples - Django, MonoEngine, Hibernate.

^[26]AKA Inductive Statistics

Plotting a *Histogram* of the data will help decide between using mean, median or mode due to being able to observe the distribution and existence of outliers.

Proposition 6.1 - Measures of Variability

Here are some measures of variability in a data set

- Range $\max - \min$.
- Variance $\text{Var}(X) := \frac{1}{n} \sum (X_i - \mu)^2$.
- Unbiased variance $\text{Var}(X) = \frac{1}{n-1} \sum (X_i - \mu)^2$.^[27]
- Standard deviation $\sigma = \sqrt{\text{Var}(X)}$.
- Quantiles - Values at which $\frac{1}{4}, \frac{1}{2}, \frac{3}{4}$ lies beneath.
- Percentiles - Values at which $N\%$ of the data lies beneath.

Definition 6.5 - Raw Moments

Let X be a random variable. The n^{th} *Raw Moment* of X is defined as

$$\mathbb{E}[X^n]$$

This means that the first moment is the expected value and the variance can be derived from the first two moments.

Definition 6.6 - Central Moment μ_n

Let X be a random variable. The n^{th} *Central Moment* of X is defined as

$$\mu_n := \mathbb{E}[(X - \mathbb{E}[X])^n]$$

The zeroth central moment is 1; first central moment is the expected value; the second central moment is the variance.^[28]

Definition 6.7 - Standardised Moments $\bar{\mu}_n$

The n^{th} *Standardised Moment* is the ratio of the n^{th} centralised moment and the standard deviation to the power of n

$$\bar{\mu}_n := \frac{\mu_n}{\sigma^n} = \frac{\mathbb{E}[(X - \mu)^n]}{(\sqrt{\mathbb{E}[(X - \mu)^2]})^n}$$

Standardised Moments differ in properties other than variability, allowing for comparison of shape of distributions.

Remark 6.3 - Cumulants

Cumulants are another alternative to moments but they seem niche so I cba to define them.

Definition 6.8 - Skewness

Skewness is a measure of distribution shape, describing how lopsided a distribution is.

$$\text{Skew}[X] := \mathbb{E} \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] = \frac{\mu_3}{\sigma^3}$$

^[27]Accounts for how sample mean varies \bar{X} from true mean μ

^[28]The third and fourth central moments are used to define skewness and kurtosis, respectively.

$\text{Skew}[X] < 0$ means the distribution has more values to the right, and $\text{Skew}[X] > 0$ it has more values to the left.

Definition 6.9 - Kurtosis

Kurtosis measures the heaviness of the tail of a distribution, compared to a normal distribution with the same variance.

$$\text{Kurt}[X] = \mathbb{E} \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right] = \frac{\mu_4}{\sigma^4}$$

Since the *Kurtosis* of a normal is $3\sigma^4$, we define *Excess Kurtosis* as $\text{Kurt}[X] - 3$.

6.1 Data Transformation

Definition 6.10 - Variance Stabilising Transforms

Variance Stabilising Transforms are a set of transforms which ensure the variance of each observation is approximately the same. This can be difficult in some cases as the variance may not be constant, but it is related to the mean^[29]

A *Box-Cox Transformation* is one such transform.

Definition 6.11 - Decorrelation

Decorrelation is the practice of removing correlation between different variables. This emphasises the differences between variables, allowing for more data to be extracted.

Principle Component Analysis is a *Decorrelation* method.

Proposition 6.2 - Decorrelation - Maths

Decorrelation transforms data so it has a diagonal *Covariance Matrix* $\Sigma = X^T X$. This is done by finding the eigenvalue of

$$\Sigma \Phi = \Phi \Lambda$$

where Λ is a diagonal matrix of eigenvalues and the columns of Φ are the eigenvectors of the covariance matrix (i.e. Φ diagonalises Σ). Equivalently

$$\Phi^T \Sigma \Phi = \Lambda$$

We can decorrelate a single vector x_i (at test time) by calculating

$$\hat{x}_i := \Phi^T x_i$$

Definition 6.12 - Whitening Data

Whitening Data is the practice of transforming data such that all variables have the same mean and variance. This makes the data more comparable.

Proposition 6.3 - Whitening Data - Maths

When *Whitening Data* we are ensuring the eigenvalues in Λ are all the same. As eigenvalues determine the length of eigenvectors, ensuring all eigenvalues are the same makes the

^[29]i.e. For a poisson distribution $\text{Var}(X) = \mathbb{E}[X]$.

covariance matrix Σ spherical (rather than elliptical).

Remark 6.4 - When to Whiten Data

There are some scenarios when it is not good to *Whiten Data*, including

- When the scale of data is important to your inference problem. In this case, eigenvalues could be a useful set of features for inference.
- The covariance matrix may be computationally intractable. And, eigenvalues decomposition takes $O(n^3)$ time.
- In machines learning you need to use the exact same transformation on training and testing data. These transformations are determined using the training data.

6.2 Density Estimation

Definition 6.13 - Density Estimation

Density Estimation is the process of computing the probability distribution which underlies the data. There are two main approaches^[30]

- i). *Parametric Approach* - Assume the histogram of the data follows some common-simple distribution^[31], this means we only need to estimate the parameters of this distribution (mean, variance). This can be extended using *Gaussian Mixture Models*^[32]

A limitation of this is that we are encoding our assumptions into the model, which may not be sound.

- ii). *Non-Parametric Approach* - Use the histogram of the data as the distribution. This does not require prior knowledge and inference is easy, but is computationally expensive. This can be extended using *Kernel Density Estimation*.

Definition 6.14 - Kernel Density Estimation

Kernel Density Estimation uses the following formula to estimate the density $\hat{f}_h(x)$ from the data $\{x_1, \dots, x_n\}$

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

where K is a smoothing kernel.

6.3 Multi-Variate Analysis

Definition 6.15 - Multivariate Analysis

Multivariate Analysis is performed on samples which consist of multiple variables. *Descriptive Statistics* are often used to describe the relationship between pairs of variables (e.g. Cross-tabulations, contingency tables & graphical representations)

^[30]There are also *Semi-Parametric Methods* such as the Dirichlet Process Mixture Model.

^[31]Often a gaussian

^[32]A collection of Gaussians each with a different weighting, fitted using the *Expectation-Maximisation Algorithm*.

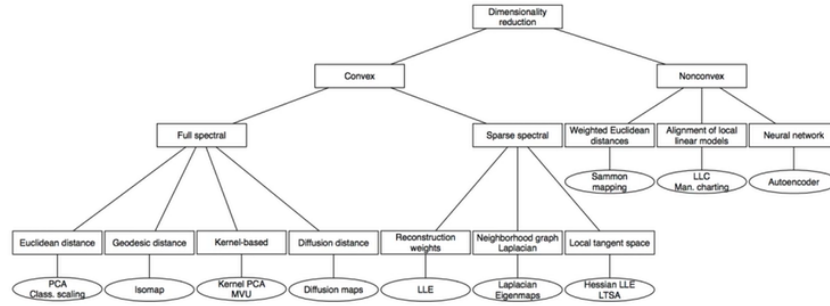


Figure 9: Taxonomy of methods for non-linear dimensionality reduction.

Definition 6.16 - Dimensionality Reduction

Dimensionality Reduction is the process of reducing the number of variables in each sample, by summarising them.^[33] We generally choose to do this due to computational limitations.

The two approaches are *Linear Dimensionality Reduction* and *Non-Linear Dimensionality Reduction* (See Figure 9).

Definition 6.17 - Principle Component Analysis

Principle Component Analysis (PCA) is a *Linear Dimensionality Reduction* which aims to keep as much variance as possible. Here is the general process to follow

- i). Standardise the data.
- ii). Obtain the eigenvectors and eigenvalues from the covariance matrix.
- iii). Sort the eigenvalues in descending order and choose the k ^[34] eigenvectors which correspond to the largest eigenvalues.
- iv). Construct the projection matrix $R \in \mathbb{R}^{k \times m}$ from these eigenvectors.
- v). Transform the original dataset $X \in \mathbb{R}^{m \times n}$ using R to obtain $P \in \mathbb{R}^{m \times k}$.

$$P = XR$$

Definition 6.18 - Random Projections

Random Projections is a simplified version of *Principle Component Analysis* which does not prioritise maximising variability. In the below equation R is chosen at random, we just hope it is good.^[35]

$$P = XR$$

where $X \in \mathbb{R}^{m \times n}$ is data matrix, $R \in \mathbb{R}^{k \times m}$ is the projection matrix and $P \in \mathbb{R}^{m \times k}$ is the resulting lower-dimensional representation.

Definition 6.19 - Gaussian Random Projections

Gaussian Random Projections is an extension of *Random Projections* which places restrictions on how R is chosen.

- i). First row is a random unit vector, chosen uniformly at random.

^[33]Reducing to 2 or 3 dimensions is popular as it is easy to visualise.

^[34]With $k \ll n$ where n is the number of dimensions in each sample.

^[35]There are extensions we aim to make this more likely.

- ii). Second row is a random unit vector chosen from the space which is orthogonal to the first row.
- iii). Third row is a random unit vector chosen from the space which is orthogonal to the first two rows.
- iv). etc.

This method means that R has *Spherical Symmetry*, *orthogonality* and *Normality*.

Theorem 6.1 - Johnson-Lindenstrauss Lemma

Given $\varepsilon \in (0, 1)$ a set X of m points in \mathbb{R}^n and a number $k > \frac{8}{\varepsilon^2} \ln(m)$, there is a linear map $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ st for all $\mathbf{u}, \mathbf{v} \in X$

$$(1 - \varepsilon) \|\mathbf{u} - \mathbf{v}\|^2 \leq \|f(\mathbf{u}) - f(\mathbf{v})\|^2 \leq (1 + \varepsilon) \|\mathbf{u} - \mathbf{v}\|^2$$

This means generating random projections as described in **Definition 6.19** (as they are orthogonal projections) will, in general, reduce the average distance between points.

Remark 6.5 - Database Friendly Random Projections

Choosing the elements of our projection matrix R as below is database friendly

$$R_{i,j} = \sqrt{3} \cdot \begin{cases} 1 & \text{with prob } 1/6 \\ 0 & \text{with prob } 2/3 \\ -1 & \text{with prob } 1/6 \end{cases}$$

This is database friendly as elements are $\{1, 0, -1\}$, making the sampling process very simple (easy for a computer).

Definition 6.20 - t -Distributed Stochastic Neighbour Embedding (t -SNE)

t -Distributed Stochastic Neighbour Embedding (t -SNE) is a *Non-Linear Dimensionality Reduction* method which tries to capture the relationship between samples and their neighbours (in the high-dimensional space), and then aims to keep these relationships the same in the lower-dimensional space (ie the closest neighbours of each sample are the same before and after transformation).

This is done by modelling the distribution of pairwise similarities as a *Student t -Distribution*. Let P_i be the distributions of pairwise similarities for input objects and Q_i for the distribution after the transformation. The problem is now to minimise the distance between the two, this is done by minimising the cost function^[36]

$$C = \sum_i KL(P_i, Q_i) = \sum_i \sum_j p_{j|i} \ln(p_{j|i} - q_{j|i})$$

Remark 6.6 - Using t -SNE

t -SNE is best used after a PCA step, it does not perform well as a pre-processing step. This is due to t -SNE focusing on pairwise relationships (which is local) and may not do well on the global scope.

Each time t -SNE is run a different visualisation may occur.

^[36]Generally using stochastic gradient descent.

7 Data Science in Production

7.1 Failure is the Rule

Remark 7.1 - *Failure is the Rule*

90% of data-science projects fail due to the complexity of the production pipeline for data-science.

A common approach is to “Fail often and Fail cheaply.”

Remark 7.2 - *Stages of Production*

There are the typical stages when applying data-science in production

- i). *Planning* - Identify the problem and justifying the need for data-science to be used. If this is successful then you need to plan a strategy, determine a budget, hire staff and outline goals (SMART goals).
- ii). *Proof-of-Concept*^[37] - Test approaches on a sample of the data. Explore the available data. Build and verify a model.
- iii). *Production* - Set up regular flows of data. Check there is no unreasonable bias in your model. Ensure the correct security and governance^[38] procedures are in place. Pay attention to when new trends arise which your model needs to be retrained for^[39]. Account for label drifting^[40].

Proposition 7.1 - *Components of the “Proof-of-Concept” Step*

A *Proof-of-Concept* can be broken down into the following steps

- Data Management:
 - Data collection.
 - Data pre-processing.
 - Data augmentation.
 - Data analysis.
- Model Learning:
 - Model selection.^[41]
 - Training.
 - Hyper-parameter selection.
- Model Verification:
 - Requirement encoding.
 - Formal verification.
 - Test-based verification.

^[37]AKA “Lighthouse Projects” as they shed light on projects.

^[38]Regular checks of the models performance and “alerts” when odds results occur.

^[39]Ideally the model would be continuously learning.

^[40]The labels in the data may change their meaning over time.

^[41]In a business setting you often need to model to be interpretable and some-what agree with preconceptions of your stakeholders.

Proposition 7.2 - Components of the Production Step

The *Production Step* can be broken down into the following steps

- Model deployment
 - Integration.
 - Monitoring.
 - Updating
- Cross-cutting aspects
 - Ethics.
 - End users' trust.
 - Security.

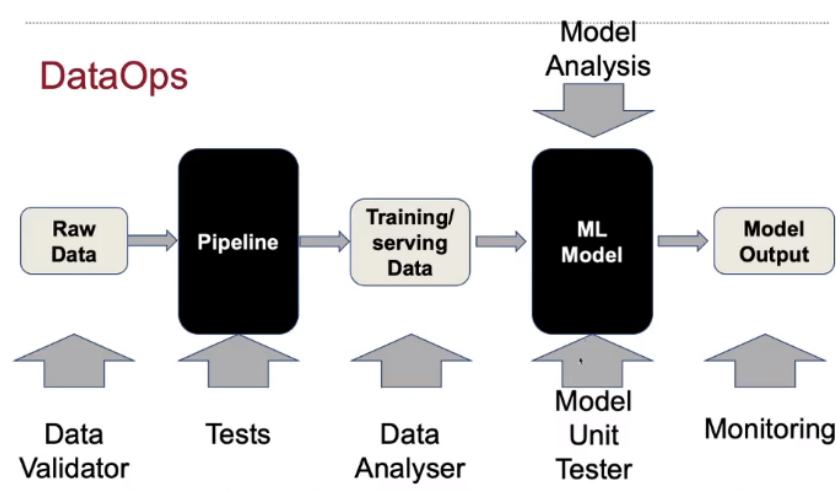
7.2 Engineering Guarantees for Data Science

Figure 10: Data Operations pipeline

Proposition 7.3 - Data Validator

A Data Validator should ensure that the data in specific fields is of the type we expect and within expected ranges/formats. This is pretty straight-forward for structured and semi-structured data, for unstructured data (e.g images and text) less validation is possible and you may need to abstract the data before validation.

Proposition 7.4 - Data Analyser

A Data Analyser should perform basic statistical procedures to check for bias, distribution skew, clusters etc.. In general, you need to be checking these features constantly in both production and deployment.

Proposition 7.5 - Model Unit Tester

A Model Unit Tester looks for errors in the training code using synthetic data. Much of this is similar to standard software testing.

They may also need to perform tests to fulfil regulatory/legal requirements (i.e. the model is explainable).

Proposition 7.6 - Tests on Features

Here are some tests which can be performed on a model (other than just visual inspection)

- *Score Decomposition* - (Not recommended these days due to non-linearity).
- *Ablation Testing* - Remove features one at a time, retrain and observe the difference.^[42]
- *Feature Permutation* - Choose a feature and augment its values to something which would never occur IRL and observe how other features compensate. This will demonstrate correlations and thus relevance.
- *Top-Bottom Analysis* - Compare distributions of inputs and outputs. This will demonstrate correlations and thus relevance.

Remark 7.3 - Feature Store

Some companies make use of a special “Feature Store” database where previously used features are stored so that their information is not lost even if they are not used in deployment.

Proposition 7.7 - Monitoring

We want to monitor the output of the model so when the output falls outside some criteria an alert is triggered and someone can investigate.

Model performance dashboards are often used to help monitoring. These dashboards typically display: model output metrics; data input metrics; and, operational telemetry.

Proposition 7.8 - Dimensions of Testing

There are often many different dimensions to consider when testing a model. Almost always you want to test the inputs and outputs, but often you will need to test several models and under different conditions (different cities, markets etc.)

Proposition 7.9 - Comparing Algorithms - Large Dataset

Here are some methods for comparing different aspects of several models, when you have access to a large dataset

- *Performance Estimation*
 - Two-way holdout method (train-test split).
 - Confidence interval via normal approximation.
- *Model Selection, hyperparameter optimisation and performance estimation*
 - Three-way holdout methods (train-validate-test split).
- *Model & algorithm comparison*
 - Multiple independent training sets + test sets
 - McNemar test.
 - Cochran’s Q + McNemar test

^[42]Popular but not recommended due to “Trigger’s Broom Paradox”.

Proposition 7.10 - *Comparing Algorithms - Small Dataset*

Here are some methods for comparing different aspects of several models, when you have access to a small dataset

- *Performance Estimation*
 - Repeated k -fold cross validation without independent test set.
 - Leave-one-out cross validation without independent test set.
 - Confidence interval via 0.632(+) bootstrap.
- *Model Selection, hyperparameter optimisation and performance estimation*
 - Repeated k -fold cross validation with independent test set.
 - Leave-one-out cross validation with independent test set.
- *Model & algorithm comparison*
 - Combined 5x2cv F test.
 - Nested cross-validation.

8 Sharing & Privacy