

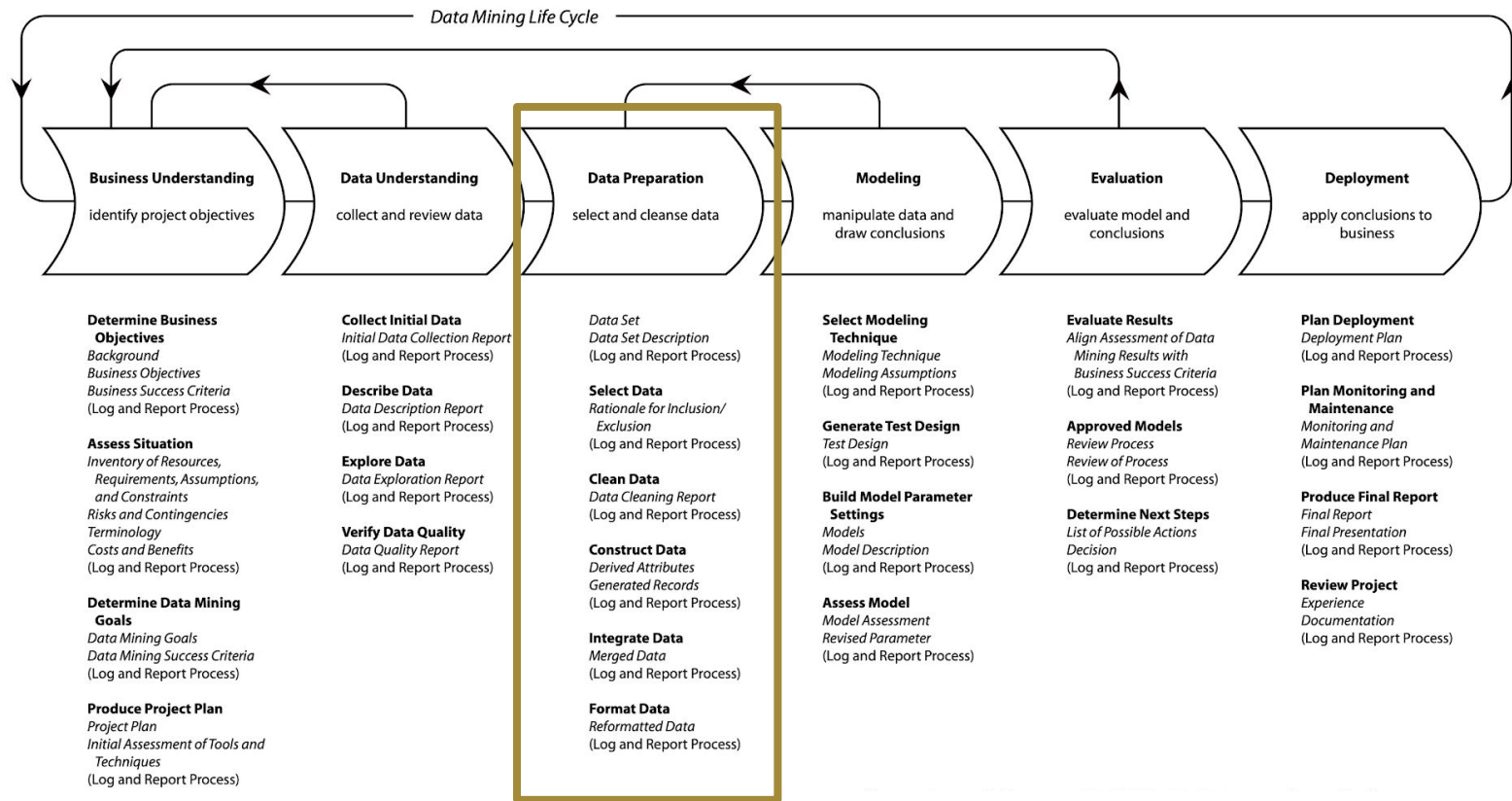
Data Preparation & Feature Engineering

Discover track | day 4

JADS

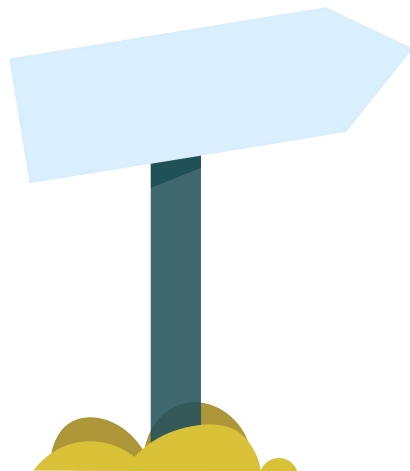
Jheronimus
Academy
of Data Science

Data preparation



Learning objectives for today

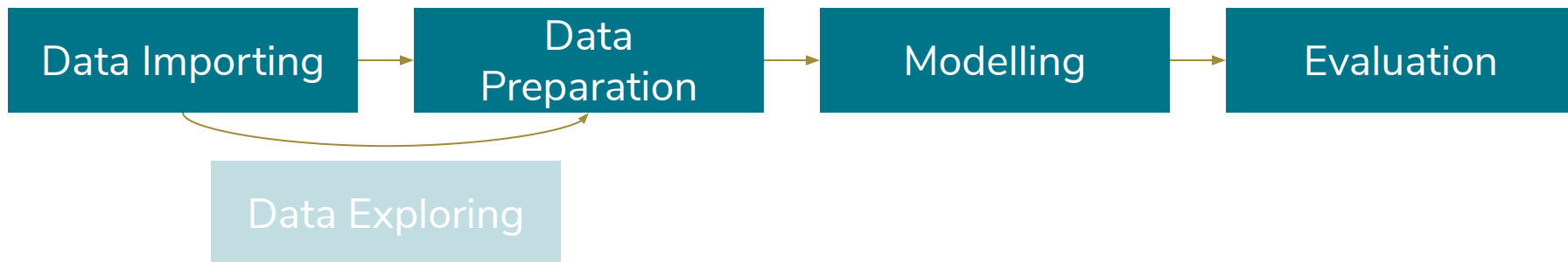
- _ Know best practices in cleaning and working with data
- _ Understand different types of data processing technologies and patterns
- _ How to transform raw data into:
 - _ technically correct,
 - _ real-world consistent, and
 - _ 'prediction-appropriate' data
- _ The importance of feature engineering
- _ Basic working knowledge of SQL for data preparation



Partly based on: [An introduction to data cleaning with R](#)

Data analysis pipeline

Training data



Test data and new data after deployment



Starting point

1. Start with the data in exactly the same shape as a newly collected dataset would have. Check whether your dataset does not already contain some non-automated pre-processing steps:
 - Selections
 - Aggregations
 - Outlier handling etc.
2. Always work from a script (SQL, R, Python)
 - never make manual adjustments to the data

Data prep steps

1. Importing and possibly combining datasets
2. Making data technically correct
 - Correct headers
 - Correct data classes (integer, string, factor, date, etc)
3. Making data real-world consistent
 - Detecting inconsistencies or errors
 - Selecting the cause of the inconsistency or error
 - Correcting the inconsistency or error
4. Making data 'prediction-appropriate'
 - Excluding variables
 - Defining Y and redefining X-variables
 - Handling correlated predictors
 - Adding domain knowledge

Work with tidy (structured) data

country	year	cases	population
Afghanistan	1999	37737	1272915272
Afghanistan	2000	3666	20495360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

variables

country	year	cases	population
Afghanistan	1999	37737	1272915272
Afghanistan	2000	3666	20495360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

observations

country	year	cases	population
Afghanistan	1999	37737	1272915272
Afghanistan	2000	3666	20495360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

values

- Each dataset much have its own **table** (dataframe)
- Each variable must have its own **column** (features)
- Each observation must have its own **row** (records, unit of analysis, grain)
- Each value must have its own **cell**

Pivoting: from wide to long format (and back)

- Use **long format** throughout your data processing pipeline
- Use **wide format** for presenting and communicating results
- Daily frustration: sources in wide format

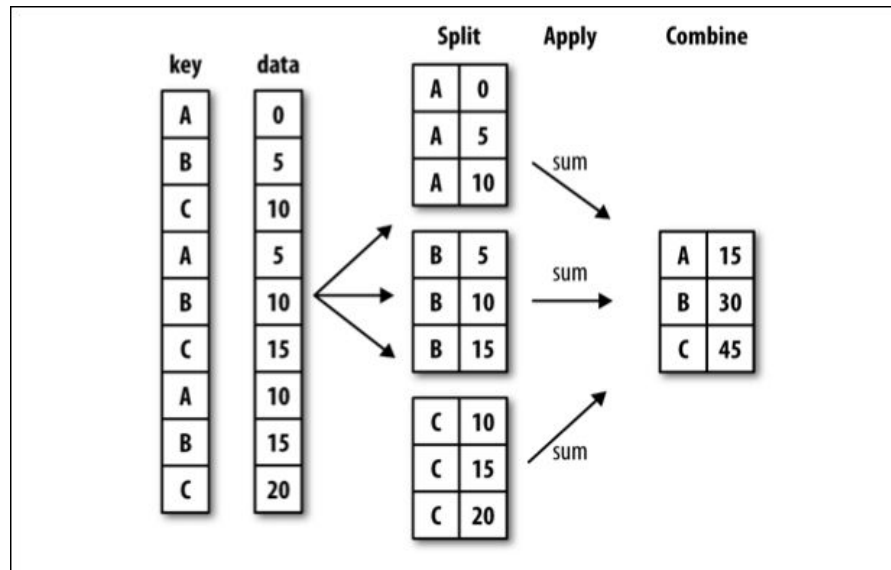
country	year	cases
Afghanistan	1999	745
Afghanistan	2000	2666
Brazil	1999	37737
Brazil	2000	80488
China	1999	212258
China	2000	213766

country	1999	2000
Afghanistan	745	2666
Brazil	37737	80488
China	212258	213766

table4

Split-apply-combine strategy for data processing

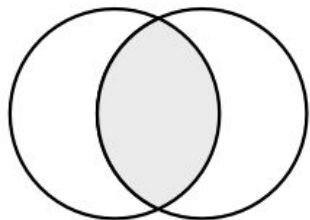
1. Extract a subset of the data for which it is easy to solve the problem.
2. Solve the problem by hand, checking results as you go.
3. Write a function that encapsulates the solution.
4. Use the appropriate plyr function to split up the original data, apply the function to each piece and join the pieces back together.



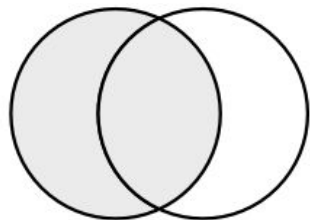
Same data processing concepts, different tools

<u>concept</u>	<u>Spreadsheet</u>	<u>SQL</u>	<u>R tidyverse</u>	<u>python pandas</u>
Sorting	<code>SORT()</code>	ORDER BY	<code>arrange()</code>	<code>sort_values()</code>
Filtering	<code>FILTER()</code>	WHERE	<code>filter()</code>	<code>filter()</code>
Split	n/a	GROUP BY	<code>group_by()</code>	<code>group_by</code>
Apply	many functions	many functions	many functions	many functions
Combine	VLOOKUP	JOIN	<code>join()</code> , <code>inner_join()</code> , <code>left_join</code> etc.	<code>join()</code> , <code>concat()</code> , <code>merge()</code> etc.

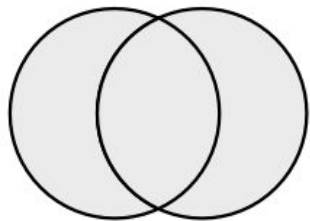
Understanding joins



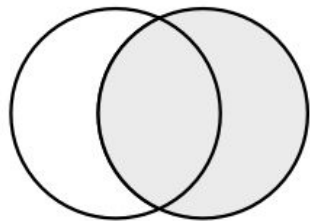
`inner_join(x, y)`



`left_join(x, y)`



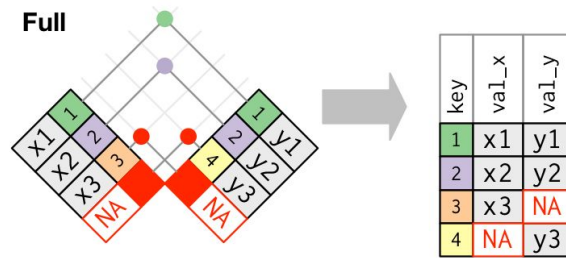
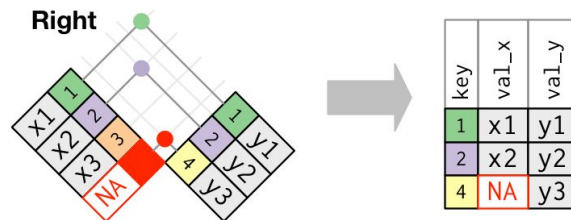
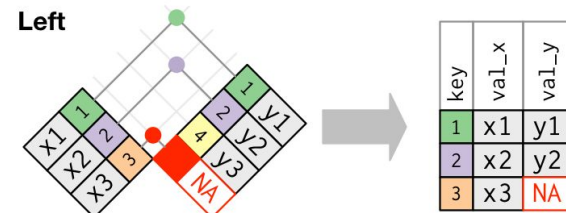
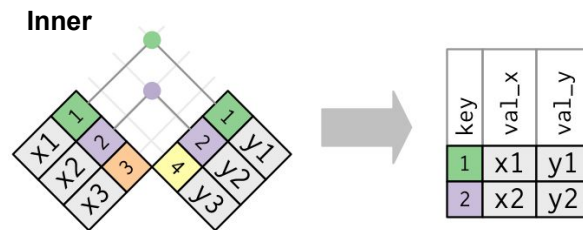
`full_join(x, y)`



`right_join(x, y)`

Understanding joins

x		y	
1	x1	1	y1
2	x2	2	y2
3	x3		



<https://r4ds.had.co.nz/relational-data.html#mutating-joins>

Making data 'prediction-appropriate'

1. Select usable variables
2. Defining the outcome Y
3. Redefining the X-variables
4. Handling correlated predictors
5. Adding domain knowledge

1. Select usable variables

Delete variables that:

- _ will not be known at the moment of modelling
- _ are (practically) constant
- _ are (practically) unique (unless they can be clustered meaningfully)

2. Defining the outcome Y

This is a crucial step in your project. Not working with the 'right' definition of Y is likely to result in either,

- a project that is of little use (as the outcome is not what users want), or
 - additional project iterations (going back in a later stage to adjust the definition because model outcomes are not satisfactory), or
 - never achieving good enough performance due to Y being of too low quality / too noisy
1. Be as sure as you can be regarding the definition and quality of Y
 2. Define the outcome Y, and
 3. Describe / visualize the newly defined Y

3. Re-defining the X-variables

It is not trivial how to approach this. Take into account that:

- Some machine learning algorithms automatically look for relevant interactions / non-linear relationships
- Some machine learning algorithms don't, and will only consider such relationships when explicitly being told to do so
- Even machine learning algorithms that do automatically look for interactions / non-linear relationships, could be helped when not having to discover those themselves

→ Finding out whether adjusting the input variables boosts model performance is a process of trial and error (model selection)

4. Handling correlated predictors

Highly correlated variables can destabilize the process of model building and tend to decrease model performance (bias-variance trade-off)

In healthcare or psychology, this typically occurs with questionnaire items

If a cluster of correlated features contains relevant information, explore different options for handling correlation:

- _ Principal Component Analysis (data-driven)
- _ Use sub scales (at least partly based on domain expertise)
- _ Use sub scales plus a few individual relevant items (at least partly based on domain expertise)
- _ Automatic feature selection using LASSO or Random Forest (data-driven)

5. Adding domain knowledge

Principles:

- _ Adding domain knowledge to your dataset is one of the most important aspects in a traditional machine learning project
- _ Your raw data are practically never structured / defined in such a way that they contain the most value to the problem you want to solve
- _ A simple model based on relevant data will perform much better than the most sophisticated model based on the raw data
- _ Feature engineering is partly data science and partly dark art
- _ Feature engineering depends largely on the availability of domain experts

Feature engineering

What variables can we create out of the variables that we have in order to increase the relevant information for modelling the outcome Y?

=> *Learning by example(s)*



Feature engineering - examples

- Converting absolute values into relative changes within the unit of analysis (e.g. change in a patient's medication use)
- Adding information about a higher unit of analysis (adding information about a clinic or therapist to a patient record)
- Categorizing numeric variables based on relevant cut-off values (e.g. converting Beck's Depression Inventory scores into mild, moderate and severe depression, or converting #minutes of daily movement into less than versus at least 30 minutes per day, changing 'last year's maternity care cost' into a yes/no dummy, or dividing age into <18, 18-45, 46-65, and 65+)

Feature engineering - examples

- Aggregating questionnaire items into relevant subscales and/or selecting specific items (e.g. converting the 45 items of the Outcome Questionnaire into subscales on Symptom Distress, Interpersonal relations and Social Role, in addition to selecting separate 'risk-items' on suicidality and aggression)
- Aggregating factor variables to decrease the number of different categories
- Describing 'agenda data' (e.g. scheduled appointments) into meaningful features (e.g. treatment intensity, continuity, etc.)
- Combining different inputs into one meaningful variable (e.g. construct 'treatment intensity' out of number of sessions and time in treatment)