



>> FAO Statistics
Division

DataLab

Paris21: Measuring Use of Food and Agricultural Statistics in Policy Making

17th June 2020



Outline

- ❖ Methodology
- ❖ Architecture, data, and workflow
- ❖ Results and web application to explore the data (ShinyApp and FAOLEX search)

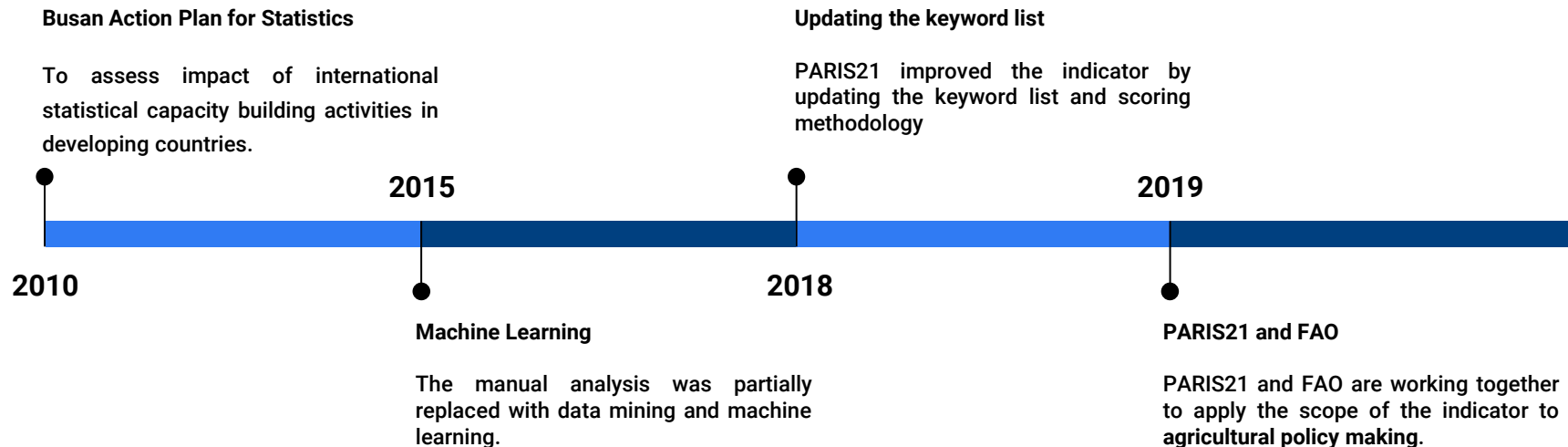


Objective

Measuring use of food and agricultural statistics in public policy documents



History of the indicator



Measuring the use of statistics

How the presence of statistical concepts is represented?



Measuring the use of statistics

Levels: sub-indexes

- The policy documents are split into sentences, and then their sentences are classified *into three mutually exclusive levels based on a set of keywords.*

COMPONENT	DESCRIPTION	EXAMPLE
Basic (Level 1)	The introduction of a measurable concept for indicating the state of and progress towards a specific development outcome	<i>The Government implemented a number of social protection interventions such as the social cash transfer, school feeding and the food security pack.</i>
Diagnosis and quantification (Level 2)	The depiction of the state of such measurable concept and the presumed causes for its variation.	<i>The mortality rate of 5 per 1 000 adults, improvement of 5% in adult literacy.</i>
Statistical analysis (Level 3)	It includes use of statistical methods to depict the state of the measurable concept and the causes of its variation	<i>The improvement of 5% in adult literacy rate is correlated to improvements in life expectancy or the income of the lower quartile.</i>

Measuring the use of statistics

Levels: sub-indexes

- In addition to the three levels, our analysis also searches for the use of *disaggregation* when a statistical term is referenced.

COMPONENT	DESCRIPTION	EXAMPLE
Disaggregation	Use of disaggregation while referencing to statistics.	---

- The structure and keywords of disaggregation originates from the stands for data disaggregation of UNSD¹.

¹<https://unstats.un.org/sdgs/files/Overview%20of%20Standards%20for%20Data%20Disaggregation.pdf>



Measuring the use of statistics

Levels: sub-indexes

Disaggregation indicators	Keywords	Disaggregation indicators	Keywords
Sex	Man, Woman, Men, Women, Girls, Boys, Female, Male	Agroecological zone	Climate variables, Type of Soil, Geomorphology
Age	Youth, Adults, infants	Water management	Irrigated, Non-irrigated
Indigenous status	Indigenous, indigenous group	Type of Conservation Facility	Medium-term, Long-term
Type of products:	Crop, livestock, mixed	Level of risk	Critical, critical-maintained, endangered, endangered-maintained
Type of products-subgroups:	wheat, rice, maize, millet, sorghum	Type of government authority	Agriculture , Industry Sector, Service Sector
Crop group, Food group	Cereals, Pulses, Fruits, Vegetables, Roots, Tubers, Oil-Bearing crops, Animals Products, Fish	Economic Sector	High income, low income, middle income
Geographical	Urban, Rural, National, Sub-national	Poor and vulnerable population	Poor, rich, poverty, wealth
Geographical Location	River basin, watershed, administrative unit	Level of Food Price Anomaly	Moderately high, High
Land Cover	Forestland, Cropland, Grassland, Wetlands, Settlements	Type of Tenure	Customary, Freehold, Leasehold
Agricultural holding type	Household, Non-household	Hydrological Units	River Basins , Aquifers
Type of Enterprise	Farming, Pastoral, Forestry, Fisheries, Hunting	Points of the Value Chain	Farm , Transport , Markets , Processors , Small-holders , large-commercial farms, firms
Size of Enterprise	Small, Medium, Large	Mountain Elevation Classes	Kapos classes

Keywords

FAO STAT

- 938 keywords
- Languages:
 - English
 - Spanish
 - French

SDGS

- 25 keywords
- Languages:
 - English
 - Spanish
 - French

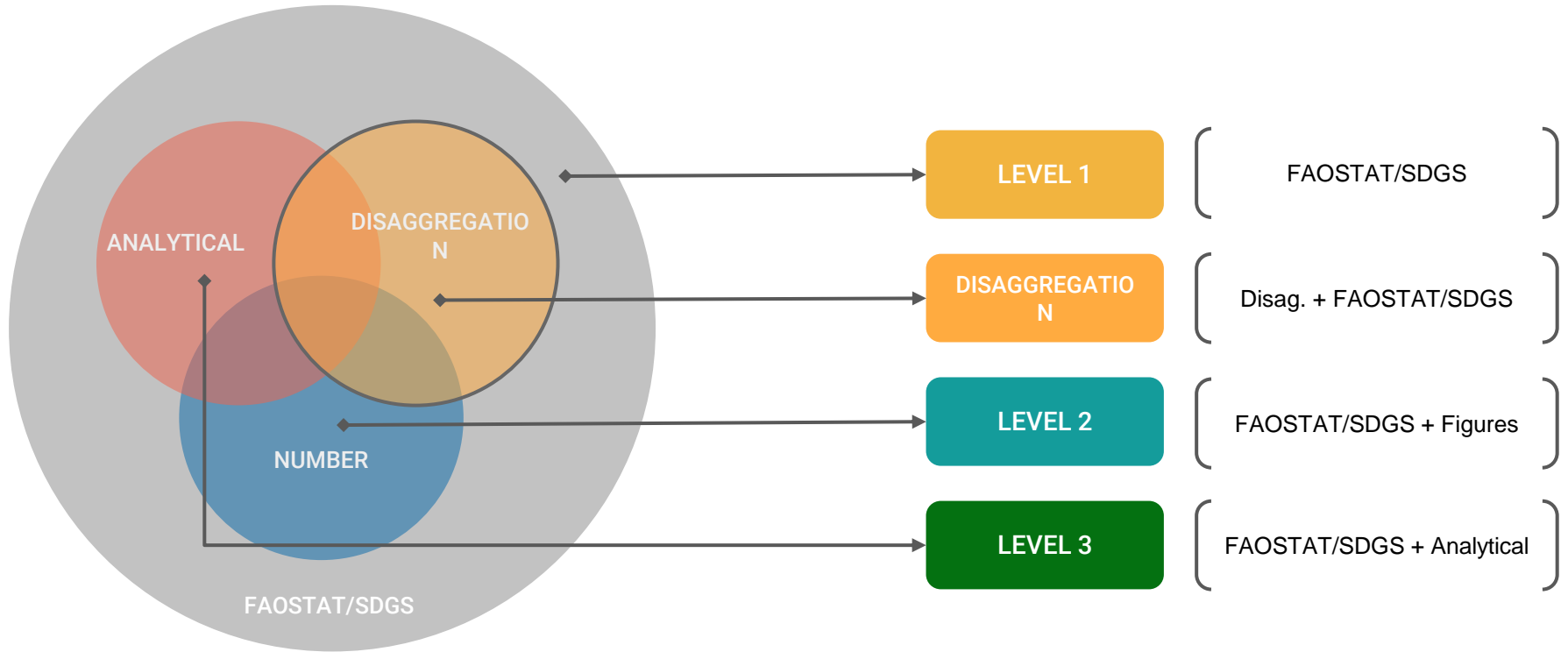
Disaggregation

- 24 indicators
- There is a list of keywords for each indicators.
- Languages:
 - English
 - Spanish
 - French

Analytical (math)

- 115 keywords
- Languages:
 - English
 - Spanish
 - French

Classification



Example

Country: Afghanistan

Title: Social Protection Sector Strategy

Year: 2008

PDF file: afg152314.pdf

Sentence classified at level 3

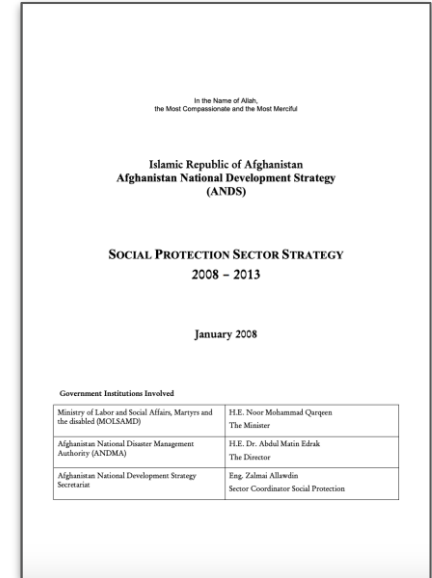
"Food **poverty** has been **estimated** to be even bigger: around **45** percent of the Afghan **population** was not able to purchase a basic food basket to provide **2100** calories consumption per day."

NUMBER

DISAGGREGATION

FAOSTAT/SDGS

ANALYTICAL



Algorithm (simple version)

LOOP OVER DOCUMENTS:

Step 1: it splits the document into sentences;

Step 2: for each sentence checks if it belongs one of the three levels;

Step 3: for each level computes the proportion of sentences;

Step 4: normalise each level to the scale min-max (0 - 1);

At this step we have the three sub-indexes: **basic**;
diagnosis and quantification; **statistical analysis**.



1. Principal Component Analysis





Weighting: Principal Component Analysis

DOCUMENT	LEVEL 1	LEVEL 2	LEVEL 3	DISAGG
doc_01	I_{11}	I_{12}	I_{13}	I_{1D}
doc_02	I_{21}	I_{22}	I_{23}	I_{2D}
.
.
.
doc_N	I_{n1}	I_{n2}	I_{n3}	I_{nD}

- Where I_{xy} is the sub-index for the level y and the document x ;
- The PCA method is applied on the matrix above to extract the weights of each level and the PCA (*Index Overall*).

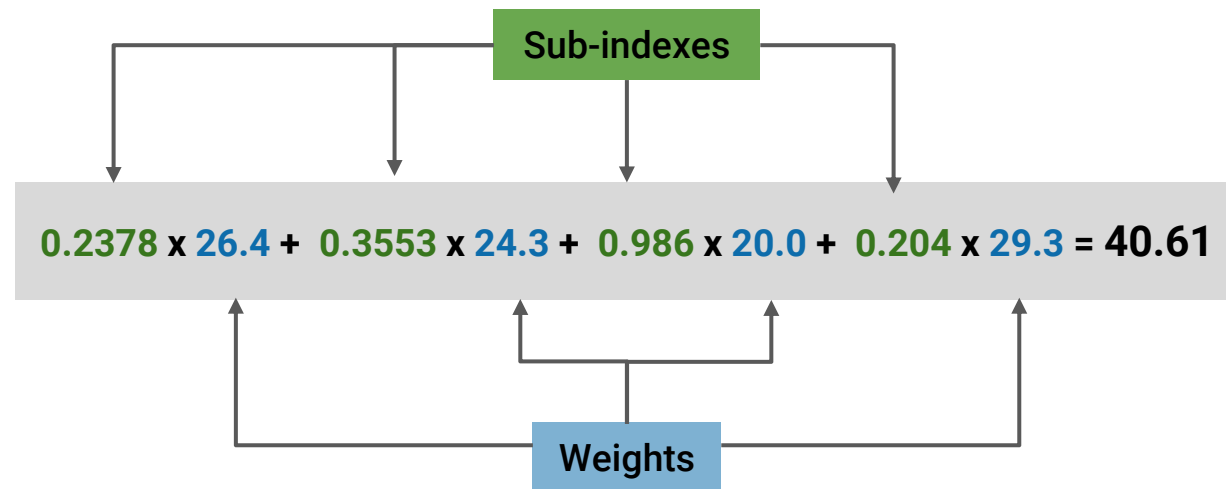


Weighting: Principal Component Analysis



SUB-INDEX	WEIGHTS PCA
Level 1	26.4
Level 2	24.3
Level 3	20.0
Disaggregation	29.3

ID	Country	Year	# sentence	Level 1	Level 2	Level 3	Level Disagg.	Score PCA
mne180387.pdf	Montenegro	2016	3873	0.2378	0.3553	0.986	0.204	40.61



2. Simple linear combination: equal weights

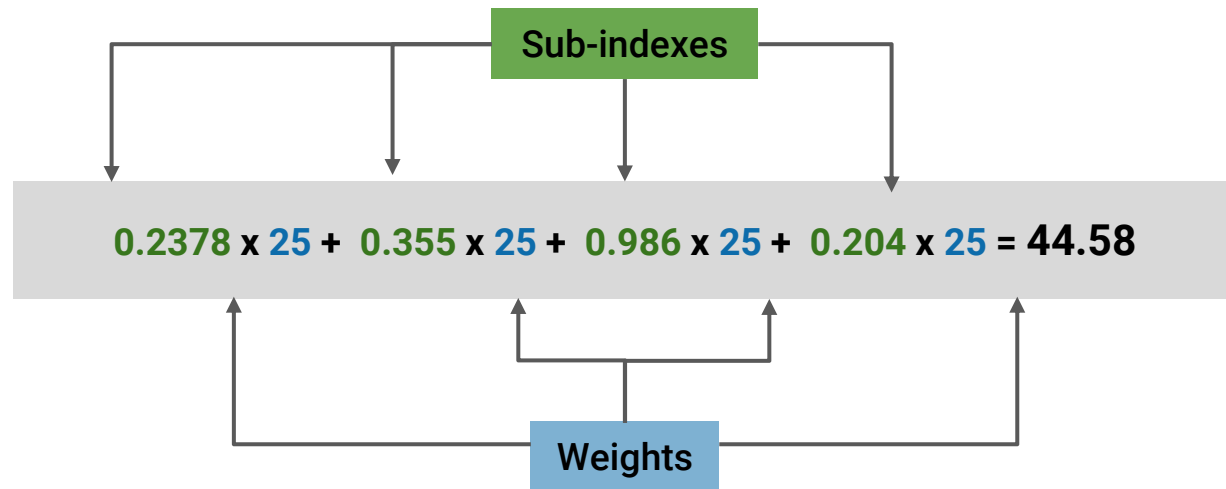




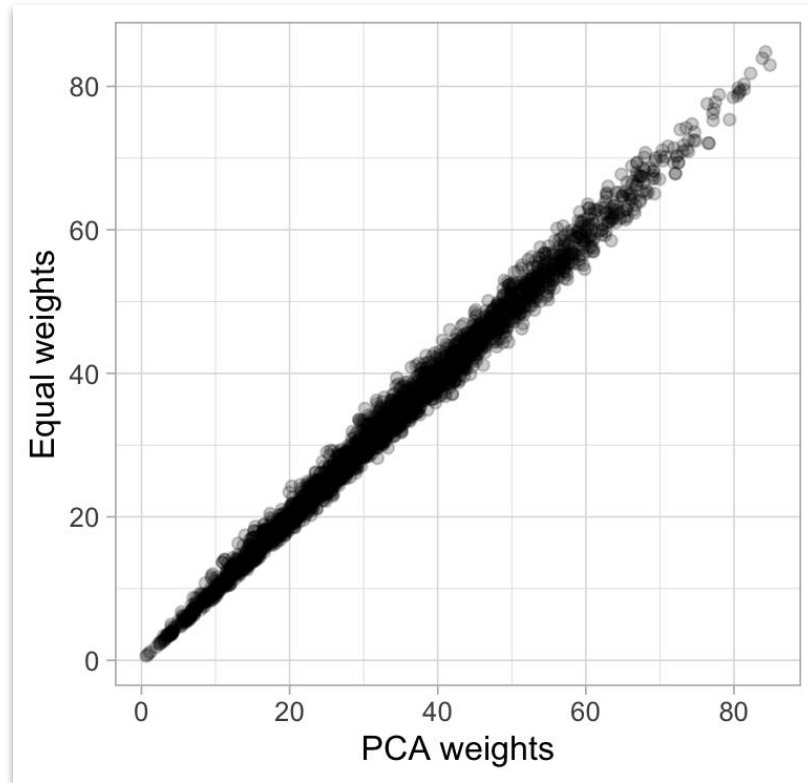
Weighting: equal weights

SUB-INDEX	WEIGHTS PCA
Level 1	25
Level 2	25
Level 3	25
Disaggregation	25

ID	Country	Year	# sentence	Level 1	Level 2	Level 3	Level Disagg.	Score Simple
mne180387.pdf	Montenegro	2016	3873	0.238	0.355	0.986	0.204	44.58



Weighting



3. Weighting by experts?



FAOLEX database

“FAOLEX is a comprehensive and up-to-date legislative and policy database, one of the world's largest online repositories of national laws, regulations and policies on food, agriculture and natural resources management.”



4295
documents



214
countries



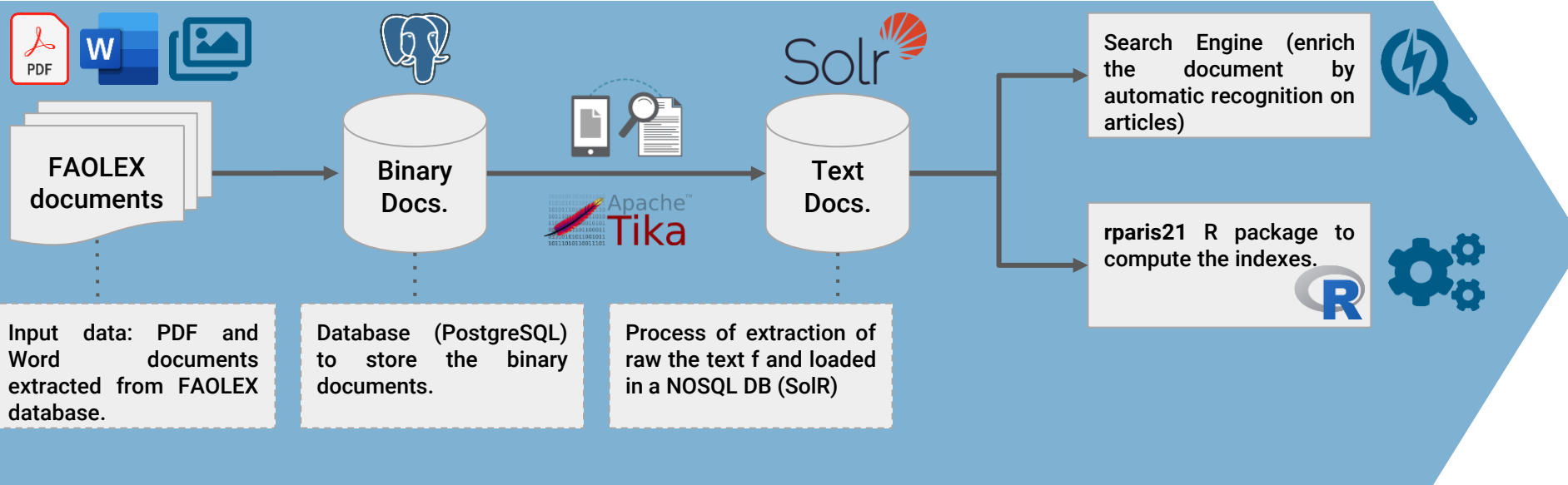
1985 - 2019
years



44
languages



Pre-processing the texts



Input data: PDF and Word documents extracted from FAOLEX database.

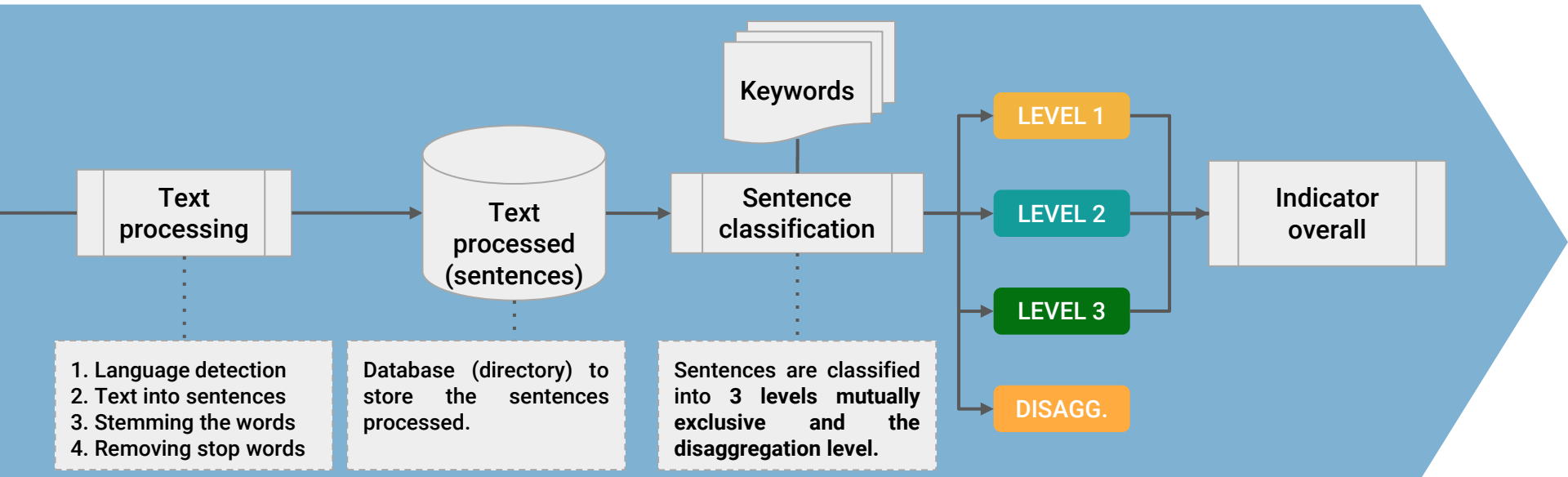
Database (PostgreSQL) to store the binary documents.

Process of extraction of raw the text f and loaded in a NOSQL DB (Solr)

Search Engine (enrich the document by automatic recognition on articles)

rparis21 R package to compute the indexes.





Search engine: user interaction



The screenshot shows a search engine interface with a search bar containing 'water management' and a 'Search' button. Below the search bar are 'Search options' including 'List', 'Entities', and 'Analyze'. A 'Sort' dropdown is set to 'Newest'. The main results area displays a document titled 'National Strategy for Agricultural Development 2016-2025' with a date of '2020-06-15T11:13:32Z' and a file name 'jor166604.pdf'. The document preview text reads: 'The National Strategy for Agricultural Development 2016-2025 is a sectoral policy aiming at achieving in ten years in the agricultural field the following results, together with economic, social and environmental consequences: (i) high agricultural productivity; (ii) efficiency in the use of irrigation water; (iii) high use of fertilizers; (iv) high partnership between'. To the right, there are filters for 'File date' (2020 (3095)) and 'Country' with a list: CHINA (81), VIET NAM (63), GHANA (57), GUATEMALA (55), COLOMBIA (52), BANGLADESH (50), CAMBODIA (49), BURKINA FASO (48). A 'Show more' link is at the bottom.

A simplified search bar with 'water management' and a 'Search' button. Below it are pagination controls: '◀ Previous Page 1 of 169 (results 1 to 10 of 1684) Next ▶'.

Search engine (Fulltext search)

Easy full text search in multiple data sources and many different file formats: Just enter a search query (which can include powerful search operators) and navigate through the results.



Solr Features

Solr is a standalone enterprise search server with a REST-like API. You put documents in it (called "indexing") via JSON, XML, CSV or binary over HTTP. You query it via HTTP GET and receive JSON, XML, CSV or binary results.

Search engine: other facilities

Thesaurus & Grammar (Semantic search)

Based on a [thesaurus](#) the multilingual semantic search engine will find [synonyms](#), [hyponyms](#) and [aliases](#), too. Using heuristics for [grammar rules](#) like [stemming](#) it finds other word forms, too.

Interactive filters (Faceted search)

Easy navigation through many results with [interactive filters](#) (faceted search) which aggregates an overview over and interactive filters for (meta) data like authors, organizations, persons, places, dates, products, tags or document types.

Exploration, browsing & preview (Exploratory search)

Explore your data or search results with an [overview of aggregated search results](#) by different facets with [named entities](#) (i.e. file paths, tags, persons, locations, organisations or products), while browsing with comfortable navigation through search results or document sets.

View previews (i.e. PDF, extracted Text, Table rows or Images).

Analyze or review document sets by preview, extracted text or [wordlists for textmining](#).

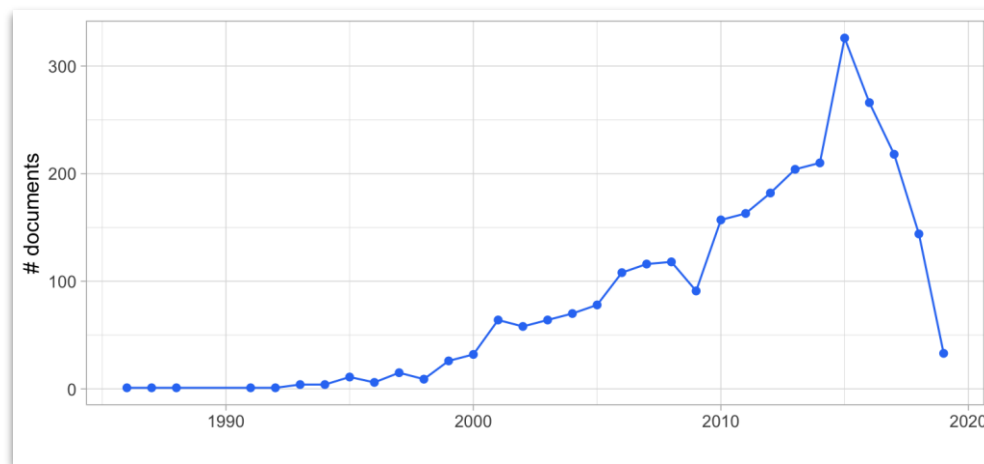
Collaborative annotation & tagging (Social search & collaborative filtering)

[Tag your documents with keywords, categories, names or text notes](#) that are not included in the original content to find them better later (document management & knowledge management) or in other research or search contexts or to be able to filter annotated or tagged documents by interactive filters (faceted search).

Or evaluate, value or assess or filter documents (i.e. for validation or collaborative filtering).



Results: data analysed



2790 / 4295
documents



200 / 214
countries

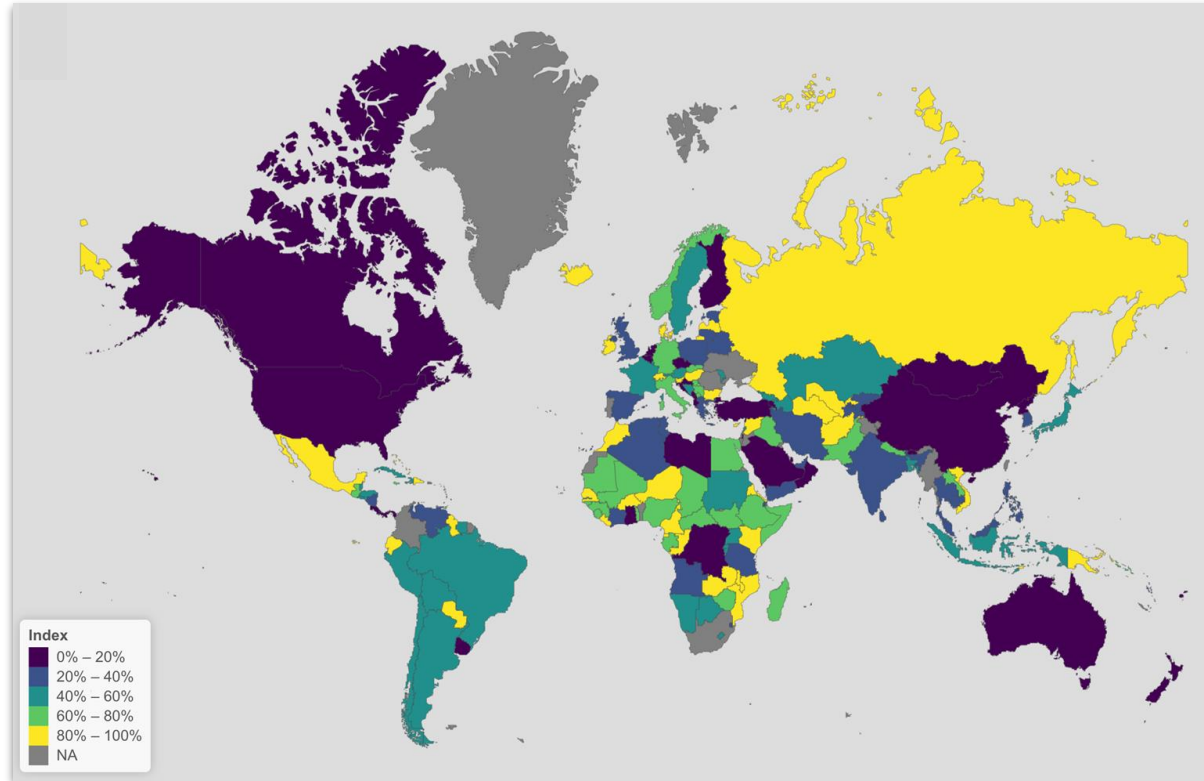


1986 - 2019
years



3 / 44
languages

Map: Index based on PCA





<https://foodandagricultureorganization.shinyapps.io/paris21/>

Next steps



- Expert review and feedback
- Improving the aggregation method for the countries, and years level.
- Improving the keywords
- Platform to explore and visualize the data and results

