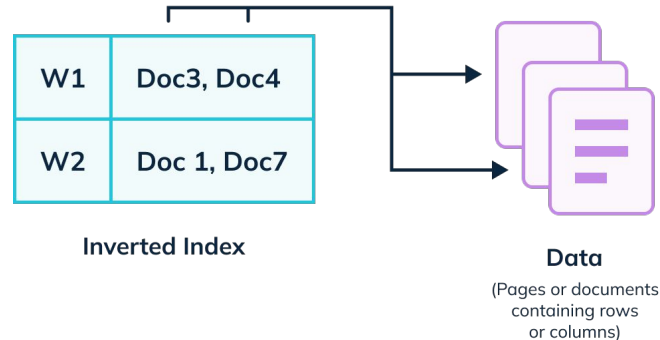# NLP 2

Inteligencia artificial avanzada para la ciencia
de datos II Modulo 5 NLP 2

Inverted Index

# Inverted Index

An **inverted index** is a data structure used in information retrieval and text search systems to efficiently store and retrieve a large amount of textual data.

It is designed to improve the speed of searching and retrieval operations by mapping terms (words or tokens) to the documents or data records in which they appear.

| W1 | Doc3, Doc4 |
|----|-----------|
| W2 | Doc 1, Doc7 |

**Inverted Index**

**Data**
(Pages or documents containing rows or columns)

# Software available

- **Lucene/Solr**: Apache Lucene is a high-performance, full-featured text search engine library written in Java. Apache Solr is an enterprise search platform built on top of Lucene. These tools are widely used for creating search engines and handling complex search operations.
- **Elasticsearch:** Elasticsearch is a distributed, RESTful search and analytics engine built on top of Lucene. It's often used for real-time search and log analytics. Elasticsearch is known for its scalability and ease of use.
- **Opensearch:** Elasticsearch fork supported by AWS, elasticsearch is an open source initiative that didn't receive earnings for AWS PaaS distribution of Elasticsearch, AWS forked the project and keep it open source to be distributed in the cloud.

# ES Architecture

Elasticsearch is a distributed, RESTful search and analytics engine designed for horizontal scalability, real-time search, and large-scale data storage and retrieval.

Its architecture is based on a cluster structure to handle vast amounts of data and to provide high availability and fault tolerance.



Elasticsearch Component Relation

# Basics

Imagine you are back in the 80's and you want to search a topic in a book.

You will go to the index , in order to find the pages in which this topic (Token) appears.

Here are some aspects you would need to know:

1. Someone read the hole book and make a map of the word and the page where it appears.
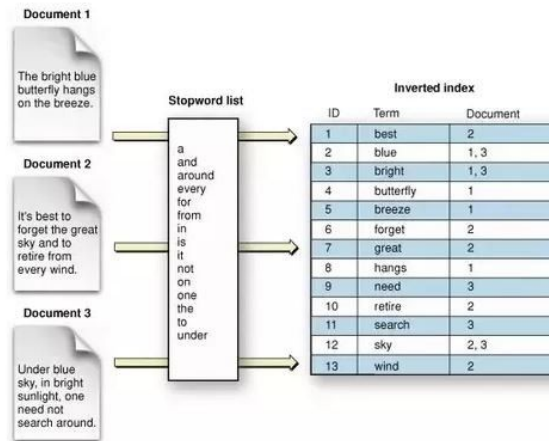2. This map is sorted and structured to be efficient to look for
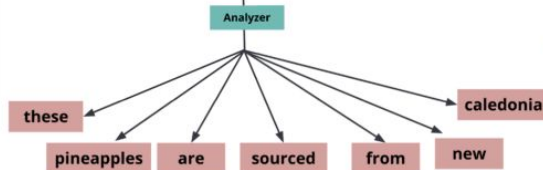
# Index time process

Reverse index software is the responsible of building that index for you, so later you can search for a word and retrieve your documents.

This means software has to:

- Process the text and tokenize
- Create the reverse index structure
- Keep some statistics to score the relevance given a sequence.
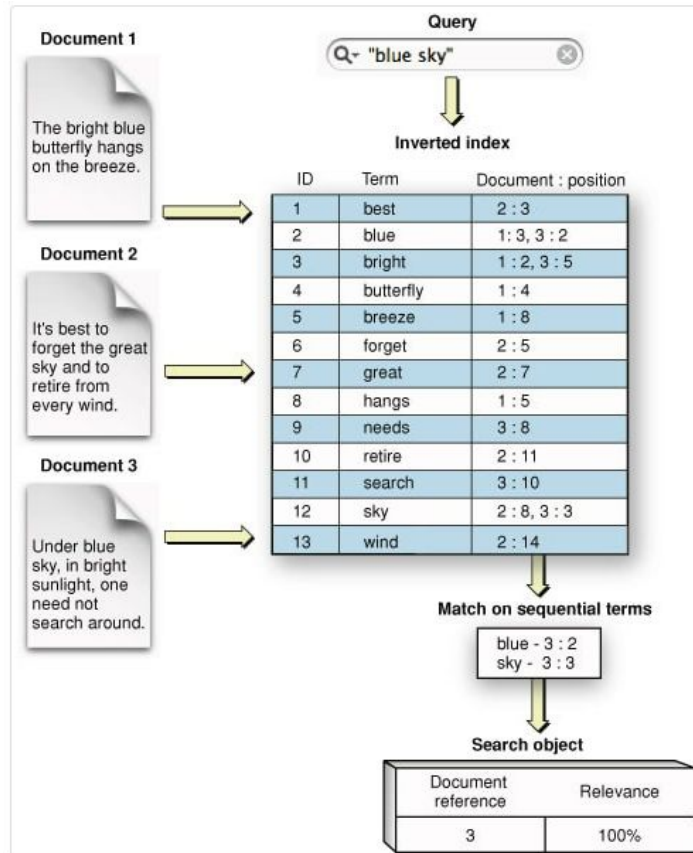
# Query time process

Also, this software will receive the input (sequence) and will be responsible for:

- Tokenize the input
- Compared with the index so it will generate some matches
- Sort results by using a score algorithm (elasticsearch uses TF/IDF as default)

# Thanks

Do you have any questions?

emmanuel.paez@tec.mx
Slack #module-5-nlp-1