

# Data Harmonization in AgroSciences

Author: Aarti Joshi



## What is Data Harmonization and its benefits?

Data harmonization refers to the process of integrating and standardizing disparate data formats into a single and unified datasets. It involves transforming data that is collected and stored in different Data Lakes, with varying naming conventions and data types, into a consistent, compatible, and consumable form.

Achieving a golden record through the process of harmonization significantly enhances the value and utility of a unified data view. Thus, data harmonization provides the business with a standardized and reliable representation of data from various sources, ensuring that all the business stakeholders within an organization or domain can rely on the same accurate information.

Large and complex datasets, prevalent in industries like business analytics, scientific research, and healthcare, greatly benefit from data harmonization. Replace the highlighted part with this. With the exponentially increasing volume of data in the AgroScience industry, the need for efficient data management and integration has become paramount. Addressing this challenge involves recognizing data harmonization as a pivotal aspect of data engineering, and it is widely acknowledged that harmonizing diverse datasets is essential to achieving a unified and coherent data landscape.



Data harmonization plays a crucial role in the field of agricultural science as well. Like other industries, agricultural science companies deal with vast amounts of data from various sources, such as crop yields, weather patterns, soil conditions, pest, and disease management, and more for field research, prototype trials etc. These datasets are often collected and stored in different formats, making it challenging to extract meaningful insights and make informed decisions.

## Why is Data Harmonization needed in the first place?

Cloud platforms have revolutionized the way organizations handle data, providing ample storage and computing capabilities. IT teams are no longer limited by infrastructure constraints when it comes to managing data. With these advancements, it has become increasingly feasible to consolidate an organization's data onto a modern data platform. However, it is crucial to ensure that the diverse data sets on this platform are interconnected and aligned through common fields.

Without integrating different data sets to obtain a comprehensive view, there is a significant risk of transforming the modern data platform into a chaotic data swamp. A data swamp occurs when data sets are simply ingested onto the platform without undergoing any data cleansing or curation techniques. To establish a well-governed modern data platform, it is imperative to apply data harmonization techniques to the data sets which helps the organization eliminate:

- **Heterogeneous Data Structures**
- **Inconsistent Data Types**
- **Missing Data**
- **Data Quality**
- **Analytical Consumption**

By harmonizing the data, organizations can ensure that the data sets are standardized, consistent, and compatible with each other. This process eliminates inconsistencies, duplicates and inaccuracies, enabling a coherent and reliable representation of the data. Through data harmonization, the modern data platform can maintain its integrity and provide a solid foundation for meaningful insights and informed decision-making.



## Data in AgroSciences

AgroSciences deal with various types of data sets that are relevant to agricultural and crop-related research, analysis, and decision-making. Some common data sets in AgroSciences include:

**Crop Yield Data:** Information on crop yields from different regions, farms, and field trials. This data helps analyze productivity, performance, and variations in crop production.

**Weather and Climate Data:** Data related to weather conditions, including temperature, precipitation, humidity, wind speed, and solar radiation. Weather data is crucial for understanding the impact of climate on crop growth, disease management, and irrigation scheduling.

**Soil Data:** Soil-related data, such as soil type, composition, fertility, pH levels, nutrient content, and organic matter. This information helps assess soil health, nutrient deficiencies, suitability for specific crops, disease management, and irrigation scheduling.

**Pest and Disease Data:** Data on pest infestations, disease outbreaks, and their impact on crops. This data helps identify and monitor pest populations, track disease incidence, and develop effective management strategies.

**Agricultural Practices Data:** Information on agricultural practices, including planting dates, irrigation methods, fertilizer application rates, crop rotation patterns, and pesticide usage. This data provides insights into farming techniques and their impact on crop performance.

**Genetic and Genomic Data:** Data related to plant genetics, genomics, and breeding. This includes genetic markers, DNA sequences, gene expression profiles, and genomic variation data. Genetic data helps in crop improvement, trait selection, and developing resistant varieties.

**Market and Economic Data:** Data on market trends, pricing, demand-supply dynamics, consumer preferences, and economic indicators relevant to the agro industry. This information assists in market analysis, forecasting, and decision-making regarding crop selection and marketing strategies.

**Remote Sensing and Satellite Imagery:** Data obtained through remote sensing technologies and satellite imagery, providing insights into vegetation indices, crop health, land cover changes, and spatial patterns. Remote sensing data aids in monitoring crop conditions, identifying stress factors, and optimizing resource allocation.



In AgroSciences, the data sets often consist of unstructured data, which presents unique challenges for data engineering and harmonization. Specifically, unstructured data in this context typically includes free-text fields found in AgroSciences data. Dealing with the abundance of textual information in AgroSciences necessitates the implementation of specialized data harmonization techniques.

## Dealing with lot of textual data

When dealing with a large number of text fields, data harmonization can indeed present challenges. Textual data can be diverse in nature, making it more complex to standardize and integrate into a cohesive framework.

To overcome the challenges associated with harmonizing a large number of text fields, organizations can utilize a range of techniques and tools in the data harmonization process. Natural Language Processing (NLP) plays a crucial role in extracting meaning from text data. NLP techniques like tokenization and named entity recognition help break down text into meaningful units, identify key entities, and understand the syntactic and semantic structures.

Machine learning algorithms can be employed to automate and streamline the harmonization process. These algorithms can be trained on existing data to learn patterns, identify similarities and differences in text elds, and make predictions or classifications based on the learned knowledge. This reduces the need for manual intervention and accelerates the harmonization process.

Text mining techniques enable organizations to extract valuable information from unstructured text data. This includes uncovering trends, patterns, and relationships within the text. Text mining algorithms, such as topic modeling, can provide deeper insights into the content and context of the text fields.



Entity recognition algorithms can identify and extract specific entities, such as names of crops, regions, geo graphics, medicines from the text elds. This allows for categorization and organization of the data based on these entities, facilitating effective integration and analysis.

Standardization methods play a vital role in ensuring consistency across text elds. This involves normalizing text by applying rules or conventions to format, spelling, capitalization, and abbreviations. Standardization enables harmonized text elds that can be easily compared, matched, and integrated.

By leveraging these techniques and tools, organizations can effectively harmonize text elds, transforming diverse textual data into a unified and structured format.

## Indexing text value variables

Once the text dataset has been harmonized, it is beneficial to undergo augmenting process to standardize and structure the text data. This involves techniques such as standardization, normalization, lemmatization, stemming, noise removal, entity recognition, and text enrichment. Standardization ensures consistent formatting, while normalization transforms the text into a common form. Lemmatization and stemming reduce words to their base forms, removing redundancy.

Noise removal eliminates irrelevant elements, and entity recognition identifies specific entities. Text enrichment involves adding additional information or metadata. Overall, augmenting enhances the uniformity, clarity, and analysis-readiness of the harmonized text data.

Upon augmenting the dataset, one can proceed with the indexing of the text value dataset. Assigning numeric keys to a text value dataset is a process known as "indexing". Indexing involves assigning unique numerical identifiers or keys to each distinct text value within a dataset. These keys serve as a reference to the corresponding text values and allow for more efficient data processing and analysis.

The indexing process typically involves creating a lookup table or index table that maps each numeric key to its respective text value. This table acts as a reference guide for translating between the numeric keys and the original text values, ensuring that the data can be interpreted correctly when needed. By using numeric keys for text values, indexing simplifies data handling, enables faster computations, and facilitates data integration across different datasets.

## Common dataset in AgroScience

Collaboration among academia, industry practitioners, policymakers, and civic society stakeholders often leads to the formation of consortiums aimed at utilizing and creating common datasets. One such example is CE-HUB.org, which offers global soil and weather data. Institutes and organizations can access this data through the consortium's open API. When integrating data from these APIs, it is crucial to ensure that the acquired dataset is harmonized and seamlessly integrated with the existing organizational data.

Seamlessly integrating the open dataset into existing datasets requires careful consideration of the formats and nomenclature used. The open dataset may have its own specific formats and naming conventions. Therefore, it is crucial to harmonize and align these formats and nomenclature with the existing datasets.

This ensures a smooth integration process, allowing for cohesive and consistent analysis across all data sources. By effectively integrating and aligning the formats and nomenclature, organizations can maximize the value and usability of the combined data sets while maintaining data integrity and accuracy.



## Data Harmonization is largely a Data Engineering problem



The advent of modern data platforms with fast computing capabilities and vast, scalable storage has revolutionized data quality improvement across organizations. These advancements enable powerful data processing tasks, such as performing Cartesian joins on datasets.

Additionally, machine learning models can be leveraged for data wrangling, allowing for the harmonization of massive datasets. By training these machine learning models on organizational data, they can automate data cleansing tasks, minimizing the need for extensive human intervention.

This means that even without deep domain knowledge, data engineering techniques can be employed to achieve data harmonization across the organization. Consequently, data harmonization can be viewed as a data engineering challenge rather than being specific to AgroSciences.

## About Modak

Modak is a solutions company that enables enterprises to manage and utilize their data landscape effectively. We provide technology, and cloud-agnostic software and services to accelerate data migration initiatives. We use Machine Learning (ML) techniques to transform how structured and unstructured data is prepared, consumed, and shared. Modak's portfolio of Data Engineering Studio provides best-in-class delivery services, managed data operations, enterprise data lake, data mesh, data fabric, augmented data preparation, data quality, and governed data lake solutions.



- **Find out more**

 [modak.com](mailto:modak.com)

- **Follow us**

 [LinkedIn.com](https://www.linkedin.com)

 [medium.com](https://medium.com)



**UAE**

Dubai Silicon Oasis, Dubai,  
UAE, 341041



**USA**

21660 W Field Pkwy,  
Deer Park, USA, IL 60010



**INDIA**

Elemental #337, Financial  
District, Nanakramguda,  
Telangana, India, 500032



**USA**

312 S 4th St, Louisville,  
USA - KY 40202

