



Data Warehousing and Its Future Role in Weather Forecast

Anushka Patil

SY BBA BA

School of Management

Maharashtra Institute of Technology WPU

Dr. Sudepta Banerjee

Faculty of Management

Maharashtra Institute of Technology WPU

Date of Submission: 05-10-2023

Date of Acceptance: 18-10-2023

Abstract

This research paper explores the implementation of a data warehouse in the field of weather forecasting. The study focuses on the use of data warehousing technology to enhance the accuracy of weather predictions and the efficiency of meteorological data management. Through a combination of historical data analysis, real-time data integration, and the development of subject-specific knowledge systems, this research demonstrates how data warehousing can significantly improve the quality of weather forecasts. This study sheds light on the potential benefits of using data warehousing technology in meteorological research and underscores the importance of data-driven approaches in the field of weather forecasting. The findings of this research have practical implications for improving weather prediction models and advancing our understanding of weather patterns.

Keywords: Data warehouse, weather forecasting, meteorology, data integration, predictive modelling, accuracy, efficiency, knowledge sharing.

I. Introduction

Access to data is the main purpose of today's forecasters' workstations. These systems typically function as graphical platforms with integrated analysis tools for visualizing both unprocessed data and derived datasets. However, these systems face the following three major difficulties:

Data underutilization: Although these systems have the capacity to store and retrieve enormous amounts of data, frequently reaching 2GB and spanning thousands of fields on a daily basis, forecasters typically only use a small portion of this

data, typically less than 1%, for operational forecasting.

Management of Historical Data: As more sophisticated storage systems are developed; a larger collection of historical data is becoming available for operational forecasting. The difficulty lies in properly maximizing the use of this previous data.

Integration of Special Observation Systems: Along with traditional data sources, new specialized observation systems like AWS, profilers, GPS/Met, storm lightning locators, and others are beginning to take shape. While these systems offer useful data for the public good, it takes adaptable data management solutions to incorporate them into operational forecasting procedures with ease.

Underutilizing the data that is already available is the most important problem with these difficulties. Implementing a data warehouse, the next-generation of decision support systems (*DSS*), could provide a potential remedy for these problems. Such a system might potentially provide a solid response to these issues and boost the effectiveness of operational forecasting procedures.

Data Warehouse – A Solution

W. H. Inmon proposes a data warehouse in his 1993 book "Building the Data Warehouse." The initial definition of a data warehouse was provided by him: "A warehouse is a subject-oriented, integrated, time-variant, and non-volatile collection of data supporting management's decision-making process."

To aid managers in their decisions, data warehouses are initially utilized in commercial enterprises. Data warehouse usage is expanding across a wider range of industries these days, including many scientific ones. But why should we



worry about a technology that was created primarily for business? Why should we employ data warehouse technologies in the field of weather forecasting?

To help short-range forecasters analyse a substantial volume of diverse data in a short amount of time, while ensuring that they use more of the available data, you can implement a combination of strategies and technologies:

Automated Data Processing:

Streamline the data collection, decoding, and quality control processes by automating as much as possible. Use intelligent algorithms to identify and correct errors, reducing the need for manual intervention.

Data Aggregation:

Aggregate data from various sources into a unified format. This simplifies data analysis by presenting forecasters with a consistent data structure.

Data Summarization:

Develop data summarization techniques that generate concise reports or visual summaries highlighting key trends, anomalies, and important data points. This allows forecasters to quickly grasp the most critical information.

Real-time Data Streaming:

Utilize real-time data streaming and processing techniques to continuously update forecasters with the latest information. Data can be analysed as it arrives, reducing the time lag between data collection and decision-making.

Data Sampling and Prioritization:

Implement smart data sampling strategies that focus on the most relevant data for the forecast. Develop prioritization algorithms to ensure that forecasters receive the most critical data first.

Machine Learning and AI Assistance:

Employ machine learning models to predict patterns and trends in the data automatically. AI algorithms can identify anomalies, highlight significant data points, and provide valuable insights to guide forecasting decisions.

Data Visualization:

Utilize interactive data visualization tools to represent data in a more digestible format. Visual representations can help forecasters quickly identify important patterns and trends.

Predictive Modelling:

Develop predictive models that can assist forecasters in making decisions based on historical data and real-time observations. These models can provide recommendations and assist with scenario analysis.

Customized Dashboards:

Design customized dashboards that present forecasters with the most relevant data, analysis tools, and decision support features, allowing them to work efficiently and access the data they need.

Data Archives and Historical Data Access:

Maintain an easily accessible archive of historical data for reference and analysis. Quick access to historical data enables forecasters to compare current conditions with past events.

Collaboration and Expert Networks:

Encourage collaboration and communication among forecasters and experts to share insights, validate data, and collectively assess the situation.

Continuous Training and Familiarization:

Regularly train forecasters on new data sources, analysis tools, and technologies. Familiarity with the full range of data resources is critical to effective decision-making.

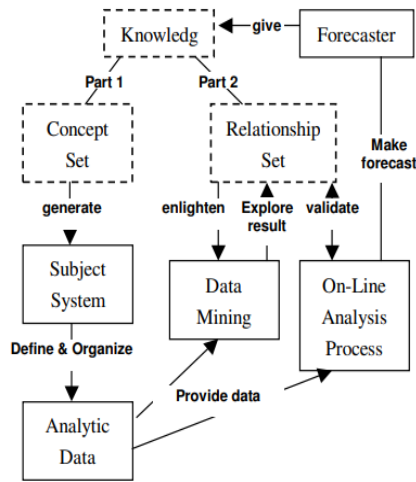
Feedback Mechanisms:

Establish feedback mechanisms for forecasters to provide input on system improvements, data quality, and the effectiveness of data utilization. Use this feedback to refine the system.

The goal is to create an environment that enables forecasters to make informed decisions quickly, even when dealing with vast datasets. By implementing automation, smart data management, predictive analytics, and user-friendly interfaces, forecasters can more effectively use a broader range of data within the time constraints of operational forecasting.

How does the system encourage forecasts to use more data? The conventional approach involves giving some substantial processing data to condense data information or developing some automatic prediction methods to save time, such as physical parameter diagnosis, conventional statistics, NWP models, MOS, MOD, etc. However, the current sophistication of the above instruments is currently insufficient to substantially lessen the labor of forecasters.

The utilization of a data warehouse can significantly benefit forecasters when it integrates the expertise of experienced forecasters into the system. By incorporating the knowledge of seasoned forecasters into the database, facilitating the accumulation of new knowledge, and offering an interactive graphical system to access and apply this knowledge effectively, forecasters can expedite data analysis and make more extensive use of data in their operational forecasting processes.



We define *Subject* in the context of forecasting as the focus of analysis, which correlates to a concept in the forecaster's understanding. The phrase "Analytic data" refers to a real value related to a topic. Based on the subject's definition, it arises from the preliminary operational data.

We use the idea that is set in the forecaster's expertise to create a subject system. For each subject in this subject system, we specify the data transformation mechanism that turns operational data into analytical data. Real-time analytical data is acquired daily during the operational process and is saved in the database. In this method, we effectively embed the forecaster's concept set by establishing a link. By doing this, we effectively incorporate the forecaster's concept set into the database by making a relationship between concept, subject, and analytical data. Forecasters can then examine the analytical evidence that supports the forecasts they have in mind.

Other terms used in Inmon's definition include: *integrated*, "which means data warehouse ought to extract data from every data source that contain useful data for the DSS mission; *time-variant*, "which means data warehouse should contain an enormous quantity of historical data to enable analysis of the climate change and rules over a period of time and *non-volatile*, "which indicates data in data warehouse are only for evaluation, no more low-level process (insert, modify, delete, etc.).

Two crucial tools for analysis are located in the right half.

DM stands for data mining. It is a tool for exploration. It automatically examines the connections between the subjects (i.e., forecaster concepts) in the analytical data set. The exploratory process is typically specified by forecasters, initiated by forecasters or automatically, and guided by their expertise. To enable forecasters swiftly build up

their expertise, the result relationship from DM will be reinforced into forecasters' understanding.

The *On-Line Analysis Process (OLAP)* is an addition. It is an interactive tool for validating. It serves as a tool for forecasters to view data, validate relationships (which may incorporate forecaster's assumptions and DM findings), and then make forecast judgments. Multidimensional analysis is the technique at its core. It will serve as forecasters' primary workstation.

In the present, DM and OLAP make the data warehouse a perfect solution. Analytical data storage with DM and OLAP together make data warehouse. For the sake of this research paper, we shall refer to Inmon's definition of data warehouse as "Data Storage".

Data Storage

Data warehouses depend on data storage. It manages these data, transforms operational data into analytical data, and gives applications, particularly DM and OLAP, accessibility to tools or other support tasks.

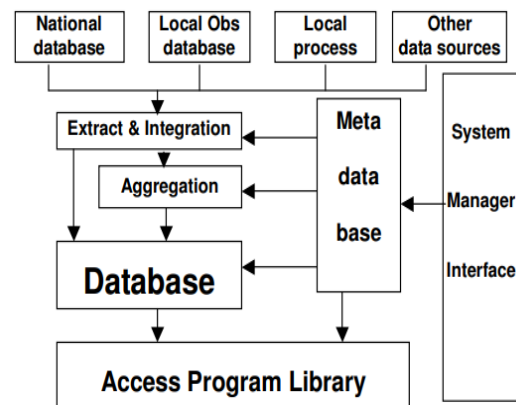


figure 2
Flowchart of data store system

Gathering forecaster knowledge is the first step in constructing a data warehouse so that the subject system may be built.

The data warehouse's subject system is its most crucial component. Subject system is the foundation for converting operational data to analytic data in data storage, as well as the foundation for a logical data model to store data. It serves as the foundation for the multi-dimensional analysis dimension set in OLAP. It serves as the foundation for the item set for association rule exploration in DM.



Currently, our subject system includes six different categories of subjects, as far as forecasters are aware:

1. The subject system must have the goal forecast component, for example, rainfall, temperature, humidity, and wind.
2. Data from an index - Some Index data are helpful; often, they serve as flags or forecasters of weather to come.
3. Numerical data-The subject system should comprise topics like the average of five stations, the greatest temperature gradient in a particular location, the minimum relative humidity over the previous five days, etc., to analyze statistical features in a region or over a time period.
4. Reconstructed data. Orthogonal function (EOF) transformation occurs to coefficients from orthogonal transformations, such as wavelet transformation. And data that has been filtered using lowpass, high pass, etc.
5. The weather. Even when analyzing model output, a forecaster's understanding must include an overall perspective. Therefore, the subject system should incorporate the weather system.
High, low, large gradient area, saddle area for scale element, convergence center, and divergence are examples of weather systems. federal database Area Observation Database Local Process Other sources of data Integrate & Extract Aggregation Database Meta database Visit the program library System Manager Interface jet stream, jet center for vector, share line, and center They need to be simpler to recognize. in computing.

6. Conceptualized models. They are passive subjects. To create conceptual models, they are made up of the aforementioned five sorts of subjects. The topics listed above are all based only on forecast knowledge. It can be expanded to cover greater ground. For instance, a system management area or a forecast service area.

The subject system in our design has a tree-like structure. The most specific idea corresponds to a leaf node. Subjects will be generalized to a wider range starting with leaf nodes. The broadest concept corresponds to the root. In accordance, the names of subjects at various tree levels may contain comparative terms like "more," "most," "less," "least," "better," "best," "bigger," "biggest," etc. When forecasters input the subject, they may choose which tree level the subject is at. The following section demonstrates how the higher-level subjects can also be generated automatically from the lower-level subjects.

We'll see in the sections that follow why such a structure is essential for DM and OLAP.

Transforming the data into analytical data

Data storage requires information from data sources, then utilizes an integration process to evaluate the information's quality, eliminate duplicates, and standardize the information so that all the information has the same meaning, unit, measurement, accuracy, etc.

The data will next be transformed and aggregated into analytical data. Translating the original data into the analytical data pertaining to one of the six type points is the first stage. Making an aggregate based on the altered data is the second stage.

According to our method, aggregation is a typical strategy used in a typical data warehouse. It involves doing computations on converted data within predetermined geographic or temporal boundaries, including counting, adding, averaging, and determining the maximum and lowest values. With grid intervals of 1, 2, 5, 10, or 20 degrees in longitude and latitude, we specifically build networks in space. Then, summaries are computed for each network box. Aggregation data is calculated for each of the time periods that time series data is broken into, which can range in length from 1 day to 30 days.

In order to provide a hierarchical picture of the data, our subject tree is often organized from shorter grid intervals or time periods (known as precise granularity) to longer ones (known as rough granularity).

To further enrich the subject tree for particular types of concepts, we have additionally created specialized algorithms designed specifically for weather forecasting in addition to these methods. These techniques add to the outcomes of the common aggregation procedures, resulting in a subject tree that is both more detailed and thorough. Data will be changed using the following manner for the six different subject kinds.

1. *Forecast goal component* - Create the subject tree by calculating the analytical data from the original data as the definition of the forecast target data. Use the fuzzy membership function to set different membership levels.
2. *Index information* - The calculation is performed in the same manner as before, but using index data definitions.
3. *Quantitative data* - The outcomes of statistics are the analytical data. Subject trees are created using the same process as aggregation; no other method is used.
4. *Reconstructed data*- Make wavelet and EOF transformations to all of the data fields, and then take a vector of the first 10 coefficients to



represent the analytical data. From keeping all 10 coefficients (leaf) to just keeping the first coefficient (root), the subject tree is created. With regard to filtered data, subject trees are created from the most angular (leaf) to the most rounded (root) fields.

5. *The weather* - The characteristics such as location, area, aspect, direction, and length of the long axis, intensity, longitude and latitude coordinate of the vertexes of the characteristic line, etc. are taken as the analytical data to create a vector and are used by the computer to identify the weather system in accordance with their definition. Similar to the spatial entities of geophysical information, weather system data can be aggregated, stored, and employed utilizing various findings from the special data warehouse. The vector components between real data and idea pattern are matched by some match function, and the subject tree is generated according to the match level. Set a few thresholds, with leaf being the highest and root being the lowest.

6. *Theoretical model* - They are made up of the first five types. According to the forecaster's expertise, analytical data is a vector made up of the attributes or values of each component and provides an idea pattern of the conceptual model. Be aware that certain complicated weather systems fall under this umbrella. The weather system for a cold front, for instance, consists of a wind shear line, a temperature high gradient area, a positive pressure changes behind the shear line, and a negative pressure change before to the shear line. It is the cold front conceptual model. The subject tree is created using the same match level as earlier.

Using metadata to manage systems

Metadata are data about data. All of the data in the data warehouse are fully specified in the metadata, including the data source, the characters, the quality, accuracy, the transformation history, the location, the access parameters, etc. It is a data warehouse terminology.

In a data warehouse, the metadata base still contains a lot of information about the system structure, such as the server's IP address or URL, the access username and password, directories, definitions for data filenames and data formats, descriptions of the algorithms, functions, and parameters for data transformation, etc.

Data storage system will be developed using this metadata basis and by offering a related data access service function. subsequently, metadata aids in managing and usage of data warehouses. Additionally registered in this information store is the subject system. All metadata is in XML. There is a DTD kernel. ENTITY is a participant in DTD.

DTD/ENTITY serves the same purpose as a pointer. Each ENTITY identifies a different part of the metadata base. Applications are able to obtain data from any component. As a demonstration, an application takes a data, utilizes the subject tree first, locates the relevant data set, establishes the server and location of the data set, locates the data by filename definition, reads the data from the file according to with the data format, and uses the data in line with the data space character and time character.

When system management applications operate data ingest, they at first examine the data moving principle in the data set, identify the source server and original data location of the data set, obtain data by filename definition, copy the data file or read the data according to the data format, and then put the data in the desired location.

Subject Graph

DTD/ENTITY

Sets of data

Filename on a Data Server Definition

Data Format Room

Time Character

According to the pre-process definition in the subject tree, pre-process changes the original data into analytical data.

Schema will take the place of DTD in the future.

Metadata

Metadata makes data management and maintenance simpler and more versatile, making it easier for applications to access historical data as well as diverse data types such special observation data and even Internet data. Several data warehouses spanning across a number of provinces may be connected as one data warehouse employing the Internet if a uniform schema is developed.

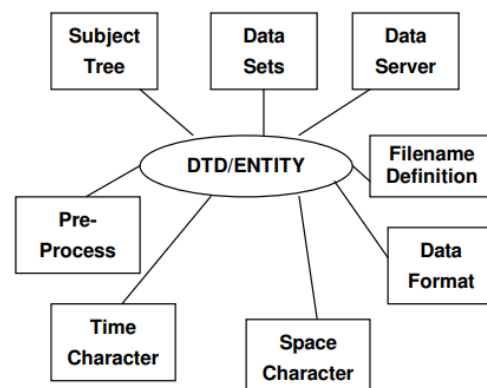


Figure 3
Structure of metadata



Data Storage

The relational database and the multi-dimensional database are the two types of data stores used in data warehouses. To make OLAP's multi-dimensional analysis more convenient, dimension tables for relational databases should be built alongside data tables. Both star mode and snowflake mode are available.

A file comprising a multi-dimensional array is known as a multi-dimensional database or cube file. When it is utilized in an application, the entire thing will be read into memory and used to create a multi-dimensional array there. Following that, OLAP will implement cube manipulations.

Typically, the OLAP data are kept in cube files, whereas the historical data are kept in relational databases. All of the data is maintained as cube files in our design.

Data Mining

Data mining investigates hidden patterns, models, and association linkages in databases. using knowledge Discovery in Databases" (KDD) technology, it explores the relationships between the forecaster's concepts and speeds the accumulation of forecaster knowledge on unused data.

Though the famous data mining methods such as multi-analysis, Bayes statistics, decision tree, classifying and clustering etc. are available usually forecasters prefer exploration of association rule method. It is a logic implication relationship.

Let C item set $I = \{i_1, i_2, \dots, i_m\}$, the item i_k ($k=1, m$) may be any weather event. Let $T = \{t_1, t_2, \dots, t_n\}$ is a historical sample set, any sample t_j ($j=1, n$) is a subset of I. If set C and D are all subset of I, and the intersection of C and D is empty. Then the implication relationship $C \rightarrow D$, or statement "If C then D" is an association rule.

Ex: let C is "There is a vortex near Place X" and D is "It will be rainy in Place X", correlation relationship requires if a vortex is near Place X then Place X will be rainy, if no vortex near Place X then no rainy in Place X. As many weather conditions may cause rain in Place X not only vortex, the second half above is wrong. Thus, we will not get a good correlation. But our implication relationship only requires the first half. So, searching implication relationship is more rational than correlation relationship for meteorological data.

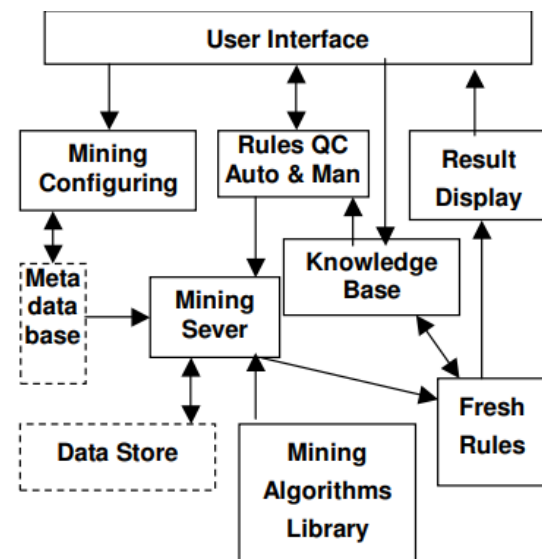


Figure 4
Flowchart for data mining

OLAP

OLAP enables forecasters to examine and analyze data from several angles.

Instead of just viewing data value or data pattern, OLAP's major goal is to view data to discover relationships or patterns in data.

Data is presented in multidimensional form in a cube which consists of dimensions and measures. Forecasters prefer to view data in graphs rather than table format. Here they can slice, dice and pivot data for analysis. They can analyze data at aggregated level as well as detail level using drill-down and roll-up operations.

Along with multi-dimensional analysis for finding more relationships, there are other 3 analysis tools based on multi-dimensional analysis: Compare Analysis (CA), Multi-Analysis (MA) and Analog Analysis (AA) in our design. **Multi-dimensional analysis** – Here selection of dimensions is prime, manipulation "Slice", "Dice" and "Pivoting" are all realized by this function.

We can divide dimensions into 3 kinds.

First is frame dimension which includes spatial, time, the coefficient of subject type. It allows the forecaster to examine the data in a spatial-temporal frame or an orthogonal base frame derived through an orthogonal transformation. The second is subject dimension which includes all subjects in subject tree. They are the content displayed in above frame. The third is element dimension which includes original data in database, such as wind, rain, temperature, height, etc.



Selection boxes for these three dimensions are displayed in OLAP.

When using graphical display principles, the sequence of the chosen dimensions is crucial.

1. If the forecaster only chooses two dimensions, the first dimension must be a frame dimension and be an X axis, while the second must be a Y axis and display the curve of Y by X. Y must not be a frame dimension.

2. If the forecaster chooses three dimensions, the X and Y axes serve as the first and second axes, respectively, and the third dimension's contour is displayed. The first two dimensions must be frame dimension, while the third dimension cannot be a frame dimension.

3. If the forecaster chooses more than three dimensions, and all of the additional dimensions are frame dimensions, a matrix of contour maps made from the first three dimensions will be displayed. Maps will show overlapped dimensions if there are more than three subject or element dimensions.

4. The scanning ruler of the forecaster can display for a few different frame dimensions. Forecasters can move the ruler mark using the mouse when scanning the ruler. It is simple to modify the dimension value. In the system interface, the rulers of space and time are constantly visible.

Multi-Analysis (MA)

Enter the Multi-Analysis mode if all the specified dimensions, including the graphics frame dimensions, are subject dimensions or element dimensions. The spatial dimensions of each object or component in this situation are all given fixed values.

The chosen dimension is referred to as "factor".

Different factors could each have a unique fixed value for the spatial dimensions.

The system displays a scattered map with the first factor as the X axis, the second factor as the Y axis, and the third factor's value (or its location in space or time, as determined by forecasters) filled in at each point in time.

Compare Analysis (CA)

A multi-dimensional study with multiple data sources is a compare analysis. For all of the compared data, the frame size should be the same. We can choose from the original data or the comparative data in the subject tree. The data will overlap in curve or contour maps if we choose multiple data points in the topic tree. The graphics of the same subjects at different moments in time will overlap in a curve or contour map.

The compare queue contains all of the compared data. In the comparison queue, the forecaster can add or remove data.

Base data is one of the data in the queue. The data will be compared to other data in the queue. The base data can be chosen by the forecaster.

Forecasters can

1. calculate the average and variance of the data in queue and highlight them

2. Can calculate the difference between each data in queue and the base data and highlight.

3. Can transform the data in the queue, such as normalization, filtering, orthogonal transformations, and compare the results.

4. Can calculate count, frequency, average and variance of the value in specified interval of data in queue.

Analog Analysis

A unique type of comparison analysis is analog analysis. Selecting the historical data that analog index with current data is greater than a threshold creates its comparison queue.

Forecasters can choose the analog index, threshold, time period and spatial range.

OLAP's goal is comparable to that of data mining.

Data mining automatically identifies relationships between subjects, whereas OLAP employs a human to carry out the same task. So, both data mining and OLAP use the same interface. Forecasters view data in OLAP to get new information, which they subsequently add to the DM knowledge base. The results from DM will be shown in the interface to provide the forecaster with hints. Online analysis mining (OLAM), a novel technology, has emerged during the past few years.

It integrates OLAP with DM.

Data is extracted and integrated from numerous data sources according to time or event. The system then transforms and aggregates the original data into analytical data in accordance with the definitions of the topics. Following that, analytical data are kept in a data storage system. DM investigates connections in a data storage database that are time- or event-driven using forecasters' setups. The DM knowledge base illuminates the exploration process and records the findings.

Building subject systems and entering relationships into knowledge bases are the initial methods used by the forecaster side to enter knowledge into the system. In his routine job, he is free to add new ideas as they come to him, or through summaries of his case studies or other research, to the subject system and knowledge base.



The forecaster uses an interactive interface to examine or modify system configurations, manage the quality of relationships in the knowledge base, validate relationships with OLAP tools, and conduct case studies or other types of research using OLAP and DM features.

In operational forecasting, the forecaster views data in analog analysis visuals from OLAP or traditional maps. He will automatically receive hints from the DM. He may choose a few subjects and the appropriate data range, give DM a task to temporally explore, and then evaluate the findings in OLAP to confirm some relationships required for the current forecast mission.

From the perspective described above, data warehouses may be seen as a novel type of Artificial Intelligence (AI) system that integrates database and weather graphics technologies. It aids forecasters in gathering, organizing, and employing their expertise in operational forecast.

Since the operational forecast mission is extremely complicated, it is obvious that the data warehouse will not become the entire system of the forecaster's workbench. However, it is a system to compile, organize, and exploit forecasters' information. It ought to be the core innovation of the next-generation forecaster's workbench.

II. Conclusion

In conclusion, a data warehouse plays a crucial role in the toolkit of the modern forecaster by integrating seamlessly meteorological data management with AI-driven analysis. Data transformation and extraction from multiple sources are the first steps in the operational process, and the analytical data that results is then stored in a unified data storage system. With the help of forecasters' knowledge and configurations, data mining (DM) reveals substantial connections within this database.

Initially, knowledge is inputted into the forecaster by building a subject system and connections within a knowledge base. Forecasters can dynamically increase their expertise in daily operations by adding new concepts and relationships as they emerge. Forecasters can manage aspects of the system with an interactive interface, verify relationships with OLAP tools, and carry out case studies with OLAP and DM features.

Forecasters use conventional maps and comparison analytical images from OLAP during operational forecasting while getting automated guidance from DM. Additionally, they can choose subjects and data ranges, ask DM to investigate relationships relevant to their present mission, and confirm the outcomes in OLAP.

From the point of view of operations, the data warehouse is positioned as a cutting-edge Artificial Intelligence (AI) system that combines database functionality with meteorological visual technology. Its main purpose is to make knowledge management, utilization, and accumulation in the field of operational weather forecasting easier. Due to its complexity, it cannot completely substitute the forecasting process, but it does serve as a foundational technology for the next-generation forecaster's workbench, providing a more adaptable and effective method of information management.

References

- [1]. Aberer, K, Hemm, K, 1996, A Methodology for Building a Data Warehouse in a Scientific Environment, First IFCIS International Conference on Cooperative Information System, Brussels, Belgium
- [2]. Ayyer, K, Data Warehousing and Knowledge Management,
- [3]. Gorawski, M, Malczok, R, 2004, Distributed Spatial Data Warehouse Indexed with Virtual Memory Aggregation Tree, Proceedings of the Second Workshop on Spatio-Temporal Database Management
- [4]. Building the Data Warehouse, John Wiley & Sons, Inc. Nestorov, S, etc. Ad-Hoc Association-Rule Mining within the Data Warehouse,
- [5]. DATA WAREHOUSING AND ITS POTENTIAL USING IN WEATHER FORECAST Xiaoguang Tan* Institute of Urban Meteorology, CMA, Beijing, China