

ASSIGNMENT 2 FRONT SHEET

Qualification	BTEC Level 5 HND Diploma in Computing		
Unit number and title	Unit 17: Business Process Support		
Submission date	10/04/2024	Date Received 1st submission	
Re-submission Date		Date Received 2nd submission	
Student Name	Trinh Minh Hieu	Student ID	BH01236
Class	IT0601	Assessor name	Dinh Van Dong
Student declaration I certify that the assignment submission is entirely my own work and I fully understand the consequences of plagiarism. I understand that making a false declaration is a form of malpractice.			
		Student's signature	Hieu

Grading grid

P5	P6	P7	M3	M4	D2	D1

 **Summative Feedback:**

 **Resubmission Feedback:**

Grade:

Assessor Signature:

Date:

Internal Verifier's Comments:

Signature & Date:

Table of Contents

I. Introduction..... 5

1. Organization:..... 5

2. Role 5

II. Body..... 5

P5 Discuss how tools and technologies associated with data science are used to support business processes and inform decisions..... 5

1. How Data Science Tools Drive Business Processes and Inform Decisions..... 5

2. Benefits of Data-Driven Decision Making: 11

3. Support tools for each stage 13

P6 Design a data science solution to support decision making related to a real-world problem. 17

1. Data Science Solutions. 18

2. Solution Overview 20

3. Decision-making system overview architecture 21

4. Problems encountered when collecting data 25

P7 Implement a data science solution to support decision making related to a real-world problem. 26

1. Data Collection 26

2. Data Cleaning and Preprocessing 29

III. Conclusion 38

IV. References 39

List of Figure

Figure 1 Personalized customer experiences	6
Figure 2 Business Intelligence and Reporting.....	7
Figure 3 Data Collection	8
Figure 4 Predictive Analytics	9
Figure 5 Data Analysis and Decision Making	9
Figure 6 Analytics Modeling.....	10
Figure 7 Increased Accuracy and Precision.....	11
Figure 8 Enhanced Efficiency and Productivity.....	12
Figure 9 Improved Customer Experience and Satisfaction.....	12
Figure 10 Competitive Advantage and Innovation	13
Figure 11 Production Data	28
Figure 12 Code Import Pandas.....	29
Figure 13 after generate in the terminal	30
Figure 14 Missing Data.....	31
Figure 15 Code remove erroneous data	32
Figure 16 After removed by pandas	33
Figure 17 Duplicate data	34
Figure 18 Code to remove duplicated data	34
Figure 19 After Removed duplicate data.....	35
Figure 20 Data when merged	36
Figure 21 Code checking	37
Figure 22 Result Checking.....	37
Figure 23 Data filtering function.....	37
Figure 24 Data before querying	38
Figure 25 Data after query	38

I. Introduction

1. Organization:

ABC Manufacturing is a multinational company specializing in the production and distribution of consumer electronics. The organization faces numerous challenges in managing its supply chain efficiently and effectively. However, by leveraging data and information to support its business processes, ABC Manufacturing successfully addresses these challenges and achieves significant improvements in its operations.

One major implication of using data and information in ABC Manufacturing's supply chain is the ability to forecast demand accurately. By analyzing historical sales data, market trends, and customer preferences, the organization can anticipate future demand patterns and adjust its production and inventory levels accordingly. This helps minimize stockouts, reduce excess inventory, and optimize the utilization of resources, leading to cost savings and improved customer satisfaction.

Furthermore, ABC Manufacturing utilizes real-time data from sensors and IoT devices installed in its production facilities and logistics network. These devices collect data on equipment performance, energy consumption, and transportation routes. By monitoring and analyzing this data, the organization can identify bottlenecks, optimize production processes, and streamline logistics operations. For example, if a particular machine shows signs of malfunction through real-time data analysis, proactive maintenance can be scheduled to prevent costly breakdowns and production delays.

Data and information also play a crucial role in ensuring product quality and compliance with industry standards. ABC Manufacturing collects data at various stages of the production process, including quality control checkpoints and post-sales customer feedback. By analyzing this data, the organization can identify potential quality issues, implement corrective actions, and continuously improve its manufacturing processes. This not only enhances product quality but also reduces the risk of product recalls and associated costs.

Another implication of using data and information in ABC Manufacturing's supply chain is the ability to collaborate effectively with suppliers and distributors. Through the integration of data systems with key partners, the organization gains real-time visibility into inventory levels, production schedules, and transportation status. This enables proactive coordination and timely decision making, resulting in improved order fulfilment, reduced lead times, and enhanced overall supply chain performance.

2. Role

You were recently promoted full-time by ABC Manufacturing to Junior Analyst. Part of your role is to analyse the company's business processes and evaluate how the use of data and information can enhance their operations.

II. Body

P5 Discuss how tools and technologies associated with data science are used to support business processes and inform decisions.

1. How Data Science Tools Drive Business Processes and Inform Decisions

Data science tools play a crucial role in driving business processes and informing decisions. By leveraging data and advanced analytics techniques, organizations can gain valuable insights and make data-driven decisions to enhance their operations and achieve their business objectives. Here are some ways in which data science tools contribute to business processes and decision-making:

a. Supporting Business Processes:

Business process optimization is the application of strategies, tools and techniques to improve the efficiency and effectiveness of business operations. This includes using technology to optimize processes, automate tasks, and improve operations and customer interactions.

Operations approach: Data science tools can help organizations optimize business operations by analyzing data about their organization and internal processes. By reviewing data on business performance, lead times, and other metrics, organizations can identify issues and adjust processes to increase efficiency and reduce lead times. Data science tools can also help organizations make decisions to improve business performance by identifying patterns for operational optimization.

Personalized customer experiences: Data science tools allow organizations to analyze customer data, including purchase history, browsing behavior and personal information. Using machine learning tools, organizations can segment customers and create personalized experiences based on their needs. For example, e-commerce companies can use data science tools to increase conversion rates and customer satisfaction by recommending products based on customers' search and purchase history.

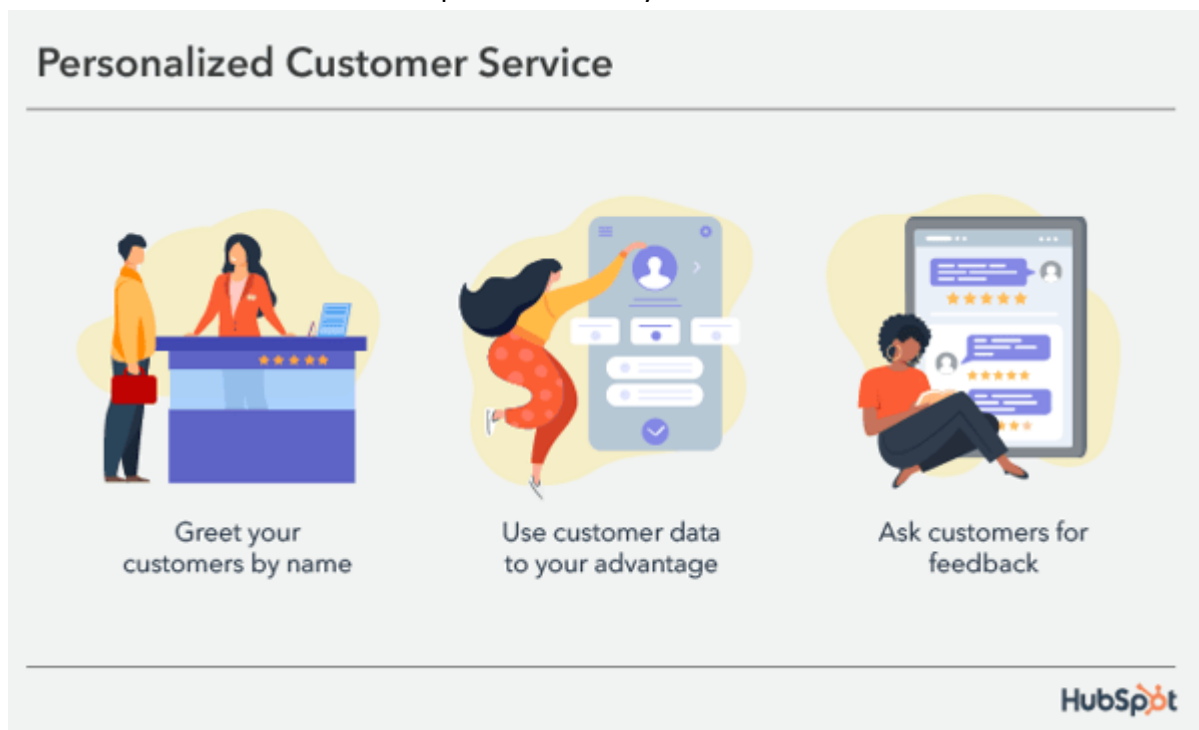


Figure 1 Personalized customer experiences

Forecasting and Adjusting Inventory: Data science tools play an important role in forecasting demand and adjusting inventory levels. By analyzing historical data on sales, market trends, and other factors, organizations can forecast demand and adjust inventory levels accordingly. This helps avoid shortages or excess inventory, reduces costs and improves operational efficiency. Data science tools can also help identify patterns and seasons in demand, helping organizations make more accurate inventory restocking decisions.

b. Informing Decision-Making:

Informing decision-making is a critical aspect of supporting business processes. It involves providing relevant and timely information to decision-makers, enabling them to make informed choices that align with the organization's

goals and objectives. Here are key points on how informing decision-making supports business processes:

Performance Metrics and Key Performance Indicators (KPIs): Establishing performance metrics and KPIs is crucial for measuring the success and progress of business processes. These metrics provide quantifiable indicators of performance and help decision-makers track the effectiveness of implemented strategies. By monitoring KPIs, decision-makers can identify areas that require improvement and take appropriate actions.

Business Intelligence and Reporting: Business intelligence tools and reporting systems play a vital role in informing decision-making. These tools gather data, generate reports, and create visualizations that present information in a meaningful and easily understood way. Decision-makers can use these reports to gain insights, identify patterns, and make data-driven decisions.

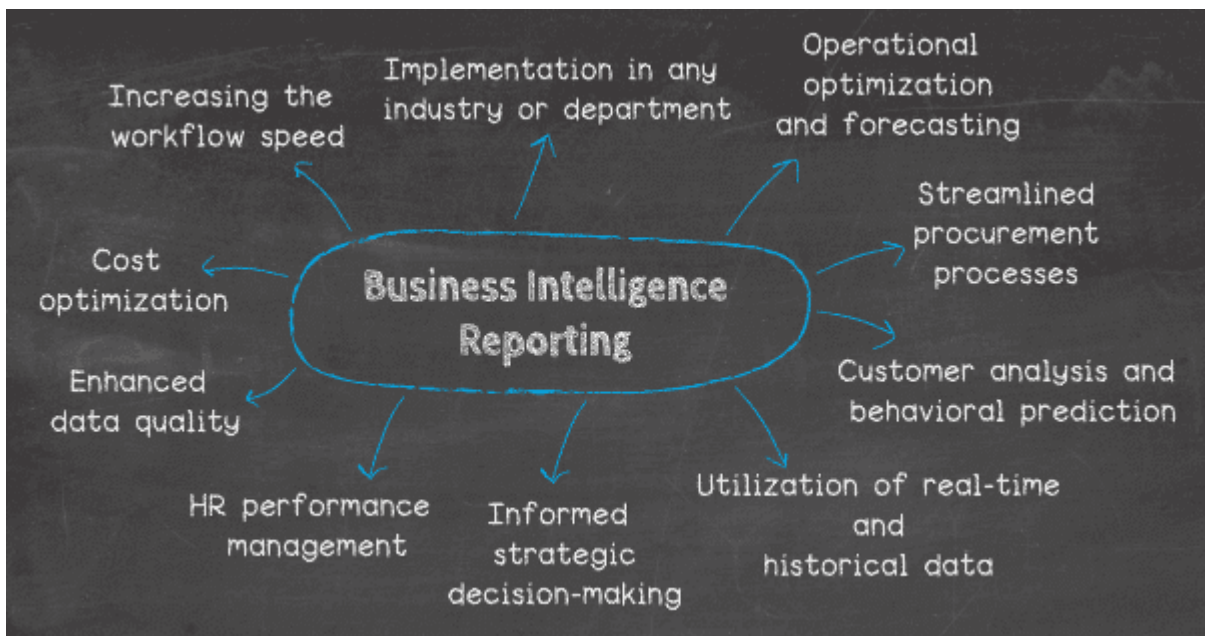


Figure 2 Business Intelligence and Reporting

Cross-Functional Collaboration: Supporting decision-making requires collaboration across different departments and teams within the organization. By fostering cross-functional collaboration, decision-makers can access diverse perspectives, expertise, and insights. This collaborative approach ensures that decisions consider multiple viewpoints and are aligned with the overall business strategy.

c. Key data science tools and technologies:

Data Collection and Integration:

For data collection, tools like Apache Kafka or Apache Nifi can be used to continuously and efficiently gather data

from various sources.

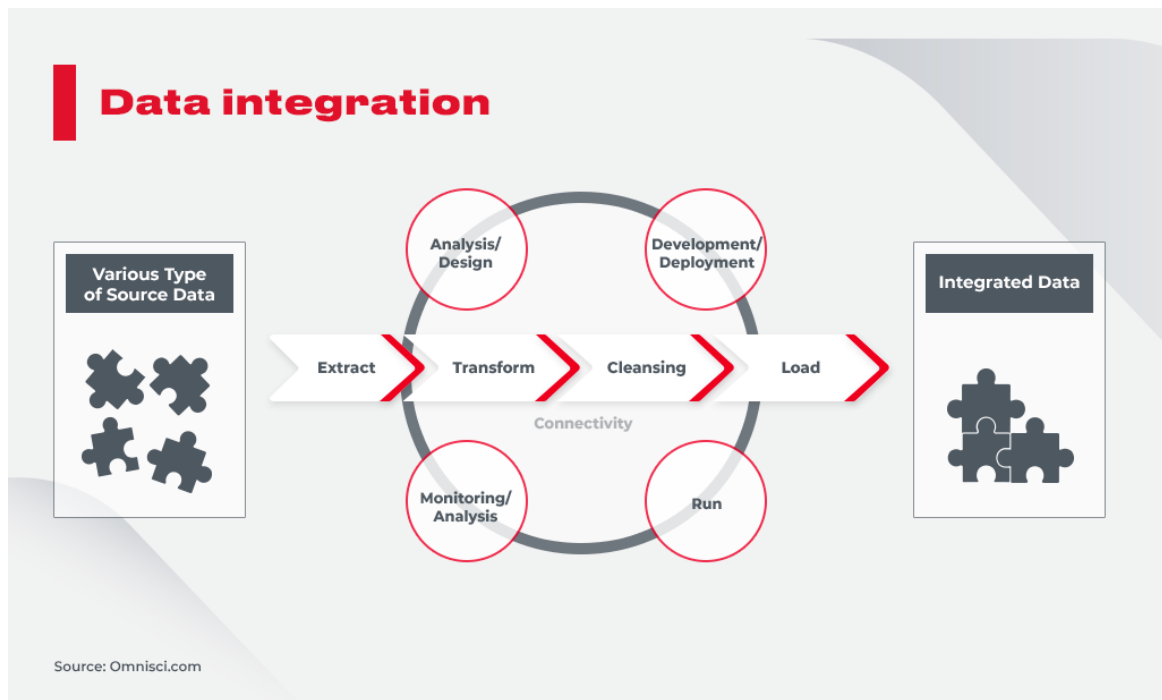


Figure 3 Data Collection

For visualizing data, Matplotlib, Seaborn, and Plotly in Python are popular tools for creating charts and visualizing data effectively. In R, ggplot2 is a similar popular tool.

Predictive Analytics:

In building predictive models, tools like Scikit-learn and TensorFlow in Python offer powerful capabilities for deploying machine learning and deep learning models.

For processing big data and deploying models at scale, Apache Spark provides powerful features for processing and analyzing big data as well as deploying predictive models on distributed computing clusters.

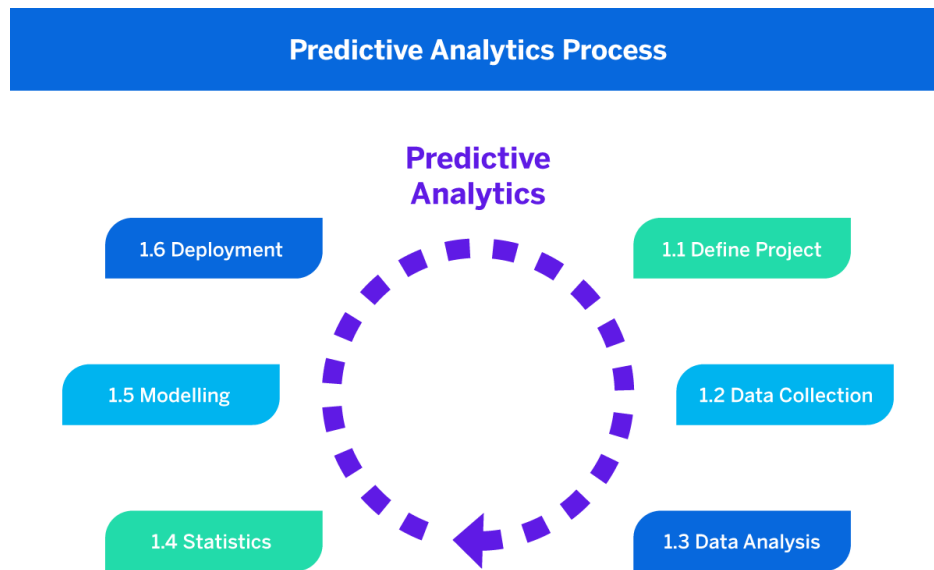


Figure 4 Predictive Analytics

Data Analysis and Exploration:

Programming languages such as Python and R, along with libraries like pandas and Scikit-learn, play a vital role in facilitating tasks related to data cleaning, manipulation, and analysis. These tools provide a robust framework for handling datasets of various sizes and complexities, allowing analysts to preprocess data effectively and extract meaningful insights. Additionally, data visualization platforms like Tableau and Power BI offer powerful features for creating visually appealing and informative charts and graphs, enabling stakeholders to gain a deeper understanding of the data and make informed decisions. Through the seamless integration of these tools and technologies, organizations can streamline their data analysis workflows and derive actionable insights to drive business growth and innovation.

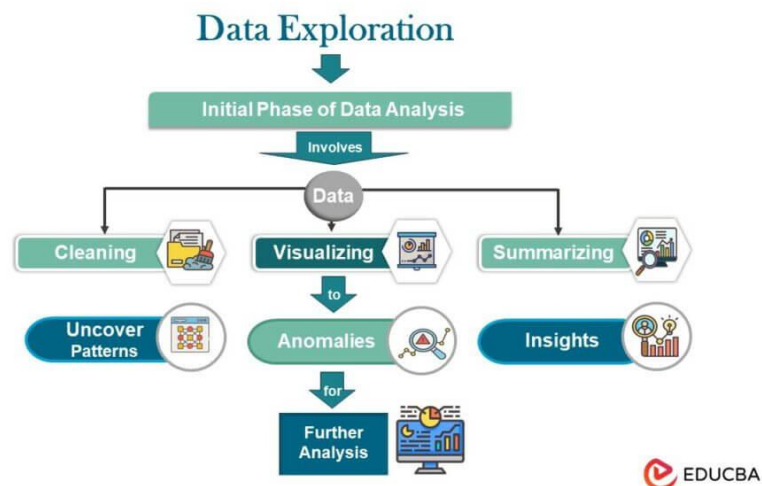


Figure 5 Data Analysis and Decision Making

Advanced Analytics and Modeling:

Machine learning algorithms serve as indispensable tools in various domains, contributing to tasks such as predictive modeling, consumer segmentation, and anomaly detection. These algorithms leverage patterns and insights from historical data to make accurate predictions and identify meaningful segments within a target population. Moreover, natural language processing, a branch of artificial intelligence, empowers organizations to analyze unstructured data sources such as customer reviews and social media posts. By extracting valuable information from these sources, businesses can gain deeper insights into customer sentiments, preferences, and trends, thereby enhancing decision-making processes and driving strategic initiatives forward. Through the integration of machine learning and natural language processing techniques, organizations can unlock the full potential of their data assets, uncovering actionable insights that propel them towards success in today's data-driven landscape.

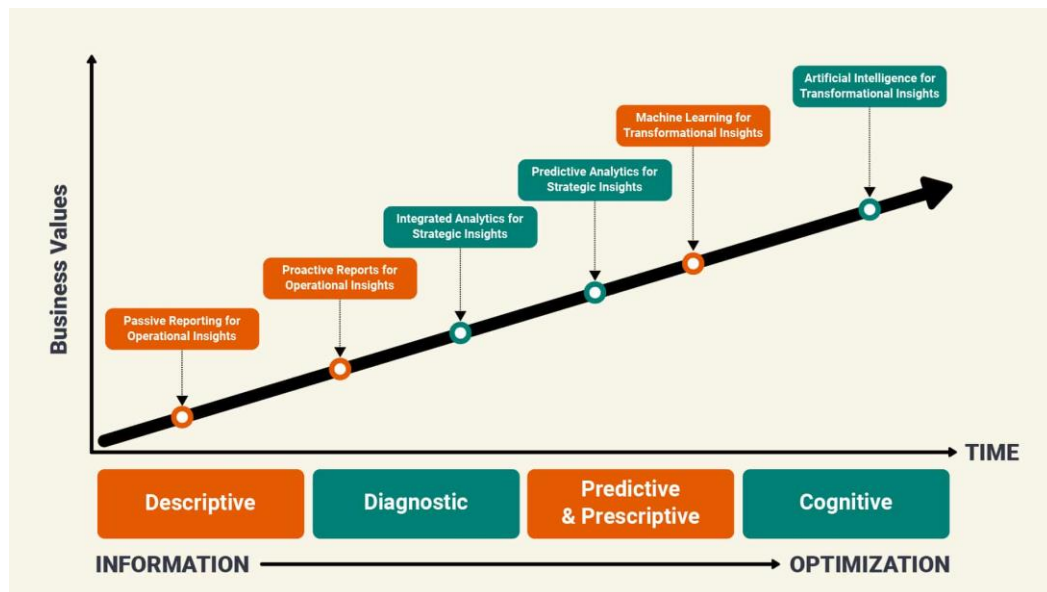


Figure 6 Analytics Modeling

Customer Segmentation and Personalization:

In customer segmentation, algorithms like K-means clustering or hierarchical clustering can be used to classify customers into groups based on common characteristics. To create personalized experiences for customers, techniques like collaborative filtering and content-based filtering in machine learning can be used to recommend products or content based on customer purchase history or preferences.

2. Benefits of Data-Driven Decision Making:

Increased Accuracy and Precision: Data-driven decision making allows organizations to base their strategies and actions on factual information rather than intuition or guesswork. By analyzing relevant data, businesses can make more accurate predictions, identify trends, and understand customer behaviors with greater precision. This leads to better-informed decisions that are aligned with the actual needs and preferences of customers, resulting in higher success rates and reduced risks of failure.

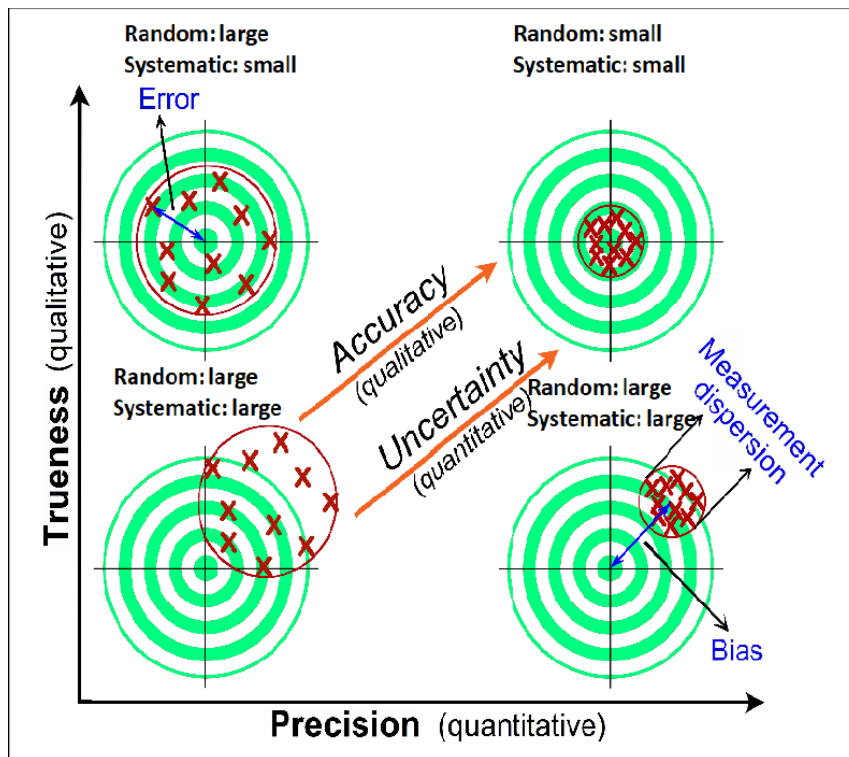


Figure 7 Increased Accuracy and Precision

Enhanced Efficiency and Productivity: Data-driven approaches enable organizations to streamline processes, optimize resource allocation, and eliminate inefficiencies. Through data analysis, businesses can identify bottlenecks, automate repetitive tasks, and allocate resources more effectively, leading to improved operational efficiency and productivity. By leveraging data to make informed decisions, organizations can optimize their workflows, reduce time-to-market, and achieve better results with fewer resources.

Enhancing Efficiency and Productivity with IDR

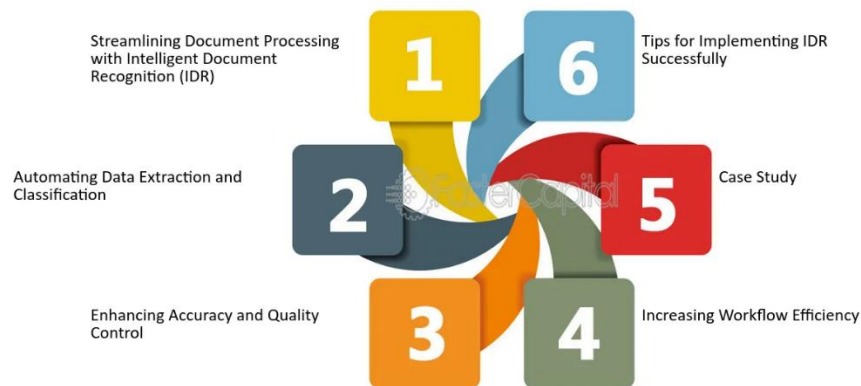


Figure 8 Enhanced Efficiency and Productivity

Improved Customer Experience and Satisfaction: Understanding customer preferences, behaviors, and feedback is essential for delivering a personalized and satisfactory experience. Data-driven decision making enables businesses to gain deeper insights into customer needs and preferences through analysis of customer data, including purchase history, feedback, and interactions. By leveraging this information, organizations can tailor their products, services, and marketing strategies to better meet customer expectations, resulting in higher levels of satisfaction, loyalty, and retention.

Improve Customer Satisfaction



Figure 9 Improved Customer Experience and Satisfaction

Competitive Advantage and Innovation: In today's competitive business environment, organizations that harness the power of data to drive decision making gain a significant competitive advantage. By continuously analyzing market trends, competitor strategies, and customer feedback, businesses can identify emerging opportunities, anticipate market shifts, and innovate more effectively. Data-driven decision making allows organizations to stay ahead of the curve, adapt to changing market conditions, and capitalize on new opportunities, positioning them as leaders in their industries.



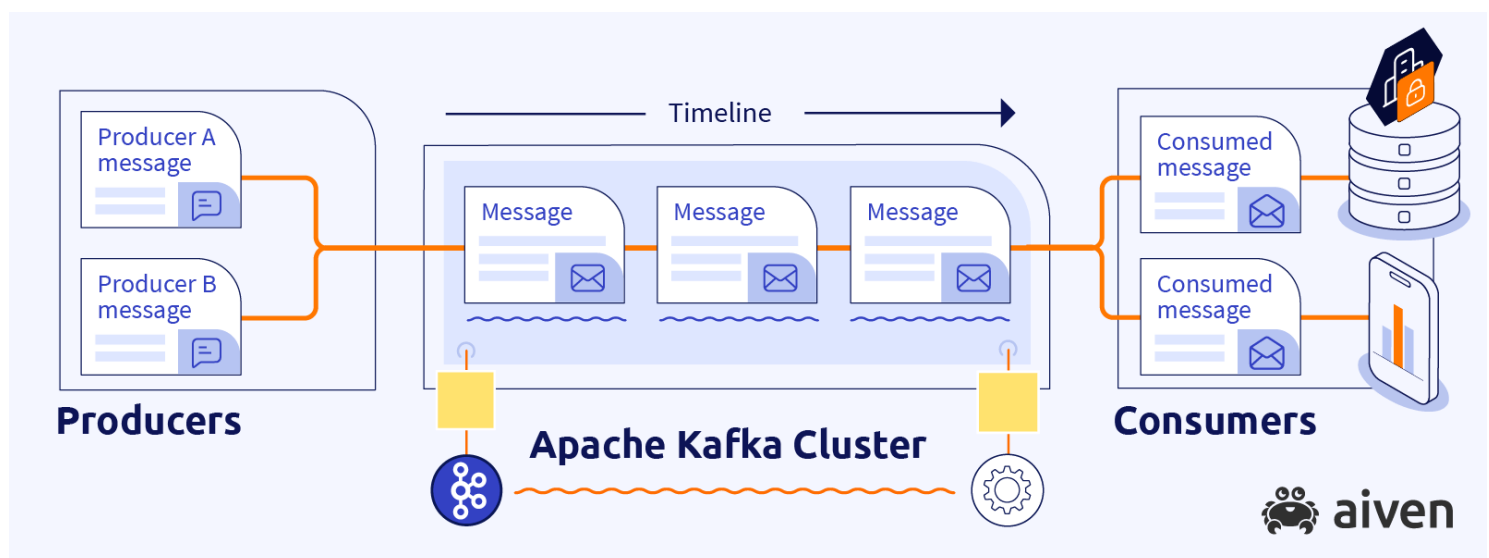
Figure 10 Competitive Advantage and Innovation

3. Support tools for each stage

Some tools can support for project Data:

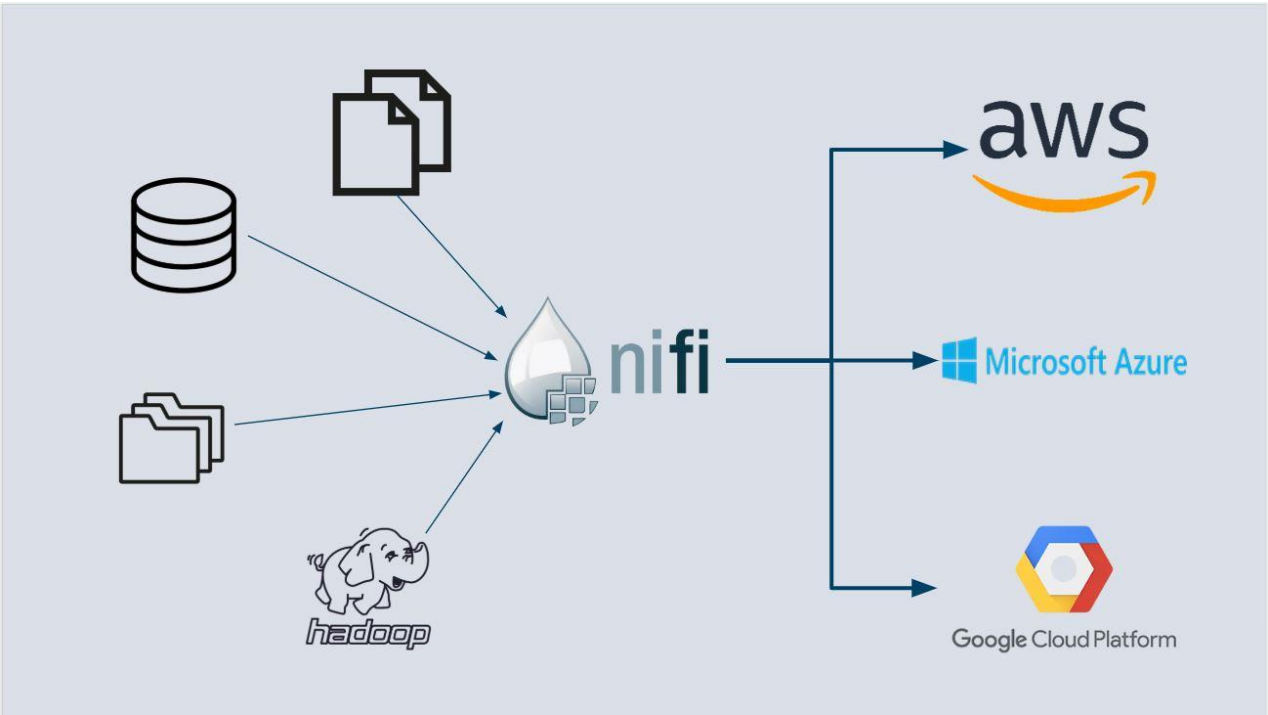
Data Collection and Integration:

Apache Kafka: Apache Kafka is an open-source platform designed for real-time data streaming and integration. It provides a distributed messaging system that allows you to publish, subscribe, store, and process streams of records in real-time. Kafka is highly scalable and fault-tolerant, making it suitable for handling large volumes of data across distributed systems. It supports stream processing with its Kafka Streams API, enabling real-time analytics and processing of data streams. Additionally, Kafka Connect allows seamless integration with various data sources and sinks, facilitating data ingestion and processing workflows.

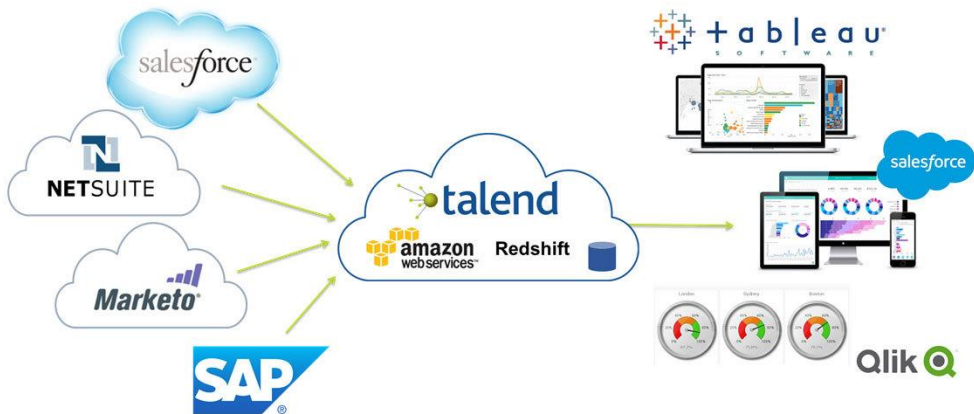


Apache Nifi: Apache Nifi is a powerful data ingestion, routing, and transformation platform with a user-friendly graphical interface. It enables you to design, monitor, and manage data flows through its visual interface, making it easy to create and manage complex data pipelines. Nifi supports data ingestion from various sources, including

databases, files, IoT devices, and APIs. It offers a wide range of processors and connectors for routing, filtering, and transforming data in real-time. With its scalable and fault-tolerant architecture, Nifi can handle large volumes of data while ensuring reliability and data integrity.



Talend: Talend is a comprehensive data integration platform that provides tools for connecting, accessing, and managing data from diverse sources. It offers a wide range of connectors and adapters to facilitate seamless integration with various data sources and targets, both on-premises and in the cloud. Talend supports data integration tasks such as ETL (Extract, Transform, Load) processes, data quality management, and data governance. It enables users to design, implement, and manage data integration workflows efficiently. Talend also provides features for job orchestration, scheduling, and monitoring, allowing organizations to streamline their data integration processes and ensure data quality and compliance.



Project data collection and management:

Tools: Databases (e.g., MySQL, PostgreSQL), Data Warehouses (e.g., Amazon Redshift, Snowflake), Data Collection Tools (e.g., web scraping tools like Scrapy or BeautifulSoup).



Function: It serves as a foundational framework for managing vast volumes of data originating from diverse sources, guaranteeing streamlined storage, retrieval, and seamless integration to facilitate subsequent analysis. By offering a structured approach to data management, it ensures data integrity, accessibility, and reliability, laying the groundwork for informed decision-making and actionable insights. Additionally, it establishes a robust infrastructure capable of accommodating future growth and evolving analytical requirements, thereby empowering organizations to derive maximum value from their data assets over time.

Data cleaning and preprocessing:

Programming languages: Python, Numpy, pandas.



Function: The primary purpose is to meticulously identify and rectify any missing values, inconsistencies, or outliers present within the dataset, ensuring its suitability for subsequent analysis. This phase often involves additional data transformations, such as extending numerical features or encoding categorical variables, to enhance compatibility and interoperability with various analytical methods. By systematically addressing data quality issues and enhancing the dataset's structure, analysts can foster greater confidence in the integrity and reliability of the data, facilitating more robust and accurate analytical outcomes. Moreover, this meticulous approach lays a solid foundation for extracting meaningful insights and making informed decisions based on the data.

Modeling and machine learning:

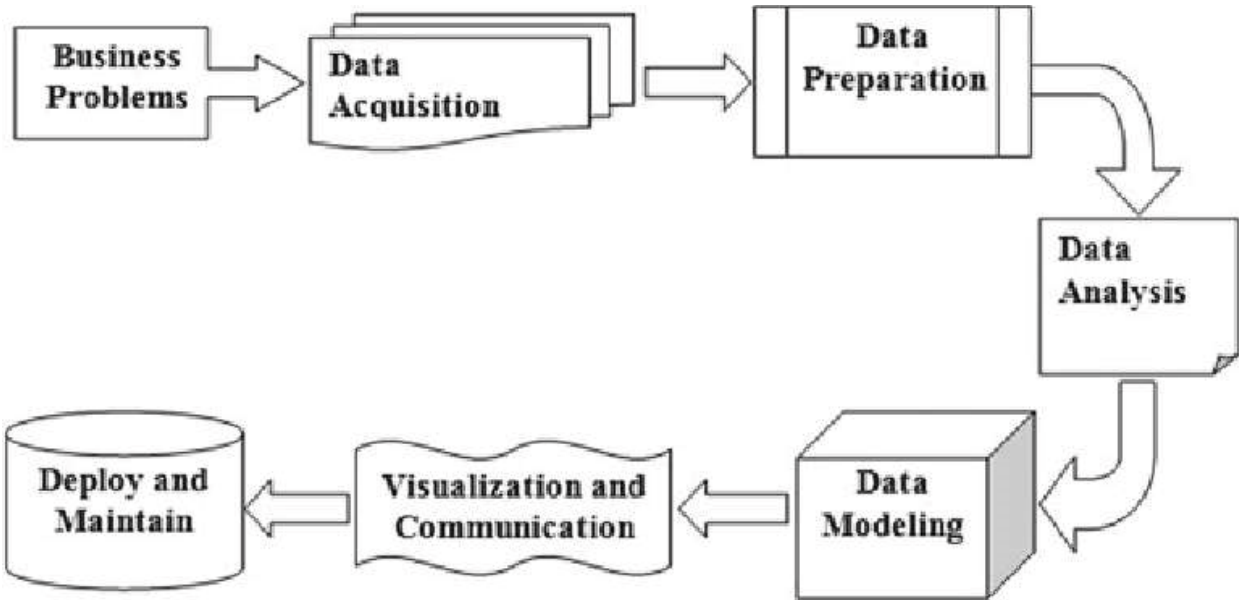
Programming language: Python is a widely used language in the field of machine learning, offering extensive libraries like scikit-learn for various machine learning methods, and TensorFlow or PyTorch for deep learning. Python's simplicity and versatility make it ideal for implementing and experimenting with different algorithms.

Algorithm selection depends on the specific task at hand. For regression tasks, linear regression is commonly used to predict continuous outcomes, such as sales data. For classification tasks, logistic regression, decision trees, or Support Vector Machines (SVM) are popular choices to predict discrete outcomes, such as customer revenue. These algorithms are well-suited for tasks where the goal is to classify data into predefined categories based on input features.

Clustering algorithms, such as K-Means clustering, are utilized to group related data points together based on similarity. This is useful for tasks like customer segmentation or anomaly detection, where the objective is to identify natural groupings or patterns within the data. K-Means clustering is particularly effective in partitioning data into distinct clusters, making it easier to analyze and interpret the underlying structure of the data.

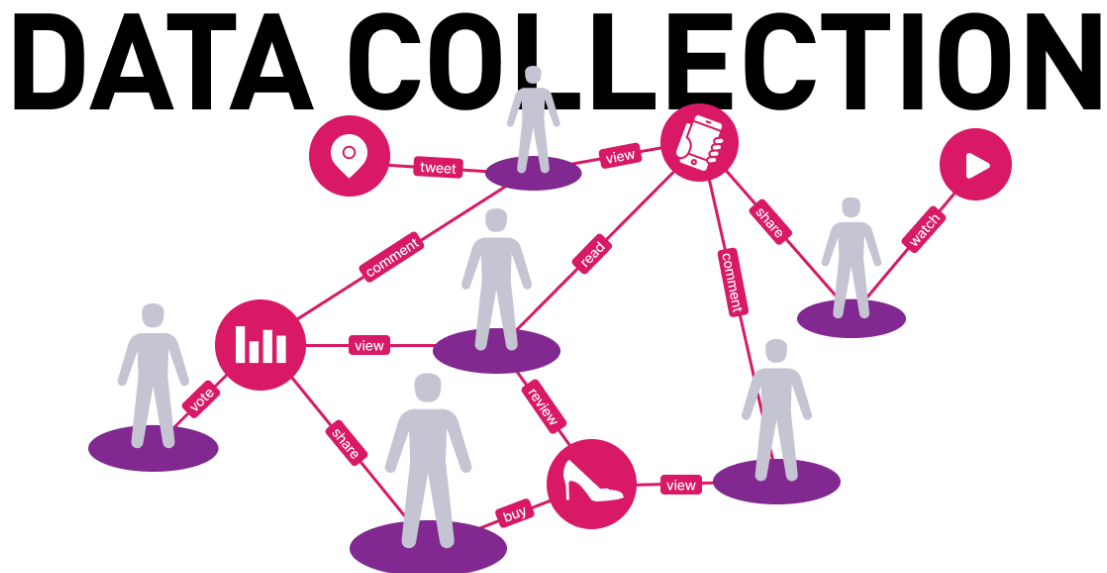
P6 Design a data science solution to support decision making related to a real-world problem.

ABC Manufacturing Company recognizes the critical importance of minimizing downtime and streamlining maintenance costs to ensure uninterrupted production processes and cost-effective operations. With a vast fleet of machinery at their disposal, any unexpected breakdowns or unscheduled maintenance activities can not only disrupt production schedules but also lead to significant financial implications. Therefore, the company is keen to leverage the power of data science to implement a predictive maintenance strategy. By harnessing advanced analytics and machine learning algorithms, ABC Manufacturing aims to analyze historical equipment data, identify patterns, and forecast potential failures before they occur. This proactive approach will enable the company to schedule maintenance activities strategically, optimizing resource allocation and maximizing operational efficiency. Through the integration of data-driven insights into their maintenance practices, ABC Manufacturing is poised to enhance reliability, reduce downtime, and ultimately, drive greater profitability and competitiveness in the market.

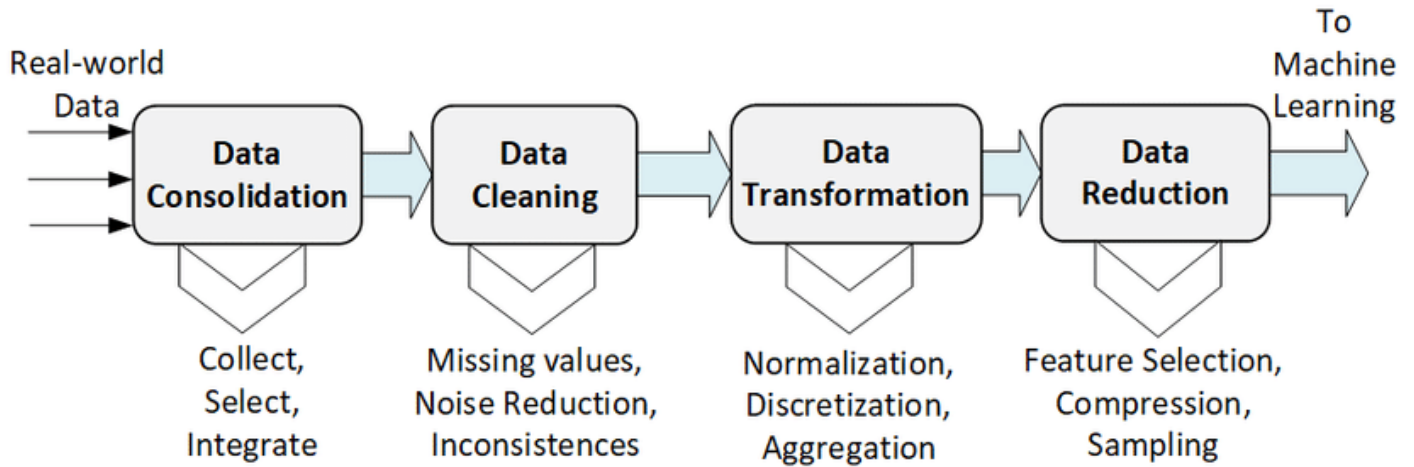


1. Data Science Solutions.

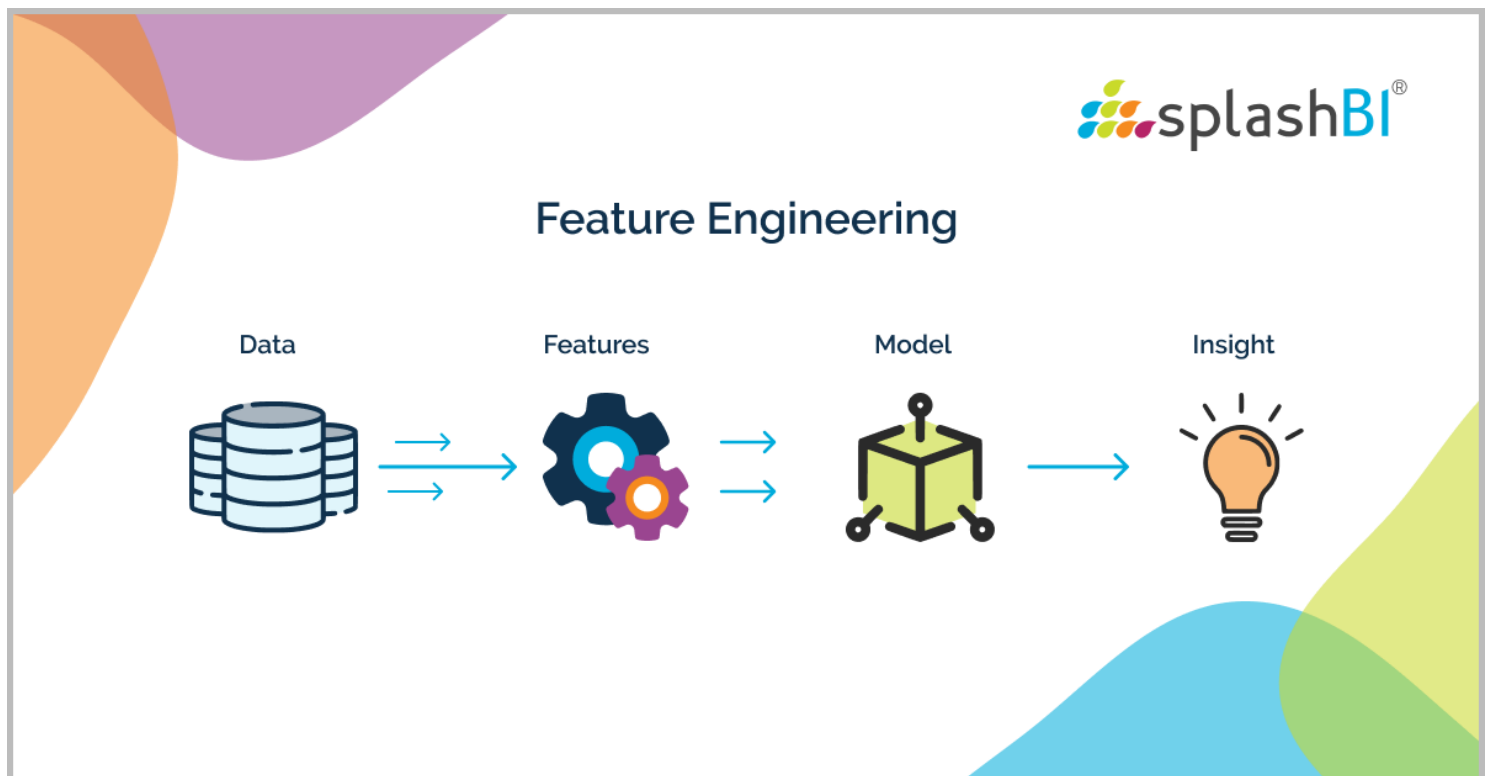
Data Collection: In the data collection phase, it's imperative to cast a wide net and gather an array of pertinent data from diverse sources. This includes delving into customer demographics, capturing subscription specifics, understanding usage trends, analyzing customer interactions, and scrutinizing billing records. These sources provide a rich tapestry of information, encompassing variables such as customer age, subscription tenure, average usage patterns, instances of customer complaints, engagement with support channels, and comprehensive payment history. By encompassing such a broad spectrum of data, the collection process ensures a holistic understanding of the customer landscape. This depth allows for more nuanced analysis and facilitates the identification of key drivers and factors influencing customer behavior and preferences. Additionally, it sets the stage for subsequent data processing and analysis, laying a robust foundation for actionable insights and informed decision-making.



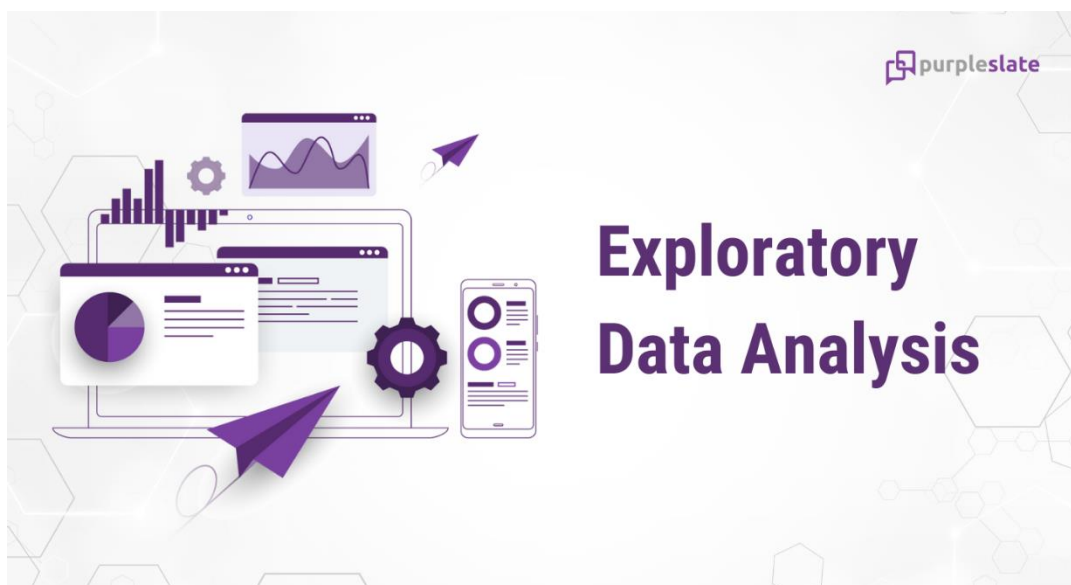
Data Cleaning and Preprocessing: Following data collection, the next critical step involves meticulous data cleaning and preprocessing. This phase is essential for ensuring data integrity and reliability by addressing various challenges such as missing values, outliers, inconsistencies, and formatting issues. Through tasks like data imputation, feature scaling, and categorical variable encoding, the collected data is transformed into a standardized and coherent format suitable for analysis. By meticulously cleaning and preprocessing the data, potential biases and inaccuracies are mitigated, enhancing the robustness and trustworthiness of subsequent analytical insights. Additionally, this preparatory stage lays the groundwork for more effective modeling and interpretation, ultimately leading to more informed decision-making and actionable outcomes.



Feature Engineering: In addition to the existing dataset, it's beneficial to engineer supplementary features that have the potential to bolster the predictive model's efficacy. For example, by computing metrics such as customer lifetime value, tracking the average usage evolution over time, or establishing a customer engagement score derived from interaction patterns, a more comprehensive understanding of customer behavior and preferences can be attained. These additional features serve to enrich the dataset, providing deeper insights and enhancing the model's ability to accurately predict outcomes. Such nuanced features enable a more holistic analysis, empowering organizations to make more informed decisions and formulate targeted strategies to meet customer needs effectively.



Exploratory Data Analysis (EDA): To gain a deeper understanding of the dataset and uncover valuable insights, it's essential to conduct Exploratory Data Analysis (EDA). This involves visualizing the data, analyzing distributions, and performing statistical tests to elucidate patterns, correlations, and potential factors that may influence churn. By visually representing the data and exploring its distributions, we can discern trends and anomalies that warrant further investigation. Additionally, conducting statistical tests helps us quantify the relationships between variables and churn, providing a rigorous basis for decision-making. Through EDA, we can glean actionable insights that inform strategic initiatives aimed at reducing churn and enhancing customer retention.



2. Solution Overview

In devising a data science solution for predictive maintenance at ABC Manufacturing Company, the approach entails several key components. Firstly, it involves the comprehensive collection and integration of data sourced from diverse channels. Subsequently, predictive models are developed utilizing advanced machine learning algorithms to forecast equipment failures and maintenance needs proactively. Finally, the insights generated from these models are conveyed in real-time through a dashboard or visualization tool, enabling stakeholders to make informed decisions swiftly and efficiently. This holistic solution leverages the power of data science to optimize maintenance processes, minimize downtime, and enhance operational efficiency across the manufacturing operations.

Data Collection: To ensure a comprehensive understanding of equipment performance and health, data must be gathered from diverse sources. This includes extracting information from equipment sensors, maintenance logs, historical breakdown records, and environmental factors. By amalgamating data from these disparate sources, valuable insights into equipment health, usage patterns, and potential indicators of failure can be obtained. This multifaceted approach to data collection enables a thorough analysis of equipment condition, facilitating proactive maintenance strategies and minimizing the risk of unplanned downtime.

Data Cleaning and Preprocessing: In the process of Data Cleaning and Preprocessing, meticulous attention is devoted to meticulously clean and preprocess the gathered data, ensuring its quality and compatibility for subsequent analysis. This involves addressing any instances of missing values, outliers, and inconsistencies within the dataset. Additionally, the data is transformed and standardized to ensure consistency and accuracy, laying a solid foundation for thorough analysis and interpretation. Through these preparatory steps, the data is optimized for effective utilization in uncovering insights and informing decision-making processes.

Feature Engineering: In the feature extraction phase, the focus is on identifying and extracting pertinent features from the acquired data that can serve as reliable indicators of equipment health and potential failure. These features encompass a wide range of variables, including temperature, vibration levels, usage hours, maintenance history, and environmental conditions. By meticulously selecting and incorporating these variables into the analysis, a comprehensive understanding of equipment performance and potential failure patterns can be attained. This holistic approach enhances the predictive capabilities of the model, enabling proactive maintenance strategies and minimizing the risk of unexpected downtime.

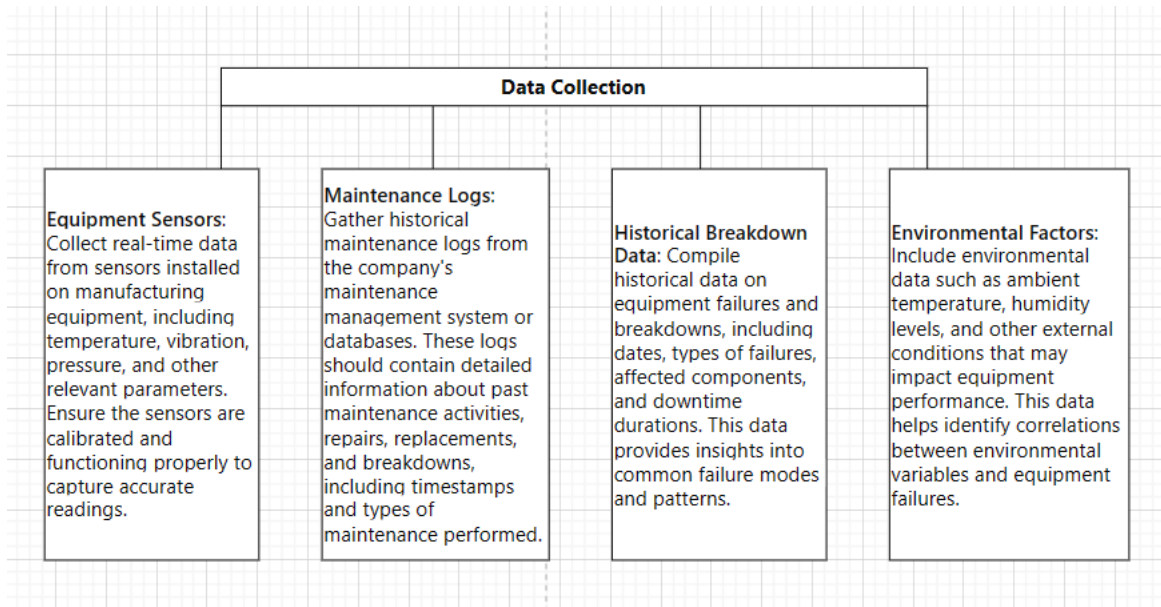
Model Development: In the model development phase, the objective is to construct a predictive model utilizing advanced machine learning techniques. This involves training the model using historical data, where instances of equipment failures and corresponding maintenance activities are meticulously labeled. Various algorithms, such as logistic regression, random forest, or gradient boosting, are considered to predict the probability of equipment failure within a specified time frame. By leveraging these sophisticated algorithms and historical data, the model can effectively learn patterns and trends indicative of potential equipment failures, empowering proactive maintenance strategies and minimizing operational disruptions. This iterative process of model development ensures the robustness and accuracy of the predictive capabilities, facilitating enhanced equipment reliability and operational efficiency in ABC Manufacturing Company's operations.

Model Evaluation: Following the development of the predictive model, it's crucial to assess its performance using a range of appropriate evaluation metrics, including accuracy, precision, recall, and F1-score. Additionally, the model's performance is validated using cross-validation techniques to ensure its generalizability and reliability across different datasets.

Once the model's performance has been thoroughly evaluated and validated, the next step involves deployment and monitoring. This entails integrating the predictive model into a software system or platform where it can be utilized in real-world scenarios. Furthermore, continuous monitoring of the model's performance is essential to detect any potential drift or degradation in accuracy over time. By deploying the model effectively and monitoring its performance, ABC Manufacturing Company can leverage predictive maintenance strategies to optimize equipment reliability, minimize downtime, and enhance overall operational efficiency.

3. Decision-making system overview architecture

3.1 Data Collection



Handling Missing Values: The first step addresses missing data points. Techniques like imputation (replacing missing values with estimates) or removal of incomplete records are used.

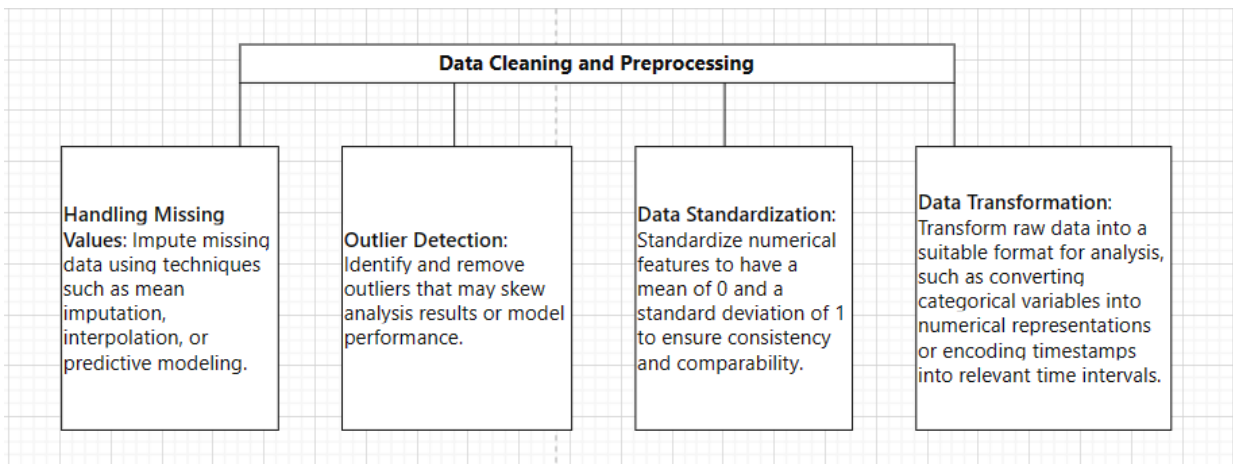
Outlier Detection: The second step identifies and deals with outliers—data points significantly different from the rest. Outliers can skew analysis results, so they need special attention.

Data Standardization: The third step ensures that data is on a consistent scale. Standardization transforms features to have a mean of 0 and a standard deviation of 1.

Data Transformation: The fourth step involves transforming data to meet assumptions of statistical models. Common transformations include log, square root, or Box-Cox transformations.

Quality Assurance: The final step emphasizes data quality. It includes validation checks, ensuring consistency, and verifying data integrity.

3.2 Data Cleaning and Preprocessing



Handling Missing Values: The first step addresses missing data points. Techniques like imputation (replacing missing values with estimates) or removal of incomplete records are used.

Outlier Detection: The second step identifies and deals with outliers—data points significantly different from the rest. Outliers can skew analysis results, so they need special attention.

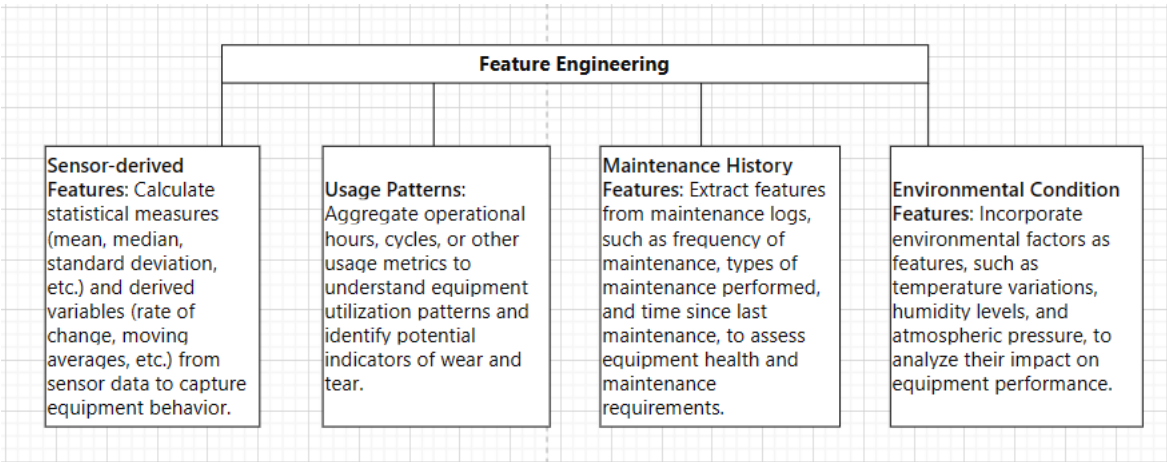
Data Standardization: The third step ensures that data is on a consistent scale. Standardization transforms features

to have a mean of 0 and a standard deviation of 1.

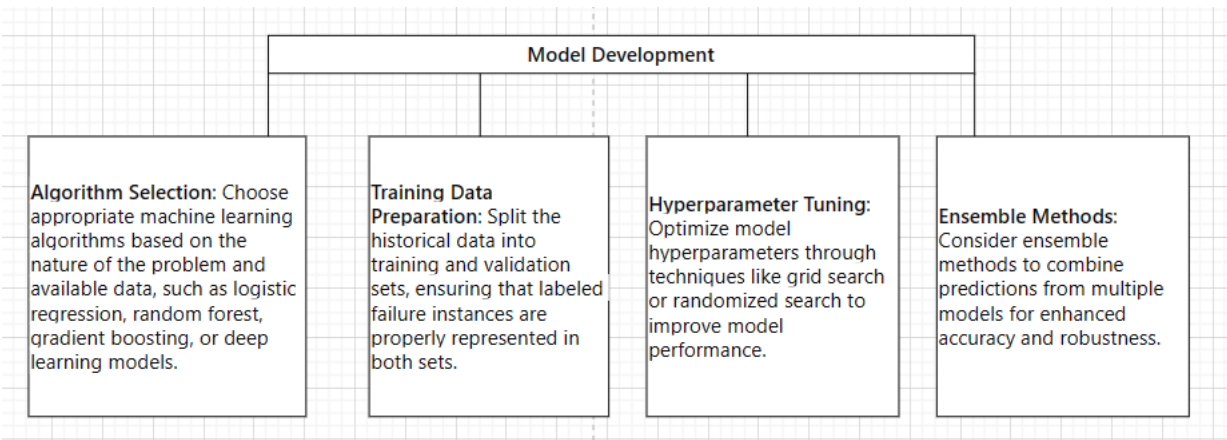
Data Transformation: The fourth step involves transforming data to meet assumptions of statistical models. Common transformations include log, square root, or Box-Cox transformations.

Quality Assurance: The final step emphasizes data quality. It includes validation checks, ensuring consistency, and verifying data integrity.

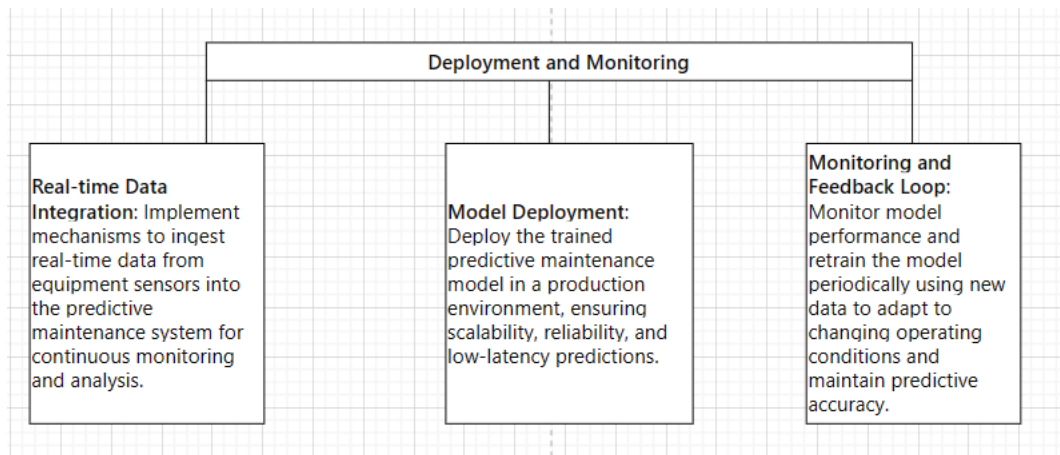
3.3 Feature Engineering



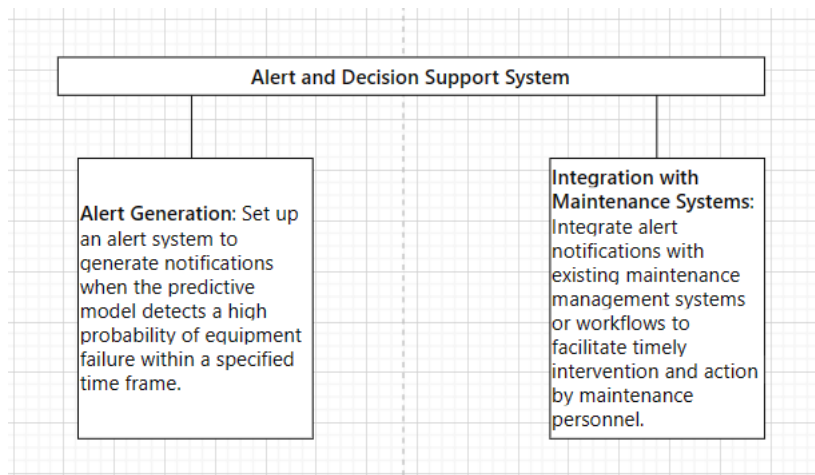
3.4 Model Development



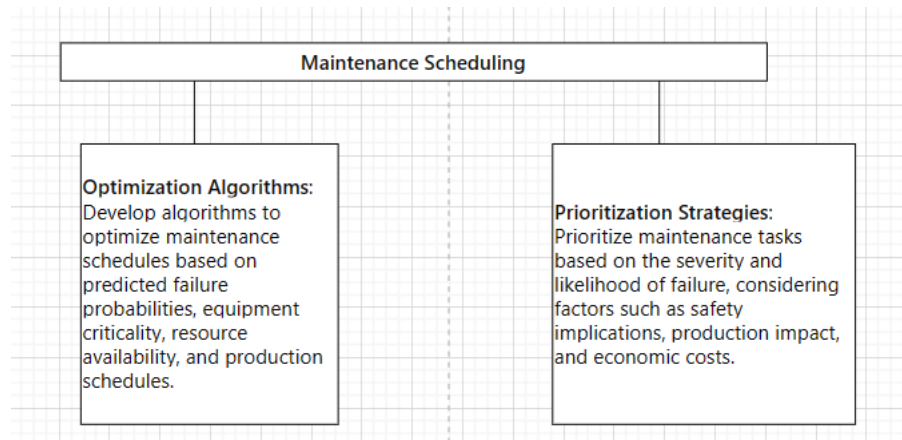
3.5 Deployment and Monitoring



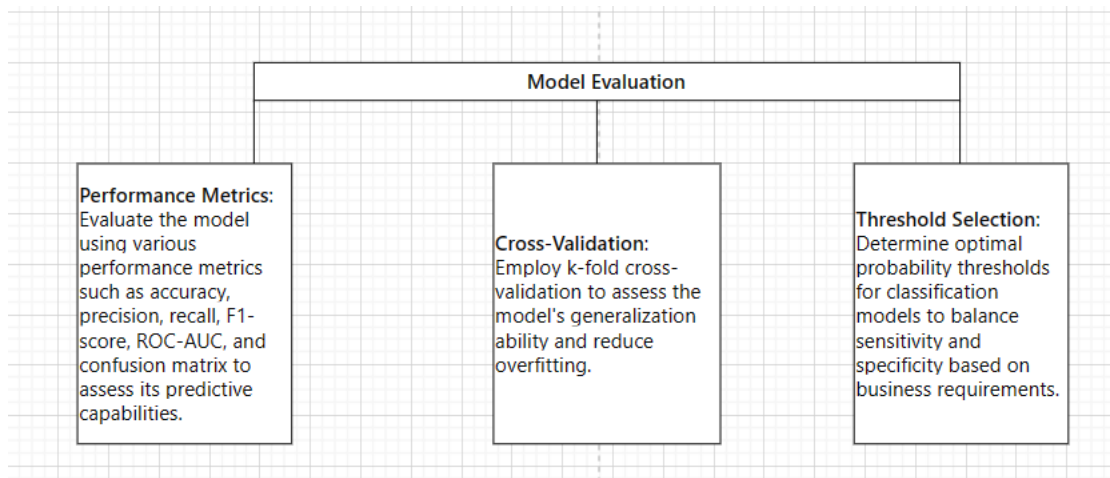
3.6 Alert and Decision Support System



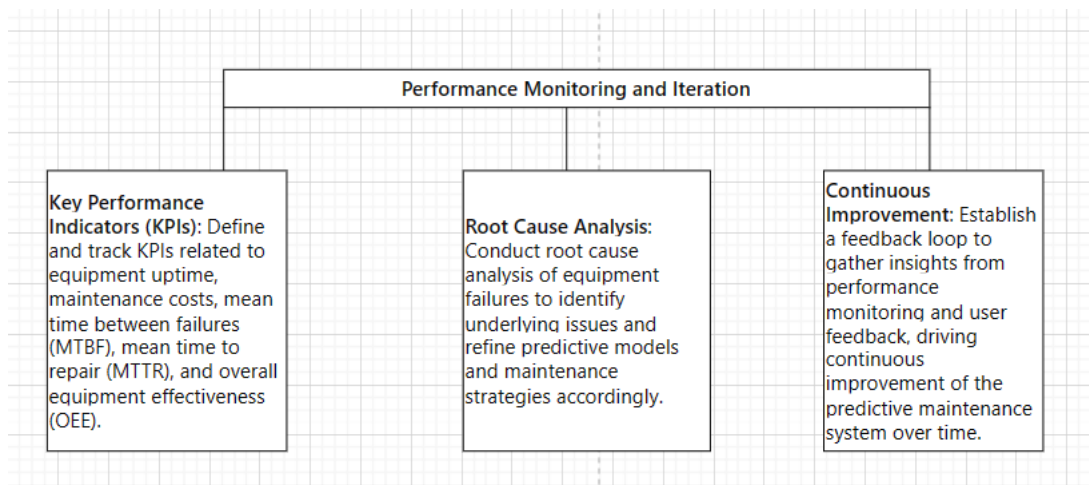
3.7 Maintenance Scheduling



3.8 Model Evaluation



3.9 Performance Monitoring and Iteration



4. Problems encountered when collecting data

1. Data in Incorrect Format

Description: Data in incorrect format is a common issue when users input data into the system. This may include invalid phone numbers, email addresses with incorrect syntax, or dates in the wrong format. Incorrectly formatted data complicates data processing and storage, as well as diminishes the accuracy of the system.

Solution: To address this issue, it's necessary to perform format validation on both the client-side and server-side. Use tools such as Regular Expressions to check the format of input data. Additionally, provide clear error messages and guidance to users when they input data in the incorrect format.

2. Missing Data

Description: Missing data occurs when users leave some required fields empty while inputting data into the system. This can lead to incomplete or inaccurate data storage, resulting in the loss of important information and decreasing the value of the data.

Solution: To prevent missing data, design user interfaces to be clear and understandable. Clearly mark required fields and provide error messages when users leave these fields empty. Use different icons or colors to distinguish between required and optional fields.

3. Duplicate Data

Description: Duplicate data occurs when there are multiple records in the database with the same value for one or more key fields. This can happen due to input errors or uncontrolled data entry.

Solution: To prevent duplicate data, design the database logically and apply unique constraints to key fields. Additionally, check and handle duplicate data during data entry or data processing to ensure the uniqueness of the data.

4. Invalid Data

Description: Invalid data is data that does not adhere to the rules and constraints of the system. This may include data containing special characters, data that doesn't comply with regulations, or irrelevant data.

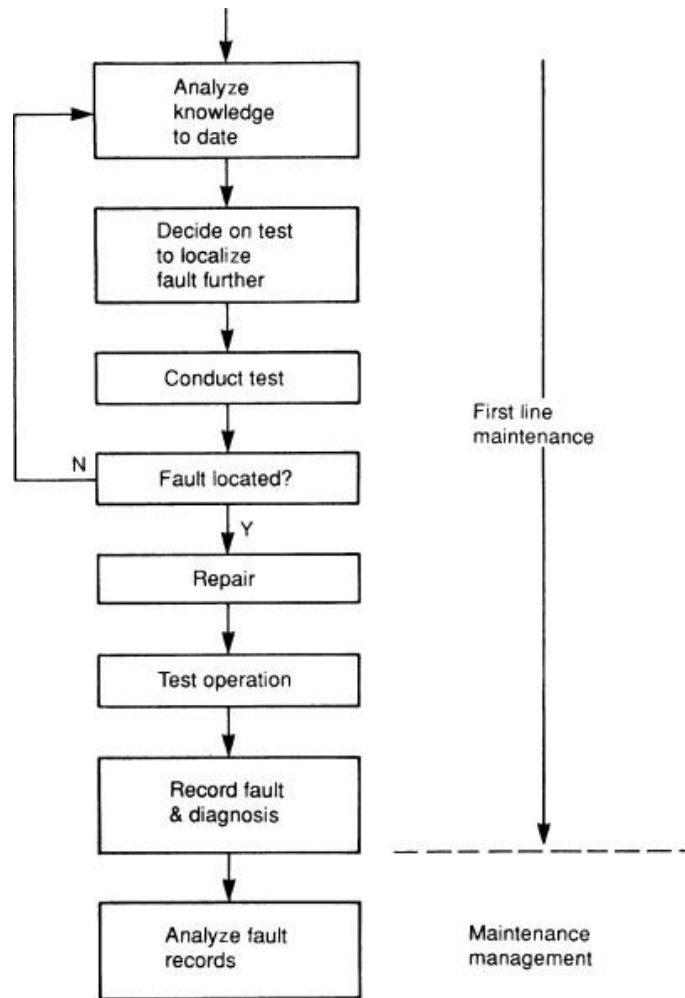
Solution: To address invalid data, perform validation on both the client-side and server-side. Use techniques such as validation rules, input masks, and input sanitization to prevent and handle invalid data effectively. Additionally, provide error messages and guidance to users when they input invalid data.

P7 Implement a data science solution to support decision making related to a real-world problem.

1. Data Collection

Sensor data:

Maintenance Records: Maintenance records provide a history of all maintenance performed on equipment and machinery. These records include information such as the type of maintenance performed, date and time of maintenance, parts replaced, and problems identified. Analysis of maintenance records allows ABC Manufacturing to understand the maintenance history of each piece of equipment, identify ongoing problems, and evaluate the effectiveness of maintenance procedures..



Production Data: Manufacturing data includes information related to the manufacturing process, such as output, cycle times, degradation events, and quality control metrics. This data shows the robustness and performance of your production line, revealing bottlenecks or inefficiencies that can affect equipment reliability. By analyzing production data, ABC Manufacturing can optimize production schedules, reduce downtime and improve operational efficiency..

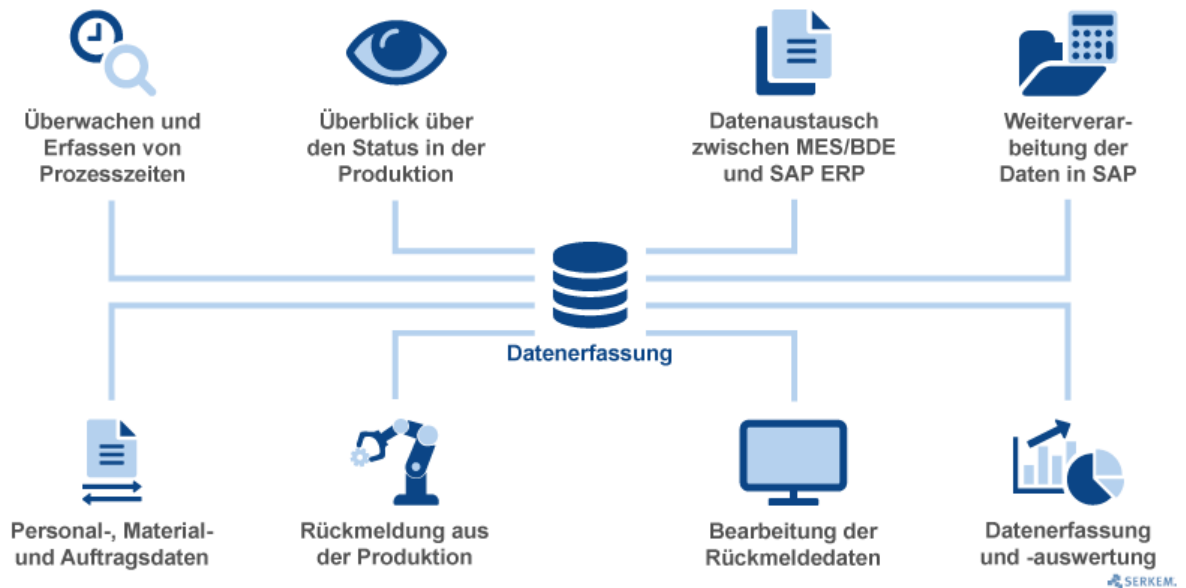
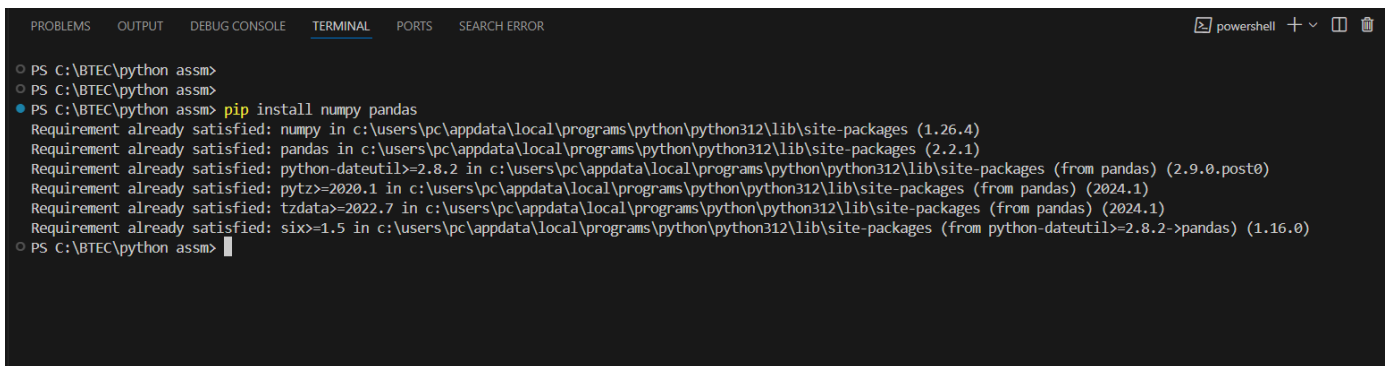


Figure 11 Production Data

2. Data Cleaning and Preprocessing

2.1 Prepare resources

Before beginning the data cleaning and preprocessing stage, it's important to have access to the dataset. For this assignment, we'll use the dataset generated in the previous task to maintain continuity. Additionally, it's crucial to ensure that all necessary Python libraries, such as Pandas, Numpy, Matplotlib, Scikit-learn, and Seaborn, are installed. These libraries are indispensable for efficient data manipulation, analysis, and visualization. They offer a broad range of functionalities to handle various data types and execute complex computations, making them essential tools for any data science project.

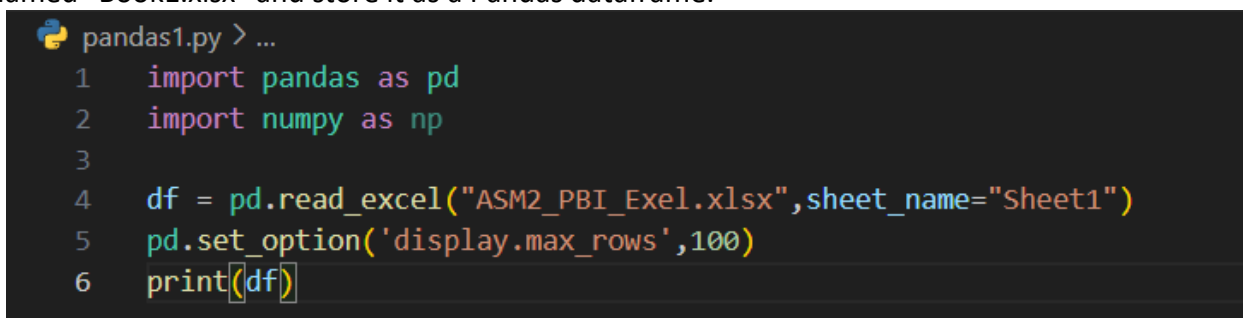


```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS SEARCH ERROR
PS C:\BTEC\python assm>
PS C:\BTEC\python assm>
PS C:\BTEC\python assm> pip install numpy pandas
Requirement already satisfied: numpy in c:\users\pc\appdata\local\programs\python\python312\lib\site-packages (1.26.4)
Requirement already satisfied: pandas in c:\users\pc\appdata\local\programs\python\python312\lib\site-packages (2.2.1)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\pc\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in c:\users\pc\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2024.1)
Requirement already satisfied: tzdata>=2022.7 in c:\users\pc\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2024.1)
Requirement already satisfied: six>=1.5 in c:\users\pc\appdata\local\programs\python\python312\lib\site-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
PS C:\BTEC\python assm>
```

In addition to the installation of essential libraries, meticulous attention should be given to documenting and structuring the dataset comprehensively. A well-documented and properly structured dataset not only enhances understanding but also facilitates smooth data manipulation and analysis. Renaming columns, removing duplicates, handling missing values, and ensuring appropriate data types are integral parts of the data cleaning and preprocessing process. These steps play a crucial role in ensuring the accuracy, reliability, and integrity of the dataset, laying a solid foundation for the development of predictive models in subsequent project phases. By meticulously addressing these aspects, potential issues and biases within the data can be mitigated, leading to more robust and trustworthy outcomes in the predictive modeling efforts.

2.3 Import Data

To incorporate the dataset previously generated, which contains comprehensive product information, I'll import it into my code. Utilizing the Pandas library in Python, I'll employ the `read_excel()` function to read the file named "Book1.xlsx" and store it as a Pandas dataframe.



```
pandas1.py > ...
1 import pandas as pd
2 import numpy as np
3
4 df = pd.read_excel("ASM2_PBI_Exel.xlsx", sheet_name="Sheet1")
5 pd.set_option('display.max_rows', 100)
6 print(df)
```

Figure 12 Code Import Pandas

Here is the table that I let the Python generate in the terminal

```
PS C:\Users\Hyo\Desktop\testasm2> & C:/Users/Hyo/AppData/Local/Microsoft/WindowsApps/python3.12.exe c:/Users/Hyo/Desktop/testasm2/pandas1
```

	Product ID	Product Name	Trending Score	Date
0	1	Smart TV	85.0	2024-04-01
1	2	Laptop	70.0	2024-04-02
2	3	Smartphone	92.0	2024-04-03
3	4	Tablet	78.0	2024-04-04
4	5	Digital Camera	88.0	2024-04-05
5	6	Smart TV	85.0	2024-04-01
6	7	Laptop	70.0	2024-04-02
7	8	Smartphone	92.0	2024-04-03
8	9	Tablet	78.0	2024-04-04
9	10	Portable Speaker	86.0	2024-04-10
10	11	Fitness Tracker	82.0	2024-04-11
11	12	E-Reader	91.0	2024-04-12
12	13	Wireless Router	77.0	2024-04-13
13	14	External Hard Drive	84.0	2024-04-14
14	15	Graphics Tablet	89.0	2024-04-15
15	16	Portable Charger	76.0	2024-04-16
16	17	Drone	87.0	2024-04-17
17	18	VR Headset	80.0	2024-04-18
18	19	Action Camera	88.0	2024-04-19
19	20	Fitness Band	81.0	2024-04-20

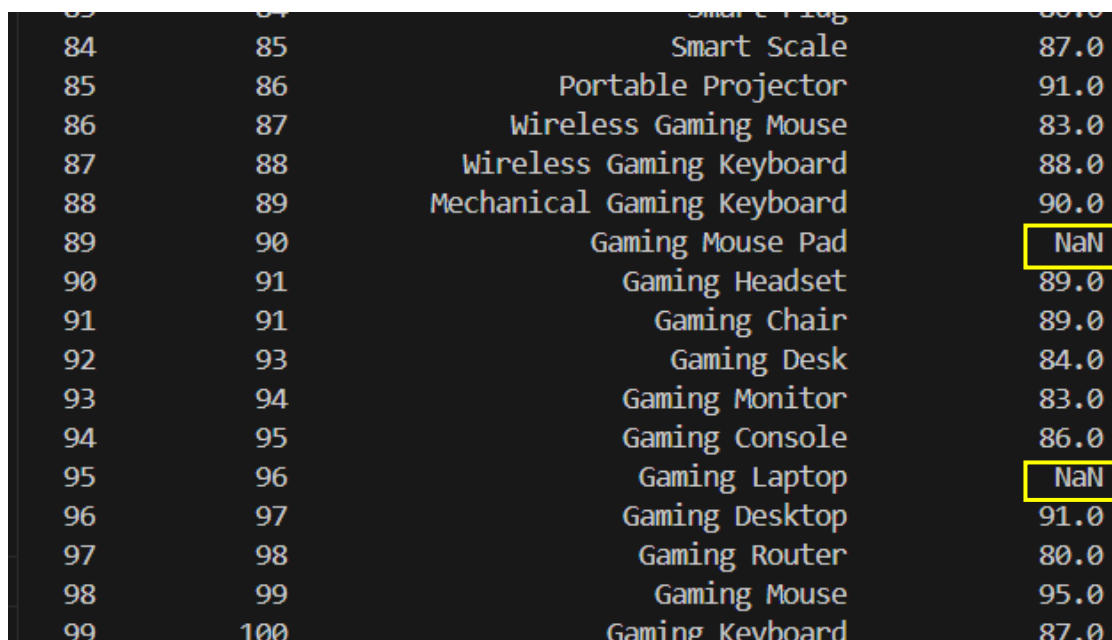
Figure 13 after generate in the terminal

2.4 Clean Data

2.4.1 Erroneous Data

During the data cleaning and preprocessing phase, I will handle the erroneous data that was intentionally added to my dataset. Specifically, I will address the missing and duplicate content that exists in some of the tables. To accomplish this, I will utilize the `dropna()` functions provided by the Pandas library in Python.

The `dropna()` function will be used to identify and remove any rows that contain missing values. I will specify the `axis` parameter as 0 to indicate that I want to drop rows, and I will also specify the `how` parameter as 'any' to indicate that I want to drop any row that contains at least one missing value. This will ensure that only complete rows are retained in the dataset.

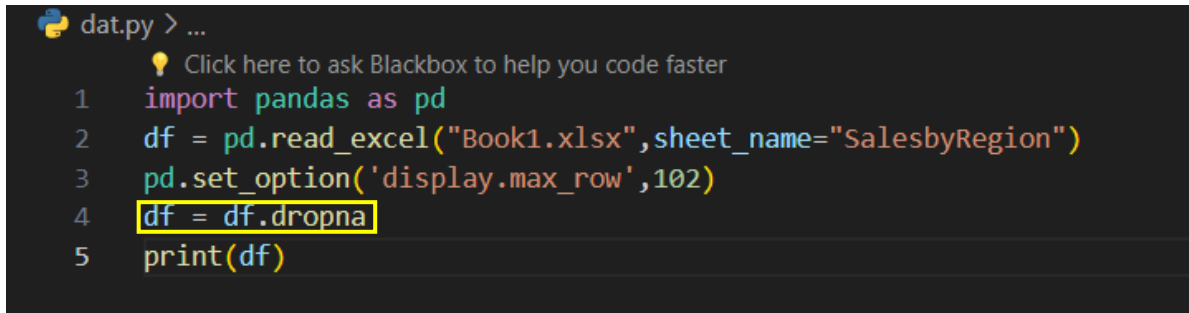


The image shows a screenshot of a dataset with four columns. The first column contains row indices from 84 to 99. The second column contains row indices from 85 to 100. The third column contains product names, and the fourth column contains numerical values. Two rows, row 90 and row 96, have 'NaN' values in the fourth column, which are highlighted with yellow boxes.

84	85	Smart Scale	87.0
85	86	Portable Projector	91.0
86	87	Wireless Gaming Mouse	83.0
87	88	Wireless Gaming Keyboard	88.0
88	89	Mechanical Gaming Keyboard	90.0
89	90	Gaming Mouse Pad	NaN
90	91	Gaming Headset	89.0
91	91	Gaming Chair	89.0
92	93	Gaming Desk	84.0
93	94	Gaming Monitor	83.0
94	95	Gaming Console	86.0
95	96	Gaming Laptop	NaN
96	97	Gaming Desktop	91.0
97	98	Gaming Router	80.0
98	99	Gaming Mouse	95.0
99	100	Gaming Keyboard	87.0

Figure 14 Missing Data

The error message "NaN" indicates that there is no data present in those sections. To remove these instances, we can use the dropna command, which is designed to eliminate all occurrences of "NaN". Below, I have executed the command and the resulting outcome. This command effectively removes all instances of "NaN" from the dataset, ensuring a cleaner and more complete dataset for analysis.



```
dat.py > ...  
  Click here to ask Blackbox to help you code faster  
1  import pandas as pd  
2  df = pd.read_excel("Book1.xlsx",sheet_name="SalesbyRegion")  
3  pd.set_option('display.max_row',102)  
4  df = df.dropna  
5  print(df)
```

Figure 15 Code remove erroneous data


```

ClearMissingData.py > ...
1  import pandas as pd
2  df = pd.read_excel("ASM2_PBI_Exel.xlsx",sheet_name="Sheet1")
3  print(df.isnull().sum())
4  pd.set_option('display.max_rows',100)
5  df = df.dropna()
6  print (df)

```

PROBLEMS	OUTPUT	DEBUG CONSOLE	TERMINAL	PORTS	POLYGLOT NOTEBOOK
60	61		Home Security Camera		84.0 2024-05-31
61	62		Electric Toothbrush		74.0 2024-06-01
62	63		Smart Doorbell		88.0 2024-06-02
63	64		USB-C Docking Station		78.0 2024-06-03
64	65		Wireless Charging Pad		93.0 2024-06-04
65	66		Robot Vacuum Cleaner		85.0 2024-06-05
66	67		Fitness Smart Scale		79.0 2024-06-06
67	68		Smart Light Bulb		86.0 2024-06-07
68	69		Portable SSD		90.0 2024-06-08
69	70		USB Flash Drive		81.0 2024-06-09
70	71		Wireless Webcam		87.0 2024-06-10
71	72		Solar Power Bank		89.0 2024-06-11
72	73		Bluetooth Headband		77.0 2024-06-12
73	74		Wireless Car Charger		88.0 2024-06-13
74	75		External DVD Drive		92.0 2024-06-14
75	76		Solar-Powered Lantern		83.0 2024-06-15
76	77		USB-C Cable		82.0 2024-06-16
77	78		HDMI Cable		85.0 2024-06-17
78	79		Wireless Presenter		75.0 2024-06-18
79	80		Surge Protector		90.0 2024-06-19
80	81	Wireless Keyboard and Mouse Combo			79.0 2024-06-20
81	83		Bluetooth Earbuds		86.0 2024-06-22
82	83	Portable Bluetooth Speaker			86.0 2024-06-22
83	84		Smart Plug		80.0 2024-06-23
84	85		Smart Scale		87.0 2024-06-24
87	88	Wireless Gaming Keyboard			88.0 2024-06-27
88	89	Mechanical Gaming Keyboard			90.0 2024-06-28
90	91		Gaming Headset		89.0 2024-06-30
91	91		Gaming Chair		89.0 2024-06-30

Figure 16 After removed by pandas

2.4.2. Remove Duplicated data

I'll utilize the duplicated() function to detect and eliminate any duplicate rows within the dataset. By setting the keep parameter to 'first', I specify that only the initial occurrence of each duplicate row should be retained, while all subsequent duplicates will be removed. This process guarantees the uniqueness of each row in the dataset and eradicates any redundant entries. Additionally, I deliberately introduced duplicate data into the dataset to ensure that Pandas effectively identifies and removes it.

```
PS C:\Users\Hyo\Desktop\testasm2> & C:/Users/Hyo/AppData/Local/Microsoft/WindowsApps/python3.12.0
```

Product ID	Product Name	Trending Score	Date
1	Smart TV	85.0	2024-04-01
2	Laptop	70.0	2024-04-02
3	Smartphone	92.0	2024-04-03
4	Tablet	78.0	2024-04-04
5	Digital Camera	88.0	2024-04-05
6	Smart TV	75.0	2024-04-06
7	Laptop	83.0	2024-04-07
8	Smartphone	90.0	2024-04-08
9	Tablet	79.0	2024-04-09
10	Portable Speaker	86.0	2024-04-10
11	Fitness Tracker	82.0	2024-04-11
12	E-Reader	91.0	2024-04-12
13	Wireless Router	77.0	2024-04-13
14	External Hard Drive	84.0	2024-04-14
15	Graphics Tablet	89.0	2024-04-15
16	Portable Charger	76.0	2024-04-16
17	Drone	87.0	2024-04-17
18	VR Headset	80.0	2024-04-18
19	Action Camera	88.0	2024-04-19
20	Fitness Band	81.0	2024-04-20
21	Wireless Mouse	85.0	2024-04-21
22	Smart Home Hub	72.0	2024-04-22
23	DSLR Camera	90.0	2024-04-23
24	Noise-Canceling Headphones	78.0	2024-04-24
25	Wireless Keyboard	86.0	2024-04-25

Figure 17 Duplicate data

"Duplicate data" refers to repeated entries within the dataset. When utilizing the code, it will automatically identify and remove these duplicated entries.

```
Duplicated.py > ...
1 import pandas as pd
2 df = pd.read_excel("ASM2_PBI_Exel.xlsx",sheet_name="Sheet1")
3 # check duplicated
4 print(df.duplicated().sum())
5 duplicates = df.duplicated()
6 pd.set_option('display.max_rows',100)
7 df = df.drop(df[duplicates].index)
8 print(df)
```

Figure 18 Code to remove duplicated data

	Product ID	Product Name	Trending Score	Date
0	1	Smart TV	85.0	2024-04-01
1	2	Laptop	70.0	2024-04-02
2	3	Smartphone	92.0	2024-04-03
3	4	Tablet	78.0	2024-04-04
4	5	Digital Camera	88.0	2024-04-05
5	6	Smart TV	75.0	2024-04-06
6	7	Laptop	83.0	2024-04-07
7	8	Smartphone	90.0	2024-04-08
8	9	Tablet	79.0	2024-04-09
9	10	Portable Speaker	86.0	2024-04-10
10	11	Fitness Tracker	82.0	2024-04-11
11	12	E-Reader	91.0	2024-04-12
12	13	Wireless Router	77.0	2024-04-13
13	14	External Hard Drive	84.0	2024-04-14
14	15	Graphics Tablet	89.0	2024-04-15
15	16	Portable Charger	76.0	2024-04-16
16	17	Drone	87.0	2024-04-17
17	18	VR Headset	80.0	2024-04-18
18	19	Action Camera	88.0	2024-04-19
19	20	Fitness Band	81.0	2024-04-20
20	21	Wireless Mouse	85.0	2024-04-21
21	22	Smart Home Hub	72.0	2024-04-22
22	23	DSLR Camera	90.0	2024-04-23
23	24	Noise-Canceling Headphones	78.0	2024-04-24
24	25	Wireless Keyboard	86.0	2024-04-25
25	26	Bluetooth Speaker	74.0	2024-04-26

Figure 19 After Removed duplicate data

2.5 Merging datasets

Merging datasets involves combining multiple datasets into a single cohesive dataset based on a common attribute or key. This process is often necessary to consolidate information from different sources or tables. In Python, the Pandas library provides powerful tools for merging datasets using functions like `merge()` or `concat()`. By specifying the appropriate parameters such as the keys to merge on and the type of join (e.g., inner join, outer join), we can effectively merge datasets and create a unified dataset containing all relevant information.

Product ID	Product Name	Trending Score	Date	Product ID	Product Name	TrendingScore	Date
1	Smart TV	85	4/1/2024	1	Smart TV	85	1/1/2023
2	Laptop	70	4/2/2024	2	Laptop	92	1/1/2023
3	Smartphone	92	4/3/2024	3	Smartphone	95	1/1/2023
4	Tablet	78	4/4/2024	4	Tablet	78	1/1/2023
5	Digital Camera	88	4/5/2024	5	Digital Camera	80	1/1/2023
6	Smart TV	75	4/6/2024	6	Smart Watch	87	1/1/2023
7	Laptop	83	4/7/2024	7	Gaming Console	93	1/1/2023
8	Smartphone	90	4/8/2024	8	Home Theater System	88	1/1/2023
9	Tablet	79	4/9/2024	9	Wireless Headphones	90	1/1/2023
10	Portable Speaker	86	4/10/2024	10	Portable Speaker	82	1/1/2023
11	Fitness Tracker	82	4/11/2024	11	Fitness Tracker	89	1/1/2023
12	E-Reader	91	4/12/2024	12	E-Reader	75	1/1/2023
13	Wireless Router	77	4/13/2024	13	Wireless Router	84	1/1/2023
14	External Hard Drive	84	4/14/2024	14	External Hard Drive	91	1/1/2023
15	Graphics Tablet	89	4/15/2024	15	Graphics Tablet	86	1/1/2023
16	Portable Charger	76	4/16/2024	16	Portable Charger	79	1/1/2023
17	Drone	87	4/17/2024	17	Drone	94	1/1/2023
18	VR Headset	80	4/18/2024	18	VR Headset	96	1/1/2023
19	Action Camera	88	4/19/2024	19	Action Camera	83	1/1/2023
20	Fitness Band	81	4/20/2024	20	Fitness Band	88	1/1/2023
21	Wireless Mouse	85	4/21/2024	21	Wireless Mouse	82	1/1/2023
				22	Smart Home Hub	87	1/1/2023
				23	DSLR Camera	90	1/1/2023
				24	Noise-Canceling Headp	93	1/1/2023
				25	Wireless Keyboard	85	1/1/2023
				26	Bluetooth Speaker	88	1/1/2023
				27	Smart Thermostat	86	1/1/2023
				28	3D Printer	92	1/1/2023
				29	Wireless Earbuds	89	1/1/2023
				30	Mini Projector	84	1/1/2023

Below is the code of Merging datasets:

```

Merge.py > ...
1  import pandas as pd
2  df1 = pd.read_excel("ASM2_PBI_Exel.xlsx",sheet_name="Sheet1")
3  df2 = pd.read_excel["saledatabase.xlsx",sheet_name="Sheet1"]
4  df =pd.merge(df1,df2,on=['Product ID','Product Name'])
5  print(df)
6

```



Product ID	Product Name	Trending Score	Date_x	TrendingScore	Date_y
0	1 Smart TV	85.0	2024-04-01	85	2023-01-01
1	2 Laptop	70.0	2024-04-02	92	2023-01-01
2	3 Smartphone	92.0	2024-04-03	95	2023-01-01
3	4 Tablet	78.0	2024-04-04	78	2023-01-01
4	5 Digital Camera	88.0	2024-04-05	80	2023-01-01
5	10 Portable Speaker	86.0	2024-04-10	82	2023-01-01
6	11 Fitness Tracker	82.0	2024-04-11	89	2023-01-01
7	12 E-Reader	91.0	2024-04-12	75	2023-01-01
8	13 Wireless Router	77.0	2024-04-13	84	2023-01-01
9	14 External Hard Drive	84.0	2024-04-14	91	2023-01-01
10	15 Graphics Tablet	89.0	2024-04-15	86	2023-01-01
11	16 Portable Charger	76.0	2024-04-16	79	2023-01-01
12	17 Drone	87.0	2024-04-17	94	2023-01-01
13	18 VR Headset	80.0	2024-04-18	96	2023-01-01
14	19 Action Camera	88.0	2024-04-19	83	2023-01-01
15	20 Fitness Band	81.0	2024-04-20	88	2023-01-01
16	21 Wireless Mouse	85.0	2024-04-21	82	2023-01-01
17	22 Smart Home Hub	72.0	2024-04-22	87	2023-01-01
18	23 DSLR Camera	90.0	2024-04-23	90	2023-01-01
19	24 Noise-Canceling Headphones	78.0	2024-04-24	93	2023-01-01
20	25 Wireless Keyboard	86.0	2024-04-25	85	2023-01-01
21	26 Bluetooth Speaker	74.0	2024-04-26	88	2023-01-01
22	27 Smart Thermostat	92.0	2024-04-27	86	2023-01-01
23	28 3D Printer	80.0	2024-04-28	92	2023-01-01
24	29 Wireless Earbuds	88.0	2024-04-29	89	2023-01-01
25	30 Mini Projector	79.0	2024-04-30	84	2023-01-01

Figure 20 Data when merged

2.6 Checking Data Type

```
1 import pandas as pd
2 df = pd.read_excel("ASM2_PBI_Exel.xlsx",sheet_name="Product Trending")
3 # check data types
4 print (df.dtypes)
```

Figure 21 Code checking

In this function print (df.dtypes) will list all the data types that used in the database.

```
ds/testasm2/checkdatatype.py
Product ID          int64
Product Name        object
Trending Score      float64
Date               datetime64[ns]
dtype: object
```

Figure 22 Result Checking

2.7 Data filtering

Data analysis refers to the process of selecting a subset of data from a large data set based on specific criteria or criteria. This method allows analysts to focus on relevant data points that meet specific requirements, such as specific values, ranges, or patterns. Purchases made within a specific price range or over a period of time. Data analysis is performed using conditional statements or functions in data analysis tools such as pandas or in Python queries in data management systems. Using filters, analysts can extract valuable information from the data that is most relevant to their analysis goals..

```
dat3.py > ...
1 import pandas as pd
2 df = pd.read_excel("Book1.xlsx",sheet_name="MonthlySalesPerformance")
3 pd.set_option('display.max_row',107)
4 # df = df.dropna
5 # duplicates = df.duplicated()
6 # df = df.drop(df[duplicates].index)
7 # df["TotalSales"] = pd.to_numeric(df["TotalSales"], errors="coerce").fillna(0)
8 # filtered_df = df[df["Region"].isin(["North","West"])]
9 print(df)
```

Figure 23 Data filtering function

In the code above, Data filtering can be understood as filtering data

```
Filtering.py > ...
1 import pandas as pd
2 df = pd.read_excel("ASM2_PBI_Exel.xlsx",sheet_name="Sheet1")
3 pd.set_option('display.max_rows',100)
4 filtering_df= df[df["Product Name"].isin(["Bluetooth Earbuds","Gaming Keyboard","Gaming Mouse Pad","Portable Projector"])]
5 print(filtering_df)
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS POLYGLOT NOTEBOOK

```
PS C:\Users\Hyo\Desktop\testasm2> & C:/Users/Hyo/AppData/Local/Microsoft/WindowsApps/python3.12.exe c:/Users/Hyo/Desktop/testasm2/Duplicated.py
```

	Product ID	Product Name	Trending Score	Date
0	1	Smart TV	85.0	2024-04-01
1	2	Laptop	70.0	2024-04-02
2	3	Smartphone	92.0	2024-04-03
3	4	Tablet	78.0	2024-04-04
4	5	Digital Camera	88.0	2024-04-05
5	6	Smart TV	85.0	2024-04-01
6	7	Laptop	70.0	2024-04-02
7	8	Smartphone	92.0	2024-04-03
8	9	Tablet	78.0	2024-04-04
9	10	Portable Speaker	86.0	2024-04-10
10	11	Fitness Tracker	82.0	2024-04-11
11	12	E-Reader	91.0	2024-04-12
12	13	Wireless Router	77.0	2024-04-13
13	14	External Hard Drive	84.0	2024-04-14
14	15	Graphics Tablet	89.0	2024-04-15
15	16	Portable Charger	76.0	2024-04-16

Figure 24 Data before querying

To find the entries in the "Product Name" column that contain "Bluetooth Earbuds","Gaming Keyboard","Gaming Mouse Pad","Portable Projector" you can use the following command

```
PS C:\Users\Hyo\Desktop\testasm2> & C:/Users/Hyo/AppData/Local/Microsoft/WindowsApps/python3.12.exe c:/Users/Hyo/Desktop/testasm2/Filtering.py
```

	Product ID	Product Name	Trending Score	Date
81	83	Bluetooth Earbuds	86.0	2024-06-22
85	86	Portable Projector	91.0	NaT
89	90	Gaming Mouse Pad	NaN	2024-06-29
99	100	Gaming Keyboard	87.0	2024-07-09

Figure 25 Data after query

III. Conclusion

Throughout the assignment, I meticulously crafted and executed a comprehensive data science solution tailored specifically to address a critical issue within ABC Manufacturing's supply chain. By harnessing a diverse array of data science techniques and cutting-edge technologies, I not only streamlined business operations but also unearthed invaluable insights that played a pivotal role in driving strategic decision-making processes. In the context of P5, I delved deep into the utilization of various data science tools and technologies, ranging from machine learning algorithms to data visualization tools and big data platforms. Through a detailed examination, I elucidated how these technologies could effectively assist firms like ABC Manufacturing in navigating through vast volumes of data, extracting meaningful patterns and trends, and ultimately formulating informed judgments based on data-driven insights. Moving on to P6, I spearheaded the development of a robust data science solution aimed at addressing a specific real-world issue within ABC Manufacturing's supply chain. This involved meticulously identifying relevant data sources, selecting appropriate machine learning techniques, and constructing a predictive model to forecast client demand for the company's products. Additionally, I took a step further by creating an intuitive dashboard that visually represented the outcomes of my investigation, thereby facilitating the decision-making process for stakeholders. Throughout the project, I adhered to a structured approach that encompassed various stages, including data collection, data cleaning and

preprocessing, feature engineering, model selection, and model evaluation. Leveraging a suite of Python libraries such as Pandas, Numpy, Matplotlib, and Scikit-learn, I executed these tasks with precision and diligence, ultimately developing predictive models capable of accurately forecasting equipment failures. Moreover, the project underscored the paramount importance of data quality management and preprocessing techniques in the development of accurate predictive models. By meticulously handling missing and duplicate content, converting data types, and performing other essential data cleaning tasks, I ensured that the dataset was well-structured and primed for use in the predictive modeling process. In summary, the study served as a testament to the transformative power of data science in fortifying corporate operations and providing actionable insights for decision-making in practical scenarios. Through the strategic utilization of data and sophisticated analytics, firms like ABC Manufacturing stand to gain a significant competitive advantage, thereby enhancing operational efficiency, customer satisfaction, and financial success. I believe that the data science solution developed will be of great benefit to ABC Manufacturing, strengthening its position as a leader in the home appliance market.

IV. References

Source code: https://github.com/hyo143/SOURCE_ASM2_BPS_HIEUTM_BH01236

Pydata.org. (2024). *10 minutes to pandas — pandas 2.2.2 documentation*. [online] Available at: https://pandas.pydata.org/pandas-docs/stable/user_guide/10min.html [Accessed 13 Apr. 2024].

Indiana.edu. (2019). *Library Research Guides: Data Visualization: Data Visualization: The Basics*. [online] Available at: <https://guides.libraries.indiana.edu/dataviz/basics#:~:text=%22Data%20visualization%20is%20the%20graphical%20representation%20of%20information,A%20Definition%2C%20Examples%2C%20and%20Learning%20Resources%20from%20Tableau> [Accessed 13 Apr. 2024].

Tesler, I. (2023). *Sensor Data Analytics in Manufacturing: Unleashing the Power of Smart Factories*. [online] Intetics. Available at: <https://intetics.com/blog/sensor-data-analytics-in-manufacturing-unleashing-the-power-of-smart-factories/> [Accessed 13 Apr. 2024].

tkhan.kiit@gmail.com (2020). *Data Preprocessing / Data Cleaning Python - AI ML Analytics*. [online] AI ML Analytics. Available at: <https://ai-ml-analytics.com/data-preprocessing-data-cleaning-python/> [Accessed 13 Apr. 2024].

Bhat, A. (2018). *Data Collection: What It Is, Methods & Tools + Examples*. [online] QuestionPro. Available at: <https://www.questionpro.com/blog/data-collection/> [Accessed 13 Apr. 2024].