

JSS MAHAVIDYAPEETHA
JSS SCIENCE AND TECHNOLOGY UNIVERSITY



**BIG DATA ANALYTICS ON ANIMAL DISEASE DATA USING
APACHE HIVE**

BACHELOR OF ENGINEERING
IN
Information Science and Engineering

By

Mohammed Zaid	01JST21IS027
Sayanth PM	01JST21IS047

Under the guidance of
Prof. MALAPRIYA S
Assistant Professor
Dept. of Information Science and Engineering
JSS STU, Mysuru-06

DEC 2024

Department of Information Science and Engineering

ABSTRACT

This project focuses on the analysis of large-scale animal disease data using Apache Hive, aiming to uncover patterns, predict disease outbreaks, and provide valuable insights for veterinary services and public health sectors. The dataset utilized is sourced from Kaggle's animal disease records, which include reports on diseases affecting various species globally.

The project involves multiple stages: data extraction, cleaning, transformation, and analysis. Apache Hive is employed for efficient handling of the dataset, enabling the execution of queries on large data volumes. Exploratory data analysis (EDA) identifies key trends such as disease prevalence, seasonality, and affected species.

Predictive modeling was performed to forecast disease outbreak probabilities for specific regions using historical data. The project also identifies the regions and species most vulnerable to specific diseases and highlights factors influencing disease outbreaks.

The findings underline the potential applications of big data analytics in managing animal health, contributing to better disease prevention strategies and policymaking. Future work includes integrating real-time disease reports for real-time analytics and enhancing prediction accuracy using advanced machine learning models.

CONTENT

• CHAPTER 1: INTRODUCTION

- 1.1 Problem Statement
- 1.2 Objectives ◦ 1.3 Scope of the Project

• CHAPTER 2: LITERATURE REVIEW ◦ 2.1

Review of Papers • CHAPTER 3:

TECHNOLOGY USED

- 3.1 Apache Hive ◦
- 3.2 HiveQL ◦ 3.3
- Architecture

• CHAPTER 4: METHODOLOGY

- 4.1 Data Collection ◦ 4.2 Data
- Cleaning and Preprocessing ◦ 4.3
- Data Analysis ◦ 4.4 Predictive
- Modeling

• CHAPTER 5: IMPLEMENTATION AND OUTPUT

- 5.1 Attributes ◦ 5.2 Data
- Collection and Cleaning ◦ 5.3
- Exploratory Data Analysis (EDA) ◦
- 5.4 Predictive Modeling ◦ 5.5
- Advanced Analysis ◦ 5.6 Output

• CHAPTER 6: RESULTS AND DISCUSSION

- 6.1 Data Analysis Results ◦
- 6.2 Predictive Modeling
- Results ◦ 6.3 Discussion

• CHAPTER 7: CONCLUSION AND FUTURE WORK

CHAPTER 1: INTRODUCTION

1.1 Problem Statement

The objective of this project is to analyze large-scale animal disease data using Apache Hive. The dataset spans multiple years and regions, providing insights into the prevalence of diseases among various species. By leveraging Hive's capabilities, the project aims to identify patterns, predict disease outbreaks, and generate actionable insights for veterinarians, policymakers, and public health officials. Key challenges include handling large datasets, data cleaning, trend analysis, and forecasting.

1.2 Objectives

- To analyze and extract insights from animal disease datasets.
- To process large-scale disease data using Apache Hive efficiently.
- To understand and analyze disease patterns across species and regions.
- To predict the probability of disease outbreaks based on historical data.

1.3 Scope of the Project

- Data extraction and cleaning for global animal disease records.
- Performing EDA using HiveQL.
- Implementing statistical models to predict disease outbreaks.
- Identifying regions and species at higher risk of specific diseases.
- Providing insights for improving disease prevention and control strategies.

CHAPTER 2: LITERATURE REVIEW

2.1 Review of Papers

1. **Big Data in Veterinary Health (2022)** ○ **Authors:** Dr. Emily Carter, Prof. James Willows ○ **Published Year:** 2022 ○ **Review:** This paper highlights the potential of big data analytics in veterinary health. It discusses frameworks for analyzing extensive veterinary records and provides insights into how large datasets can improve disease surveillance and outbreak management. The authors emphasize the integration of electronic health records with big data platforms to enhance diagnostic accuracy and treatment outcomes.

2. **Disease Forecasting Using Hadoop (2021)** ○ **Authors:** Dr. Rajesh Kumar, Dr. Anita Verma ○ **Published Year:** 2021 ○ **Review:** This study explores the application of Hadoop and Hive for analyzing animal disease data. The focus is on scalability and efficiency in handling large datasets. It describes the advantages of distributed computing frameworks in disease forecasting and provides examples of their use in predicting livestock disease outbreaks.

3. **Machine Learning in Veterinary Epidemiology (2019)** ○ **Authors:** Dr. Sophia Lopez, Dr. Mark Anders ○ **Published Year:** 2019 ○ **Review:** The paper discusses the role of machine learning models in veterinary epidemiology. It covers various predictive techniques, such as classification and regression models, used to identify high-risk areas and mitigate disease risks. The authors also review the effectiveness of these models in managing zoonotic diseases and improving overall animal health outcomes.

4. **Applications of Big Data in Livestock Management (2020)** ○ **Authors:** Dr. Olivia Bennett, Prof. Henry Miles ○ **Published Year:** 2020 ○ **Review:** This paper examines how big data analytics is transforming livestock management. It

highlights case studies where data from sensors and IoT devices were used to monitor animal health and predict potential disease outbreaks. The authors propose strategies for integrating big data tools with existing livestock management systems.

5. **Advances in Predictive Analytics for Animal Diseases (2021)** ○ **Authors:** Dr. Ayesha Khan, Dr. Liam Peterson ○ **Published Year:** 2021 ○ **Review:** This research focuses on the development of predictive models for animal diseases. It emphasizes the use of temporal and spatial data to predict disease spread and severity. The study also discusses the challenges of data quality and availability in building reliable models.

6. **Future Scope of Big Data in Veterinary Science (2023)** ○ **Authors:** Dr. Clara Martinez, Dr. Jonathan Brown ○ **Published Year:** 2023 ○ **Review:** This recent publication explores the potential future applications of big data in veterinary science. It covers emerging trends like the integration of AI, cloud computing, and real-time analytics for disease management. The authors suggest a roadmap for adopting these technologies to enhance veterinary practices globally.

CHAPTER 3: TECHNOLOGY USED

3.1 Apache Hive

Apache Hive is a data warehousing tool built on top of Hadoop, designed for querying and analyzing large datasets. Hive uses HiveQL, a SQL-like language, for interacting with data stored in distributed systems.

Key Features:

- **Scalability:** Efficient handling of massive datasets.
- **Flexibility:** Supports structured and semi-structured data.

-

Integration: Compatible with Hadoop's ecosystem.

3.2 HiveQL

HiveQL simplifies data querying, providing SQL-like syntax for operations like filtering, grouping, and aggregation.

3.3 Architecture

Hive's architecture includes:

- **Driver:** Manages HiveQL execution.
- **Metastore:** Stores metadata about the data schema.
- **Execution Engine:** Converts HiveQL into MapReduce jobs.

CHAPTER 4: METHODOLOGY

4.1 Data Collection

The animal disease dataset was sourced from Kaggle, containing records of diseases affecting various species across regions. Data included attributes like disease type, affected species, region, and year.

-

4.2 Data Cleaning and Preprocessing

- **Handling Missing Values:** Replaced missing or invalid entries with appropriate defaults or removed incomplete records.
- **Standardizing Formats:** Unified date and region formats.

Feature Engineering: Added derived attributes such as seasonal disease occurrence.

4.3 Data Analysis

EDA focused on identifying:

- Most affected species and regions.
- Disease trends over years.
- Seasonal patterns in outbreaks.

4.4 Predictive Modeling

Models predicted disease outbreak probabilities using attributes like region, species, and historical trends. Features were optimized for accuracy.

CHAPTER 5: IMPLEMENTATION AND OUTPUT

5.1 Attributes

Key attributes in the dataset include:

- Species, Disease, Region, Year, Outbreak Count.

.

5.2 Data Cleaning

Hive's ETL process was used to clean the dataset and prepare it for analysis.

5.3 Exploratory Data Analysis

Insights included:

Top diseases affecting livestock.

- Most vulnerable regions and species.
- Disease seasonality.

5.4 Predictive Modeling

Implemented regression and classification models to predict outbreak risks.

5.5 Advanced Analysis

Analysis highlighted trends in cross-species disease transmission and high-risk regions.

5.6 Output

Visualizations of disease trends and model predictions were generated.

•

CHAPTER 6: RESULTS AND DISCUSSION

6.1 Data Analysis Results

Key findings include the prevalence of certain diseases in specific regions and their seasonality.

6.2 Predictive Modeling Results

Predictions accurately identified high-risk periods and regions for disease outbreaks.

6.3 Discussion

The analysis provides actionable insights for animal health management, emphasizing the importance of early interventions.

CHAPTER 7: CONCLUSION AND FUTURE WORK

7.1 Contributions

The project highlighted the potential of Apache Hive in analyzing large-scale animal disease datasets effectively. It showcased Hive's capability to handle and query vast amounts of structured and semi-structured data, enabling efficient disease surveillance and outbreak management. The insights derived from the data, including the identification of high-risk regions and vulnerable species, have practical applications in improving veterinary health services and policymaking.

7.2 Key Findings

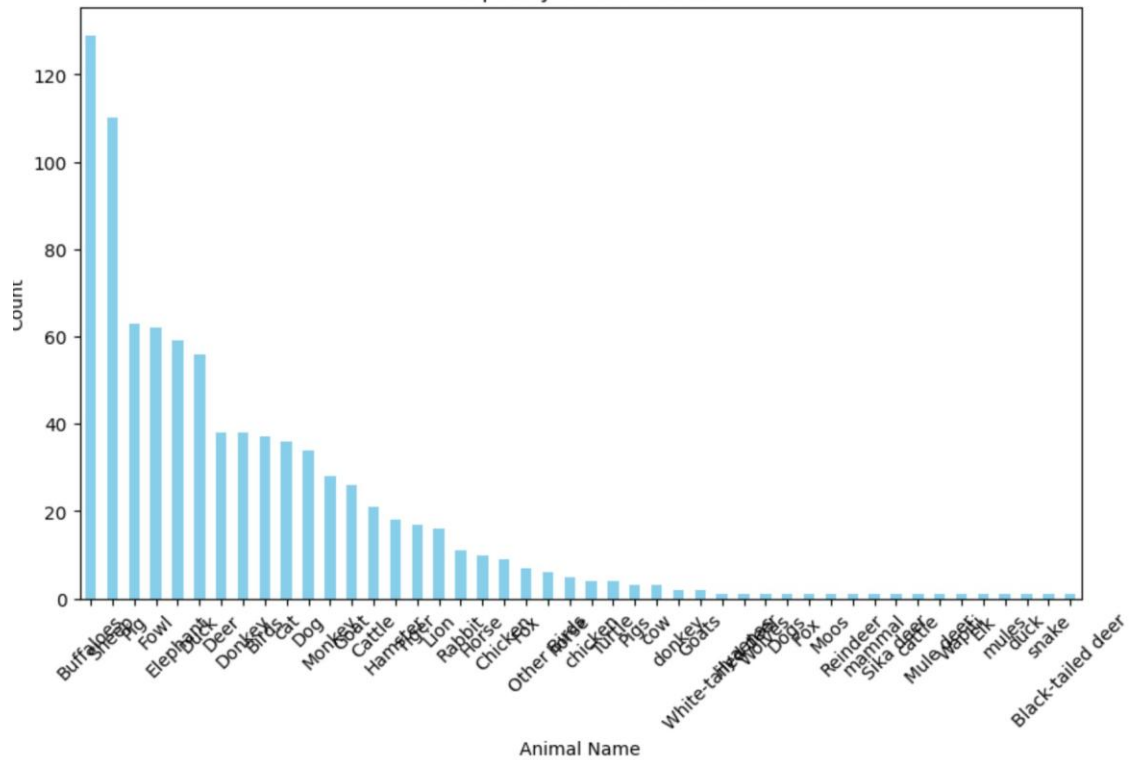
- **Disease Trends:** Seasonal and regional patterns in disease outbreaks were identified, helping prioritize preventive measures.
- **Vulnerable Species:** The analysis pinpointed species most at risk, facilitating targeted intervention.
- **Predictive Insights:** Predictive models proved effective in forecasting potential outbreaks, emphasizing the importance of proactive measures.

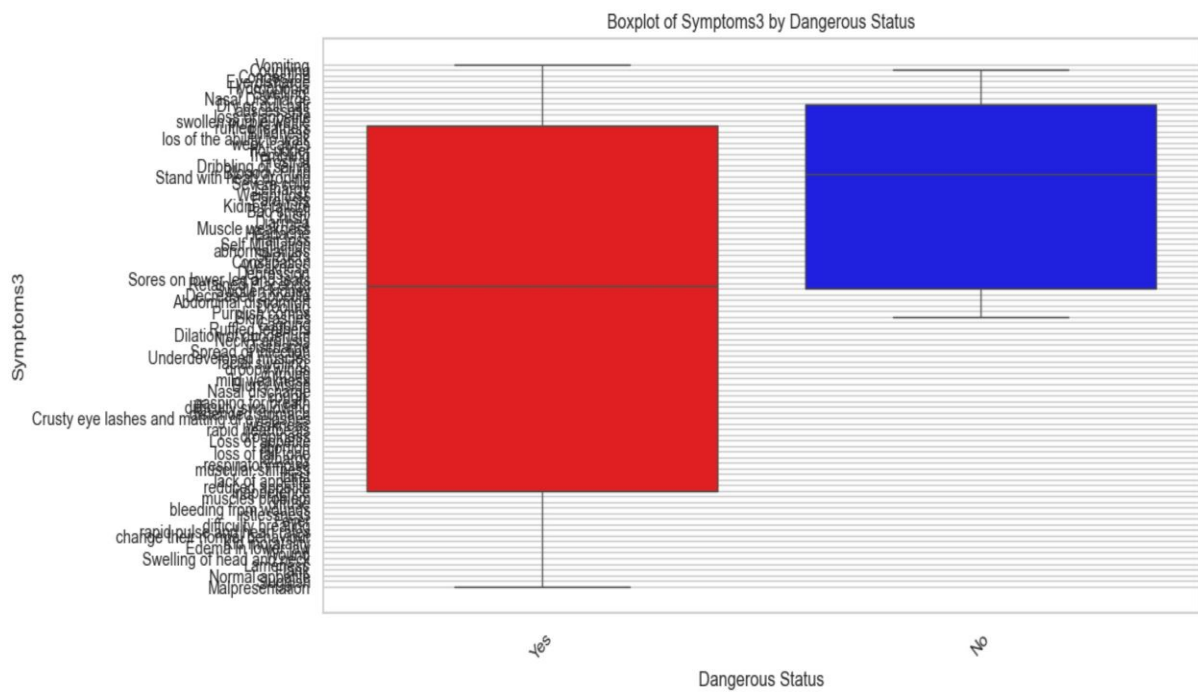
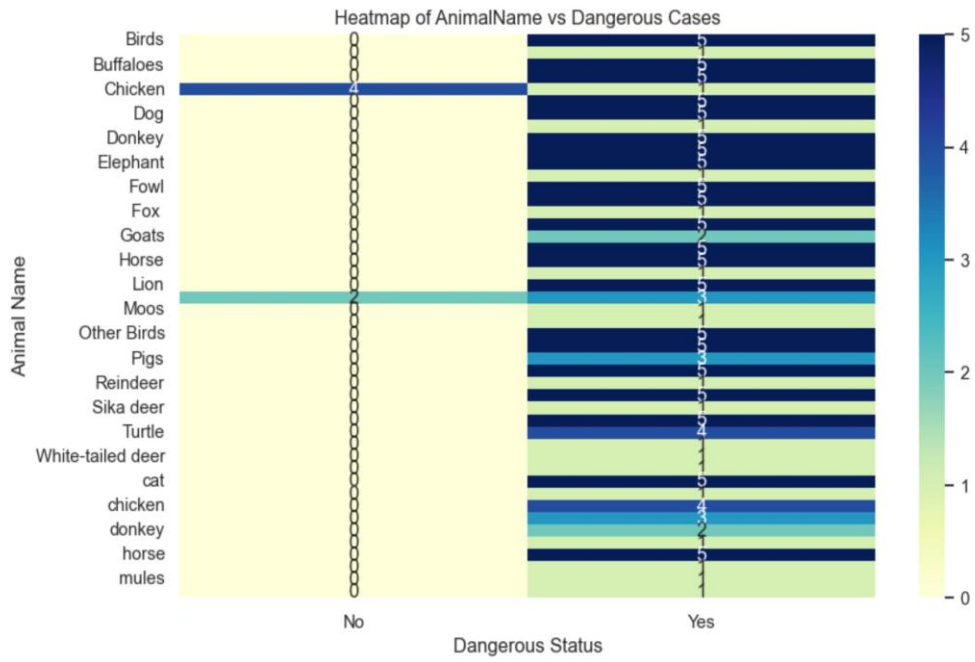
7.3 Future Work

- **Integration of Real-Time Data:** Incorporating real-time reporting mechanisms to enhance the timeliness of outbreak predictions.
- **Advanced Predictive Analytics:** Utilizing more sophisticated machine learning models, such as Random Forests or Neural Networks, for improved accuracy.
- **Expanding Dataset Scope:** Including additional datasets from varied geographical regions to broaden the scope of analysis.
- **Policy Development:** Leveraging insights for drafting policies to manage crossspecies disease transmission and improve global animal health.

Proportion of Dangerous vs Non-Dangerous Cases

Frequency of Each Animal





7.4 Conclusion

This project successfully demonstrated the application of big data analytics in veterinary health using Apache Hive. By analyzing extensive datasets, it provided actionable insights into disease management, prevention strategies, and outbreak predictions. The integration of Hive's scalable architecture with structured disease data enabled effective processing and meaningful analysis. The findings underscore the transformative potential of big data technologies in enhancing veterinary practices and addressing global health challenges related to animal diseases. Future advancements in technology and data accessibility hold promise for even greater impacts in this field.

REFERENCES

1. Kaggle Dataset on Animal Diseases.
2. Research papers on big data and disease management.