

Statistical Approaches to Machine Learning in a Big Data Environment

Zhuoyong Liu*

The High School Attached to Northwest Normal University, Lanzhou, 730030, China

* Corresponding Author Email: 1908216389@qq.com

Abstract. This thesis explores statistical approaches to machine learning in a big data environment. Firstly, the connections and differences between machine learning and statistics in a big data environment are introduced, as well as the statistical foundations in machine learning models. Secondly, the application of statistical methods in big data analysis is discussed, including the combination of traditional data analysis and machine learning. Then, the challenges and limitations of statistical methods in the big data environment, such as high dimensionality and huge amount of data, are discussed. Then, common statistical methods in the big data environment, including linear regression, decision trees, and support vector machines, are described in detail. Finally, the research findings are summarised and future directions and research trends are outlined. Through the research in this paper, a deeper understanding of statistical methods for machine learning in big data environment is provided, which provides an important reference for big data analysis and application.

Keywords: machine learning; statistical methods; big data environment; data analysis.

1. Introduction

With the rapid development of information technology and the popularity of the Internet, the era of big data has arrived, and large-scale data integration and processing has become a reality. In such a context, machine learning, as a powerful data analysis tool, has gradually received widespread attention. However, machine learning in the big data environment still faces many challenges, including the huge amount of data, the uneven quality of data, and the limitation of computational resources.

The aim of this research is to explore in depth the combination of machine learning and statistical methods in big data environments in order to address the challenges and problems in large-scale data analysis. Specifically, we will focus on the application of statistical methods in machine learning models and explore their role in improving model interpretability, reducing overfitting, and handling high-dimensional data. The introductory section will introduce the background and motivation of the research, outlining the purpose and importance of this study. Next, an overview of machine learning in big data environments will be explored in depth, including the characteristics of big data, the role of machine learning in big data processing, and the challenges faced.

2. Overview of machine learning in a big data environment

2.1. Big data concepts and characteristics

Big data refers to data sets that are huge in size, diverse in type, and complex in processing, and is characterised by three main aspects: Volume (large volume of data), Velocity (fast data processing speed), and Variety (diverse data types). Firstly, the volume of big data usually exceeds the processing capacity of traditional data processing tools, and needs to be processed using technologies such as distributed computing. Secondly, the processing speed of big data is demanding, requiring data to be analysed and processed in real-time or near real-time to meet business needs. In the era of big data, effective use of big data resources and mining the potential value in the data are of great significance for decision-making, marketing, product innovation and other aspects of enterprises.

2.2. The role of machine learning methods in big data processing

Machine learning methods play a key role in big data processing. Firstly, machine learning algorithms are able to extract useful information and patterns from huge amounts of data and help users understand the patterns and trends behind the data^[1]. Secondly, machine learning models are able to automate tasks such as data classification, clustering, prediction and optimisation, greatly improving the efficiency and accuracy of data processing. In addition, machine learning is able to discover non-linear relationships and hidden patterns between data, helping users discover new insights and findings. In the big data environment, machine learning algorithms are able to handle data with high dimensionality, high complexity and high noise, with strong robustness and generalisation capabilities.

2.3. Challenges of machine learning in a big data environment

The sheer volume of data may lead to high computational and storage costs, which need to be handled using distributed computing and storage techniques^[2]. Big data tends to have high dimensionality and sparsity, which can increase model complexity and training time, as well as easily trigger dimensionality catastrophe and overfitting problems. In addition, big data often contains a large amount of noise and redundant information, which can reduce the accuracy and interpretability of models.

3. The place of statistical methods in big data analysis

3.1. Fundamentals of statistical methods

Statistical methods are the basis of data analysis, and their fundamentals include both descriptive and inferential statistics. Descriptive statistics reveal the basic characteristics and distribution of data, such as mean, variance and frequency distribution, by summarising, organising and presenting the data^[3]. Inferential statistics, on the other hand, infer the characteristics of the whole by analysing the sample data, including methods such as parameter estimation, hypothesis testing and regression analysis. Among them, parameter estimation is to estimate the value of the overall parameter through sample data, hypothesis testing is to judge whether the hypothesis of the overall parameter is valid or not through sample data, and regression analysis is to study the relationship between the independent variable and the dependent variable^[4]. The basic principles of statistical methods are widely used in practical applications in various fields, such as medicine, finance, social science, etc., providing powerful tools and methods for data analysis and decision-making. In the big data environment, the fundamentals of statistical methods are still applicable, but they need to be combined with advanced technologies such as machine learning to meet the challenges of large data volume, high dimensionality and complexity.

3.2. Application of statistical methods to traditional data analysis

Statistical methods are commonly used in exploratory analysis of data, revealing the distributional characteristics of data and the relationship between variables through descriptive statistical means, providing a basis for further analysis. Statistical methods are widely used in hypothesis testing to verify the validity or otherwise of research hypotheses through statistical inference of sample data^[5]. In addition, statistical methods play an important role in parameter estimation and regression analysis, helping researchers to estimate the values of overall parameters and explore the causal relationships between variables. Statistical methods are widely used in data analysis, decision making and scientific research in fields such as medicine, social science and market research.

3.3. Applicability and limitations of statistical methods in a big data environment

Statistical methods can still play an important role in big data analysis, especially in the exploratory analysis of data, hypothesis testing and parameter estimation, etc., with rich theoretical basis and

practical experience. Statistical methods have certain advantages in dealing with high dimensionality and high complexity data, such as cluster analysis and principal component analysis, which can effectively reduce the dimensionality of data and extract the main features of data. However, statistical methods also have some limitations in the big data environment, such as the computational and storage pressure that may be faced when dealing with large-scale datasets, and traditional statistical methods may not be able to meet the requirements of real-time and efficiency.

4. Integration of machine learning and statistical methods

4.1. Connections and differences between machine learning and statistics

Machine learning and statistics are both important branches in the field of data science, and they are closely related, but there are some differences. Firstly, they have slightly different goals: statistics is mainly concerned with making inferences and testing hypotheses on data to understand the process of data generation and general characteristics, while machine learning focuses more on building models from data to achieve prediction and classification of unknown data^[6]. Secondly, they differ in methods and techniques: statistics usually uses a frequentist approach, focusing on statistical theories such as probability distributions and parameter estimation; whereas machine learning is more inclined to use data-based methods, such as neural networks and support vector machines, focusing on pattern recognition and model optimisation. In addition, machine learning typically focuses more on the efficiency and scalability of algorithms, as they need to deal with large-scale datasets; whereas statistics focuses more on the interpretability of models and the accuracy of statistical inference. Despite these differences, there is also an increasing amount of crossover and integration between machine learning and statistics, which draw on and promote each other, and together they are driving the development of data science.

4.2. Statistical foundations in machine learning models

Many basic principles and methods of statistics are incorporated in machine learning models. Firstly, probability theory is an important foundation in machine learning, and many machine learning models are constructed based on probabilistic models, such as plain Bayesian classifiers and Gaussian mixture models. Second, parameter estimation methods in statistics are widely used in the training process of machine learning models, such as maximum likelihood estimation and least squares for estimating the parameter values of models. In addition, methods such as cross-validation and hypothesis testing in machine learning also originate from statistics and are used to assess the generalisation ability and statistical significance of models^[7]. In addition, the principle of variance-bias trade-off in statistics is widely used in machine learning to adjust the complexity of a model to avoid overfitting or underfitting. In conclusion, the statistical foundation in machine learning models not only provides theoretical support for the construction of the model, but also provides an important guarantee for the model's interpretability, reliability and generalisation ability.

4.3. The role and significance of statistical methods in machine learning

Statistical methods play a crucial role in machine learning, and their significance is mainly reflected in the following aspects. Firstly, statistical methods provide a theoretical foundation and methodological guidance for machine learning, helping to construct reasonable models and algorithms. Secondly, statistical methods can help machine learning extract useful information and patterns from data, and achieve in-depth analysis and understanding of data through parameter estimation, hypothesis testing and other methods. In addition, statistical methods can help machine learning to assess the accuracy and generalisation ability of the model, and verify the reliability and stability of the model through cross-validation, hypothesis testing and other methods. Finally, statistical methods can also help machine learning to solve problems such as data imbalance, missing values and outliers, and improve the robustness and adaptability of models. Therefore, statistical

limitations, such as easy overfitting, sensitivity to noise, etc., which require parameter tuning and model optimisation to improve performance and stability.

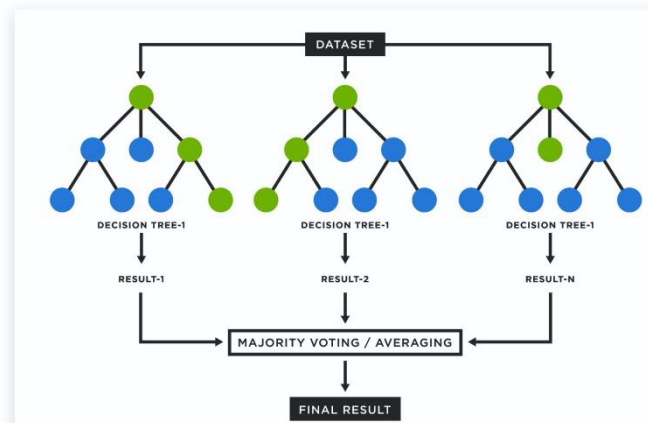


Figure 2. Decision trees and random forests

5.3. Support vector machines

Support Vector Machine (SVM) is a classical supervised learning algorithm commonly used for classification and regression tasks. The basic principle is to find a hyperplane to separate different categories of data and maximise the distance from the hyperplane to the nearest data point. SVM can map linearly indistinguishable data into a high-dimensional space by introducing a kernel function, thus achieving non-linear classification. In big data environment, support vector machine is favoured for its better optimization ability and generalisation in high dimensional space. It has better robustness and generalisation ability, and is well adapted to high dimensional data and sparse data. However, the training process of support vector machines consumes a lot of computational resources and time, and may face performance and efficiency challenges when dealing with large-scale datasets. Therefore, in practical applications, data size and algorithm complexity need to be carefully considered and appropriate parameters and techniques need to be selected to improve the performance and efficiency of support vector machines in big data environments.

5.4. Cluster analysis

Cluster analysis is a commonly used unsupervised learning method that aims to divide the samples in a dataset into a number of categories, such that the samples within the same category have a high degree of similarity and the samples between different categories have a low degree of similarity. In a big data environment, cluster analysis can help to discover underlying patterns and group structures in the data, thus enabling structured and inductive analysis of the data. Common clustering algorithms include K-mean clustering, hierarchical clustering and DBSCAN, etc. K-mean clustering is a distance-based clustering method that assigns samples to the nearest cluster centres by iteratively optimising the distance from the samples to the cluster centres. Hierarchical clustering, on the other hand, is a tree-based clustering method that gradually aggregates samples into different categories through a bottom-up or top-down aggregation process. DBSCAN is a density-based clustering method that identifies cluster clusters by searching for high-density regions, and is suitable for discovering cluster clusters of arbitrary shapes. In practical applications, clustering analysis can help identify user behavioural patterns, market segmentation, anomaly detection and other tasks, providing important references for enterprise decision-making and business optimization.

5.5. Principal component analysis

Principal Component Analysis (PCA) is a commonly used dimensionality reduction technique designed to transform high dimensional data into low dimensional data through linear transformation while retaining most of the information in the data set. PCA achieves dimensionality reduction by

finding the principal components in the data, i.e., the directions in which the data has the highest variance. Principal components are linear combinations of original features with maximum variance and therefore retain the most important information in the data. In a big data environment, principal component analysis can help reduce the dimensionality of data, reduce the storage space and computational cost of the data set, and help discover hidden patterns and laws in the data. The application scenarios of principal component analysis include feature extraction, data visualisation, data compression and so on. Although principal component analysis may lose some information during the dimensionality reduction process, it is still a powerful tool that can effectively handle large-scale data sets and provide useful support for subsequent data analysis and modelling.

6. Summary

6.1. Summarising research findings

Through the research in this paper, we delve into statistical methods for machine learning in big data environments and analyse the related concepts, principles and applications in detail. We found that statistical methods are still significant in big data analysis, especially playing a key role in data exploration, model evaluation and interpretation. At the same time, we also identified challenges and limitations of statistical methods in the big data environment, such as high computational and storage costs, and difficulties in handling high-dimensional and sparse data. In addition, we have found that the combination of machine learning and statistical methods is important for solving complex problems in big data analytics, which can improve the accuracy and robustness of models.

6.2. Challenges and limitations of statistical methods in a big data environment

The high dimensionality and complexity of big data makes it possible that traditional statistical methods may not be able to deal with them effectively, as these methods tend to assume that the data conforms to specific probability distributions, which is often violated by big data. The high growth rate of big data leads to a dramatic increase in computational and storage pressure, and traditional statistical methods may not be able to meet the requirements of real-time and efficiency, requiring the use of parallel computing and distributed storage and other technologies to accelerate the computational process.

References

- [1] Mayhew M , Atighetchi M , Adler A ,et al.Use of machine learning in big data analytics for insider threat detection[C]//MILCOM 2015.IEEE, 2015.
- [2] Alam M , Amjad M .A precipitation forecasting model using machine learning on big data in clouds environment[J].Mausam: Journal of the Meteorological Department of India, 2021(4):72.
- [3] Hullman J , Kapoor S , Nanayakkara P ,et al.The worst of both worlds: A comparative analysis of errors in learning from data in psychology and machine learning[J].arXiv e-prints, 2022.
- [4] Iniesta R , Stahl D , Mcguffin P .Machine learning, statistical learning and the future of biological research in psychiatry[J].Psychological Medicine, 2016, 46(12):2455-2465.
- [5] Nakhaeizadeh G , Taylor C C .Machine learning and statistics : the interface[J].Journal of the American Statistical Association, 1997, 93(442).
- [6] Li-Pang C .Statistical Inference and Machine Learning for Big Data[J].Biometrics, 2023(4):4.
- [7] Franke B , Plante J F , Roscher R ,et al.Statistical Inference, Learning and Models in Big Data[J].International Statistical Review, 2016, 84.