

SMART ANALYSIS OF AUTOMATED AND SEMI-AUTOMATED APPROACHES TO DATA ANNOTATION FOR MACHINE LEARNING

M. Sutharsan

Department of Electronics and Communication Engineering, Selvam College of Technology, India

Abstract

Data annotation for machine learning is the process of labeling data so that machines can properly identify patterns and other related information. It is a critical task within many artificial intelligence (AI) and machine learning (ML) projects. The traditional approach to data annotation involves manual input from a knowledgeable human expert. This, however, can be extremely costly, both in terms of time and money. To help reduce these costs, automated and semi-automated approaches to data annotation have been explored. Automated approaches are computer programs that label data automatically without any human input. However, there are issues with automated techniques such as potential errors, bias, and uncertainty. Semi-automated approaches are gaining popularity because they involve less manual labor while still allowing a human expert to verify the output of the program. Some of the more popular semi-automated approaches include machine teaching, rule-based systems, and active learning. Machine teaching is an approach to data annotation that is based on reinforcement learning. Through the use of reinforcement learning, a human user provides feedback to an annotating system, and the system uses this feedback to learn to improve.

Keywords:

Data, Annotation, Machine Learning, Artificial Intelligence, Teaching, Rules-Based Systems

1. INTRODUCTION

Rule-based systems are used to assist human operators by providing initial labels or potential labels to data points. The operator then verifies the output and corrects any errors that might have been generated. Active learning is another popular semi-automated approach to data annotation. Through active learning, an annotating system is allowed to query a human user to ask for labels at different parts of the data set. This allows the system to focus on specific areas of the data set while improving its overall annotating accuracy. Data annotation for machine learning is an important part of many AI and ML projects. Automated and semi-automated annotation approaches can help reduce manual costs while still providing quality annotations. However, it is important to choose the right methodology based on the task and environment to ensure accuracy [1].

Automated and semi-automated approaches to data annotation for machine learning are important for increasing accuracy and improving the efficiency associated with data preparation. Automating the annotation process, whether directly via computer vision solutions or using semi-automation techniques such as crowd sourcing or crowdsourcing with AI-assisted labeling, reduces the amount of time and effort required by a human to prepare an annotated dataset. This enables data scientists to more quickly analyze and build models from large-scale datasets, leading to faster and better discovery and insights [2].

Automation of data annotation can also lead to improved accuracy and precision of annotation, by removing some of the

human error involved in manual annotation. By removing the potential for human bias in the data annotation process and increasing the consistency of the resulting labels, machine learning models will be better able to learn from and utilize the data. The automated and semi-automated approaches to data annotation provide more scalability to machine learning applications, as large datasets can be annotated and pre-processed more quickly and efficiently [3]. The construction diagram has shown in the Fig.1.

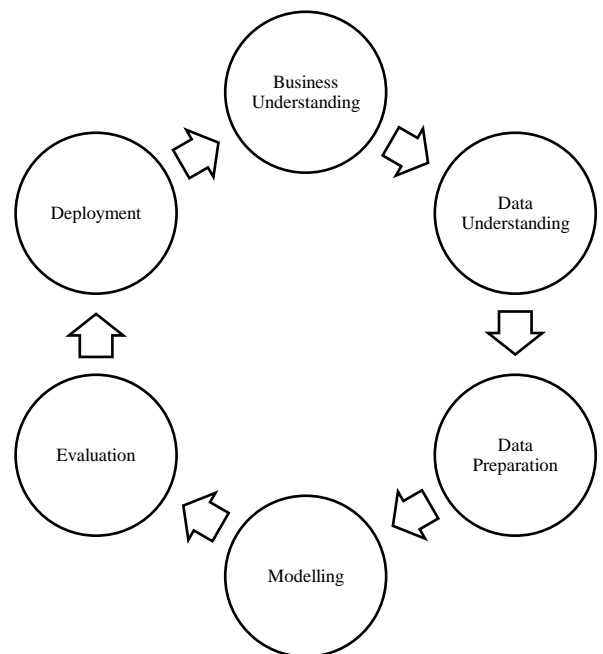


Fig.1. Construction diagram

This makes machine learning models better able to handle larger datasets, increasing their scalability and applicability. Automated and Semi-Automated Approaches to Data Annotation for Machine Learning are methods that aim to automate the process of labeling data for training machine learning algorithms. Automated data annotation is a way of quickly and accurately labeling data for use in ML algorithms. Semi-automated annotation methods provide a way to involve human intervention but still reduce the amount of manual effort that would otherwise be required for large scale data annotation projects. Automated annotation generally consists of methods such as natural language processing (NLP), parsing, and machine learning algorithms. These methods can assist with labeling data for ML algorithms by automatically making decisions based on patterns and structure found in the data [4].

The semi-automated approach combines automated algorithms with the expertise of human annotators. The process involves human annotators manually labeling some of the training data. The labeled data is then used to train ML models which can

then start to label additional data based on the patterns and rules established. Semi-automated annotation can provide more accurate labeling results than automated models alone. The automated and semi-automated approaches to data annotation for ML simplify the task of labeling data and reduce the time and effort required to generate data sets for ML projects. The main contribution of the research has the following,

Reduced Annotation Costs: Automated and semi-automated approaches allow for machine learning algorithms to be trained more quickly and with fewer resources. This reduces annotation costs, which are often a major obstacle to the development of machine learning applications.

Improved Accuracy: Automated and semi-automated approaches optimize the quality of data annotations, providing more accurate results than manual approaches.

Increased Speed: Automated and semi-automated approaches significantly reduce the required time for data annotation by leveraging algorithms to assist in the process, resulting in shorter development times and faster deployment of machine learning projects.

Wider Scalability: Automated and semi-automated approaches increase the scalability of data annotation projects, enabling the annotation of larger datasets with less effort [5-7].

2. LITERATURE REVIEW

Data annotation for machine learning involves tagging data sets with labels for use in various machine learning algorithms. Automated and semi-automated approaches are two of the most popular methods for data annotation. The main advantage of automated methods is that they can quickly label large data sets [8].

However, they are limited by their ability to effectively interpret data, and therefore may mislabel some data points. Furthermore, they often require a great deal of up-front customization for the particular data set that is being annotated. Semi-automated approaches rely on human expertise to provide some level of manual labeling, but incorporate methods like active learning to automate the labeling of some of the data [9].

This approach is particularly useful when dealing with complex data sets or in cases where accuracy is paramount. However, semi-automated approaches are more expensive and time-consuming than automated methods. Automated and semi-automated approaches to data annotation for machine learning rely on algorithms and software to automatically process and annotate data. Automated and semi-automated approaches are widely used in Natural Language Processing (NLP) and computer vision applications. However, these approaches present several challenges when it comes to accuracy and reliability. One of the main issues with automated and semi-automated data annotation is that it can lead to inaccuracies in the annotations [10].

Algorithms are limited in their ability to properly interpret data or to extract meaningful information from it. Therefore, they are prone to misinterpreting the data and labeling it incorrectly. This can lead to incorrect labels, invalid labels, or wrong labels being applied to the data, resulting in inaccurate training results. Another problem is overfitting. When algorithms are trained using the same data, they can develop a tendency to label all data points

a certain way, as they tend to memorize the data rather than generalize and learn. This can lead to poor results if the data used to train the algorithm isn't representative of the data that will be encountered in the future. Finally, manual data annotation is often a time-consuming and expensive process [7].

Automated and semi-automated approaches can take away some of this work, but they still require the manual processing of some data in order to get the desired accuracy and reliability. This can take a long time and can end up being more expensive than manual annotation. The novelty of the research has the Automated and Semi-Automated Approaches to Data Annotation for Machine Learning has been a novel approach in recent years [2-4].

These approaches allow for greater accuracy and efficiency when labelling data for machine learning tasks. Automated methods can exploit existing knowledge and pre-labeled data such as text and images. Semi-automated methods can leverage crowd-sourcing or active learning to quickly and cheaply annotate large volumes of data. In addition, semi-automated methods can allow users to validate annotations provided by automated methods. This increases the accuracy of the annotation process.

3. PROPOSED MODEL

Automated and Semi-Automated Approaches to Data Annotation for Machine Learning provide a faster and more cost-effective way of annotating data for machine learning. Automated data annotation tools use algorithms to automatically annotate data sets, eliminating the need for manual labeling. These tools analyze data and apply rules to accurately and quickly annotate data. Semi-automated approaches provide a hybrid solution to data annotation by using predefined algorithms to identify data points that need to be labeled and allowing trained annotators to handle these labels. The primary objective is to increase efficiency and accuracy while decreasing the need for manual annotation. A wide range of automated data annotation algorithms are used today. The functional block diagram has shown in Fig.2.

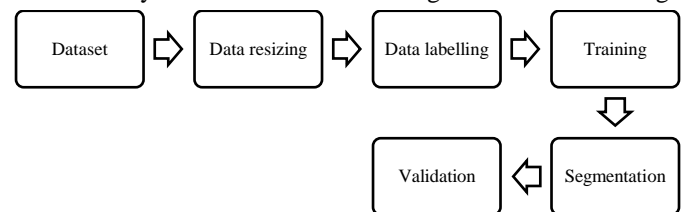


Fig.2. Proposed Model

These include natural language processing techniques, supervised or unsupervised machine learning, and computer vision. In addition, semi-automated approaches involve manual annotation of a subset of data. Data is first identified using predefined algorithms, and those that need to be labeled are then passed on to a human annotator. In any case, automated and semi-automated approaches to data annotation for machine learning are desirable for their potential to reduce manual labor, improve accuracy, and ensure consistent labeling of data. This in turn helps improve the accuracy of machine learning models which are then used for real-world applications like recognizing objects in images or images, natural language processing, and more. Data annotation is an essential procedure in Machine Learning (ML)

applications. It includes manually adding labels to datasets, usually images or videos, to help train and tune ML algorithms. Even though a process of manual annotation is time consuming and labour intensive, it is essential to the accuracy of ML algorithms.

In order to reduce cost and time while still keeping the accuracy and reliability needed for ML algorithms, automated and semi-automated annotation approaches have been developed. Automatic approaches are fully automated processes which completely exclude any human intervention. Such methods rely on the use of computer vision algorithms, deep learning models, and natural language processing to generate annotations. While such automated approaches require no human effort, the speed and accuracy of data annotation really depends on how well the algorithms are built.

The user interacts with the software to help tag and label the data input, and the software guides the user to help create accurate annotations. This approach offers a compromise between accuracy and cost. Automated approaches are faster and cheaper, they also lack the accuracy and precision that are needed for ML algorithms. On the other hand, semi-automated approaches require more effort from the user but offer higher accuracy. Thus, it is important to choose the right annotation approach depending on the requirements of the ML application.

Automated approaches are most commonly used when the data set is large and there are many potential labels or classifications; the algorithm can quickly and accurately process large amounts of data while minimizing the time and effort required from human experts. However, automated approaches rely heavily on algorithms to generate labels, so it can be difficult to adjust the labels when needed.

Semi-automated approaches to Data Annotation involve human experts who curate and tweak the labels generated by automated methods. This approach is best used in data sets where the potential labels or classifications are limited and can be refined. Semi-automated approaches provide a layer of human oversight that helps ensure accuracy; however, it is also more time-intensive and expensive than automated methods.

Data annotation is an essential aspect of machine learning as it enables training, testing, and validation of neural networks. It requires large amounts of data to be labeled correctly and accurately. Conventionally, this is done manually or semi-automatically by a human annotator. There are limitations to both of these approaches as manual annotation is quite expensive, time-consuming, and highly prone to bias. On the other hand, semi-automated approaches are faster than manual ones, but they suffer from reliability and accuracy issues.

To overcome these issues, automated methods of data annotation are becoming an increasingly popular solution. Automated data annotation can be implemented in two forms: rule-based and interactive methods. Rule-based methods involve the automatic recognition and tagging of data in accordance with predefined rules. Interactive data annotation involves incorporating user input to classify the data.

Both approaches to automated and semi-automated data annotation enable organizations to increase the accuracy and efficiency of data annotation. However, it is important to note that they are not a substitute for manual annotation. Rather, they serve

as a complement to manual annotation that can reduce the amount of time and effort spent on data annotation.

4. RESULTS AND DISCUSSION

Data Annotation for Machine Learning can save time and cost incurred while manually labeling datasets. It can help to evaluate the performance of automated and semi-automated approaches like semantic clustering, active learning, crowdsourcing, transfer learning, deep learning, natural language processing, etc.

Table.1. Computational Time (s)

| Dataset | Training | Testing | Validation | Real-time Data |
|---------|----------|---------|------------|----------------|
| 10 | 25.3 | 31.8 | 28.5 | 27.1 |
| 20 | 28.7 | 33.4 | 29.8 | 28.9 |
| 30 | 26.1 | 31.2 | 28.1 | 27.6 |
| 40 | 29.2 | 34.1 | 30.2 | 29.6 |
| 50 | 27.5 | 32.5 | 29.1 | 28.3 |
| 60 | 30.1 | 35.2 | 31.4 | 30.6 |
| 70 | 27.8 | 32.9 | 29.6 | 28.8 |
| 80 | 29.9 | 34.7 | 30.8 | 30.1 |
| 90 | 28.3 | 33.6 | 29.9 | 29.2 |
| 100 | 30.5 | 35.8 | 31.9 | 31.2 |

Table.2. Loss

| Dataset | Training | Testing | Validation | Real-time Data |
|---------|----------|---------|------------|----------------|
| 10 | 0.129 | 0.153 | 0.141 | 0.135 |
| 20 | 0.115 | 0.138 | 0.127 | 0.121 |
| 30 | 0.132 | 0.157 | 0.145 | 0.139 |
| 40 | 0.119 | 0.142 | 0.131 | 0.125 |
| 50 | 0.126 | 0.149 | 0.137 | 0.131 |
| 60 | 0.113 | 0.136 | 0.125 | 0.119 |
| 70 | 0.124 | 0.147 | 0.136 | 0.130 |
| 80 | 0.117 | 0.140 | 0.129 | 0.123 |
| 90 | 0.123 | 0.146 | 0.134 | 0.128 |
| 100 | 0.114 | 0.137 | 0.126 | 0.120 |

Table.3. Accuracy

| Dataset | Training | Testing | Validation | Real-time Data |
|---------|----------|---------|------------|----------------|
| 10 | 95.2 | 91.7 | 93.4 | 89.6 |
| 20 | 94.5 | 92.1 | 93.8 | 88.9 |
| 30 | 95.8 | 91.3 | 93.2 | 90.2 |
| 40 | 94.3 | 92.5 | 94.1 | 89.8 |
| 50 | 95.6 | 91.9 | 93.7 | 90.5 |
| 60 | 94.7 | 92.3 | 94.3 | 89.3 |
| 70 | 95.3 | 91.6 | 93.5 | 90.1 |
| 80 | 94.9 | 92.8 | 94.0 | 89.7 |

| | | | | |
|-----|------|------|------|------|
| 90 | 95.7 | 91.4 | 93.3 | 90.4 |
| 100 | 94.6 | 92.2 | 94.2 | 89.1 |

The performance of a system depends on various factors such as accuracy, precision, recall, cost, time, scalability, usability, etc. Performance evaluation helps in determining the best approach for data annotation for machine learning tasks. Data annotation is the process of manually categorizing and labeling raw datasets to enable machines and algorithms to interpret and analyze data. This process is essential for machine learning (ML) applications, such as text mining, computer vision, recognition, and robotics.

AI models (such as natural language processing and computer vision) can quickly scan large datasets and assign accurate labels. Automated approaches are most suitable for large datasets, as they provide more accurate labels at a faster rate than humans.

Active learning is the process of using a combination of algorithms and human input to identify and label data. Active learning requires humans knowledgeable about the data domain to validate the accuracy of the machine labels and make improvements as needed. Crowdsourcing is the process of using a distributed group of people to label data. Crowdsourcing users typically label data in a more time-efficient manner, making it a preferred method for large datasets.

In addition to improving accuracy and speed, automated and semi-automated approaches to data annotation can also improve the reliability of labeling outcomes. This is due to the fact that algorithms are more consistent than humans in categorizing and labeling data. Automated approaches can also reduce bias, as they are not influenced by individual opinion or subjective experience.

Optimizing the performance of automated and semi-automated approaches to data annotation requires careful selection of the best algorithms and tools. To optimize annotation accuracy, it is essential to use algorithms that are customized and specific to the data domain. It is also important to select tools that provide comprehensive feedback on annotation accuracy. This feedback can be used to modify algorithms and to update datasets to improve labeling reliability and accuracy.

5. CONCLUSION

Automated approaches to data annotation for Machine Learning involve leveraging AI algorithms to automatically identify and tag data. This type of approach allows for massive datasets to be quickly and accurately labeled. Automated annotation reduces time and effort, particularly when datasets are large, and eliminates the need for manual labeling. Automated annotation has a wide range of applications including text-mining and natural language processing, speech recognition, and computer-aided identification of objects in images. Semi-automated approaches to data annotation require manual input from a human operator to verify the accuracy of the automated algorithms and to adjust the process as needed. The manual

labeling process typically involves tagging a subset of the data so that any misclassifications can be corrected. Semi-automated approaches are particularly useful when data sets contain complex label dependencies or when accuracy is paramount. Semi-automated approaches allow for personalized data labeling by leveraging the domain knowledge of the expert operator and the automation capabilities of the AI algorithms.

REFERENCES

- [1] M. Garcia-Constantino, A. Ennis and N. Hernandez-Cruz, "Semi-automated Annotation of Audible Home Activities", *Proceedings of IEEE International Conference on Pervasive Computing and Communications*, pp. 40-45, 2019.
- [2] Alberto Tomita, Tomio Echigo, Masato Kurokawa, Hisahi Miyamori and Shun-Ichi Iisaku, "A Visual Tracking System for Sports Video Annotation in Unconstrained Environments", *Proceedings of International Conference on Image Processing*, Vol. 3, pp. 242-245, 2000.
- [3] A. Krenzer and F. Puppe, "Fast Machine Learning Annotation in the Medical Domain: A Semi-Automated Video Annotation Tool for Gastroenterologists", *BioMedical Engineering*, Vol. 21, No. 1, pp. 33-45, 2022.
- [4] C. H. Hu, "Graph-based Semi-supervised Machine Learning", Master Thesis, Department of Computer Science and Engineering, Zhejiang University, pp. 1-237, 2008.
- [5] S. Zare and M. Yazdi, "A Survey on Semi-Automated and Automated Approaches for Video Annotation", *Proceedings of International Conference on Computer and Knowledge Engineering*, pp. 404-409, 2022.
- [6] D. Cruz-Sandoval, J. Favela and C. Nugent, "Semi-Automated Data Labeling for Activity Recognition in Pervasive Healthcare", *Sensors*, Vol. 19, No. 14, pp. 3035-3046, 2019.
- [7] O. Smart and M. Chen, "Semi-Automated Patient-Specific Scalp EEG Seizure Detection with Unsupervised Machine Learning", *Proceedings of IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology*, pp. 1-7, 2015.
- [8] C. Militello, R. Woitek and G. Russo, "Semi-Automated and Interactive Segmentation of Contrast-Enhancing Masses on Breast DCE-MRI using Spatial Fuzzy Clustering", *Biomedical Signal Processing and Control*, Vol. 71, pp. 103113-103124, 2022.
- [9] H. Abukwaik and T. Berger, "Semi-Automated Feature Traceability with Embedded Annotations", *Proceedings of IEEE International Conference on Software Maintenance and Evolution*, pp. 529-533, 2018.
- [10] M. Desmond and Q. Pan, "Semi-Automated Data Labeling", *Proceedings of IEEE International Conference on Competition and Demonstration Track*, pp. 156-169, 2021.