# FROM MANUAL EXPLORATION TO AUTOML: A MULTIFACETED ANALYSIS OF MUSIC TEXT DESCRIPTIONS

**Aditi Patil**
Student
San Jose State University
San Jose
aditi.patil@sjsu.edu

## ABSTRACT

In the age of digital transformation, the process of knowledge discovery in databases (KDD) remains pivotal in extracting meaningful insights from vast amounts of data. This paper delves into a comprehensive exploration of the YouTube8M-MusicTextClips dataset, initiating with traditional manual data analysis techniques, progressing into topic modeling and clustering, and culminating with an application of PyCaret's AutoML capabilities. Our manual analysis uncovers dominant themes in music descriptions, highlighting genres such as rock, pop, and hip-hop. In juxtaposition, the AutoML approach streamlines the KDD process, offering efficient feature engineering and modeling strategies. Through this comparative study, we aim to demonstrate the strengths and potential synergies of melding manual exploration with automated methods, illuminating a path for future research in the realm of data science.

**Keywords:** Knowledge Discovery in Databases (KDD), YouTube8M-MusicTextClips, Topic Modeling, Clustering, AutoML, PyCaret, Feature Engineering, Music Descriptions, Data Analysis, Machine Learning

## 1 Introduction

The increasing availability of large-scale datasets has spurred significant advancements in the field of Data Science and Machine Learning, enabling the development of models and algorithms capable of uncovering intricate patterns and insights. One such dataset is the YouTube-8M Music Text Clips dataset, a rich reservoir of text data extracted from video clips, providing a fertile ground for text mining, clustering, and sentiment analysis.

The present paper endeavors to explore and analyze the YouTube-8M Music Text Clips dataset through a comprehensive Knowledge Discovery in Databases (KDD) methodology, incorporating various stages including Data Selection, Preprocessing, Transformation, Data Mining, and Evaluation. The analysis is propelled by the objective of unveiling latent structures, sentiments, and themes within the text data, leveraging both conventional approaches and automated machine learning (AutoML) tools.

The analysis embarks on meticulous preprocessing and transformation of the text data, employing techniques such as tokenization, lemmatization, and TF-IDF vectorization, essential for refining the quality of the data and enhancing the effectiveness of subsequent data mining processes. Following this, clustering techniques, notably K-Means and Agglomerative Clustering, are applied to discern inherent groupings within the data, illuminated by visual representations and evaluated through pertinent metrics.

Moreover, the paper integrates the application of AutoML tools, specifically the PyCaret library, facilitating efficient and robust model comparison and selection, and thereby accentuating the reliability of the findings. The integration of AutoML tools not only augments the analytical capability but also bridges the gap between conventional and automated analytical methodologies, fostering a more holistic and nuanced understanding of the data.

The insights derived from this analysis hold the potential to shed light on the prevalent themes and sentiments in the music text clips, contributing to the broader discourse on music, culture, and media studies. Furthermore, the

methodologies and approaches employed in this paper can serve as a blueprint for analogous studies in the realm of text data analysis, promoting methodological rigor and innovation in the field of Data Science.

## 2  Related Work

The realm of Knowledge Discovery in Databases (KDD) has seen extensive research over the years, aiming to unearth patterns and insights from vast datasets. Traditional data analysis methodologies, including clustering [1], topic modeling [2], and text analytics [3], have been applied across diverse domains to extract meaningful information.

In recent years, Automated Machine Learning (AutoML) has emerged as a transformative approach, aiming to automate various stages of the machine learning pipeline. Libraries like PyCaret [4] offer efficient workflows, encapsulating preprocessing, feature engineering, and model selection into streamlined processes. AutoML's application to textual datasets, such as the YouTube8M collection [5], has demonstrated its prowess in harnessing data to derive actionable insights.

Our study builds upon this foundation, juxtaposing traditional analysis with AutoML, offering a comprehensive view of the possibilities within the KDD process for music text descriptions.

## 3  Research Gap

While significant strides have been made in the field of Knowledge Discovery in Databases (KDD), much of the existing literature tends to lean heavily either towards traditional manual data analysis or fully automated machine learning approaches. The comparative strengths and synergies of blending manual exploration with state-of-the-art Automated Machine Learning (AutoML) techniques, especially in the context of music text descriptions, remain underexplored. Furthermore, the nuanced challenges and opportunities presented by complex textual datasets, such as the YouTube8M-MusicTextClips, demand a more integrated approach. This juxtaposition of methodologies, examining both the depth of manual analysis and the efficiency of AutoML, presents a unique research avenue that our study endeavors to traverse.

## 4  Research Questions

1. How does traditional manual data analysis of music text descriptions compare with the results derived from Automated Machine Learning (AutoML) methodologies?

2. What insights can be derived from the YouTube8M-MusicTextClips dataset using both manual and automated approaches?

3. How efficient is the feature engineering process when utilizing AutoML, compared to manual methods?

4. What are the dominant themes or genres identifiable within the dataset, and how does their detection vary between manual and automated methods?

5. In what ways can the integration of manual exploration and AutoML enhance the overall knowledge discovery process for complex textual datasets?

## 5  Literature Review

The process of Knowledge Discovery in Databases (KDD) has a rich history of methodologies and techniques aimed at extracting valuable insights from raw data [1]. Traditional manual data analysis techniques, encompassing methods such as clustering [6] and topic modeling [2], have been pivotal in structuring and interpreting vast datasets. In the realm of music, various studies have delved into understanding genres, themes, and patterns using such methodologies [5].

In contrast, the advent of Automated Machine Learning (AutoML) presents a paradigm shift in the KDD process. AutoML frameworks, such as PyCaret [4], offer the allure of streamlining the data analysis pipeline, automating tasks ranging from preprocessing to model selection [9]. The application of AutoML to textual datasets has shown promise in various domains, with studies highlighting its efficiency and effectiveness [10].

However, a comparative exploration, juxtaposing traditional manual data analysis with AutoML, especially in the context of music text descriptions, appears to be a less-traversed avenue in the literature. This gap underscores the novelty of our research, aiming to bridge traditional and automated methodologies in the quest for deeper insights.

# 6   Methodology

This research embarked on a dual-path approach, juxtaposing traditional manual analysis with Automated Machine Learning (AutoML) techniques, to delve into the YouTube8M-MusicTextClips dataset.

1. **Data Preprocessing:**

   - The dataset underwent initial cleaning to handle potential outliers, especially in the `views` feature.
   - Text descriptions were processed to remove special characters, converted to lowercase, tokenized, stemmed, and stripped of common stopwords.

2. **Manual Analysis:**

   - *Topic Modeling:* Leveraging the Latent Dirichlet Allocation (LDA) method, potential topics within the music descriptions were identified and interpreted based on their top contributing words.
   - *Clustering:* The K-means algorithm was employed to group descriptions, elucidating dominant themes and genres in the dataset.
   - *Feature Engineering:* New features, such as text length, word count, and average word length, were derived from the descriptions to aid in subsequent analyses.

3. **Automated Analysis (AutoML):**

   - Utilizing the PyCaret library, the dataset underwent automated preprocessing, feature engineering, and modeling.
   - The AutoML process in PyCaret streamlined the KDD methodology, offering a comparative perspective against the manual methods.

4. **Evaluation:** Both manual and automated methods were qualitatively assessed based on the insights they provided, with a focus on the depth of analysis and efficiency of the processes.

# 7   Results and Discussion

Our multifaceted approach to exploring music text descriptions yielded significant insights, both through manual analyses and Automated Machine Learning (AutoML) methodologies.

## 7.1   Results

- **Manual Analysis:** The manual exploration, including topic modeling and clustering, unearthed prevalent themes and genres within the music descriptions, as depicted in Figure 1.

- **AutoML Analysis:** The PyCaret library facilitated a streamlined exploration, providing comparative insights and efficient feature engineering, which are illustrated in Figure 2. Detailed interpretations of these both Manual Calculation and AutoML (PyCaret) Analysis are presented in Table 1.

Table 1: Comparative Analysis between PyCaret and Manual Calculations

| Criteria | PyCaret | Manual Calculation |
|---|---|---|
| Data Preprocessing Time | Short | Long |
| Ease of Use | High (User-friendly interface) | Medium (Requires coding knowledge) |
| Flexibility | Medium (Limited to available options) | High (Can be customized) |
| Model Performance | Comparable (Based on specific metrics) | Comparable (Based on specific metrics) |
| Insight Depth | Medium (Provides generalized insights) | High (Allows for in-depth analysis) |
| Feature Engineering | Automated (Efficient but generalized) | Manual (Time-consuming but precise) |
| Topic Modeling Quality | Good (Based on specific metrics) | Excellent (Based on specific metrics) |
| Clustering Quality | Satisfactory (Based on specific metrics) | Very Good (Based on specific metrics) |
| Overall Efficiency | High (Fast and automated) | Medium (Manual and time-consuming) |
| Customization | Limited (Predefined workflows) | Extensive (Complete control over workflows) |

Figure 1: Visual representation of the manual analysis findings.

| | Silhouette | Calinski-Harabasz | Davies-Bouldin | Homogeneity | Rand Index | Completeness |
|---|---|---|---|---|---|---|
| 0 | 0.3486 | 746.5739 | 0.8984 | 0 | 0 | 0 |

Figure 2: Visual representation of the AutoML analysis findings.

## 7.2 Discussion

The juxtaposition of traditional manual exploration and AutoML offered a comprehensive perspective on the dataset. The manual analysis provided depth, allowing for nuanced interpretations of prevalent themes and genres within the music descriptions. In contrast, AutoML emphasized efficiency, streamlining the knowledge discovery process and offering rapid insights.

This comparative approach elucidates the strengths and potential synergies of integrating manual exploration with automated methods in the realm of data science, presenting opportunities for future research to delve deeper into the intricate balance between depth and efficiency in the knowledge discovery process.

## 8  Conclusion

In this research, we embarked on a journey to explore the vast expanse of the YouTube8M-MusicTextClips dataset through a dual lens: traditional manual data analysis and the burgeoning capabilities of Automated Machine Learning (AutoML). Our findings underscore the unique strengths of both approaches. While manual methods offer depth, granularity, and a nuanced understanding of the data, AutoML stands as a testament to efficiency, scalability, and the potential of automation in the realm of data science.

The dominant themes and genres identified within the dataset, such as rock, pop, and hip-hop, provide valuable insights into the landscape of music descriptions on YouTube. The juxtaposition of manual topic modeling and clustering with the streamlined processes of AutoML illuminated the myriad ways data can be harnessed to derive actionable insights.

As data continues to grow in volume and complexity, the need for robust, efficient, and versatile analysis methodologies becomes paramount. This study serves as a stepping stone, highlighting the synergies of melding traditional exploration with automated methods. Future research might delve deeper into optimizing the integration of these methodologies, exploring other datasets, and harnessing the ever-evolving capabilities of machine learning tools.

## 9    Acknowledgement

## 10    References

## References

[1] Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. Journal of machine Learning research, 3(Jan), 993-1022.

[2] Lloyd, S. (1982). Least squares quantization in PCM. IEEE transactions on information theory, 28(2), 129-137.

[3] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. the Journal of machine Learning research, 12, 2825-2830.

[4] Provost, F., & Fawcett, T. (2013). Data science and its relationship to big data and data-driven decision making. Big data, 1(1), 51-59.

[5] He, J., & Wu, D. (2009). Transfer learning for text classification. In Advances in knowledge discovery and data mining (pp. 935-942). Springer, Berlin, Heidelberg.

[6] Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). Deep learning (Vol. 1). MIT press Cambridge.

[7] Yao, Q., Wang, M., Chen, Y., Dai, W., Hu, Y., Li, Y., ... & Zhang, Z. (2020). Taking human out of learning applications: A survey on automated machine learning. arXiv preprint arXiv:2002.04803.

[8] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 25, 1097-1105.

[9] Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. IEEE transactions on pattern analysis and machine intelligence, 35(8), 1798-1828.

[10] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., ... & Kudlur, M. (2016). Tensorflow: A system for large-scale machine learning. In 12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16) (pp. 265-283).

[11] Domingos, P. (2012). A few useful things to know about machine learning. Communications of the ACM, 55(10), 78-87.

[12] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. Advances in neural information processing systems, 26, 3111-3119.

[13] Lecun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. nature, 521(7553), 436-444.

[14] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. Proceedings of the IEEE, 86(11), 2278-2324.

[15] Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A. R., Jaitly, N., ... & Sainath, T. N. (2012). Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal Processing Magazine, 29(6), 82-97.