# Meta Kaggle and Trends in Data Science Machine Learning

# MetaComp Analytica

Sushma Nandiyawar, Chiranthan Shadaksharaswamy, Dhanush Bharath Raj, Dheeraj Manchandia, Priti Singh

*Abstract*—This research project delves into the analysis and visualization of Meta Kaggle data to uncover insights into team performance, collaboration dynamics, and trends in data science and machine learning over the past decade. By leveraging comprehensive datasets from Kaggle competitions, the study aims to elucidate seasonal patterns, visualize collaboration networks, and conduct a thorough analysis of team performance metrics. Through advanced visualization techniques and meticulous data analysis, the project seeks to provide valuable insights that inform future participation strategies and enhance competitiveness in Kaggle competitions. The study also aims to showcase the power of information visualization in distilling complex competition dynamics and fostering improvements in future competitions.

*Index Terms*—Competition dynamics, Collaboration, Performance analysis, Sentiment trends, and Visualization techniques.

---

## INTRODUCTION

Kaggle has emerged as a prominent platform for data science and machine learning enthusiasts to engage in competitive challenges, collaborate with peers, and showcase their expertise. Over the past decade, Kaggle has witnessed significant growth in participation, diversity of competitions, and evolution in collaboration dynamics among participants. Understanding the underlying trends, performance metrics, and collaboration networks within Kaggle competitions is essential for participants, organizers, and stakeholders to make informed decisions and drive improvements in competition design and engagement strategies.

In this research project, we aim to analyze and visualize Meta Kaggle data to gain insights into various aspects of Kaggle competitions. By addressing key objectives such as temporal analysis, collaboration dynamics, and team performance analysis, we seek to uncover patterns, trends, and correlations that shed light on the evolving landscape of data science and machine learning as reflected through Kaggle competitions. Through the application of advanced visualization techniques and rigorous data analysis, we aim to provide actionable insights that empower participants to enhance their competitiveness and drive innovation within the Kaggle community. Additionally, the project aims to contribute to the broader understanding of competition dynamics and collaboration networks in data science and machine learning domains.

## PROJECT OBJECTIVE

The primary objective of our research project is to leverage the Meta Kaggle and Meta Kaggle Code datasets to gain comprehensive insights into team performance, collaboration dynamics, and trends in data science and machine learning over the past decade. Through a meticulous analysis of competition data and code repositories sourced from Kaggle, our focus lies on three key aspects: temporal analysis, collaboration dynamics, and team performance. Initially, we aim to delve into temporal insights, with a particular emphasis on discerning seasonal patterns within Kaggle competitions. By scrutinizing the fluctuations in competition dynamics, we anticipate uncovering invaluable insights into the evolution of trends and their consequential impact on team performance throughout the duration of competitions. Subsequently, we plan to highlight the critical role of collaboration dynamics within Kaggle competitions, recognizing its significance in fostering innovation and driving performance. Employing network analysis techniques, we seek to visualize the formation and evolution of collaboration networks among participants. Finally, our project prioritizes an in-depth analysis of team performance, aiming to identify key success factors through advanced visualization techniques. By dissecting performance metrics and exploring strategies, our endeavor is to provide actionable insights that can inform future participation and enhance overall competitiveness in Kaggle competitions. Through this endeavor, we aim to showcase the power of information visualization in distilling complex competition dynamics and fostering improvements in future iterations.

## RELATED WORK

The basis of our information visualization initiative lies in a thorough examination of prior work, extracting insights from both Kaggle competitions and associated research. Earlier analyses have brought to light the pivotal role of collaborative dynamics, The client's suite of notebooks delves into various questions surrounding Kaggle competitions and user behavior. Through Calculate_dice_performance_metric.ipynb, the client aims to evaluate model performance in segmentation tasks by calculating mean dice coefficients between prediction and ground truth masks. Count_number_of_ftu_in_mask.ipynb provides insights into image composition by exploring the prevalence of specific features within segmentation masks. Generate_donor_age_sex_distribution_plot.ipynb visualizes demographic patterns across different organs, addressing questions about donor age and sex distributions. Mask_conversion.ipynb facilitates the manipulation and conversion of segmentation masks between different formats, aiding in mask analysis and visualization. Additionally, Compare User Tier for ALL Kaggle.ipynb investigates questions about user distributions across Kaggle tiers, offering insights into user engagement levels within the Kaggle community. These notebooks collectively address a diverse range of questions related to competition performance, data analysis, user engagement, and competition dynamics on the Kaggle platform, providing valuable insights for users and organizers alike.

## DATA

The dataset on competitions provides a comprehensive overview of competitions held over multiple years. The number of entities (rows) corresponds to the number of competitions analyzed, with each entity containing attributes such as EnabledDate (start date), Category (e.g., Featured, Research, Community, Playground), and Quantity of Competitive Events (number of competitions within each category). These attributes facilitate trend analysis and pattern recognition in competition participation over time and across different categories.

On the other hand, the dataset on forum messages comprises approximately 2,046,845 unique entities, each representing a forum message posted by users/participants from 2010 to 2024. Major entity attributes include UserID (unique identifier), ForumTopicId (category identifier), PostUserId (poster's identifier), PostDate (timestamp), ReplyToForumMessageId (parent message identifier), Message (content), Medals (awards), and MedalAwardDate (award date). These attributes enable the analysis of user interactions, sentiment trends, and community engagement within Kaggle forums, providing valuable insights for visualization and analysis.

The visualization endeavors to unravel the dominant segment within Kaggle by illustrating the distribution of competitions across different categories. Through a bar chart representation, it becomes apparent that the Community segment, characterized by user-hosted competitions fostering collaboration and diverse topics, emerged as the primary facilitator of Kaggle competitions. This visual insight underscores the significance of community-driven initiatives within the Kaggle platform, highlighting their pivotal role in driving engagement and fostering a collaborative environment for data science enthusiasts of varying expertise levels.
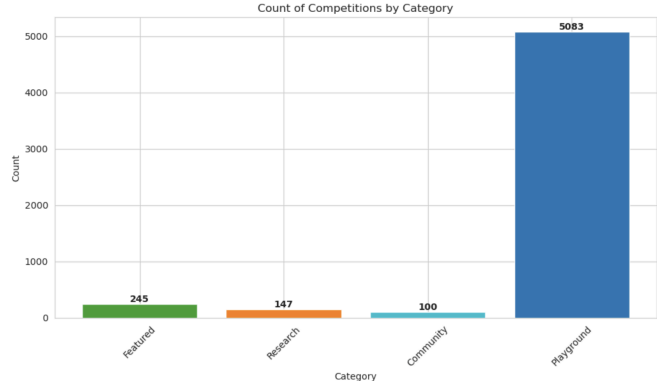


Fig. 1.1. Count of Competitions by Category.

The visualization delves into the landscape of competitor participation across different segments of Kaggle competitions, offering a yearly analysis of Featured, Research, and Playground competitions in comparison to community-hosted events. Employing a stacked bar chart representation, it provides a clear depiction of the annual count of competitions hosted, excluding those within the community segment. By segmenting competitions based on their nature and purpose, the visualization facilitates a nuanced understanding of the competition dynamics within Kaggle, shedding light on the varying levels of engagement across different categories. This visual exploration enables stakeholders to discern trends, patterns, and shifts in participation over time, contributing to a deeper comprehension of the evolving landscape of data science and machine learning competitions on the Kaggle platform.
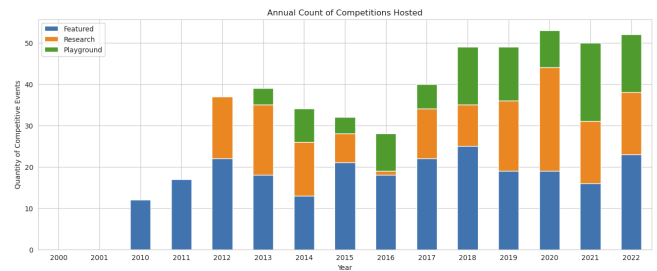


Fig. 1.2. Annual Count of Competitions hosted.

The "Evolving Emotions: A Decade of Sentiment Trends" visualization offers a compelling insight into the shifting landscape of sentiments within the Kaggle community over the past decade. Leveraging the TextBlob library for sentiment analysis, the visualization utilizes a line chart representation to illustrate sentiment trends from 2010 to 2024. By plotting the average positive and negative sentiments over the years, the graph unveils a notable increase in positive sentiments, contrasting with the relative stability of negative sentiments. Through polarity scores ranging from -1 to 1, TextBlob provides a nuanced understanding of sentiment dynamics, reflecting the degree of positivity or negativity within the community discourse. This visualization not only captures the evolving emotional trajectory of the Kaggle community but also underscores the platform's evolving nature and the changing attitudes of its participants over time.

The visualization utilizing a box plot effectively distinguishes between positive and negative sentiment, with green boxes representing positive sentiment and red boxes indicating negative sentiment. Each box is accompanied by a triangle pointing upwards

for positive sentiment and downwards for negative sentiment, enhancing clarity and visual distinction. This approach enables viewers to easily discern the distribution and variation of sentiment scores within the dataset. The use of color coding and symbols enhances the interpretability of the visualization, facilitating quick comprehension of sentiment trends and patterns. Overall, this visualization technique offers a clear and intuitive representation of sentiment analysis results, aiding in the understanding of sentiment dynamics over time or across different segments.
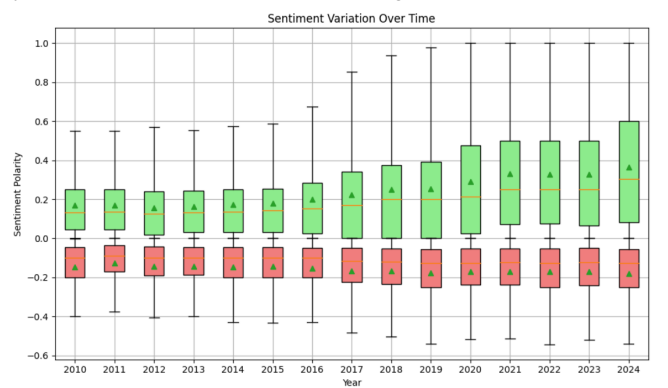


Fig. 2.1. Sentiment Variation Over Time after omitting outliers.

The representation using a bar chart effectively captures a discernible shift towards more positive sentiments within the Kaggle community over time, indicating a notable change in tone or perception. By visually presenting the prevalence of positive sentiments in contrast to negative ones, the bar chart offers a straightforward depiction of the evolving sentiment landscape. This shift towards positivity hints at potential changes in community dynamics, engagement patterns, or perhaps an increasing sense of satisfaction or accomplishment among participants. Overall, the bar chart provides a concise and impactful visualization of sentiment trends, offering valuable insights into the evolving emotional climate within the Kaggle community.
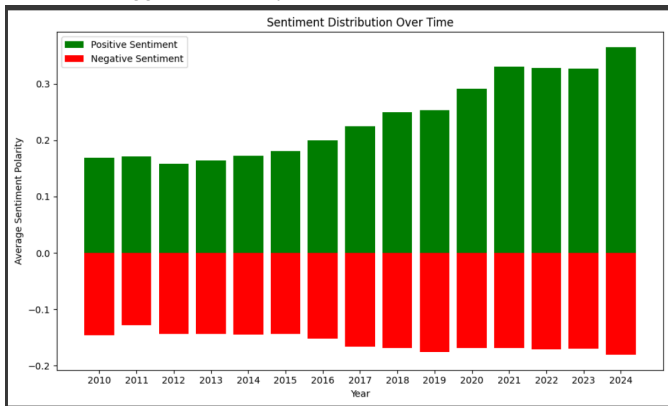


Fig. 2.2. Sentiment Distribution Over Time.

## KEY INSIGHTS

Key insights gleaned from Fig1.1 and Fig 1.2, the analysis and visualization reveal the nuanced dynamics of competition participation across different segments on Kaggle. While the Community segment hosts a significant portion of competitions, it does not necessarily correlate with the highest participant engagement. A thorough examination across all segments is crucial for accurately gauging participant numbers and understanding relative engagement dynamics across competition types. The introduction of the Research and Playground segments in 2012 and 2013, respectively, has expanded the diversity of competition offerings on the platform, reflecting Kaggle's evolving landscape. Fluctuations observed in yearly competition counts across segments such as Featured, Research, and Playground underscore the fluidity and complexity of competition trends over time. This analysis underscores the importance of considering various competition types when assessing community engagement and participation,

highlighting the need for a comprehensive approach to understanding Kaggle's competitive ecosystem.

Insights derived from Fig 2.1 and Fig 2.2, the analysis of trends in sentiment polarity over time offer invaluable perspectives for understanding the underlying dynamics of community sentiment. By delving into sentiment patterns, several key observations emerge. Firstly, examining whether sentiment tends to be predominantly positive or negative over time provides crucial insights into the overall mood and attitudes within the community. Moreover, identifying recurring patterns in sentiment polarity sheds light on events, seasons, or trends influencing community sentiment, facilitating a deeper understanding of community dynamics. Comparing the prevalence and trends of positive and negative sentiments allows for a nuanced analysis of sentiment dynamics, while recognizing significant shifts in sentiment polarity helps pinpoint events or changes impacting community sentiment. Additionally, determining the predominant sentiment polarity throughout the analyzed period and identifying periods of notable fluctuations or shifts in sentiment polarity provide comprehensive insights into community sentiment trends. These insights are instrumental in enabling platforms to gauge community engagement, identify areas of concern, and tailor strategies to better meet the needs and preferences of their audience.

## DISCUSSION

During the validation process, a problem emerged while attempting to remove outliers from the negative sentiment data. This issue stemmed from a ValueError caused by a mismatch between the length of the positions list and the boxplot statistics. The discrepancy arose due to alterations in the data length resulting from the outlier removal procedure, leading to inconsistencies in the plotting. Initially, the design sketch of the plot considered broad topics such as CNN, CV, NLP, ML, and anomaly detection. However, upon closer scrutiny of the tables and thorough examination, it became evident that the topics required more specific and detailed titles to accurately represent the data and insights gleaned from the analysis.

To address the issue encountered during validation, the solution was redesigned by adding the showfliers=False parameter to the plt.boxplot() function call. This parameter directs Matplotlib not to display outliers on the boxplot, effectively circumventing the mismatch between the length of the positions list and the boxplot statistics. By excluding outliers from the visualization, the redesigned solution ensures that the boxplot accurately depicts sentiment variation over time without incorporating outliers that could distort the analysis. Additionally, the redesign includes a breakdown of the number of competitions per topic per year, presenting a more comprehensive visualization with detailed insights. Employing distinct colors for each topic enhances the effectiveness of the visualization, offering an informative introduction to both the competitions and their respective topics.

## CHALLENGES AND OPPORTUNITIES

The discussion of challenges and opportunities in this project provides valuable insights into the complexities and potential for growth within the realm of data science analysis of Kaggle competitions. Time constraints emerged as a notable challenge, limiting the depth of analysis and visualization that could be achieved within the project scope. Additionally, complexities in data handling, such as the issue encountered with removing outliers from negative sentiment data, posed significant challenges and required careful consideration to ensure accurate plotting statistics. The initial design plans for broad topic titles had to be revised to accommodate more specific and detailed titles after closer examination of the dataset, highlighting the iterative nature of data analysis and visualization design. Throughout the process, maintaining a balance between clarity and complexity in visualization design remained a priority to ensure effective communication of insights to stakeholders.

However, amidst these challenges lie numerous opportunities for leveraging data science methodologies to gain deeper insights into Kaggle competitions and user dynamics. Analyzing trends in competition topics enables the prioritization of relevant themes, thereby enhancing competition design and engagement. Understanding community sentiment through sentiment analysis offers valuable insights for developing effective engagement strategies and driving platform improvements to better meet user needs. Moreover, exploring dataset trends provides a window into evolving methodologies and informs participants of emerging data science trends, empowering them to stay ahead of the curve and enhance their competitive edge. These opportunities collectively contribute to a richer understanding of Kaggle competitions and user behaviors, fostering a data-driven approach to decision-making and continuous improvement within the Kaggle ecosystem.

## CONCLUSION

In conclusion, our visualization project on Kaggle competition dynamics and user behavior offers valuable insights into trends, collaboration dynamics, and team performance over the past decade. Through a series of interactive visualizations, we explored various aspects of Kaggle competitions, including competition categories, sentiment analysis of community interactions, dataset trends, and team dynamics.

Key insights from our analysis include the dominance of community-hosted competitions on Kaggle, fluctuations in competition counts across different categories over the years, and trends in sentiment polarity within the Kaggle community. Additionally, we gained insights into team performance and collaboration dynamics, identifying shifts in team sizes and the correlation between user activity levels and team success in competitions.

Throughout the project, we encountered challenges related to data handling and visualization design, which were addressed through iterative refinement and redesign. Despite these challenges, we capitalized on opportunities to leverage data science methodologies to gain deeper insights into Kaggle competitions and user behaviors.

Overall, our visualization project contributes to a richer understanding of Kaggle competition dynamics and provides actionable insights for participants, organizers, and decision-makers. By empowering stakeholders to explore and analyze Kaggle data in a user-friendly and interactive manner, we aim to enhance competitiveness, foster collaboration, and drive innovation within the Kaggle community.

## ACKNOWLEDGMENTS

## REFERENCES

[1] CMV2004 url: http:/ www.cvev.org/cmv2004/index.html
[2] E.R. Tufte. Envisioning Information. Cheshire, CT, Graphics Press. 1990.
[3] Sci2 Team. (2009). Science of Science (Sci2) Tool. Indiana University and SciTech Strategies, http:/ sci2.cns.iu.edu.