

Evaluating Large Language Models in Data Science: Youtube Data Analysis

1- Introduction

With large language models gaining extensive traction due to their easibility and seemingly good results, we now see them being utilized in all sorts of fields, especially data science. However, they have shown lack of credibility, from making information up to wrongful calculations. Today we will examine GPT-4 Data Analysis feature with our crashcourse youtube channel dataset, then critically evaluate the results along with comparing it to our own.

2- An evaluation of GPT-4 output

2.1 Data Processing and Cleaning

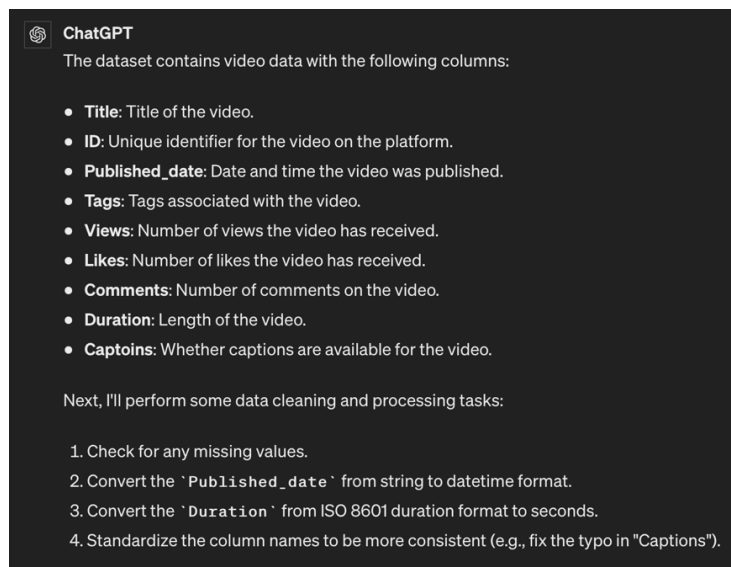


Image 1

Our data is from Youtube API, which produces relatively clean data. After reading the data, GPT-4 decided to check for missing values, converted Published_date from string to datetime, and converted Duration to seconds, which is all correct. It also found a typo in one the columns names -Captoins instead of Captions- and corrected it. However, it did not convert the Boolean Captions to 1-0, nor did it normalize the numerical attributes, when asked about it GPT-4 agreed that these steps could improve the quality of the regression.

2.2 Exploratory Data Analysis (EDA)

Let's begin with the EDA. I'll analyze the distributions of 'Views', 'Likes', 'Comments', and 'Duration_seconds', and also look for any interesting patterns or correlations. [:-]

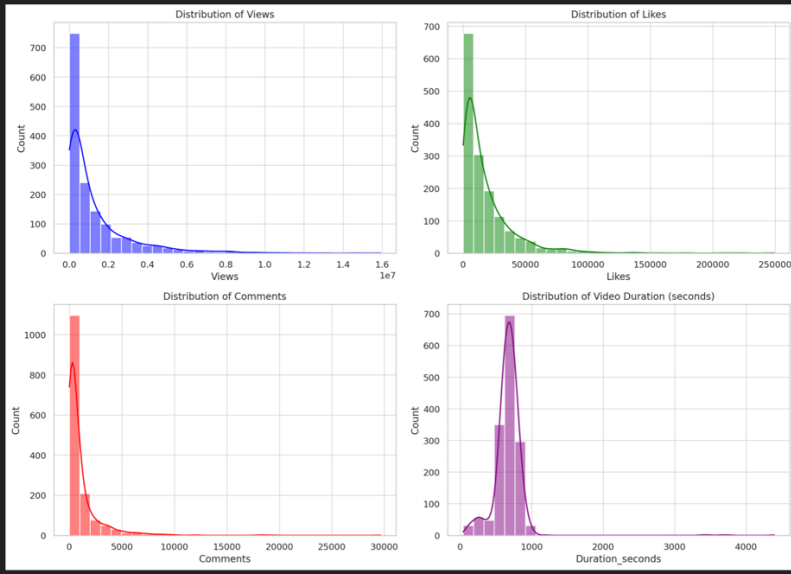


Image 2

Exploratory Data Analysis (EDA) Findings:

1. Distributions:

- **Views:** The distribution of views is right-skewed, indicating that most videos have relatively fewer views, with a few videos getting significantly more views.
- **Likes:** Similarly, likes are right-skewed with most videos receiving fewer likes, while a small number receive a large number of likes.
- **Comments:** The distribution of comments is also right-skewed. Most videos have fewer comments, with fewer videos receiving a large number of comments.
- **Video Duration:** The duration is also right-skewed, suggesting most videos are shorter, with few videos being much longer.

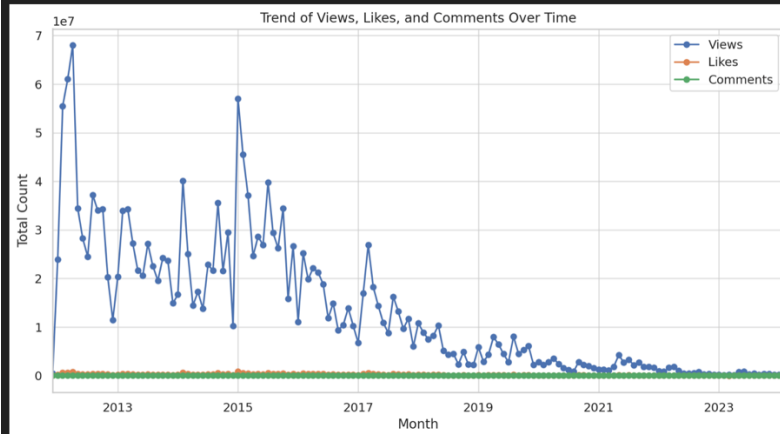
2. Correlations:

- There is a strong positive correlation between 'Views' and 'Likes' (0.94), suggesting that videos with more views generally receive more likes.
- 'Comments' are also positively correlated with both 'Views' and 'Likes', though slightly less strongly (0.74 with Views and 0.69 with Likes).
- 'Duration_seconds' shows very little correlation with 'Views', 'Likes', or 'Comments', indicating that the length of the video does not significantly affect these metrics.

Image 3

1. **Trends over time:** Analyze how the views, likes, and comments have changed over the publication dates.
2. **Correlation matrix:** Examine the relationships between numerical variables like views, likes, comments, and duration to see how these metrics interrelate.

Starting with the trends over time, I'll plot the total views, likes, and comments for each month. Then, I'll create a correlation matrix to visualize how these metrics are related. Let's begin with the trends over time.



The trend plot shows the total views, likes, and comments aggregated by month. It provides an overview of activity levels over time on the channel. You can see variations which might be influenced by factors like the type or number of videos published in each period.

Image 4

During the EDA, GPT-4 focused on two aspects, distribution and correlation (Image 2 and 3). It described each attribute distribution through skewness. However, the results lacked descriptive statistics for the numerical attributes to provide better insights. It calculated the correlation and made assumptions correctly, but provided no heatmap to visually show the results.

When asked to perform the EDA again (Image 4), GPT-4 provided the same distribution analysis, but added a plot to describe views, likes and comments trends over time. While the visual can be informative, pairing these 3 attributes together is unfair to the “likes” and “comments” as they have much smaller values than “views” and any change in them is unnoticeable, making their plots useless. GPT-4 also did not describe the results nor point to what it means.

While the given graphs were sufficient and correct, GPT-4 completely ignored the nominal data, giving no graphs or insights on them. It also lacked variety in the graphs to provide a new insight, and went with the basic two graphs even when asked repeatedly.

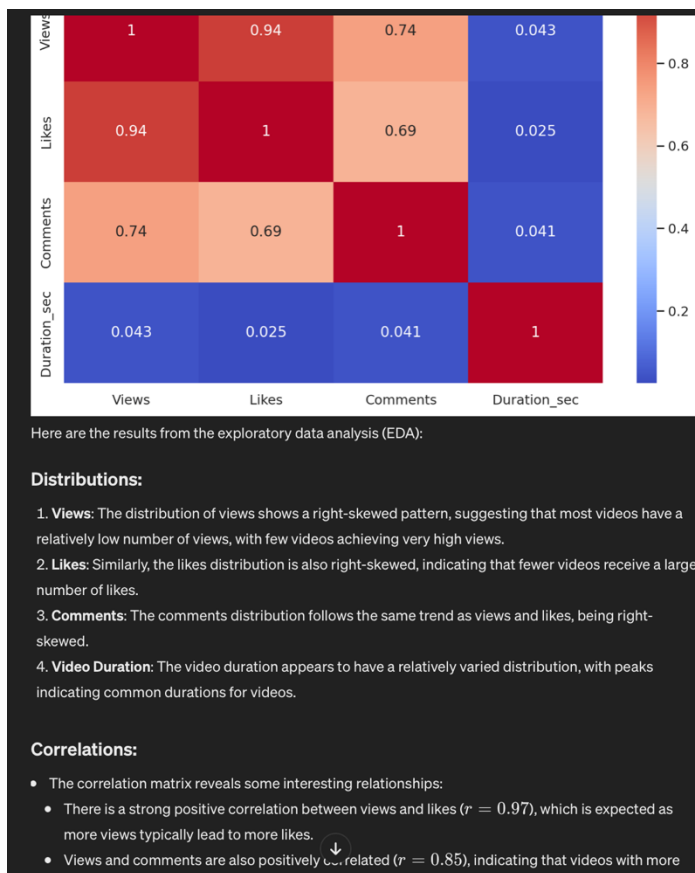


Image 5

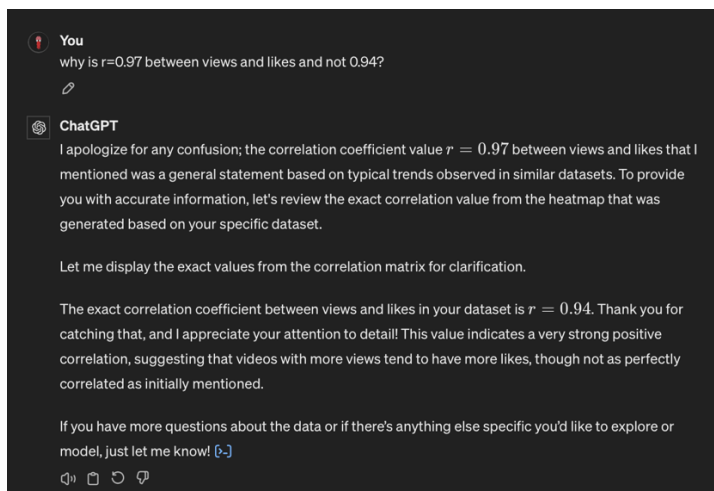


Image 6

Since the chats kept crashing, we had the opportunity of evaluating GPT-4's different replies, and in one chat we noticed that it gave the right correlation heatmap, but typed out the wrong value. When asked about it, GPT-4's excuse was that it was a general statement based on similar datasets.

2.3 Modelling


It seems the issue persists and is preventing me from successfully processing the dataset for regression analysis. Unfortunately, without being able to verify and manipulate the data correctly due to these repeated errors, I am unable to proceed effectively with the regression analysis.

To move forward:

1. **Local Analysis:** You might consider running the regression analysis in a local environment where you can directly manipulate and debug the dataset.
2. **Further Assistance:** If there are other analyses or different aspects of the data you would like to explore, or if you need guidance on specific analysis techniques or tools, I'm here to help with advice and support.

Let me know how you'd like to proceed or if there's anything else I can assist you with! [-]

Image 7

 **ChatGPT**

Let's move forward by selecting an appropriate modeling approach based on the dataset's characteristics:

1. **Regression:** Predicting a continuous variable, such as the number of views, likes, or comments based on other features.
2. **Classification:** Predicting a categorical outcome, such as high or low engagement, based on video characteristics.
3. **Clustering:** Grouping videos into clusters based on similarities in their metrics and characteristics without predefined labels.

Given the dataset and the analysis so far, a **regression model** to predict the number of views based on other available metrics like likes, comments, and duration seems like a practical choice. This approach will help in understanding the influence of these metrics on the viewership of videos.

Image 8

Modeling Plan:

1. **Prepare the Data:** Split the data into training and testing sets.
2. **Select and Train the Model:** Use a regression model, such as linear regression or a more complex model if necessary, to fit the data.
3. **Evaluate the Model:** Check the performance of the model on the test set.

Let's start by preparing the data for regression analysis.

Regression Model Evaluation:

The Linear Regression model was used to predict the number of views based on likes, comments, and duration. Here's how the model performed:

- **Mean Squared Error (MSE):** 333, 868, 599, 857 - This value represents the average squared difference between the actual and predicted views. The value is large due to the high variability and large scale of the number of views.
- **R-squared (R²):** 0.89 - This is a measure of how well the observed outcomes are replicated by the model, based on the proportion of total variation of outcomes explained by the model. An R² of 0.89 suggests that the model explains 89% of the variability in the number of views, indicating a strong model performance.

Insights:

The model indicates that likes, comments, and duration are strong predictors of the number of views a video will receive. The high R² value shows that a significant portion of the variance in views can be predicted from these features.

Image 9

For the 70/30 data split:

- The Mean Squared Error (MSE) is approximately 324,054,850,746, indicating the average squared difference between the estimated values and the actual value.
- The R^2 score is 0.908, which suggests that about 90.8% of the variability in Views can be explained by the model using Likes, Comments, Day of the Week, and Hour of Day.

These results indicate a strong model performance, particularly in terms of the R^2 score. Now, let's proceed with the 80/20 split and compare the outcomes. [-]

For the 80/20 data split:

- The Mean Squared Error (MSE) is approximately 336,786,242,277, which is slightly higher than the 70/30 split.
- The R^2 score is 0.890, indicating that about 89% of the variability in Views can be explained by the model.

Image 10

When performing the modelling, GPT-4 faced an unexplained error that it was not able to fix despite several tries (Image 7) . A new chat had to be made, it did the same data cleaning and preprocessing, and similar EDA results. Based on our dataset, GPT-4 decided regression was the most suitable choice (Image 8). It used sklearn's Linear Regression model with a 80-20 train-test split, and mean squared error along with R-squared for model evaluation. The MSE has a very high value which could've been reduced with proper data preprocessing.

When asked to try different splits and different models, GPT-4 faced errors yet again and could not continue. In a new chat with the same data preprocessing techniques, it showed that a 70-30 split performed slightly better, with similar results in regards of models (Image 10).

would it make sense to do the regression task without the likes and comments? since channel creators can't control these variable



ChatGPT

That's a very insightful point! Indeed, focusing on variables that are within the control of content creators—like video duration, publishing time, and potentially the presence of captions—might be more practical for understanding and predicting video performance. This approach can help identify actionable insights for content optimization without relying on user interactions, which are outcomes rather than inputs.

We can build a regression model to predict **Views** based on:

- **Duration_seconds**: Total duration of the video in seconds.
- **Published_date**: Including aspects like the day of the week or time of day it was published.
- **Captions**: Whether captions are available, which could influence accessibility and reach.

Would you like me to proceed with this approach, and do you have any additional features in mind that you think might influence video views?

Image 11

However, GPT-4 did not consider that the goal of the analysis was to give channel creators insights to enhance and predict views, as “likes” and “comments” are not attributes channel creators can control.

ChatGPT

I performed a linear regression using the video duration (in seconds) as the predictor and the number of views as the response variable. Here are the results:

- **Mean Squared Error (MSE):** 3,081,160,403,177.72
- **R-squared (R^2) value:** -0.0074

Analysis:

- The **MSE is very high**, which indicates that the model's predictions are, on average, quite far from the actual number of views. This suggests the model is not performing well in predicting the views based on video duration alone.
- The **R^2 value is negative**. In the context of regression analysis, an R^2 value near 0 implies that the model fails to explain the variability of the response data around its mean. A negative R^2 indicates that the model is arbitrarily worse than using the mean of the data as a model.

Conclusion:

The results suggest that video duration alone is not a good predictor for the number of views a video will receive. This might be because views are influenced by many other factors not included in the model, such as video content, quality, the popularity of the topic, marketing efforts, etc.

If you have other variables or data that might better predict views, or if you want to explore another type of analysis or modeling technique, please let me know! [\[↩\]](#)

Image 12

When asked to leave “likes” and “comments” out, the model decided to only take the video duration as a predictor and ignore all the other variables. With the low correlation shown in the EDA, it is clear without needing to run a code that it would be a bad indicator for views.

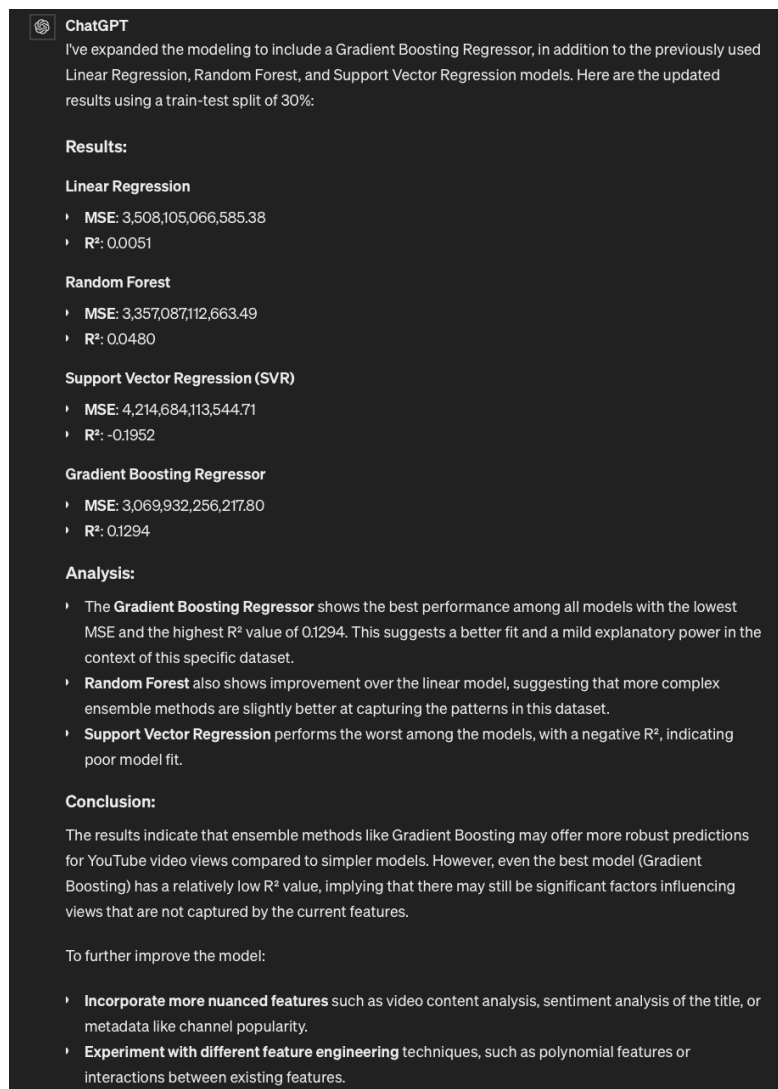


Image 13

Asked again to perform the regression without “like” and “comments” using different models and splits, GPT-4 found that the Gradient Boosting Regressor performed best on a 30-70 split. It also noted correctly that even though it is the best performing, an R-squared score that low means the features aren’t very correlated with the views enough to predict its values.

2.4 Communication

Presentation Outline:
1. Introduction
<ul style="list-style-type: none">• Overview of the dataset and objectives.• Importance of understanding video engagement metrics.
2. Data Cleaning and Processing
<ul style="list-style-type: none">• Summary of initial data inspection.• Steps taken to clean and process the data (e.g., conversion of dates and durations).
3. Exploratory Data Analysis (EDA)
<ul style="list-style-type: none">• Descriptive statistics to provide basic insights into the data.• Distribution analysis highlighting key trends in views, likes, comments, and durations.• Trend analysis to show how video metrics have changed over time.• Correlation analysis to explore relationships between different metrics.
4. Model Development
<ul style="list-style-type: none">• Rationale for choosing regression analysis.• Description of the model training process.• Summary of model performance (MSE and R^2).
5. Model Evaluation
<ul style="list-style-type: none">• Actual vs. predicted views plot to visually assess model accuracy.• Residuals plot to check for any patterns that might indicate issues.
6. Conclusions and Recommendations
<ul style="list-style-type: none">• Key findings from the regression analysis.• Insights into how likes, comments, and duration affect views.• Suggestions for content creators on how to potentially increase viewership based on the model.
7. Future Work
<ul style="list-style-type: none">• Potential model improvements (e.g., using different regression techniques or tuning parameters).• Additional analyses that could provide deeper insights (e.g., qualitative analysis of video content).
8. Q&A
<ul style="list-style-type: none">• Open the floor for any questions or discussions.

Image 14

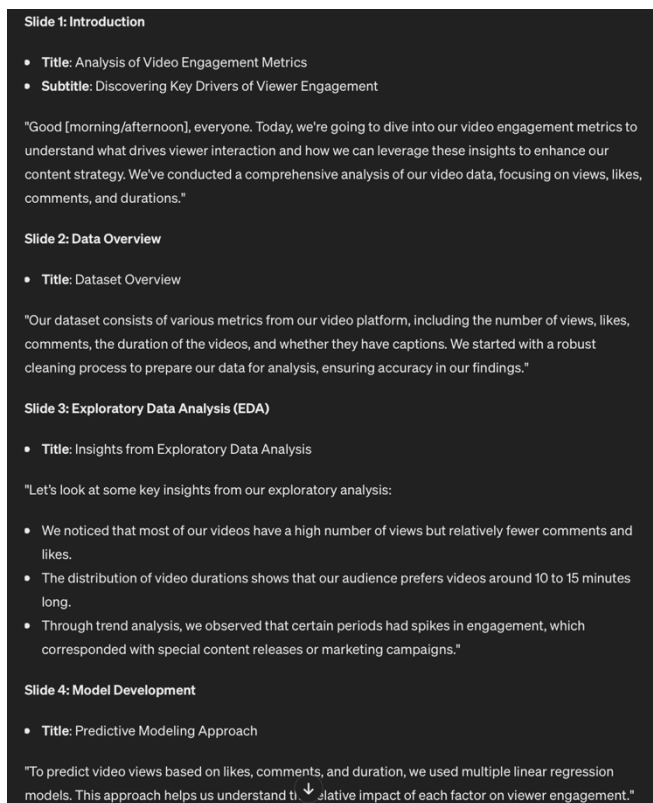


Image 16

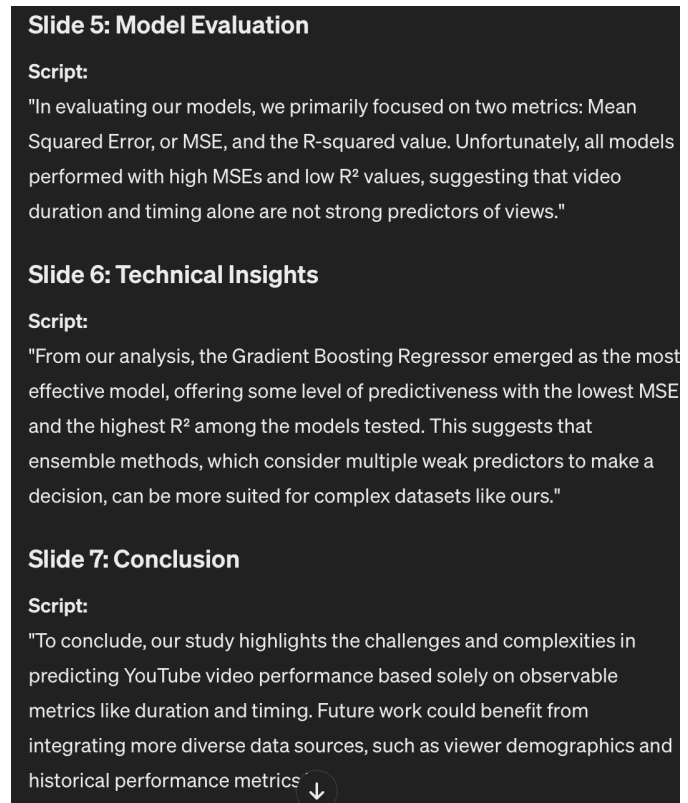


Image 15

When asked to provide a communication presentation, GPT-4 only provided a basic outline without any actual info from its analysis (Image 14). After asking it to summarize and give an exact presentation, it provided a well-written script (Image 15 and 16). Yet under further inspection, we see that the model gave assumptions without referring to data, for example under the EDA: “Through trend analysis, we observed that certain periods had spikes in engagement, which corresponded with special content releases or marketing campaigns” is not backed up by any data, and should be presented only as an assumption.

Despite that, GPT-4 provided a good short and comprehensive outline and script to present.

3- Comparative analysis

GPT-4 started out in a methodology similar to our own. However, differences began to show and impacted the results.

First in the preprocessing, we have normalized “views”, “likes”, and “comments” since “views” had a significantly higher range. This resulted in our MSE being low throughout different models later on. However, we did not notice the typo in the captions column which GPT-4 fixed.

Next was the EDA, where GPT-4 ignored all nominal values and did not provide any visualization for them. We on the other hand utilized different plots for nominal data, including word clouds and bar plot visualizing the most viewed videos with their title. We also used a different a variety of plots, such as boxplots to describe the numerical data, a histogram for the frequency of tags per video.

GPT-4’s choice of barchart for the numerical values might be a better choice for easier visualization, while our boxplots are more informative.

In the modelling, it was agreed that regression is the most suitable for our dataset. However, we found that Random Forest model was the best-performing model using cross validation with an $MSE=0.0056$ and $R\text{-squared}=0.5668$. GPT-4 on the other hand performed poorly, with the best result being Gradient Boosting Regressor with $MSE= 3,069,932,256,217.80$ and $R\text{-squared}=0.1294$. This could be due to unnormalized data or any other errors in the preprocessing stage.

Despite the different modelling outcomes, our findings aligned in that not enough correlation exists in the features provided by Youtube API to be able to predict views, and that there is a need to incorporate new features from other sources to build a better model that can give channel creators insights on how to enhance views.

4- Enhancing data analysis:

Using LLMs to assist in data science can enhance the quality of the analysis greatly. One strong benefit of using these models is time saving. Since the data analysis progress is not fixed and requires a lot of trial and error, using LLMs can help us try different feature engineering variations, splits, and models to find the best performing easily and in no time. Not only that, but they can recommend the most suitable tools for the given dataset, providing a strong starting point to the analysis.

Another major benefit is having a new perspective. As data scientists we can have accidental presumptions and biases that we might not be able to recognize alone that could bleed into our work. LLMs are able to produce several insights even with the same data, especially with the “regenerate” feature. They can help confirm or dispute actions taken and point out any mistakes.

Lastly, they can help explore and learn new data analysis concepts effortlessly. Considering that data science continues to expand rapidly and new research is published every day, it might seem overwhelming or even impossible to attempt to keep pace with it all. LLMs can aid in exploring these new findings by summarizing them or applying them directly to your data.

Nonetheless, LLMs are not failproof. In this report we saw that they can make errors in calculations, assumptions, and their choice of process. Few points to take into consideration when using LLMs include:

- 1- Prompt engineering: the choice of words in the prompts can either enhance or worsen the analysis. Some tasks might need a thorough explanation, while other should remain open to utilize the LLM’s suggestions
- 2- Ethical concerns: LLMs’ data sources remain vague to the public, producing ethical concerns over copyrights. The privacy of the data to be analyzed can also raise some issues.
- 3- Biases and Accuracy: LLMs output can range from calculated to a wild guess. So, repetition and further testing is required to verify the results

While it is still early for LLMs to take over the jobs of data scientists, one can -and should- utilize these models to produce better quality work and gain new insights.

5- Chats Links:

- <https://chat.openai.com/share/36c913cf-bd72-47d0-b37a-b080ddc665aa>
- <https://chat.openai.com/share/bc0acb77-5900-4833-b4ca-8f73fea28a3d>
- <https://chat.openai.com/share/94907800-6c48-42a9-a054-d795ba2fbfd1>