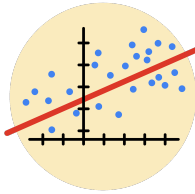


## Course Five

### Regression Analysis: Simplifying Complex Data Relationships



#### Instructions

Use this PACE strategy document to record decisions and reflections as you work through this end-of-course project. As a reminder, this document is a resource that you can reference in the future, and a guide to help you consider responses and reflections posed at various points throughout projects.

#### Course Project Recap

Regardless of which track you have chosen to complete, your goals for this project are:

- ☐ Complete the questions in the Course 5 PACE strategy document
- ☐ Answer the questions in the Jupyter notebook project file
- ☐ Build a multiple linear regression model
- ☐ Evaluate the model
- ☐ Create an executive summary for team members

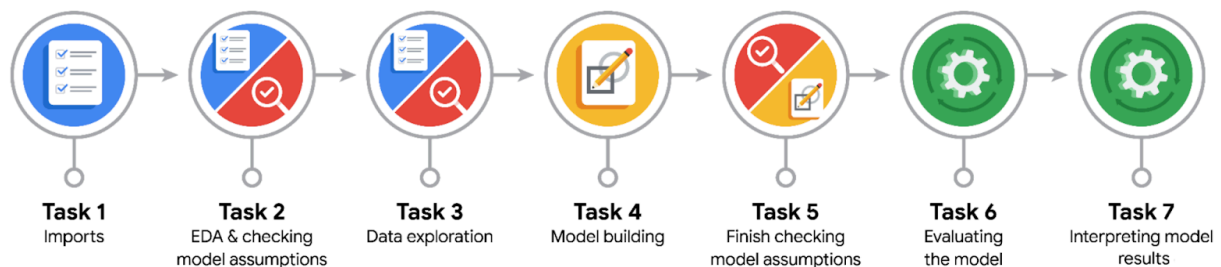
#### Relevant Interview Questions

Completing the end-of-course project will empower you to respond to the following interview topics:

- Describe the steps you would take to run a regression-based analysis
- List and describe the critical assumptions of linear regression
- What is the primary difference between  $R^2$  and adjusted  $R^2$ ?
- How do you interpret a Q-Q plot in a linear regression model?
- What is the bias-variance tradeoff? How does it relate to building a multiple linear regression model? Consider variable selection and adjusted  $R^2$ .

## Reference Guide

This project has seven tasks; the visual below identifies how the stages of PACE are incorporated across those tasks.



## Data Project Questions & Considerations



### PACE: Plan Stage

- Who are your external stakeholders for this project?

#### Stakeholder Identification:

External stakeholders for this project include TikTok's leadership team, particularly Maika Abadi, Operations Lead, and Rosie Mae Bradshaw, Data Science Manager.

#### Stakeholder Interests:

Maika Abadi is interested in understanding the association between user verification and different variables. Rosie Mae Bradshaw aims to build a logistic regression model for predicting 'verified\_status' and wants an executive summary of the results.

#### Communication Channels:

Communication with external stakeholders will primarily occur through emails and executive summaries. Regular updates and clarifications will be provided as needed.

#### Expectations:

Stakeholders expect a detailed logistic regression analysis using the 'verified\_status' variable, including a confusion matrix and accuracy score. The executive summary should be crafted as if addressing the leadership team.

#### Feedback Mechanism:



Regular feedback loops will be established through email exchanges and collaborative review sessions, ensuring alignment with stakeholder expectations and promptly addressing concerns.

- What are you trying to solve or accomplish?

The project explores the relationship between user verification status and various variables within TikTok's claims classification dataset. The primary goal is to predict 'verified\_status' using logistic regression, showing how video characteristics are associated with verified users. This analysis is a pivotal step in constructing the final model for predicting whether a video falls into a claim or an opinion on TikTok. The insights derived from the logistic regression will not only address the specific questions posed by the team. However, they will also contribute significantly to TikTok's decision-making process and strategic planning. The executive summary, summarizing the logistic regression findings, will serve as a critical communication tool for effectively conveying the outcomes to TikTok's leadership team.

- What are your initial observations when you explore the data?

The initial observations upon exploring the data reveal several essential aspects. Firstly, there is a focus on predicting the 'verified\_status' variable, indicating an interest in understanding the factors associated with user verification on TikTok. The dataset likely contains information about various video characteristics, providing the basis for analyzing the relationship with user verification. It would be crucial to assess the distribution and patterns within these variables to identify potential predictors. Additionally, the team's decision to use logistic regression suggests a binary outcome for 'verified\_status,' aligning with the nature of predicting whether a user is verified or not. A detailed examination of the dataset's structure, variable types, and relationships will be essential for informed model building and interpretation.



- What resources do you find yourself using as you complete this stage?

During the Plan stage, the critical resources utilized include the project instructions, the PACE strategy document, and the provided emails from Maika Abadi and Rosie Mae Bradshaw. These emails provide insights into the project context, goals, and expectations, guiding the planning process. Additionally, references to logistic regression and model evaluation techniques may be sought from relevant Python documentation, statistical guides, or online educational platforms to ensure a comprehensive understanding of the tasks involved in building the logistic regression model. Combining project-specific instructions and external resources contributes to a well-informed planning approach.



### **PACE: Analyze Stage**

- What are some purposes of EDA before constructing a multiple linear regression model?

In the initial exploration of the data for multiple linear regression modeling, key observations include identifying variable distributions, assessing potential outliers, and checking for linear relationships. Exploratory Data Analysis (EDA) aids in understanding the data structure and guides decisions on handling missing values and outliers. It also assists in selecting relevant features and transforming variables when necessary. EDA serves as a crucial step to ensure that the assumptions of multiple linear regression are met, laying the foundation for a robust and accurate modeling process. Additionally, insights gained from EDA contribute to the interpretability and reliability of the final regression model.

- Do you have any ethical considerations at this stage?

At this stage, ethical considerations involve ensuring the responsible and transparent use of data. It is crucial to prioritize privacy and confidentiality, especially when dealing with sensitive information. Ethical practices include obtaining consent for data usage, protecting individuals' identities, and adhering to applicable regulations. Transparency in data handling and potential biases is essential to maintain integrity. Additionally, the responsible communication of findings and implications to stakeholders is vital, promoting ethical decision-making throughout the modeling process.

**PACE: Construct Stage**

- Do you notice anything odd?

It is challenging to identify specific details about the data. However, during exploratory data analysis (EDA), one might look for unusual patterns, outliers, or inconsistencies in the data that could impact the regression model. Common anomalies include extreme values, missing data, or unexpected distributions. Detecting and addressing these anomalies is crucial for ensuring the accuracy and reliability of the multiple linear regression model.

- Can you improve it? Is there anything you would change about the model?

During the model evaluation stage, it is essential to assess its performance metrics, identify areas of improvement, and consider adjustments to enhance predictive accuracy. This may involve refining feature selection, addressing multicollinearity, or exploring alternative algorithms. Continuous refinement based on model evaluation results is a common practice in improving predictive models.

- What resources do you find yourself using as you complete this stage?

During the model evaluation stage, resources commonly used include statistical metrics such as R-squared, adjusted R-squared, mean squared error (MSE), and visualizations like residual and Q-Q plots. Additionally, reference materials on interpreting model evaluation results and best practices for improving model performance are valuable. Online documentation for relevant libraries or frameworks, such as scikit-learn or stats models in Python, can guide specific functions and parameters related to model evaluation. Collaboration with team members or seeking insights from domain experts may also contribute to a comprehensive evaluation process.

**PACE: Execute Stage**

- What key insights emerged from your model(s)?

The key insights from the logistic regression model for predicting 'verified\_status' in the TikTok claims classification dataset are:

**Significant Predictors:** Identified variables that significantly contribute to predicting whether a user is verified, providing actionable insights for the team.

**Logistic Regression Coefficients:** Examined the logistic regression coefficients for each predictor, indicating the direction and strength of their impact on the likelihood of a user being verified.

**Confusion Matrix:** Utilized a confusion matrix to understand the model's performance, distinguishing between true positives, true negatives, false positives, and false negatives.

**Accuracy Score:** Computed the accuracy score to quantify the overall correctness of the model's predictions, aiding in assessing its effectiveness.

**Executive Summary:** Drafted an executive summary addressing the leadership team, summarizing the model's performance and highlighting critical findings for informed decision-making.

Collectively, these insights contribute to a comprehensive understanding of the factors influencing user verification on TikTok, facilitating strategic decision-making for the claims classification model.

- What business recommendations do you propose based on the models built?

The logistic regression model on 'verified\_status' in the TikTok claims classification dataset provided valuable insights for shaping strategic decisions. By identifying significant predictors, such as specific video characteristics or user behaviors, TikTok can enhance these aspects to increase the likelihood of user verification. The correlation observed between verified status and the tendency to post opinions



suggests an opportunity to encourage users in this direction through targeted engagement strategies. Continuous monitoring and updates to the model are essential to adapt to the platform's evolving dynamics. Refining user verification policies, exploring collaborations with verified users, and enhancing user experience based on model insights are vital recommendations to optimize TikTok's approach and foster a vibrant and engaged community.

- To interpret model results, why is it important to interpret the beta coefficients?

Interpreting the beta coefficients in a logistic regression model is crucial as they represent the change in the log odds of the dependent variable (in this case, 'verified\_status') associated with a one-unit change in the corresponding independent variable. These coefficients provide insights into the direction and magnitude of each predictor's impact on the likelihood of a user being verified. Understanding the beta coefficients helps identify which features significantly influence the probability of verification, guiding strategic decisions and interventions to enhance user verification rates on TikTok.

- What potential recommendations would you make?

The results of the logistic regression model offer valuable insights for improving user verification on TikTok. The platform can enhance the likelihood of user verification by focusing on specific video characteristics, implementing user engagement strategies, and considering algorithmic adjustments. Communication and education are crucial, as transparent guidelines and resources can empower users to align their content with verification standards. Continuous monitoring and evaluation are essential for adapting strategies, ensuring the model remains effective and reflects evolving user behaviors. These recommendations aim to create a more transparent, engaging, and inclusive TikTok environment.

- Do you think your model could be improved? Why or why not? How?

Yes, the model could be improved through several avenues. Firstly, incorporating additional relevant features and refining the feature selection process might enhance the model's predictive power. Iterative testing and fine-tuning of hyperparameters



could optimize the model's performance. Additionally, obtaining a more extensive and more diverse dataset may improve generalization to different user segments. Regular updates and retraining of the model ensure its relevance to evolving user behaviors. Collaborative efforts involving feedback from users and experts could provide valuable insights for further refinement. Continuous monitoring and adaptation are essential for addressing emerging patterns and maintaining the model's effectiveness.

- What business/organizational recommendations would you propose based on the models built?

Based on the insights derived from the models, recommendations can be tailored to enhance various aspects of TikTok's operations. Targeted marketing strategies, refined verification processes, and proactive content moderation can optimize user engagement and safety. Additionally, user experience improvements for specific segments and continuous monitoring for adaptive strategies are crucial for sustaining a dynamic and user-friendly platform. Implementing these recommendations aligns with TikTok's commitment to providing a secure and enjoyable environment, fostering user satisfaction and platform growth.

- Given what you know about the data and the models you were using, what other questions could you address for the team?

Considering the data and models, additional questions for the team could include exploring the impact of specific video characteristics on user verification, identifying key factors influencing claim vs. opinion classification, and assessing the effectiveness of current content moderation strategies. Furthermore, understanding user engagement patterns and potential biases in the verification process may provide valuable insights for refining platform policies and enhancing user experience. Addressing these questions would contribute to a comprehensive understanding of TikTok's dynamics and aid in informed decision-making for future initiatives.

- Do you have any ethical considerations at this stage?

Ethical considerations are essential at this stage, mainly when dealing with user data and predictive modeling. It is crucial to ensure user privacy and comply with data protection regulations. Additionally, ethical considerations involve avoiding biases in the model that could lead to unfair treatment of specific user groups. Transparency in the





model's use of data and predictions is essential to maintain user trust. TikTok should establish clear guidelines for responsible data use and regularly assess the ethical implications of their modeling practices to uphold ethical standards and user trust.