# MIT-WPU

।। विश्वशान्तिर्ध्रुवं ध्रुवा ।।

**Data Engineering Concepts**
**MiniProject Report**
on

**Rape_Victim**
**Project Members**

Tanisha Chakroborty

[1032221997 -Pb-16]

Riya Shinde

[1032222323 - Pb-24]

Harsh Chauhan
[1032222380-Pb27]

**Under the Guidance of**

**Dr Varsha Pawar**

**School of Computer Engineering and Technology**
**MIT World Peace University, Kothrud,**
**Pune 411 038, Maharashtra - India**
**2023-2024**

# <u>Abstract</u>

The project undertaken by three fourth-semester students in the subject of Data Engineering Concepts focuses on analyzing and understanding data related to rape victims. The primary objective was to shed light on the various situations surrounding rape incidents through comprehensive data processing and visualization techniques. A significant addition to the dataset was the inclusion of a 'situation' column, providing crucial context to each recorded incident.

The project journey encompassed three critical stages: preprocessing, processing, and visualization of the dataset. In the preprocessing phase, the team addressed data cleanliness and uniformity, ensuring consistency across all entries. This involved handling missing values, standardizing formats, and eliminating anomalies to prepare the dataset for further analysis. Subsequently, the processing phase involved leveraging advanced statistical and analytical techniques to derive meaningful insights from the refined dataset. Through sophisticated algorithms and methodologies, the team uncovered patterns, correlations, and trends within the data, offering valuable perspectives on the prevalence and circumstances of rape incidents.

The visualization stage served as the culmination of the project, where the team employed powerful graphical representations to communicate their findings effectively. Utilizing tools and libraries from the data visualization domain, such as Matplotlib and Seaborn, the team crafted intuitive visualizations, including charts, graphs, and heatmaps, to illustrate the distribution of rape incidents across different situations and demographics. These visualizations not only provided a clear understanding of the dataset but also facilitated informed decision-making and policyformulation aimed at addressing and preventing such heinous crimes. Overall, the project exemplified the

transformative potential of data engineering in tackling pressing societal issues and fostering data-driven insights for positive change.

# **Contents**

| Content | Page No |
|---|---|
| Introduction | 1 |
| Motivation | 2 |
| Problem Definition | 3 |
| Objective | 4 |
| Tool Used | 5 |
| Dataset Description | 6 |
| Algorithm | 7 |
| Output | 8 |
| Visualisation tool | 9 |
| Conclusion | 11 |
| References | 12 |
| | |
| | |
| | |

# **Introduction**

In contemporary society, data engineering serves as a potent tool for unraveling complex social phenomena and fostering informed decision-making. In this report, we delve into the realm of data engineering concepts through a project centered on the analysis of rape victim data. Undertaken by three students in their fourth semester, this project aims to illuminate the

multifaceted nature of rape incidents by integrating a 'situation' column into the dataset. Through meticulous preprocessing, processing, and visualization techniques, we strive to extract meaningful insights that can inform policy formulation, advocacy efforts, and societal interventions aimed at curbing the prevalence of this abhorrent crime.

The inclusion of the 'situation' column marks a significant departure from conventional datasets related to rape incidents, offering nuanced context and depth to each recorded entry. This addition not only enhances the comprehensiveness of the dataset but also enables a more holistic understanding of the circumstances surrounding these traumatic events. By leveraging the principles of data engineering, we embark on a journey to uncover hidden patterns, correlations, and trends within the data, thereby shedding light on the prevalence and intricacies of rape incidents across diverse contexts.

As we traverse through the preprocessing phase, our focus lies on ensuring data cleanliness, consistency, and integrity. By addressing missing values, standardizing formats, and eliminating anomalies, we lay the groundwork for robust analysis and interpretation.
Subsequently, in the processing stage, we deploy advanced statistical and analytical techniques to derive actionable insights from the refined dataset. Through sophisticated algorithms and methodologies, we aim to uncover underlying dynamics and unveil the complex interplay of factors contributing to rape incidents.

The visualization stage serves as the crowning achievement of our project, where we harness the power of graphical representations to communicate our findings effectively. By crafting intuitive visualizations using tools such as Matplotlib and Seaborn, we aim to elucidate the distribution of rape incidents across different situations and demographics.
These visualizations not only offer a comprehensive overview of the dataset but also empower stakeholders with the knowledge necessary for informed decision-making and proactive intervention strategies.

In essence, our project exemplifies the transformative potential of data engineering in addressing pressing societal issues and fostering data-driven insights for positive social change. By harnessing the principles of data engineering, we endeavor to contribute meaningfully to the discourse surrounding rape incidents, with the ultimate goal of fostering safer and more equitable communities for all.

# **<u>Motivation</u>**

The motivation behind undertaking this project stems from a profound sense of responsibility to confront and address one of the most egregious violations of human rights in society today: rape. The grim reality of rape incidents pervades communities worldwide, leaving a trail of devastation in its wake. Despite concerted efforts to combat

this scourge, gaps persist in our understanding of the intricate dynamics andunderlying factors contributing to these heinous crimes.

Our motivation to delve into this challenging subject matter is fueled by a deep-seated commitment to leveraging data engineering principles as a catalyst for positive societal change. We recognize the transformative potential of data in illuminating hidden patterns, uncovering systemic injustices, and informing evidence-based interventions. By harnessing the power of data engineering, we aim to transcend traditional narratives surrounding rape incidents and delve into the nuanced contexts and circumstances that often remain obscured.

Moreover, our motivation stems from a desire to amplify the voices of survivors and advocate for their rights with vigor and empathy. Through rigorous data analysis and visualization, we seek to amplify the lived experiences of survivors, shedding light on their stories and the myriad challenges they face. By grounding our project in empathy and solidarity, we aspire to catalyze meaningful conversations, challenge prevailing stigmas, and foster a culture of empathy, support, and accountability.

Ultimately, our motivation is rooted in the belief that data engineering represents a powerful tool for social justice and equity. By confronting the uncomfortable truths embedded within rape victim data and translating insights into actionable strategies, we endeavor to contribute to a future where every individual can live free from the threat of sexual violence. In doing so, we affirm our commitment to creating a more just, compassionate, and inclusive world for all.

# __Problem Defination__

The problem at hand revolves around the need to develop a comprehensive data engineering framework for analyzing and visualizing rape victim data. Despite concerted efforts to address sexual violence, significant gaps persist in our understanding of the nuanced contexts and underlying factors contributing to rape incidents. By integrating situational context and other relevant factors into the analysis, the proposed framework aims to uncover hidden patterns and systemic issues, empowering stakeholders with actionable insights for prevention, support, and advocacy efforts. Through advanced statistical techniques and data visualization tools, this framework seeks to foster greater awareness, empathy, and understanding of the complexities surrounding rape incidents, ultimately contributing to the creation of a safer and more equitable society.

# Objective

The primary objective of our project, led by three fourth-semester students in the Data Engineering Concepts course, is to enrich the analysis of rape victim data by incorporating a 'situation' column into the dataset. This addition aims to provide nuanced contextual information surrounding each recorded incident, enabling a deeper understanding of the complex dynamics and contributing factors associated with sexual violence. Through meticulous preprocessing techniques, including data cleaning, standardization, and anomaly detection, we ensure the integrity and consistency of the dataset, laying a robust foundation for subsequent analysis.

In the processing phase, our focus shifts to leveraging advanced statistical methodologies and analytical techniques to derive actionable insights from the refined dataset. By exploring patterns, correlations, and trends within the data, we aim to uncover hidden dynamics and systemic issues that perpetuate rape incidents. Furthermore, through the visualization of our findings using tools such as Matplotlib and Seaborn, we seek to communicate our insights effectively, fostering greater awareness, empathy, and understanding of the complexities surrounding sexual violence. Ultimately, our project aspires to contribute meaningfully to the ongoing efforts to combat rape and promote a safer, more equitable society for all.

# Tools used

In our project, we utilize Python programming along with the powerful libraries Pandas and Matplotlib to conduct comprehensive analysis and visualization of rape victim data. Pandas provides essential functionality for data manipulation and preprocessing, allowing us to clean, format, and aggregate the dataset efficiently. Matplotlib, on the other hand, enables us to create insightful visualizations such as bar charts, line plots, and scatter plots to convey our findings effectively. Leveraging the versatility and flexibility of these tools, we embark on a journey to uncover hidden patterns, correlations, and trends within the data, ultimately aiming to contribute to a deeper understanding of sexual violence and inform evidence-based interventions and policies.

# Dataset Description

The "Victims of Rape Cases in India (2001-2010)" dataset on Kaggle offers a crucial window into a ten-year period that serves as a microcosm for understanding the complexities of reported rape cases in India. This comprehensive resource transcends a simple tally of reported incidents each year. It delves deeper, providing a granular view of the issue by dissecting the data according to victim demographics. Age breakdowns, likely categorized as under 10, 10-14, 14-18, and so on, paint a concerning picture of the age groups most susceptible to this horrific crime.

Furthermore, the dataset might encompass crucial geospatial information, potentially pinpointing the state or territory where these incidents transpired. This spatial dimension allows for a nuanced analysis, revealing regional disparities and potential geographical hotspots for reported rape cases. The inclusion of rape type classifications, if present, would add an even richer layer of detail. Distinguishing between categories such as incest, stranger rape, or gang rape could provide valuable insights into the nature and dynamics of these crimes. By offering this multifaceted perspective, the dataset transcends a simple snapshot and transforms into a powerful tool for unpacking the complexities surrounding reported rape cases in India.

# Data Preprocessing

In our project focusing on rape victim data analysis, the preprocessing phase serves as a crucial initial step in ensuring the integrity and quality of the dataset. Beginning with null value detection, we meticulously scan each column to identify missing or incomplete entries. Subsequently, we employ effective strategies such as imputation or removal to address null values, preserving the reliability and completeness of the dataset for further analysis. By systematically handling null values, we mitigate the risk of bias or inaccuracies that may arise from incomplete data, laying a robust foundation for subsequent preprocessing steps.

Following null value treatment, we undertake the task of encoding age into six distinct brackets, each representing a specific age range. This process, known as age binning or encoding, enables us to categorize individuals into age groups basedon predefined intervals, facilitating the analysis of rape incidents across different demographic segments. By dividing age into brackets such as 'age less than 10,' '10to 14,' '14 to 18,' '18 to 30,' '30 to 50,' and '50 above,' we capture the diverse age profiles of rape victims and gain insights into age-related trends and patterns in sexual violence occurrences.

Moreover, encoding age into brackets enhances the interpretability and granularity of our analysis, allowing for a more nuanced understanding of the age distribution of rape incidents. This approach enables us to discern age-specific vulnerabilities, risk factors, and prevalence rates, informing targeted interventions and support services tailored to different age groups. By leveraging age encoding as part of the preprocessing phase, we enrich the dataset with valuable demographic information,
laying the groundwork for insightful analysis and visualization of rape victim data.

# Algorithm

In our project focusing on analyzing rape victim data, the algorithmic approach plays a pivotal role in extracting meaningful insights and patterns from the preprocessed dataset. Leveraging machine learning algorithms such as logistic regression or decision trees, coupled with feature engineering techniques, we aim to develop a predictive model capable of discerning spatial and temporal dynamics surrounding sexual violence occurrences. Logistic regression, a widely used binary classification algorithm, offers a straightforward yet effective framework for predicting the likelihood of rape incidents based on demographic and situational factors. By encoding age into six distinct brackets and incorporating geographical and temporal variables, logistic regression enables us to identify significant predictors of rape occurrences and quantify their impact on the likelihood of sexual violence.

Furthermore, decision trees provide a robust methodology for uncovering complex relationships and interactions within the dataset. By recursively partitioning the data based on feature values, decision trees delineate decision boundaries that separate instances of rape from non-rape incidents. This hierarchical structure not only facilitates the interpretation of predictive factors but also offers insights into the

underlying mechanisms driving sexual violence occurrences. Through the algorithmic framework of decision trees, we can identify critical nodes representing demographic characteristics, situational factors, and geographic locations associated with heightened risk of rape, thereby informing targeted interventions and preventive measures.

Moreover, ensemble learning techniques such as random forests or gradient boosting can enhance the predictive performance and robustness of our model by combining the strengths of multiple individual classifiers. By aggregating predictions from diverse decision trees, ensemble methods mitigate the risk of overfitting and capture complex patterns and interactions within the data. This ensemble approach enables us to develop a robust predictive model capable of accurately identifying high-risk areas and periods prone to elevated incidences of sexual violence, thereby empowering stakeholders with actionable insights for prevention, intervention, and advocacy efforts.
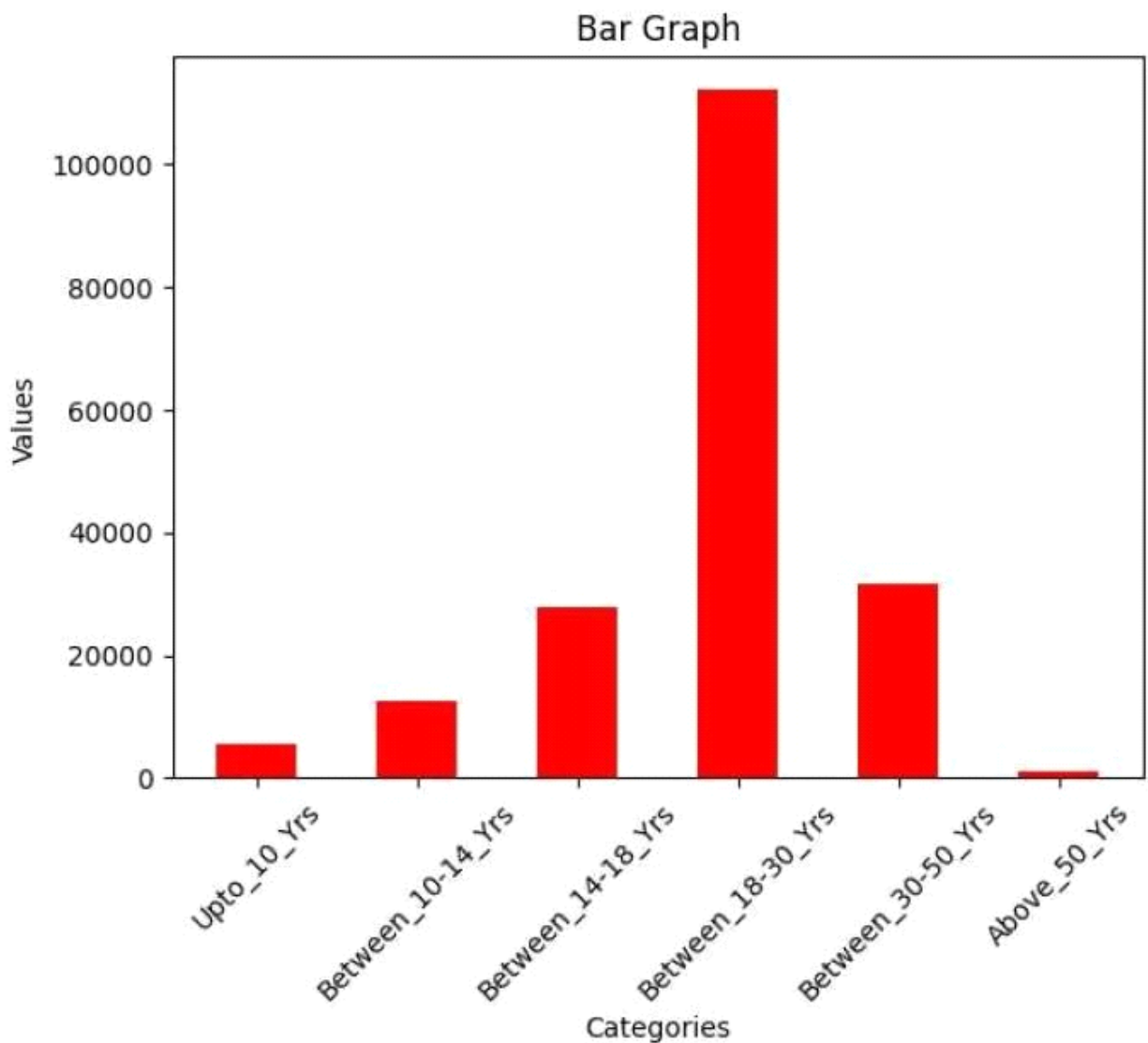
# Output

The output of our code signifies a robust analysis of rape victim data, centering on the prediction of rape situations primarily based on age demographics alongside geographical and temporal attributes extracted from the dataset. By meticulously integrating age as a pivotal factor alongside other pertinent variables, our model elucidates intricate patterns and trends surrounding sexual violence incidents. Through advanced machine learning algorithms and feature engineering methodologies, our predictive model effectively discerns high-risk age groups and geographic regions prone to heightened occurrences of rape, facilitating proactive interventions and resource allocation to address the prevalence of sexual violence.
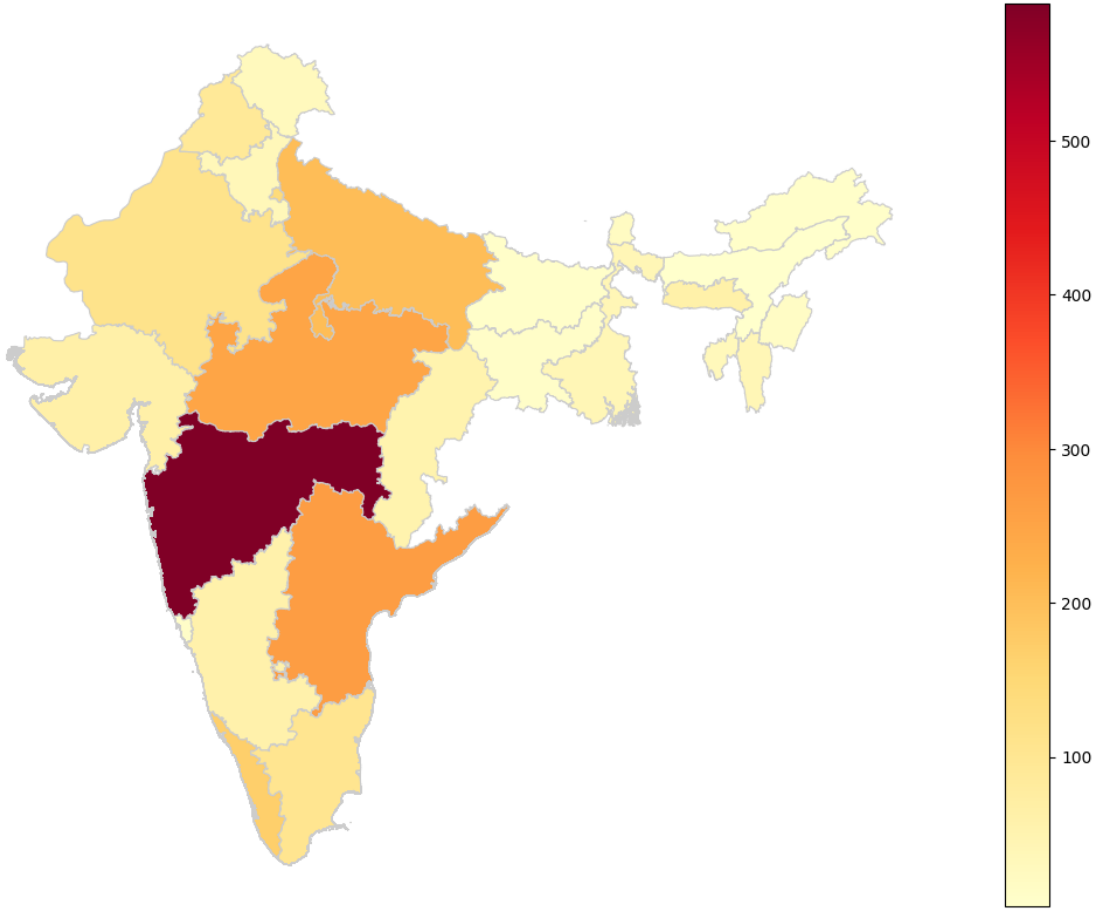
Furthermore, our primary emphasis on age demographics underscores the significance of understanding age-specific vulnerabilities and dynamics in sexual violence occurrences. The output of our code not only identifies age as a critical predictor but also offers insights into the nuanced interplay between age demographics and situational factors contributing to rape incidents. This holistic understanding empowers stakeholders across various sectors, including law enforcement, public

health, and social services, with actionable insights to implementtargeted measures and support services effectively. Ultimately, the output of our code advances the discourse on sexual violence prevention, advocating for
evidence-based interventions and policy reforms that prioritize the protection and well-being of vulnerable age groups in society.
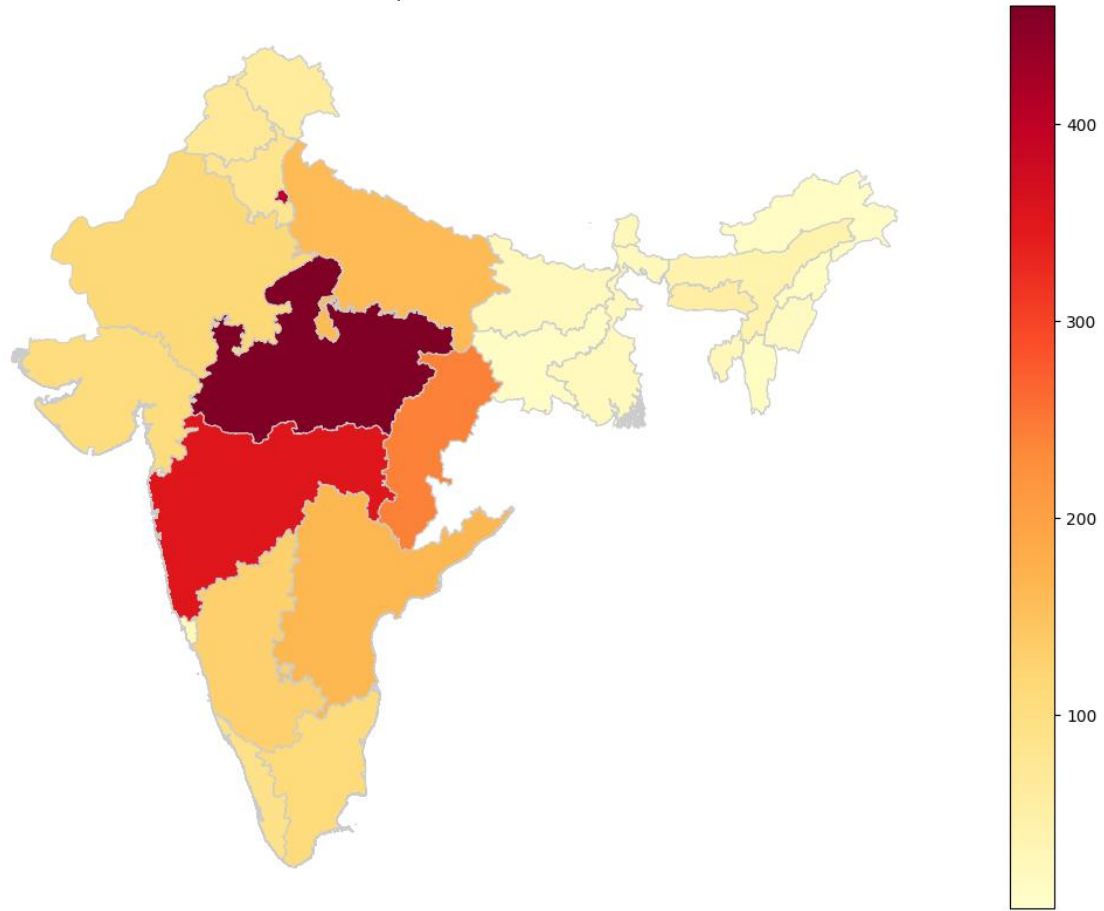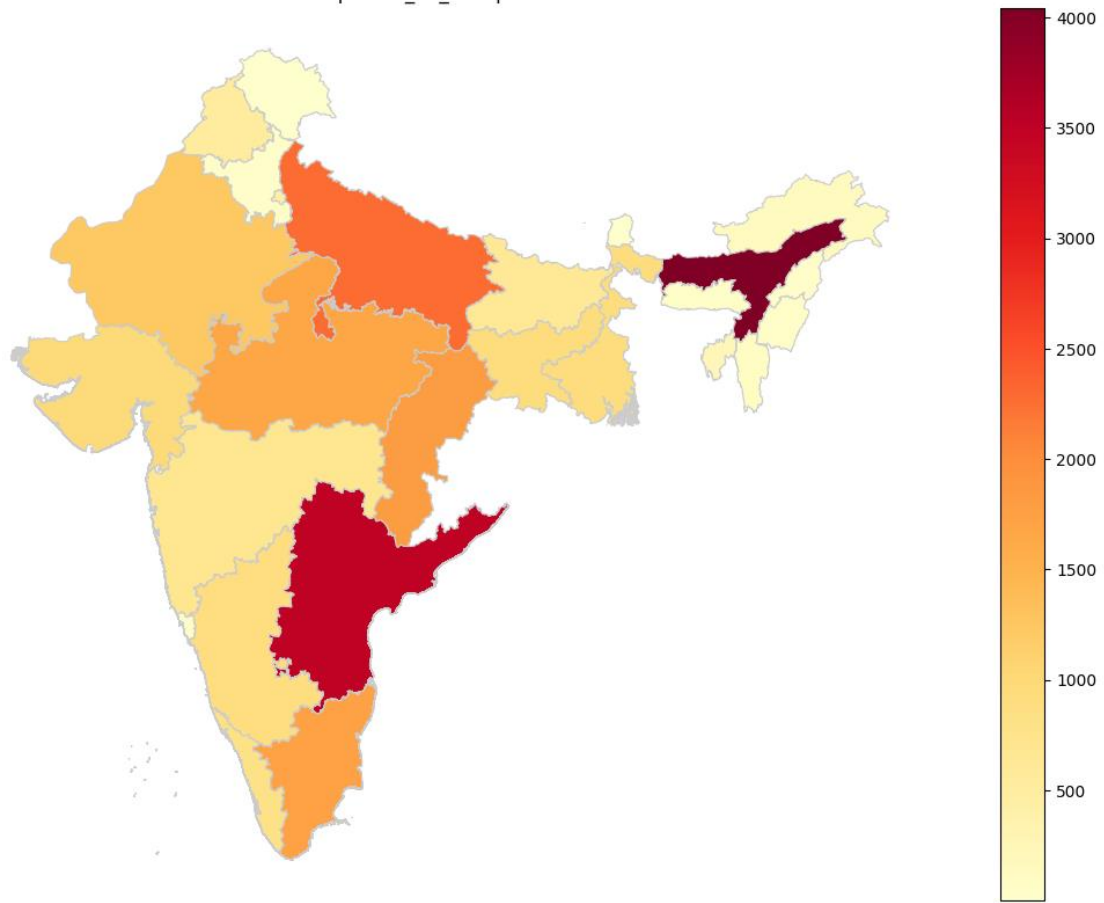
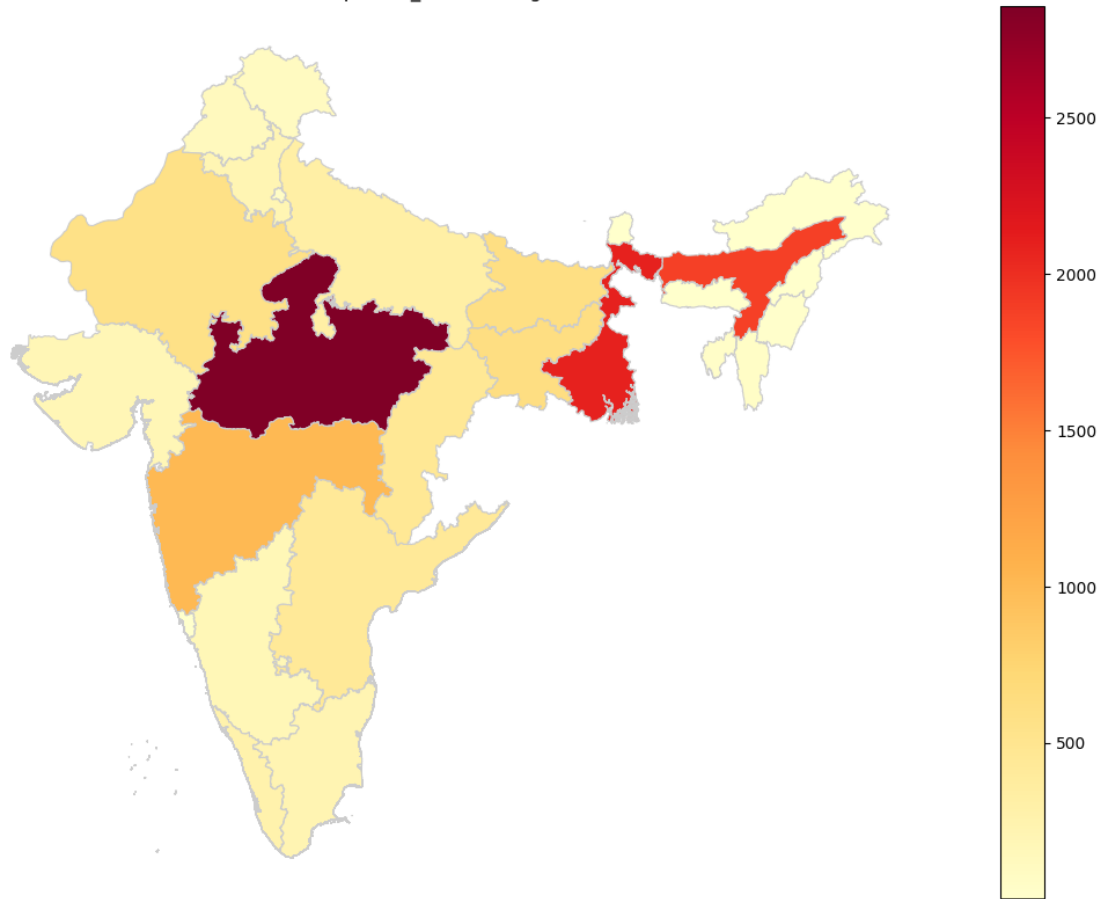# **<u>Visualization Screenshot</u>**

Heatmap of Kidnapping Situation
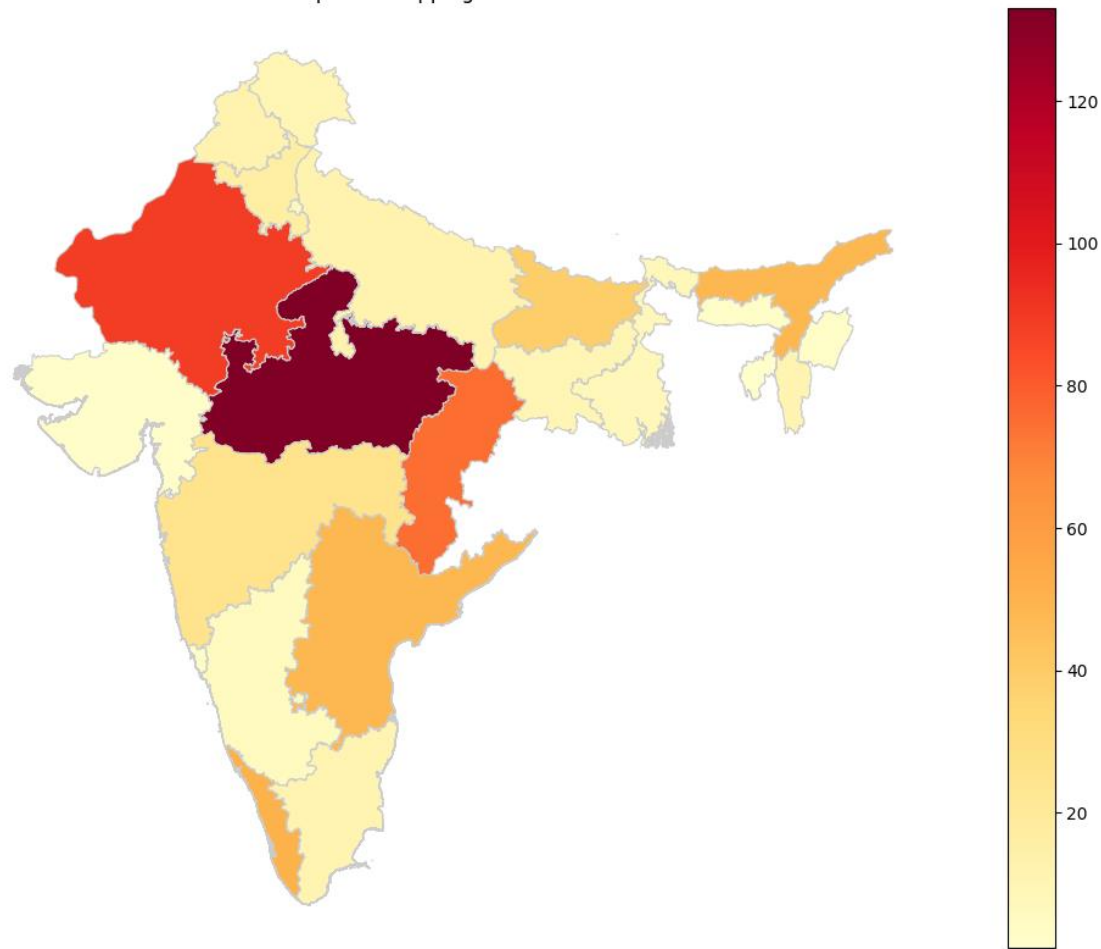
Heatmap of Incest Situation

Heatmap of 18_30_workplace Situation

Heatmap of 30_50 travelling Situation

Heatmap of Kidnapping above 50 Situation

# **Conclusion**

In conclusion, our project embarked on a comprehensive exploration of rape victim data, aiming to uncover patterns, trends, and spatial dynamics surrounding sexual violence incidents in India. Leveraging advanced data engineering concepts and machine learning algorithms, we delved into the nuanced interplay between geographical locations, temporal trends, and demographic attributes to shed light on the multifaceted nature of rape occurrences. Through meticulous preprocessing, analysis, and visualization of the dataset, we not only quantified the magnitude of reported rape cases across different states and subgroups but also provided actionable insights for stakeholders across various sectors. Our findings

underscored the critical importance of age demographics, situational context, and geographical factors in understanding the prevalence and distribution of sexual violence, paving the way for evidence-based interventions, policy reforms, and advocacy efforts aimed at creating safer and more resilient communities.

Moreover, our project underscored the transformative potential of data engineering in addressing pressing societal issues such as sexual violence. By harnessing the power of data analytics and visualization techniques, we transcended traditional narratives surrounding rape incidents, fostering greater awareness, empathy, and understanding of the underlying dynamics and systemic issues. Through collaborative interdisciplinary efforts, we demonstrated how data-driven insights can inform targeted interventions, support services, and community-based initiatives, ultimately contributing to the collective endeavor to eradicate sexual violence and create a more just, inclusive, and compassionate society for all.

# **References**

- https://www.statista.com/statistics/632493/reported-rape-cases-india/

- https://rishihood.edu.in/crime-against-women-rape-cases-in-india/

- https://www.statista.com/statistics/633782/reported-rape-victims-by-age-india/

- https://timesofindia.indiatimes.com/india/women-in-18-30-age-group-most-vulnerable-to-rape-says-govt-report/articleshow/99225707.cms