# GISELLE'S DATA ADVENTURE

## DATA ANALYSIS AND MACHINE LEARNING

**By Giselle Halim**

# HELLO! I'M
# GISELLE.

Fueled by a passion for data and machine learning, I'm dedicated to harnessing the power of analytics to drive actionable insights. My journey has been enriched by experiences at Kalbis University and the prestigious Bangkit Academy, where I've honed my skills in data analysis, visualization, and machine learning. I've further solidified my expertise by earning the Google Data Analytics Certification.

Equipped with a solid foundation in data analysis tools like Excel, SQL, Power BI, and Tableau, I've refined my skills through practical projects. My technical expertise extends to Python, Scikit-learn, and TensorFlow, where I've explored the potential of machine learning algorithms. I'm eager to contribute my expertise to a dynamic organization as a data analyst or data scientist, where I can leverage data-driven insights to innovate and drive positive outcomes.

# EDUCATION



## Kalbis University    (2020 - 2024)

*Information Systems – Big Data Analytics*

**GPA: 3.94 / 4.00**

**Achievements:**
- National Finalist of ASEAN Data Science Explorers 2024
- 2nd place - Internal competition for Information System Analysis and Design course

**Organization:**
Kalbis University Information Systems Student Association (HIMSI GALAKSI)

# EXPERIENCES



## Bangkit Academy 2023
## Machine Learning Student

**Graduated with a score of 93.5/100.** For the capstone team project, **conceptualized and developed an Android-based educational application for traditional fabric motif identification using Machine Learning** techniques, reaching 93% accuracy in recognizing diverse motifs.

## Orbit Future Academy
## AI Mastery Student

**Graduated with a score of 88/100.** For the capstone team project, **spearheaded the development of an AI model for a web application designed to detect chili plant diseases**, successfully training the model on 500 images and achieving a detection accuracy rate of 82%.

# EXPERIENCES

## Data Scientist Virtual Intern
Nov - Dec 2023

### Home Credit Indonesia x Rakamin Academy

Developed a credit risk prediction model using company-provided loan data using Random Forest and reached 90% accuracy.

## Data Scientist Virtual Intern
Sep - Oct 2023

### ID/X Partners x Rakamin Academy

Worked on building a credit risk prediction model for a lending company using loan data using Random Forest and reached 90% accuracy.

## Data Scientist Virtual Intern
Aug - Sep 2023

### Kalbe Nutritionals x Rakamin Academy

Enhanced business strategies by creating a Tableau dashboard and developing predictive models for customer segmentation using K-Means Clustering.

## Big Data Analytics Virtual Intern
Jul - Aug 2023

### Kimia Farma x Rakamin Academy

Analyzed data using SQL to generate insights and visualized findings using Google Data Studio. Created a comprehensive dashboard to track medicine sales from raw data.

# TECHNICAL SKILLS



- **Languages:** Python, SQL
- **Libraries:** Pandas, Numpy, Scikit-Learn, TensorFlow, Keras, Matplotlib, Seaborn
- **Web Technologies:** HTML, Flask, Streamlit
- **Tools:** Power BI, Tableau, Looker, SAP Analytics Cloud, MySQL, Microsoft Excel
- **Data Analysis:** Data Cleaning, Exploratory Data Analysis (EDA), Predictive Analysis, Cluster Analysis, Sentiment Analysis
- **Machine Learning:** Predictive Modeling, Image Classification, NLP, Recommender System

# CERTIFICATES



Machine Learning Specialization



TensorFlow Developer by DeepLearning.AI



TensorFlow: Data and Deployment



TensorFlow: Advanced Techniques



Google Data Analytics Professional Certificate



Data Analysis with Python

# CERTIFICATES



Machine Learning Implementation



Machine Learning Operations (MLOps)



Accenture Data Analytics Job Simulation



PwC Switzerland Power BI Job Simulation



ASEAN DSE 2024 Enablement Session



ASEAN DSE 2024 National Finalist

# PROJECTS

Showcasing past works related to **data analysis, data science, and machine learning**

# DIABETES IN ASEAN ANALYSIS



A comprehensive analysis of diabetes cases in ASEAN was conducted using SAP Analytics Cloud. This project focused on understanding the rising prevalence of diabetes, its associated complications, mortality rates, risk factors, and the significant economic burden it imposes on the region. By aligning with Sustainable Development Goal 3 (Good Health and Wellbeing), the analysis sought to identify effective solutions. In addition to sector-specific recommendations, a gamified app was conceptualized to promote reduced sugar consumption and contribute to mitigating the diabetes epidemic. This project was selected to represent our team in the national finals of the ASEAN Data Science Explorers 2024 Competition, underscoring its significance and potential impact.

**FULL SLIDE**

# MEDICINE SALES DASHBOARD



Leveraged SQL to cleanse and structure data for a comprehensive dashboard analysis of medicine sales at Kimia Farma. Developed a dashboard with Looker to visualize sales trends, total revenue, and product-level performance over a two-week period. The dashboard facilitates quick analysis and data-driven decision-making for optimizing inventory and sales strategies.

**FULL PROJECT HERE**

**LOOKER DASHBOARD**

# CALL CENTER DASHBOARD



In today's saturated telecom market, where providers bombard customers with claims of "better price" and "best service," a clear understanding of customer needs is crucial. This Power BI dashboard empowers a major telecom company to cut through the noise. By tracking KPIs like overall customer satisfaction, call answer rates, and call duration, the dashboard provides actionable insights to improve customer experience and optimize call center operations.

**FULL PROJECT HERE**

**FULL SLIDE**

# SOCIAL MEDIA CONTENT ANALYSIS



This data analysis project, conducted for SocialBuzz, a leading social media and content creation firm, aimed to optimize their content strategy using Power BI. By analyzing vast amounts of social media data, the project delivered actionable insights into content performance, audience engagement, and trends. The analysis focused on identifying the top 5 content categories driving the most engagement, providing recommendations for content optimization, and uncovering opportunities for further growth.
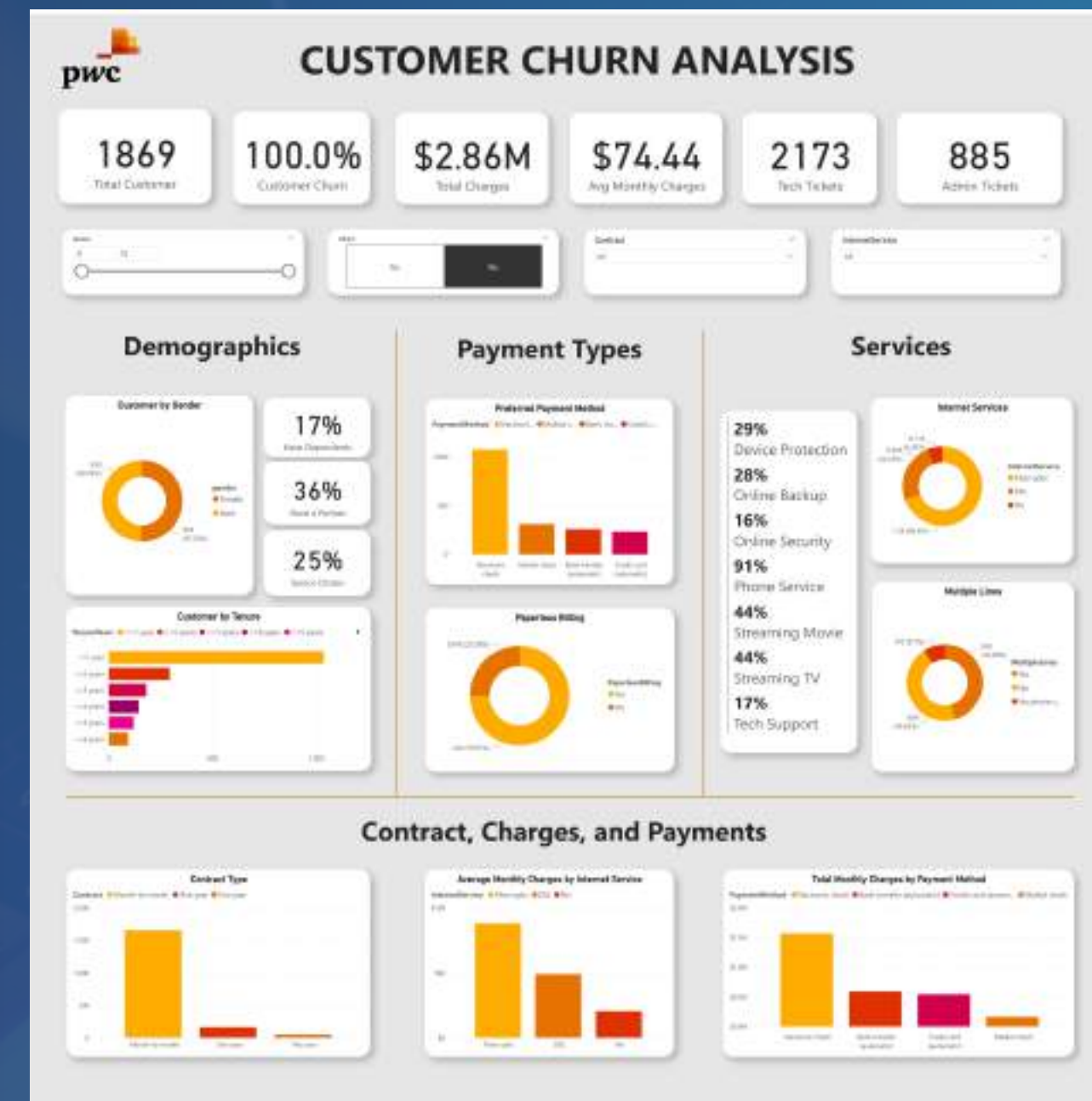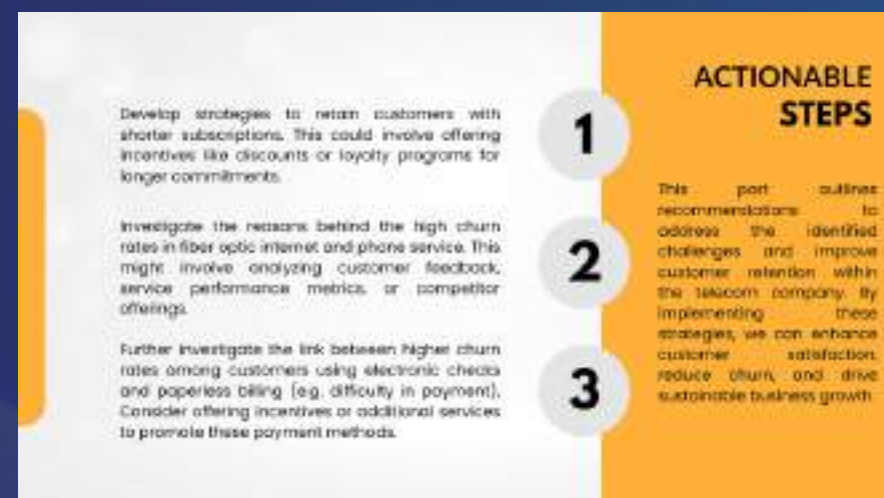
**FULL PROJECT HERE**

**FULL SLIDE**

# CUSTOMER CHURN DASHBOARD

This data analysis project focused on understanding customer churn within a telecom company. By leveraging Power BI, a comprehensive dashboard was developed to visualize customer demographics, service usage patterns, and other relevant factors contributing to churn. While the overall churn rate of 26.5% fell within industry standards, it highlighted the need for proactive retention strategies. The dashboard provided valuable insights to identify at-risk customers and implement targeted interventions to mitigate churn and enhance customer satisfaction.



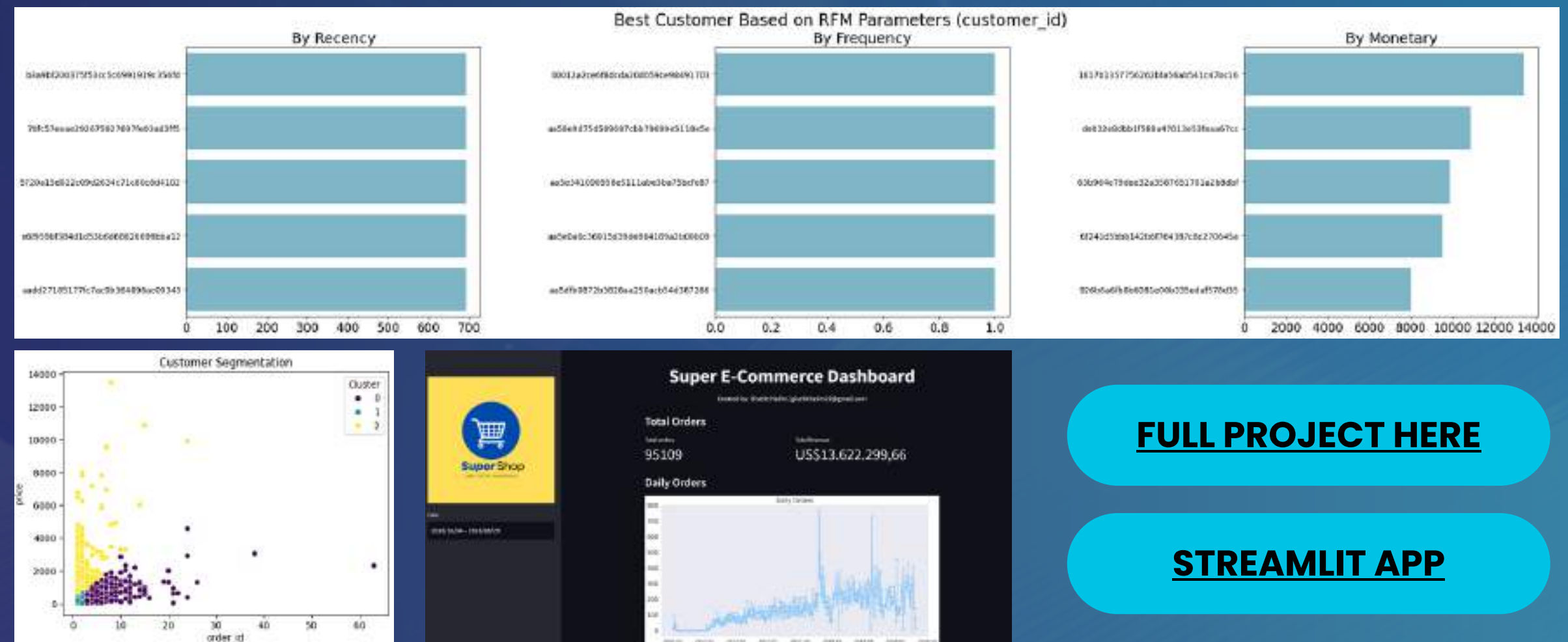**FULL PROJECT HERE**

**FULL SLIDE**

# E-COMMERCE SALES ANALYSIS

Conducted data analysis for an e-commerce platform using Python to address key business questions. The analysis included calculating total sales and profit, evaluating monthly sales trends, identifying preferred payment methods, and performing RFM (Recency, Frequency, Monetary) analysis. K-Means Clustering was implemented to segment customers, enabling targeted marketing strategies. This provided insights into sales performance, customer behavior, and preferences, aiding in strategic decision-making. A web dashboard built with Streamlit provided a simple summary of the transaction data.
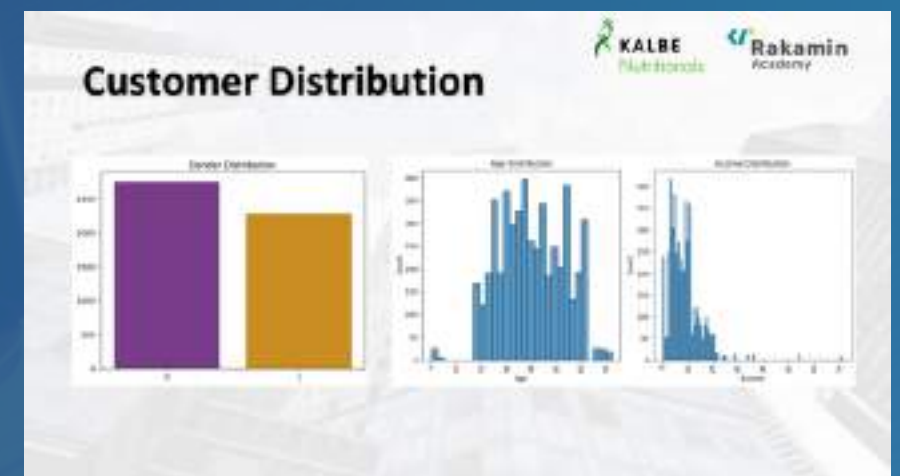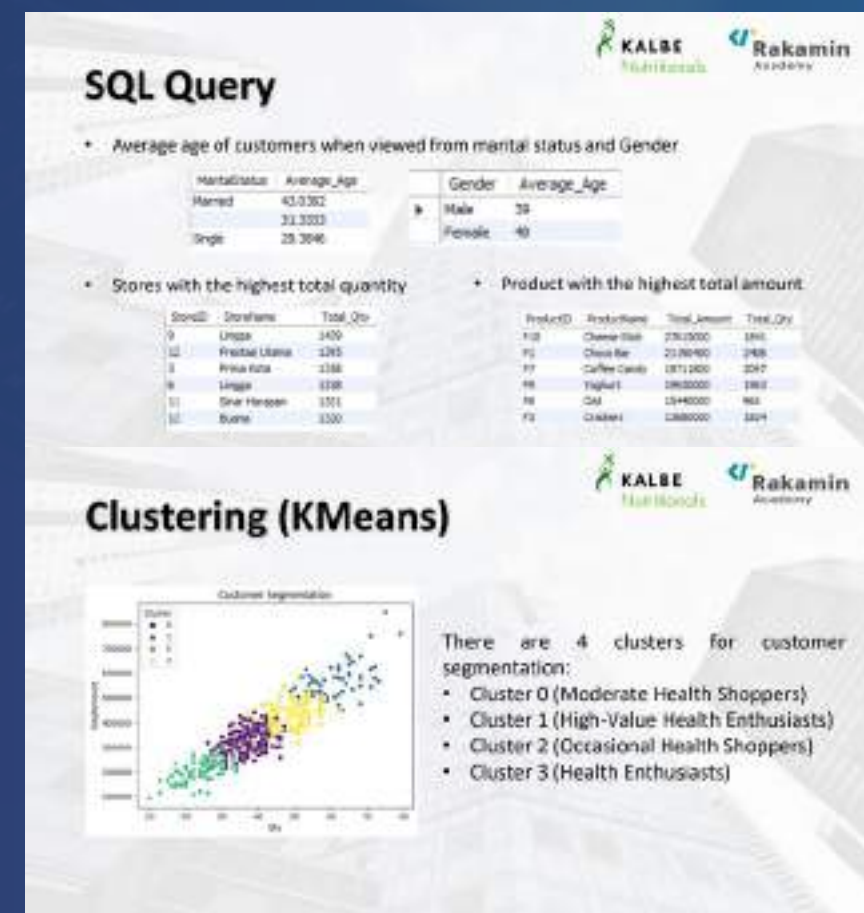


**FULL PROJECT HERE**

**STREAMLIT APP**

# HEALTH PRODUCT SALES ANALYSIS

Leveraged SQL to analyze customer and store data, gaining valuable insights for business improvement. Python was employed for exploratory data analysis, while a Tableau dashboard provided a comprehensive overview of health product sales performance, including metrics like total sales, revenue, and sales trends. The dashboard featured visualizations of sales patterns, product popularity, and regional distribution, enabling swift identification of top-performing areas and products. K-Means Clustering was implemented to segment customers, enabling targeted marketing strategies.



**FULL PROJECT HERE**

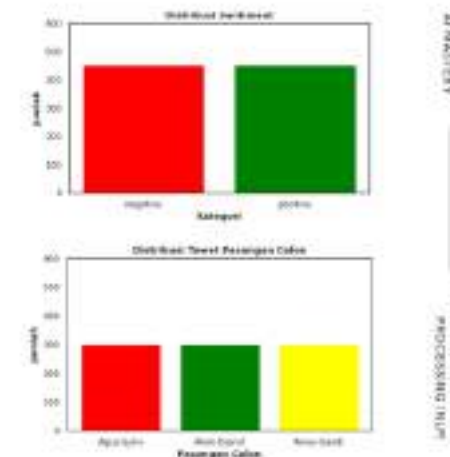**TABLEAU DASHBOARD**

**FULL SLIDE**

# 2017 JAKARTA LOCAL ELECTION SENTIMENT ANALYSIS

This project employed a Bernoulli Naive Bayes classifier to conduct sentiment analysis on tweets related to the 2017 DKI Jakarta Local Leader Election. Data preprocessing like case folding, stemming, and removing stopwords were done. Feature extraction techniques were implemented to optimize model performance. The model effectively analyzed public sentiment towards the three candidates, providing valuable insights into voter preferences during the election campaign. With an accuracy rate of 88%, the model demonstrated its ability to accurately gauge public opinion based on social media data.



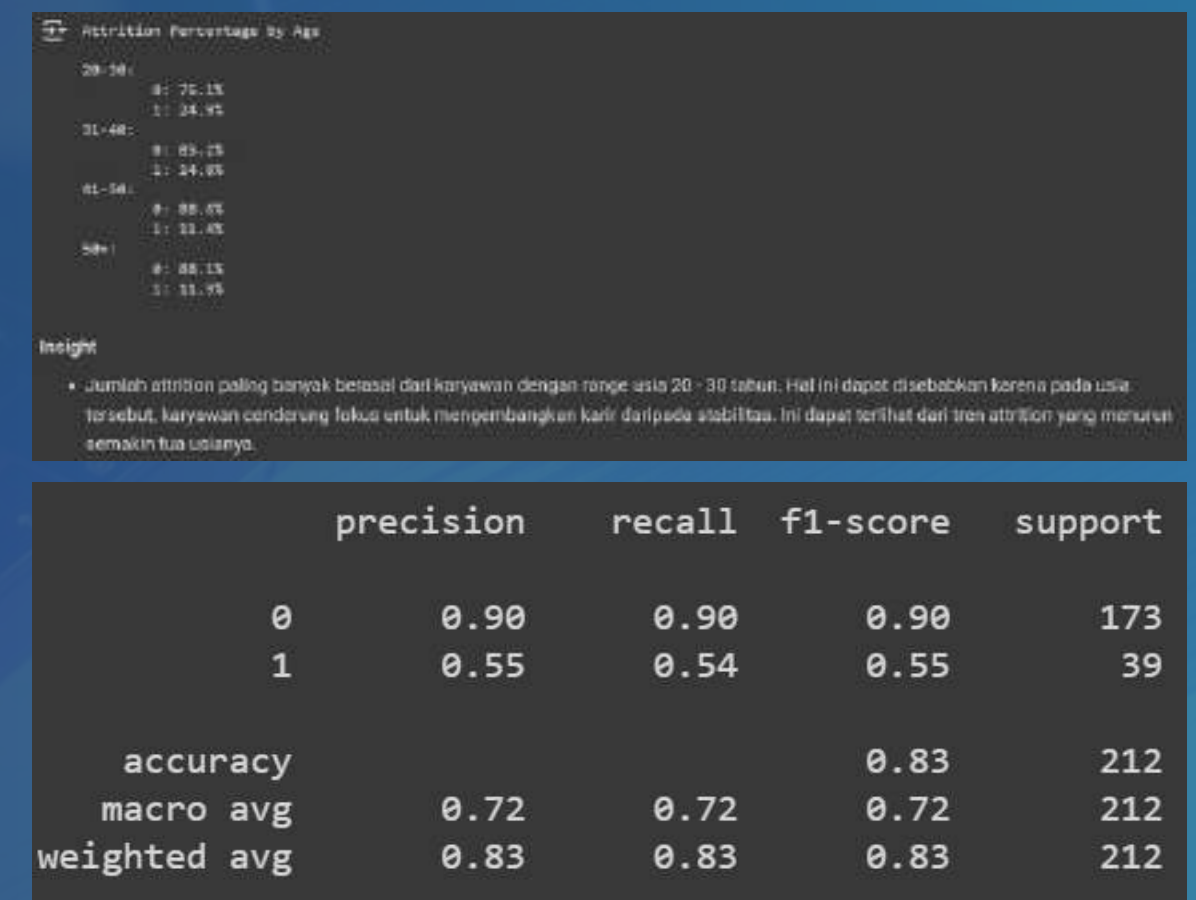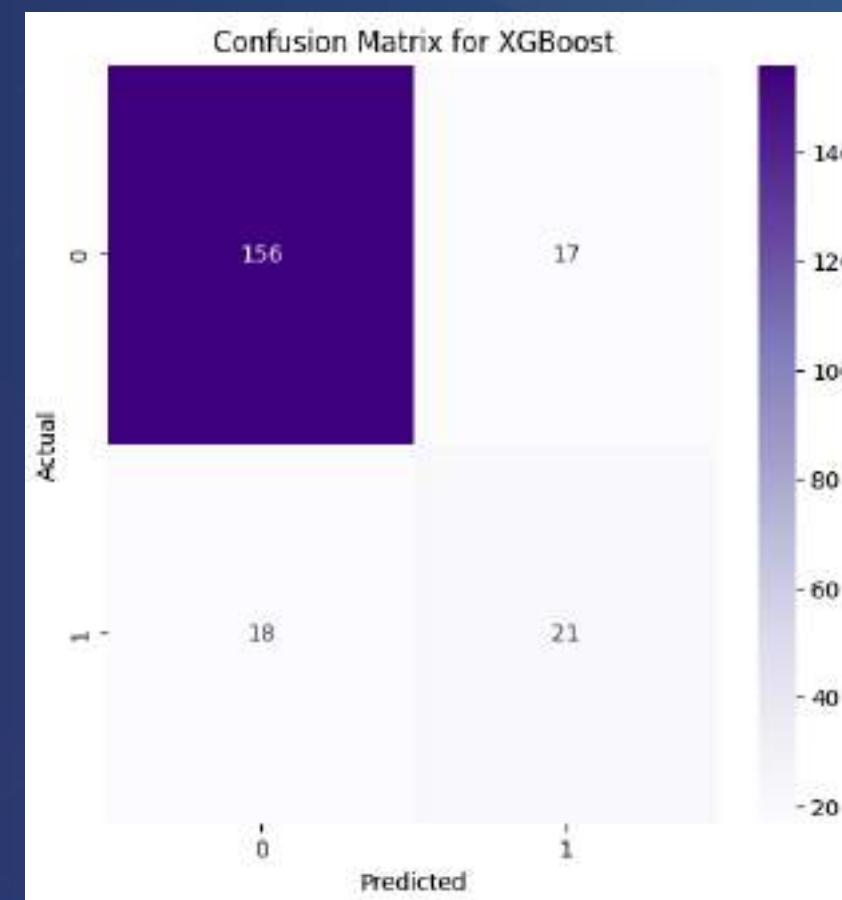**FULL PROJECT HERE**

**FULL SLIDE**

# EMPLOYEE ATTRITION ANALYSIS

Leveraged Python and machine learning to analyze over 1,000 employee data and predict attrition rates by examining factors like job role, department, business travel, and satisfaction. A predictive model built with XGBoost achieved 83% accuracy in predicting employee departures. PowerBI dashboards visualized the findings to inform data-driven retention strategies.
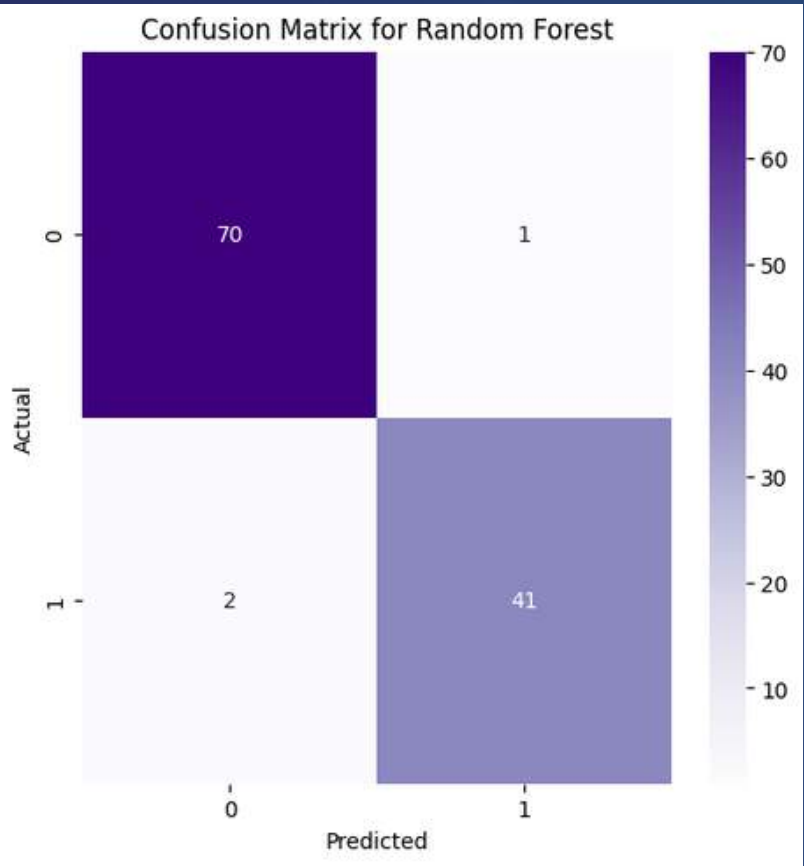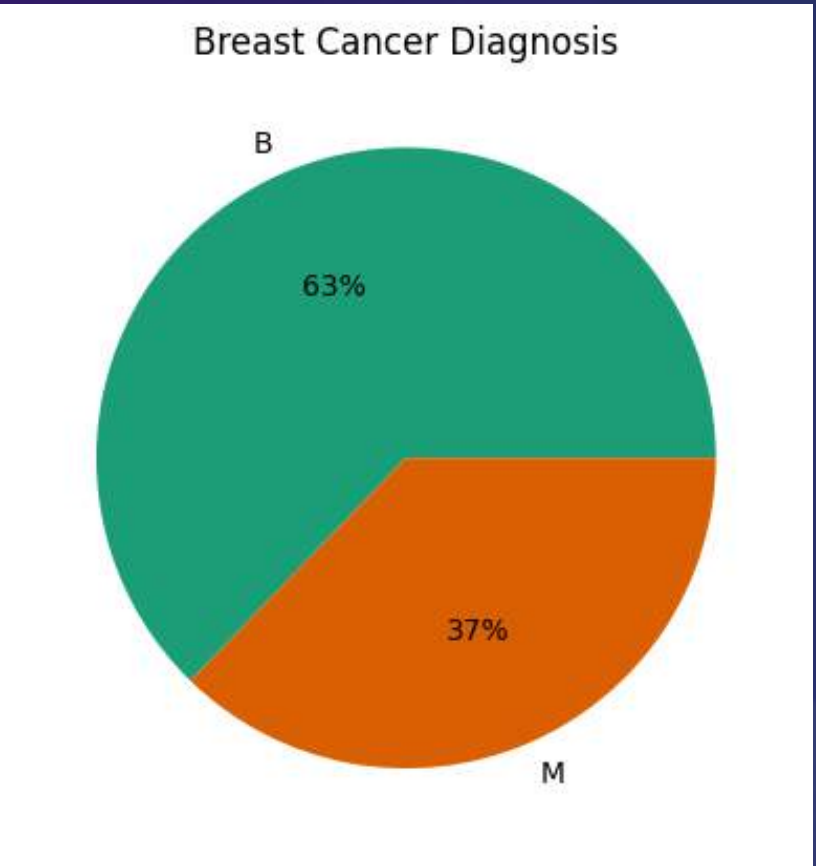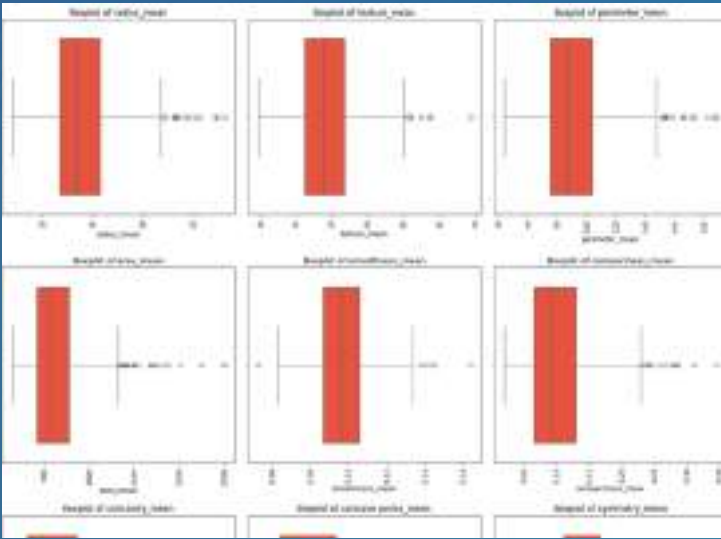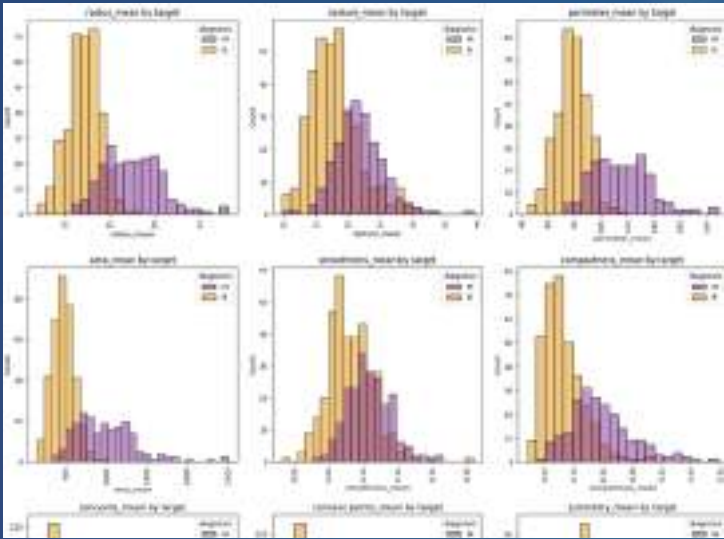


**FULL PROJECT HERE**

# BREAST CANCER PREDICTION

Employed data analysis to identify patterns within the breast cancer dataset. Developed a breast cancer prediction model utilizing the Random Forest algorithm to enhance accuracy and reliability. Two additional models, Gradient Boosting and Stacking, were compared to evaluate performance. Compared against the two models, Random Forest emerged as the top performer with 97% testing accuracy, providing a robust tool for early detection and improved patient outcomes.


Breast Cancer Diagnosis


Confusion Matrix for Random Forest

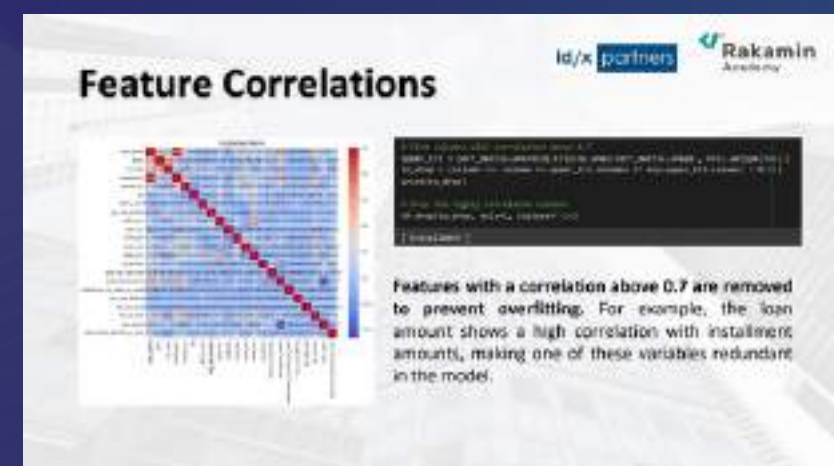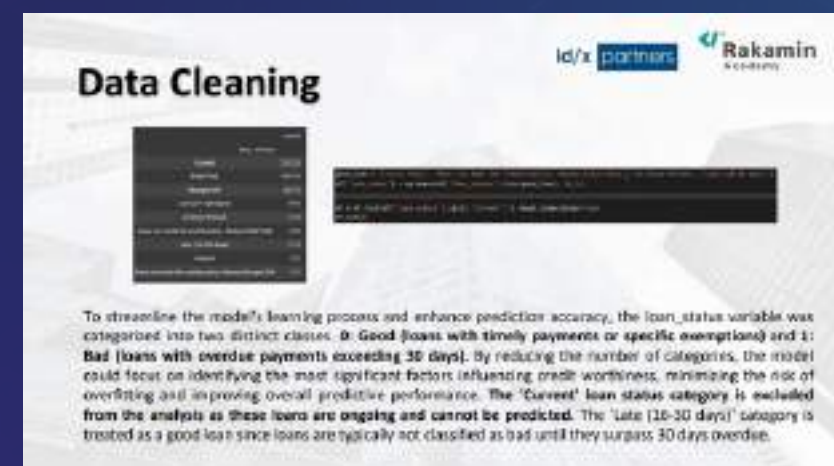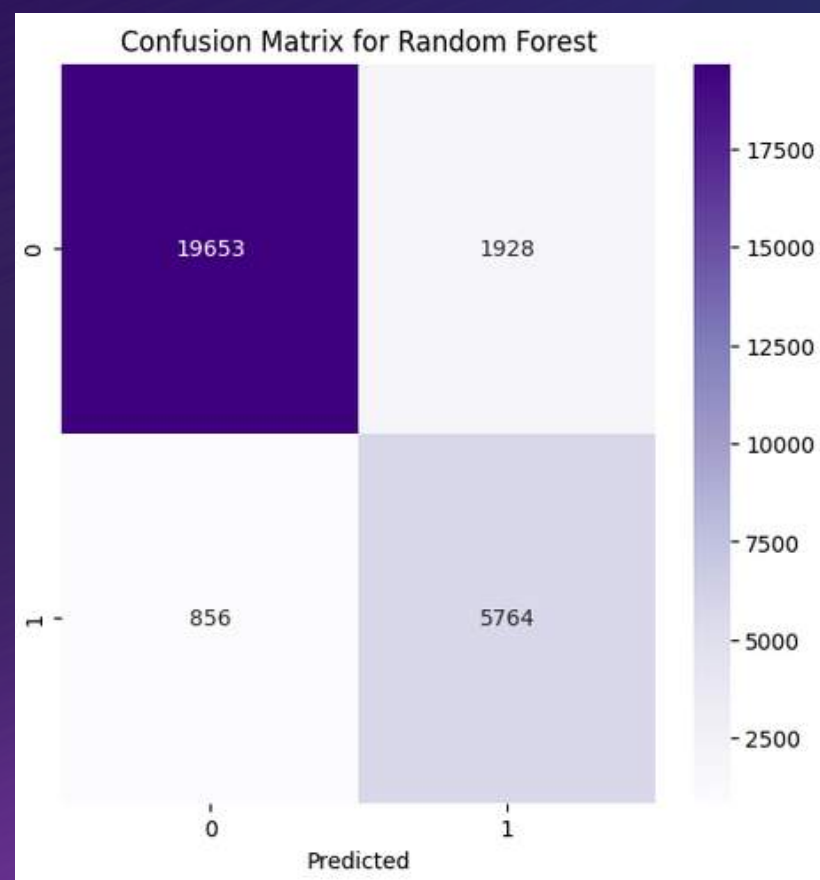| Model | Train Acc | Test Acc | Precision | Recall | Specificity | F1-score | ROC-AUC score |
|---|---|---|---|---|---|---|---|
| Gradient Boosting | 1.00 | 0.96 | 0.96 | 0.95 | 0.93 | 0.95 | 0.95 |
| Random Forest | 1.00 | 0.97 | 0.97 | 0.97 | 0.95 | 0.97 | 0.97 |
| Stacking Model | 1.00 | 0.96 | 0.96 | 0.96 | 0.95 | 0.96 | 0.96 |





**FULL PROJECT HERE**

# CREDIT RISK PREDICTION

Leveraged Python for exploratory data analysis (EDA) on over 400,000 credit risk data rows spanning 7 years, uncovering actionable insights to enhance company operations. The project focused on analyzing patterns of bad loans and predicting credit risk, with the goal of reducing the instance of bad credit due to a high default rate compared to industry standards. Addressed data imbalance using SMOTE oversampling and developed a Random Forest model that achieved 90% accuracy and an AUC of 89%, enabling more informed, data-driven decisions and improving credit risk management.



**FULL PROJECT HERE**

**FULL SLIDE**