# Predict Clicked Ads Customer Classification by using Machine Learning

**Created by:**
**Bagus Ariobimo**
bagusariobimo7@gmail.com
https://www.linkedin.com/in/bagus-ariobimo
https://github.com/riyouuyt

"A data enthusiast who has completed a course in this field and is ready to start his career. Have a excellent understanding in statistics, programming, and data processing. Proficient in using tools such as Python, R, and SQL. Able to collect, clean and analyze data with the necessary techniques. Have skills in data visualization and simple statistical modeling. Creative problem solver and always thrive for learning. Ready to contribute to data-driven projects and collaborate in teams. Committed to further developing data science skills and achieving significant results in data analysis."

**Project Overview** ✿

A company in Indonesia wants to evaluate the effectiveness of the advertisements they broadcast. This is crucial for the company to assess the success of their marketing efforts in attracting customers to view the advertisements. By analyzing historical advertising data and identifying insights and patterns, this project aims to assist the company in defining their marketing targets. The primary objective of this project is to create a machine learning classification model that can determine the right customer targets.

**Goals** 🎯

- Create Effective Model: Build a machine learning model to predict ad clicks accurately.
- Enhance Marketing Strategy: Use data insights to improve the company's marketing campaigns.
- Optimize Targeting: Help the company target the right audience for better ad performance.

**Objectives** 📋

- Model Development: Develop a robust ad click prediction model.
- Insight Generation: Identify actionable insights from historical ad data.
- Data-Driven Decisions: Promote data-driven marketing decisions within the company.

## 📚 Python Libraries for The Project:

**pandas** 🐼: for data manipulation.
**numpy** 🔢: for numerical operations.
**seaborn** 🌊: for data visualization.
**matplotlib** 📊: for plotting.
**scikit-learn × :** for machine learning tasks.

---

## 🤘 Machine Learning Goals:

- The goal is to predict whether a user will click on an ad based on their behavior and demographics.
- Feature engineering and model selection will be used to achieve this goal.

## 📊 Dataset Overview:

**Total Rows: 1000**
**Total Columns: 11**

### Columns Explanation:

- **'Daily Time Spent on Site'**: Average time spent by users on the website.
- **'Age'**: Age of users.
- **'Area Income'**: Income level of users' areas.
- **'Daily Internet Usage'**: Daily internet usage by users.
- **'Male'**: Gender of users (binary: "Male" or "Female").
- **'Timestamp'**: Timestamps or dates of user data.
- **'Clicked on Ad'**: Indicator of whether a user clicked on an ad (binary: "Yes" or "No").
- **'City'**: City where users are located.
- **'Province'**: Province or region of users.
- **'Category'**: Categorical variable or category.

## NUMERICAL FEATURES ▦

☑ KEY INSIGHTS BASED ON HISTOGRAM:

1. **DAILY TIME SPENT ON SITE:**
   - NEARLY SYMMETRICAL DISTRIBUTION.
   - SLIGHT NEGATIVE KURTOSIS INDICATES A FLATTER DISTRIBUTION.
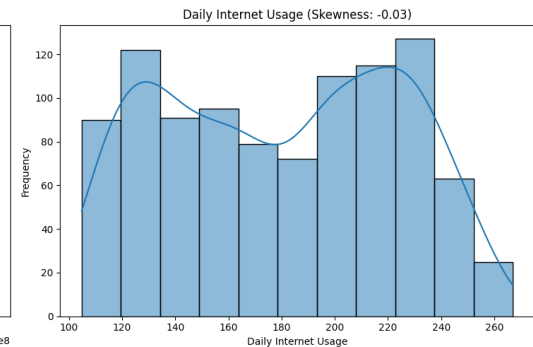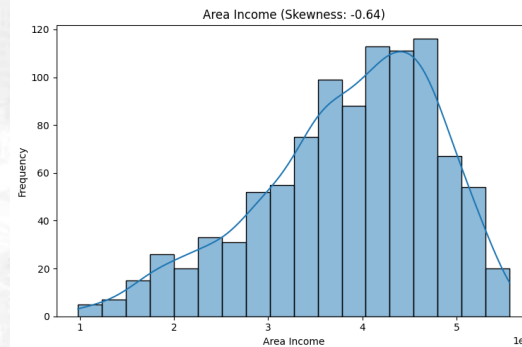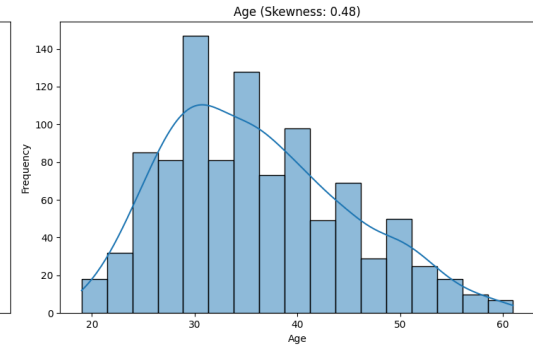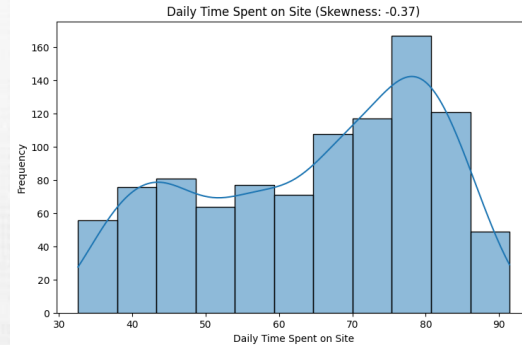
2. **AGE:**
   - SLIGHTLY POSITIVELY SKEWED DISTRIBUTION.
   - NEGATIVE KURTOSIS SUGGESTS A RELATIVELY FLAT PEAK.

3. **AREA INCOME:**
   - MODERATELY NEGATIVELY SKEWED DISTRIBUTION.
   - KURTOSIS CLOSE TO ZERO INDICATES A NORMAL-LIKE DISTRIBUTION.

4. **DAILY INTERNET USAGE:**
   - NEARLY SYMMETRICAL DISTRIBUTION.
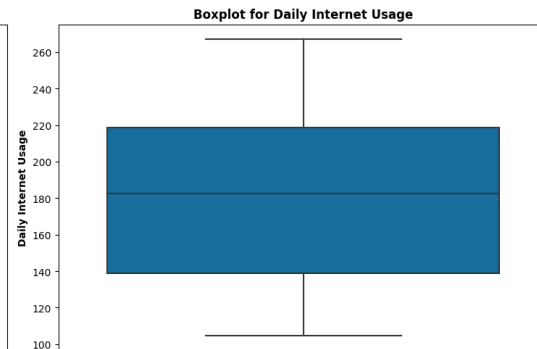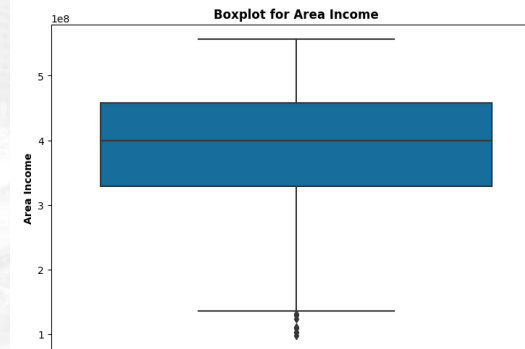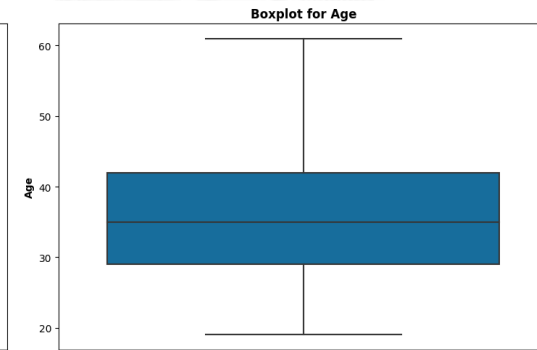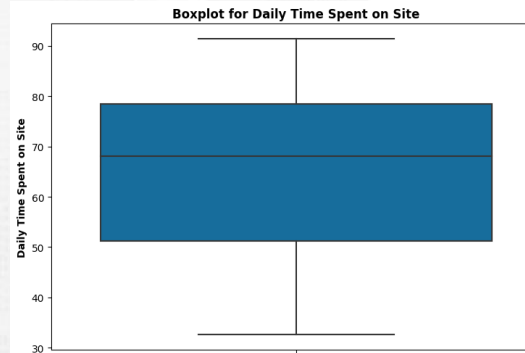   - NEGATIVE KURTOSIS INDICATES A FLAT PEAK.



**For more info of the project, click here**

## NUMERICAL FEATURES 🔢

🔍 KEY TAKEAWAYS:

- FOR THE FEATURE DAILY TIME SPENT ON SITE, THE DATA EXHIBITS A SPREAD WITH AN INTERQUARTILE RANGE (IQR) OF 27.19. USERS TEND TO SPEND VARYING AMOUNTS OF TIME ON THE SITE, WITH THE MIDDLE 50% SPENDING BETWEEN 51.27 AND 78.46 UNITS.

- IN TERMS OF AGE, THERE'S A NOTABLE SPREAD IN USER AGES, AS INDICATED BY AN IQR OF 13.00. THE MIDDLE 50% OF USERS FALL BETWEEN THE AGES OF 29 AND 42

- THE FEATURE AREA INCOME DISPLAYS SIGNIFICANT VARIABILITY, WITH AN IQR OF 129,722,500.00. USERS' INCOMES IN THE MIDDLE 50% RANGE FROM 328,633,000.00 TO 458,355,400.00, INDICATING DIVERSE INCOME LEVELS.

- DAILY INTERNET USAGE DEMONSTRATES A WIDE RANGE OF USAGE, WITH AN IQR OF 80.08. USERS IN THE MIDDLE 50% CONSUME DAILY INTERNET SERVICES BETWEEN 138.71 AND 218.79 UNITS.



**For more info of the project, click here**

## CATEGORICAL FEATURE

### 🏙 PROVINCE DISTRIBUTION:

THE TOP 5 PROVINCES WITH THE HIGHEST REPRESENTATION AMONG USERS ARE:

- **DAERAH KHUSUS IBUKOTA JAKARTA** (253 USERS)
- **JAWA BARAT** (210 USERS)
- **JAWA TIMUR** (90 USERS)
- **BANTEN** (76 USERS)
- **JAWA TENGAH** (53 USERS)

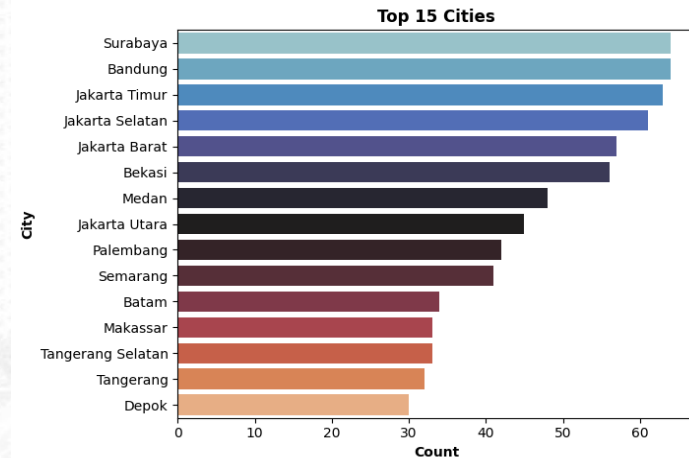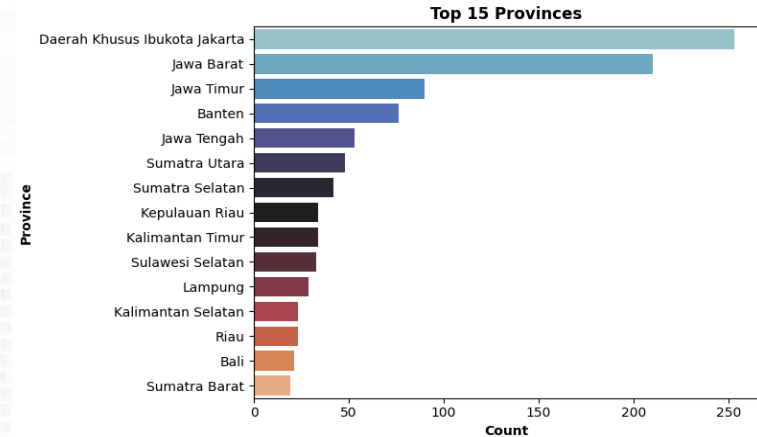JAKARTA'S CAPITAL REGION, DAERAH KHUSUS IBUKOTA JAKARTA, HAS THE HIGHEST USER COUNT.

### 🏢 CITY DISTRIBUTION:

THE TOP 5 CITIES WITH THE HIGHEST REPRESENTATION AMONG USERS ARE:

- **SURABAYA** (64 USERS)
- **BANDUNG** (64 USERS)
- **JAKARTA TIMUR** (63 USERS)
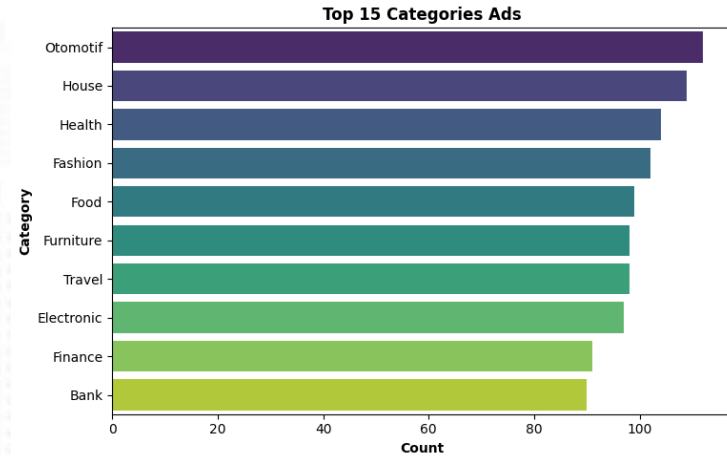- **JAKARTA SELATAN** (61 USERS)
- **JAKARTA BARAT** (57 USERS)

SURABAYA AND BANDUNG ARE THE MOST PROMINENT CITIES AMONG USERS.



Top 15 Provinces



Top 15 Cities

# Exploratory Data Analysis
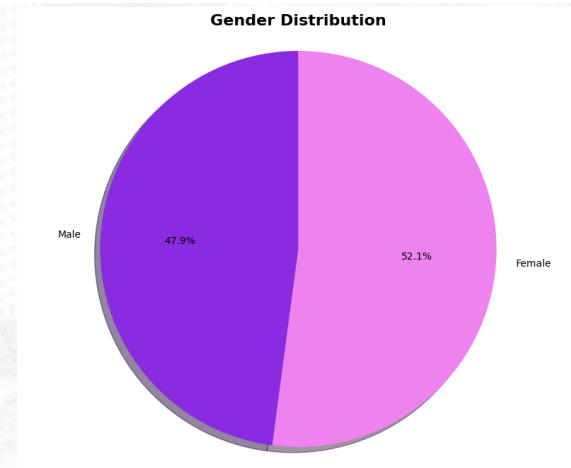
## CATEGORICAL FEATURE

### ADS CATEGORIES 📢

WE EXAMINED THE DISTRIBUTION OF USERS AMONG DIFFERENT AD CATEGORIES. THE TOP AD CATEGORIES INCLUDE **'OTOMOTIF,' 'HOUSE,' 'HEALTH,' 'FASHION,' 'FOOD,' 'FURNITURE,' 'TRAVEL,' 'ELECTRONIC,' 'FINANCE,'** AND **'BANK.'** UNDERSTANDING THESE CATEGORIES HELPS US TAILOR OUR MARKETING STRATEGIES TO TARGET SPECIFIC INTERESTS AND PREFERENCES OF OUR USERS.

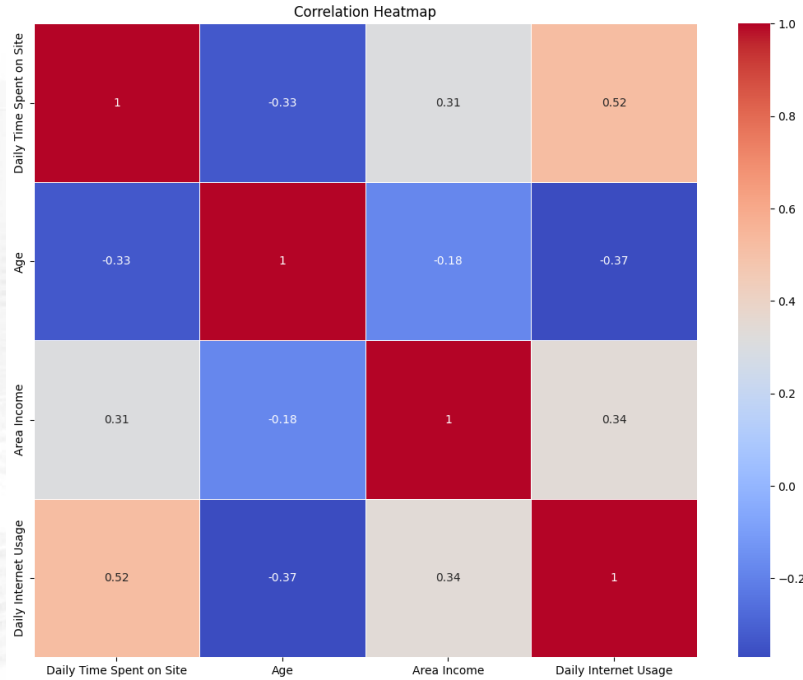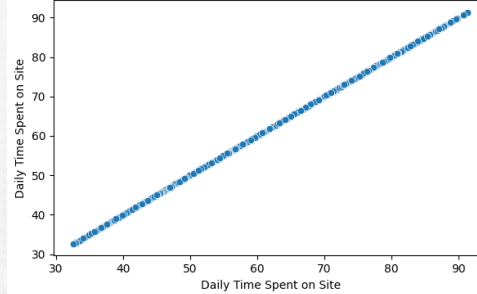### GENDER DISTRIBUTION 👥

WE ANALYZED THE GENDER DISTRIBUTION AMONG OUR USERS. APPROXIMATELY 47.9% OF USERS IDENTIFY AS **MALE**, WHILE 52.1% IDENTIFY AS **FEMALE**. ADDITIONALLY, THERE WERE 3 CASES (0.3%) WITH MISSING GENDER INFORMATION. THIS GENDER BREAKDOWN ALLOWS US TO CUSTOMIZE OUR MARKETING CAMPAIGNS TO RESONATE WITH DIFFERENT SEGMENTS OF OUR AUDIENCE.
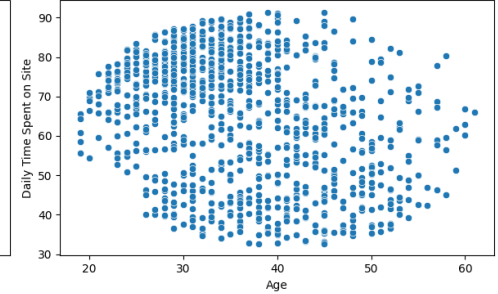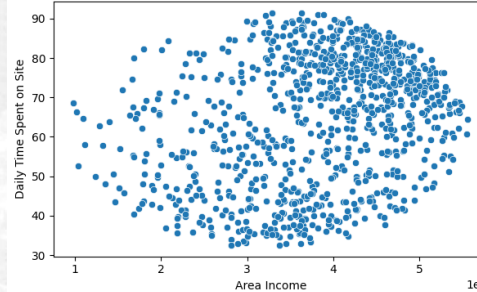


Top 15 Categories Ads



Gender Distribution

## BIAVARIATE ANALYSIS

SUMMARY FOR THE ALL THE CORRELATION COLUMNS ABOVE:

- **DAILY TIME SPENT ON SITE VS. AGE:** NEGATIVE CORRELATION (-0.3314). YOUNGER USERS TEND TO SPEND MORE TIME ON THE SITE.

- **DAILY TIME SPENT ON SITE VS. AREA INCOME:** POSITIVE CORRELATION (0.3083). USERS FROM HIGHER-INCOME AREAS SPEND SLIGHTLY MORE TIME ON THE SITE.

- **DAILY TIME SPENT ON SITE VS. DAILY INTERNET USAGE:** STRONG POSITIVE CORRELATION (0.5183). MORE TIME ON THE SITE IS ASSOCIATED WITH HIGHER DAILY INTERNET USAGE.

- **AGE VS. AREA INCOME:** NEGATIVE CORRELATION (-0.1793). YOUNGER USERS OFTEN RESIDE IN LOWER-INCOME AREAS.

- **AGE VS. DAILY INTERNET USAGE:** NEGATIVE CORRELATION (-0.3705). YOUNGER USERS HAVE HIGHER DAILY INTERNET USAGE.

- **AREA INCOME VS. DAILY INTERNET USAGE:** POSITIVE CORRELATION (0.3381). HIGHER AREA INCOME RELATES TO HIGHER DAILY INTERNET USAGE.

- **STRONGEST CORRELATION:** DAILY TIME SPENT ON SITE AND DAILY INTERNET USAGE (0.5183). A STRONG POSITIVE RELATIONSHIP EXISTS.

- **AGE VS. DAILY TIME SPENT ON SITE:** NEGATIVE CORRELATION (-0.3314). OLDER USERS SPEND LESS TIME ON THE SITE.

IN SUMMARY, USER ENGAGEMENT IS POSITIVELY LINKED TO INTERNET USAGE, WITH YOUNGER USERS BEING MORE ACTIVE. ADDITIONALLY, USERS FROM HIGHER-INCOME AREAS TEND TO ENGAGE MORE WITH THE SITE. HOWEVER, AGE AND AREA INCOME ARE NEGATIVELY RELATED, SUGGESTING THAT YOUNGER USERS OFTEN RESIDE IN LOWER-INCOME AREAS. THESE INSIGHTS CAN INFORM MARKETING STRATEGIES AND TARGET DEMOGRAPHIC

☑ **Customer Type and Behavior Analysis on Advertisement** ☑

In our comprehensive analysis, we delved into customer behavior patterns and demographics, shedding light on their interaction with digital advertisements.

- Ads Categories 📣: We examined the distribution of users across various ad categories, helping us understand their preferences and interests. The top categories include **'Otomotif,' 'House,' 'Health,' 'Fashion,' 'Food,' 'Furniture,' 'Travel,' 'Electronic,' 'Finance,' and 'Bank.'** This knowledge empowers us to tailor our marketing strategies to align with specific customer interests.

- Gender Distribution 👫: Our analysis revealed a balanced gender distribution, with approximately 47.9% of users identifying as Male and 52.1% as Female. There were 3 cases (0.3%) with missing gender information. This insight enables us to create targeted and inclusive marketing campaigns.

- Customer Type 🎯: Through behavioral analysis, we identified different customer types based on their interaction with advertisements. By categorizing users into various segments, we gain a deeper understanding of their preferences, needs, and responsiveness to our marketing efforts.

- Age Groups 🎂: We explored age intervals to determine which age ranges are more likely to engage with advertisements. This information aids in crafting age-specific marketing strategies for better results.

**For more info of the project, click <u>here</u>**
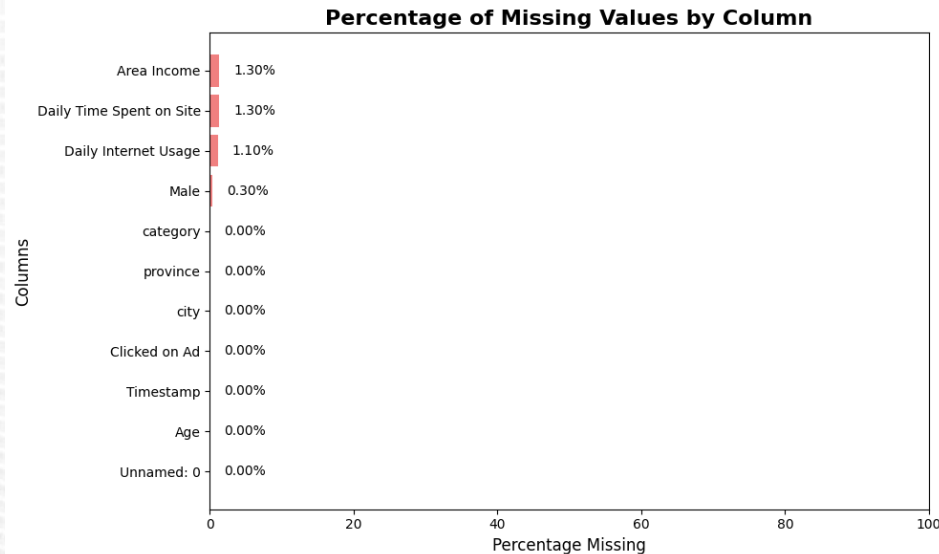
**Checking Null Values in Dataset** 🕵️

In our dataset, we meticulously examined the presence of missing values. Here's a quick overview:

📊 Columns with Null Values:
- **'Daily Time Spent on Site' and 'Area Income'** each had a 1.3% null rate.
- **'Daily Internet Usage'** had a 1.1% null rate.
- **'Male'** exhibited a mere 0.3% of missing values.
- **'Unnamed: 0,' 'Age,' 'Timestamp,' 'Clicked on Ad,' 'city,' 'province,'** and **'category'** were entirely free from null values.

Ensuring data integrity is essential for reliable analysis. We'll handle these missing values appropriately to maintain data quality. 🛠️

**Percentage of Missing Values by Column**

| Column | Percentage Missing |
|---|---|
| Area Income | 1.30% |
| Daily Time Spent on Site | 1.30% |
| Daily Internet Usage | 1.10% |
| Male | 0.30% |
| category | 0.00% |
| province | 0.00% |
| city | 0.00% |
| Clicked on Ad | 0.00% |
| Timestamp | 0.00% |
| Age | 0.00% |
| Unnamed: 0 | 0.00% |

**For more information about the project, click here**

☐ **Data Cleaning and Preprocessing Summary** ☐

Our data underwent a comprehensive cleaning and preprocessing journey, ensuring its readiness for analysis. Here's what we accomplished:

- 🚹 **Handling Missing Values:** We meticulously addressed missing values, ensuring minimal data loss: **'Daily Time Spent on Site,' 'Area Income,' 'Daily Internet Usage,'** and **'Male'** columns had missing values, which were handled appropriately. **'Unnamed: 0,' 'Age,' 'Timestamp,' 'Clicked on Ad,' 'city,' 'province,' and 'category'** columns remained free from null values.
- ☐ **Duplicate Data:** We identified and successfully removed duplicate rows from the dataset, leaving us with a pristine dataset, completely devoid of duplicates.
- 🔁 **Feature Encoding:** We encoded categorical columns **'Gender'** and **'Clicked on Ad'** using one-hot encoding, making them ready for machine learning models.
- 😛 **Column Renaming:** We renamed columns for clarity and consistency, including renaming **'Gender'** columns for better interpretation.
- ✏️ **Column Removal:** We removed the **'Unnamed: 0'** column, which didn't contribute to our analysis.

Now, with our dataset clean, structured, and encoded, we're primed for in-depth analysis and predictive modeling! 🚀🔍

In our data preprocessing journey, we took crucial steps to prepare the dataset for analysis:

- 🎯 **Splitting Data:** We divided our dataset into features (X) and the target variable (y). The **'Clicked on Ad_Yes'** column serves as our target, while the remaining columns comprise our features.

- 🕰 **Extracting Date Information:** To unlock valuable insights, we extracted date-related details from the **'Timestamp'** column:
  - We converted the **'Timestamp'** column to datetime format to ensure accurate time-related calculations.
  - We created new columns, including **'Year,' 'Month,' 'Week,' and 'Day,'** to dissect the timestamp data into various time components.

These meticulous preparations set the stage for our exploratory data analysis and predictive modeling, empowering us to uncover meaningful patterns and make informed decisions! 🚀☑

Machine Learning Experiments:

1. **Experiment 1 (Before Normalization/Standardization):**
- Low accuracy (49%) and poor performance.
- Top Features: Daily Time Spent on Site, Daily Internet Usage.

2. **Experiment 2 (After Normalization/Standardization):**
- High accuracy (96%) and improved performance.
- Top Features: Daily Time Spent on Site, Daily Internet Usage.



Experiment 1 Confusion Matrix



Experiment 2 Confusion Matrix

**For more information about the project, Click Here**

# Data Modeling

Top Features by Coefficient Magnitude:

- Daily Time Spent on Site: -2.83
- Daily Internet Usage: -2.76
- Area Income: -1.60
- Age: 1.40
- category_Furniture: -0.56
- city_Jakarta Pusat: -0.51
- province_Kalimantan Timur: -0.50
- city_Semarang: 0.36
- category_Travel: -0.35
- category_Electronic: -0.32



These features have the highest absolute coefficient magnitudes in your logistic regression model. Features with negative coefficients have a negative impact on the prediction of clicking on an ad, while features with positive coefficients have a positive impact.

**For more information about the project, Click Here**

Rakamin
Academy

☑ Interpretation of Model Results ☐

- Experiment 1 lacked data normalization and performed poorly in predicting ad clicks. The top features, Daily Time Spent on Site and Daily Internet Usage, were not sufficient for accurate predictions.
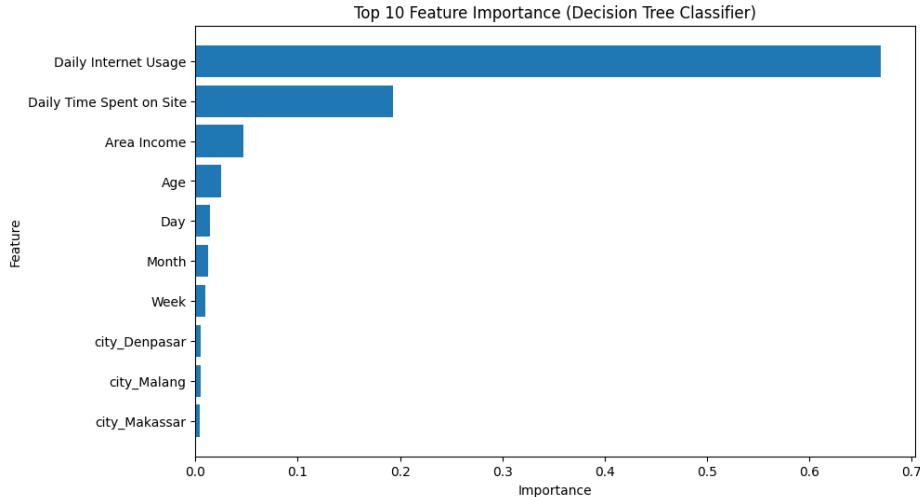- Experiment 2, with data normalization, significantly improved accuracy and precision. It successfully identified potential ad clickers.
- The top features affecting predictions are **Daily Time Spent on Site** and **Daily Internet Usage**, highlighting their importance in understanding user behavior.
- Data normalization plays a crucial role in enhancing model performance, making it a recommended practice.

Use these insights to refine marketing strategies and target potential ad clickers effectively. 🚀

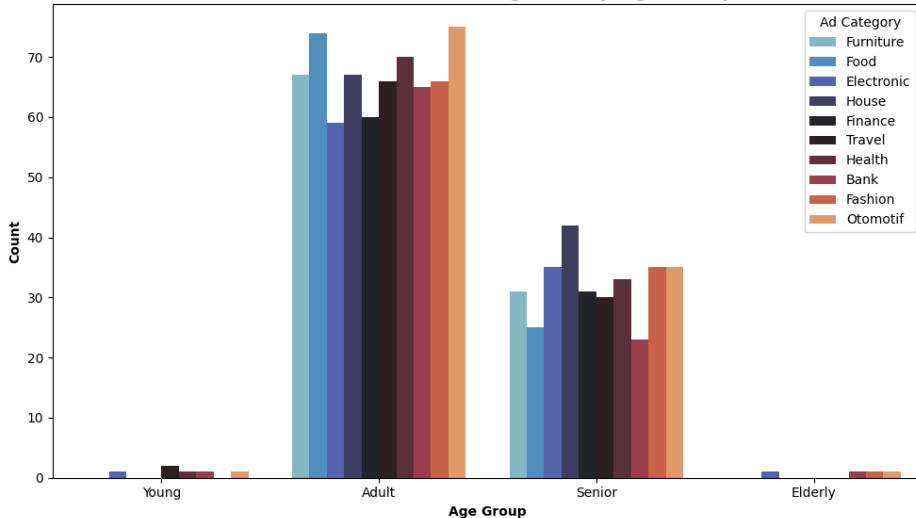Top 10 Feature Importance (Decision Tree Classifier)

The feature importance analysis has identified that **"Daily Internet Usage"** and **"Daily Time on Site"** are the two most influential factors affecting user behavior and ad click-through rates. 🚀

- **Daily Internet Usage**: This feature is a critical determinant of user engagement. Users who spend more time online tend to be more active and attentive. This is a strong indicator of their responsiveness to ads.

- **Daily Time on Site**: The time users spend on the site directly impacts their exposure to ads. Longer visits provide more opportunities for users to interact with advertisements.

These two factors, **"Daily Internet Usage"** and **"Daily Time on Site,"** are highly valuable in understanding user behavior and optimizing ad campaigns. 📊💡

**For more information about the project, Click** Here

**Interaction with Ad Categories by Age Group**

## Age Group Distribution 📊

Based on the feature importance, The age is pretty importance and can add some recommendation into the business distribution is predominantly concentrated in the "Adult" and "Senior" categories, with these two groups making up nearly the entire user base. The "Young" and "Elderly" groups represent only a small portion of the users.

## Business Recommendation 🚀

- Target Adult and Senior Groups: 🎯 The marketing strategy should be primarily designed to target the "Adult" and "Senior" age groups since they make up the majority of users. Tailor advertisements, promotions, and content to resonate with the interests and preferences of these age groups.

- Explore Opportunities: 🔍 While the "Young" and "Elderly" groups are smaller in number, they still represent potential market segments. It may be worth exploring specific strategies to attract and engage these demographics, considering their unique interests and needs.

This business recommendation aims to optimize marketing efforts by focusing on the most significant user segments while exploring potential growth opportunities. ☑

**For more information about the project, Click Here**

# BUSINESS SIMULATION

Rakamin Academy

## 💰 *Business Schemas* 💰

### Company Profile: CarMend Pro 🚗🔧

CarMend Pro is a customer-focused automotive services company specializing in car maintenance and repair. We connect car owners with trusted mechanics, offering convenient and reliable solutions to keep their vehicles in top condition. Our transparent pricing, user-friendly platform, and dedication to innovation make us a leader in the automotive industry.

### Marketing Team Target: 🎯

- Objective:

To drive traffic to the CarMend Pro website, generate interest in our services, and increase sales. 🚗💻📈

- Target Audience: 🎯

Car owners looking for automotive maintenance and repair services. 🚗🔧💡

**For more information about the project, Click Here**

## Assumed Clicked Customers 🖱

Before the machine learning for the classification customers exist, the marketing teams are assuming the prospect customer which there are:

- **Interested Prospects (Assumed 30%):** ☐
  - Clicked on the ad and visited the website.
  - Spent more time on the site, viewed multiple pages.
  - Engaged with interactive tools like the "Price Estimator."

- **Non-Interested Prospects (Assumed 70%):** 😕
  - Clicked on the ad but didn't engage much.
  - Bounced quickly from the landing page.

- **Marketing Team's Assumed Clicked Customers:** 📊
  - Total Clicks: 193
  - Interested Prospects (30% of clicks): 58
  - Non-Interested Prospects (70% of clicks): 135

**For more information about the project, Click Here**

## Expected Revenue Sales and CVR (Conversion Rate): 💰☑

- **Marketing Team's Expected Results:**
  - Total Conversions: 58 (assuming all 58 interested prospects converted).
  - Total Revenue Sales: $11,600 (58 conversions x $200 average transaction value).
  - Conversion Rate (CVR): 100% (since all expected interested prospects are assumed to convert).

- **Non-Interested Prospects:** 💰
  - Assuming a 2% conversion rate for non-interested prospects.
  - Expected Conversions: 3 (135 x 2%).
  - Average Transaction Value: $200
  - Revenue from Non-Interested Prospects: $600 (3 x $200).

- **Total Expected Results:**
  - Total Conversions: 61 (58 from interested + 3 from non-interested)
  - Total Revenue Sales: $12,200 ($11,600 from interested + $600 from non-interested)
  - Conversion Rate (CVR): Approximately 31.6% (61 conversions out of 193 clicks)

**For more information about the project, Click Here**

|  | Not Clicked The Ads | Clicked The Ads |
|---|---|---|
| **Not Clicked The Ads** | 88 | 6 |
| **Clicked The Ads** | 5 | 94 |

Certainly, let's explain that the model was created using the **DecisionTreeClassifier** and compare the revenue and Conversion Rate (CVR) with and without the machine learning model

It achieved an impressive accuracy of **94.3%**. The confusion matrix revealed the following:

- **True Positives (TP)**: 94
- **True Negatives (TN)**: 88
- **False Positives (FP)**: 6
- **False Negatives (FN)**: 5

**For more information about the project, Click Here**

## Comparison Revenue with Machine Learning and without Machine Learning

- **Marketing Team's Expected Results:**
  - Total Conversions: 58 (assuming all 58 interested prospects converted).
  - Total Revenue Sales: $11,600 (58 conversions x $200 average transaction value).
  - Conversion Rate (CVR): 100% (since all expected interested prospects are assumed to convert).

- **Machine Learning Model's Results:**
  - Total Conversions: 182 (as previously mentioned).
  - Total Revenue Sales: $36,400 (as previously mentioned).
  - Conversion Rate (CVR): Approximately 94.3% (as previously mentioned).

**Comparison Summary:**

With Machine Learning (ML) 🤘 :

The model predicts significantly higher conversions and revenue than the marketing team's initial expectations. ML delivers a Conversion Rate (CVR) of around 94.3%, which is lower than the assumed 100% but still quite high. ☑ 💰 Without Machine Learning (ML) The marketing team's expected results would fall short, achieving only 58 conversions and $11,600 in revenue. The CVR remains high, but it doesn't reach 100%. 🙅‍♂️🤷‍♂️

**For more information about the project, Click Here**

Thank You