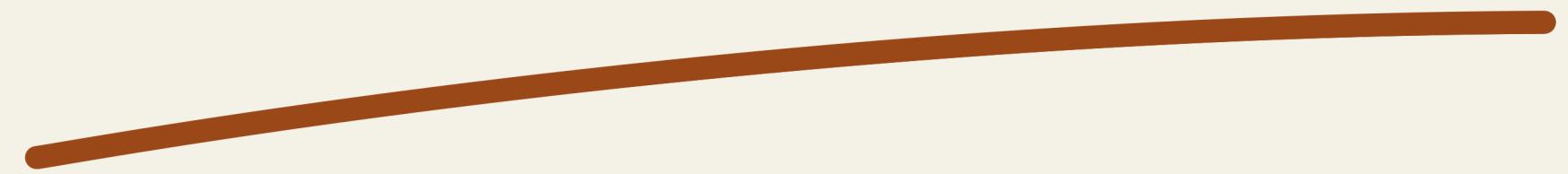




*Lhutfia Ichsan's*

# Portfolio Presentation



# Content

O1 ABOUT AUTHOR

O2 MINI PROJECT 1 - ANALYZING E-COMMERCE  
BUSINESS PERFORMANCE WITH SQL

O3 MINI PROJECT 2 - INVESTIGATE BUSINESS HOTEL  
USING DATA VISUALIZATION

O4 MINI PROJECT 3 - PREDICT CUSTOMER  
PERSONALITY TO BOOST MARKETING CAMPAIGN

O5 MINI PROJECT 4 - PREDICT CUSTOMER CLICKED  
ADS BY USING MACHINE LEARNING

O6 MINI PROJECT 5 - IMPROVING EMPLOYEE  
RETENTION BY PREDICTING EMPLOYEE RETENTION  
USING MACHINE LEARNING



# Hello, I'm Lhutfia

I am a recent Information Systems Student who has interested in Data Processing. I enjoy taking on responsibilities, contributing to a team, and finishing tasks on my own.

---

Honest and on time is my work ethic. And I'am exited for my upcoming adventure!

[BACK TO CONTENT PAGE](#)



# Mini Project 04

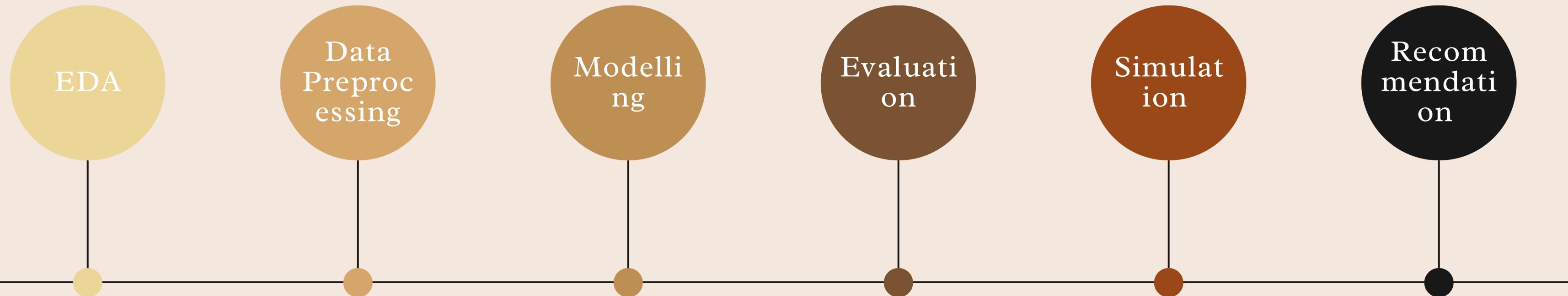
---

## *Predict Customer Clicked Ads by Using Machine Learning*



[BACK TO CONTENT PAGE](#)

# Workflow >>



Melakukan eksplorasi data saat ini

Mempersiapkan data sebelum masuk kedalam machine learning

- Implementasi terhadap data dengan menggunakan dua schema pemodellan
- Menggunakan business metric untuk menentukan model

Melakukan simulasi perbedaan sebelum dan sesudah penggunaan machine learning

Memberikan rekomendasi bisnis yang relevant untuk

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
 #   Column           Non-Null Count   Dtype  
--- 
 0   daily_time_spent_on_site    987 non-null    float64
 1   age                          1000 non-null   int64  
 2   area_income                  987 non-null    float64
 3   daily_internet_usage        989 non-null    float64
 4   male                         997 non-null    object  
 5   timestamp                    1000 non-null   object  
 6   clicked_on_ad                1000 non-null   object  
 7   city                         1000 non-null   object  
 8   province                     1000 non-null   object  
 9   category                     1000 non-null   object  
dtypes: float64(3), int64(1), object(6)
memory usage: 78.2+ KB
```

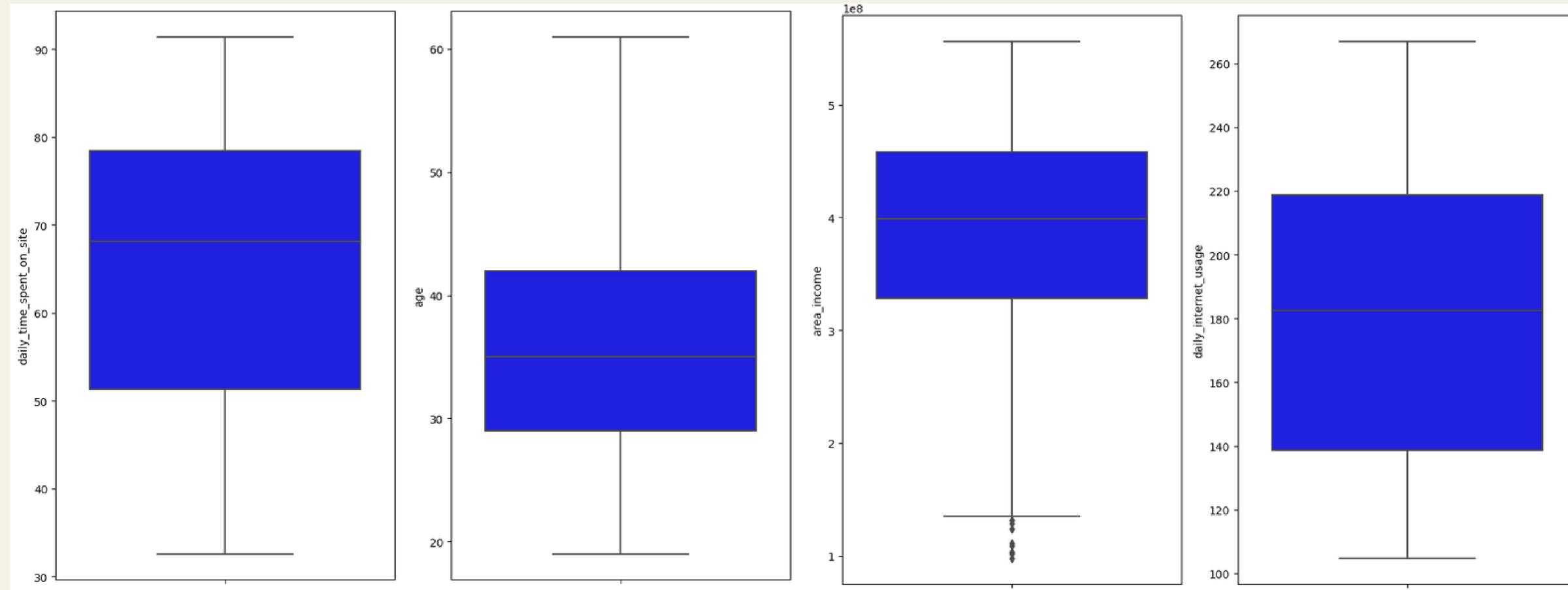
EDA

## How was the Dataset?

- 100 rows, 10 columns
- Terdapat missing values
- Terdapat unnormal data type
- Tidak terdapat data yang aneh untuk masing-masing columns

# Univariate Analysis - EDA

L<sub>I</sub>

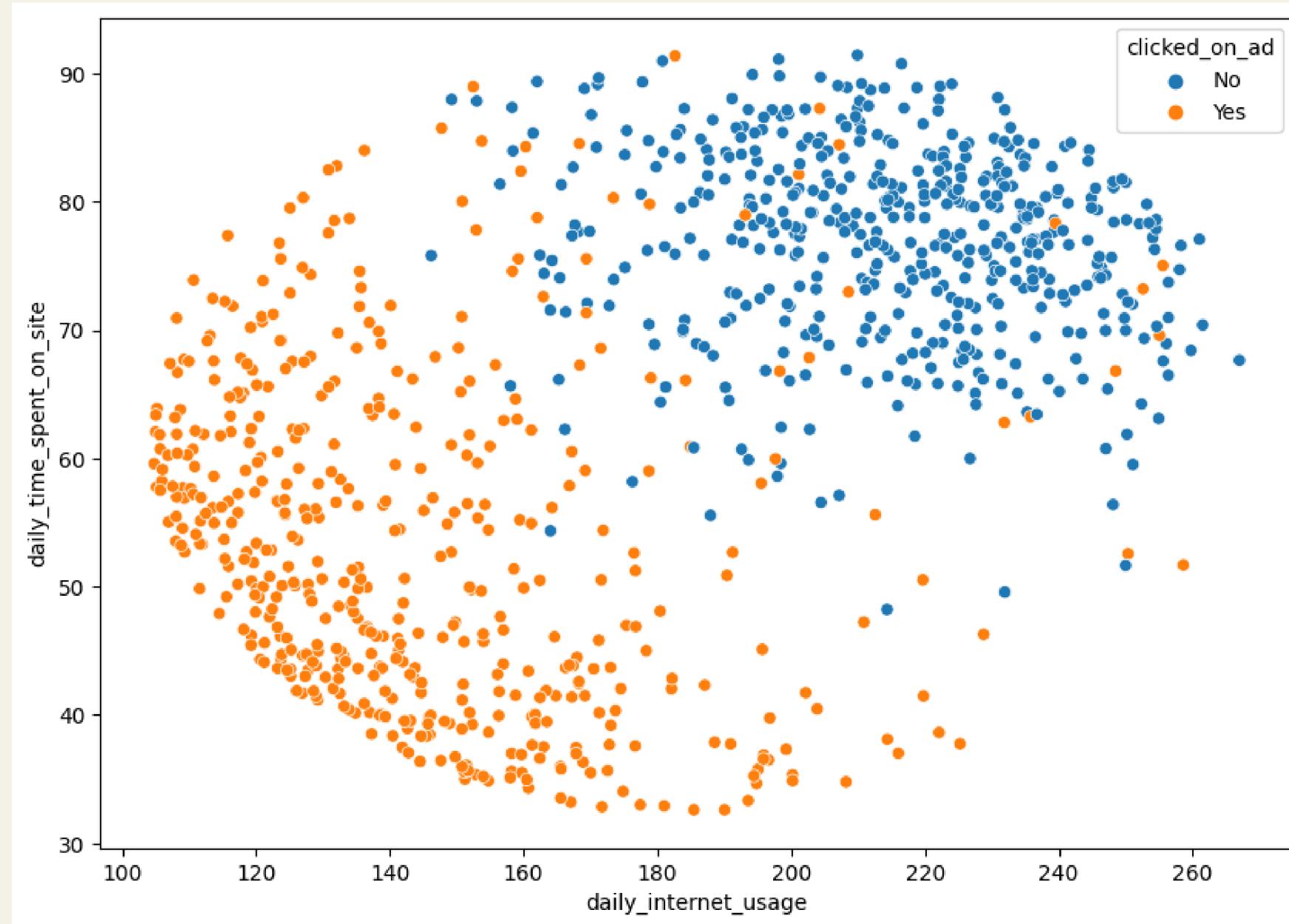


EDA

See Skewed Data  
Using Boxplot

## Bivariate Analysis - EDA

$L_I$

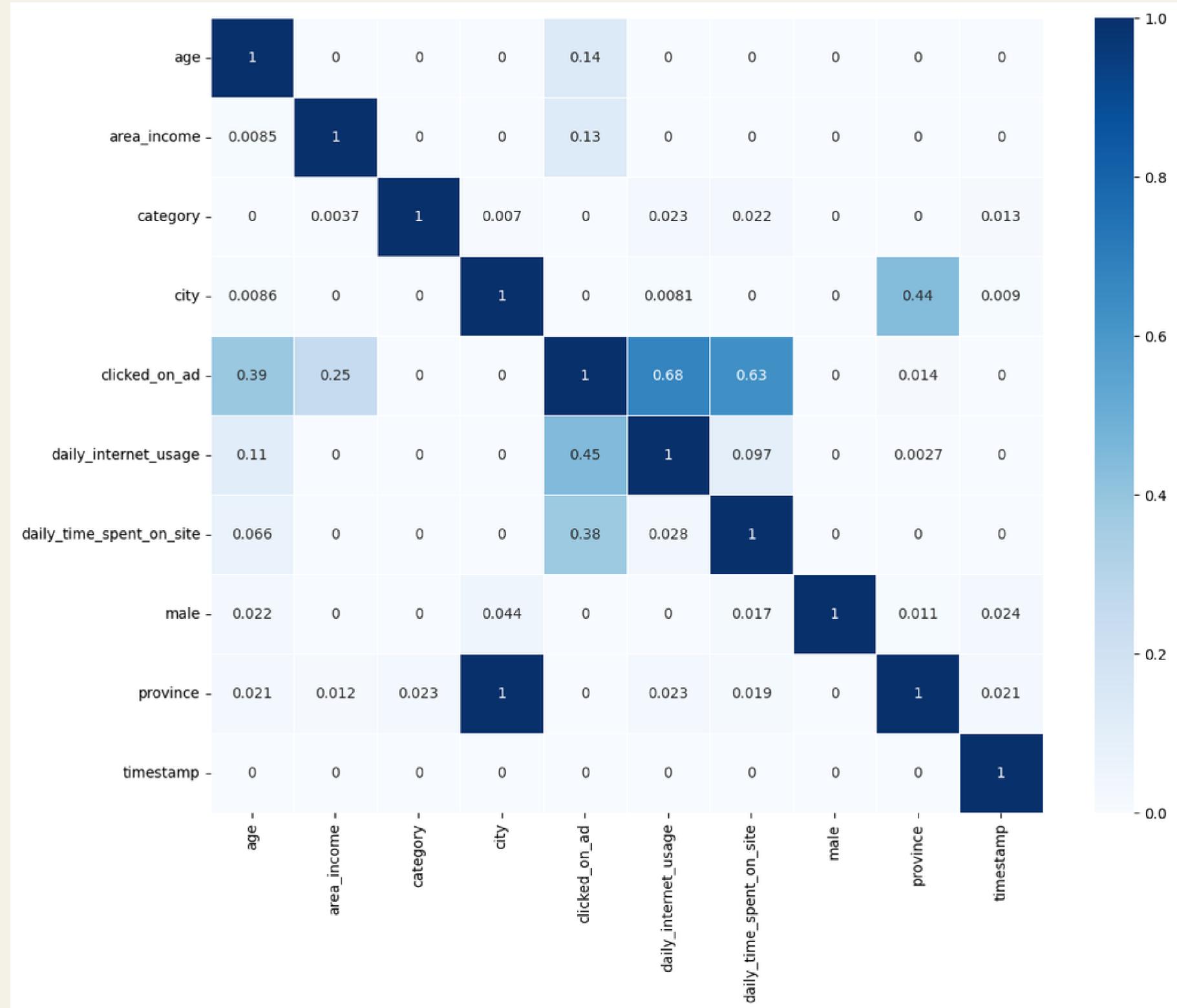


EDA

Using Scatter plot for  
*Crowd Zero*

# Multivariate Analysis - EDA

L  
I



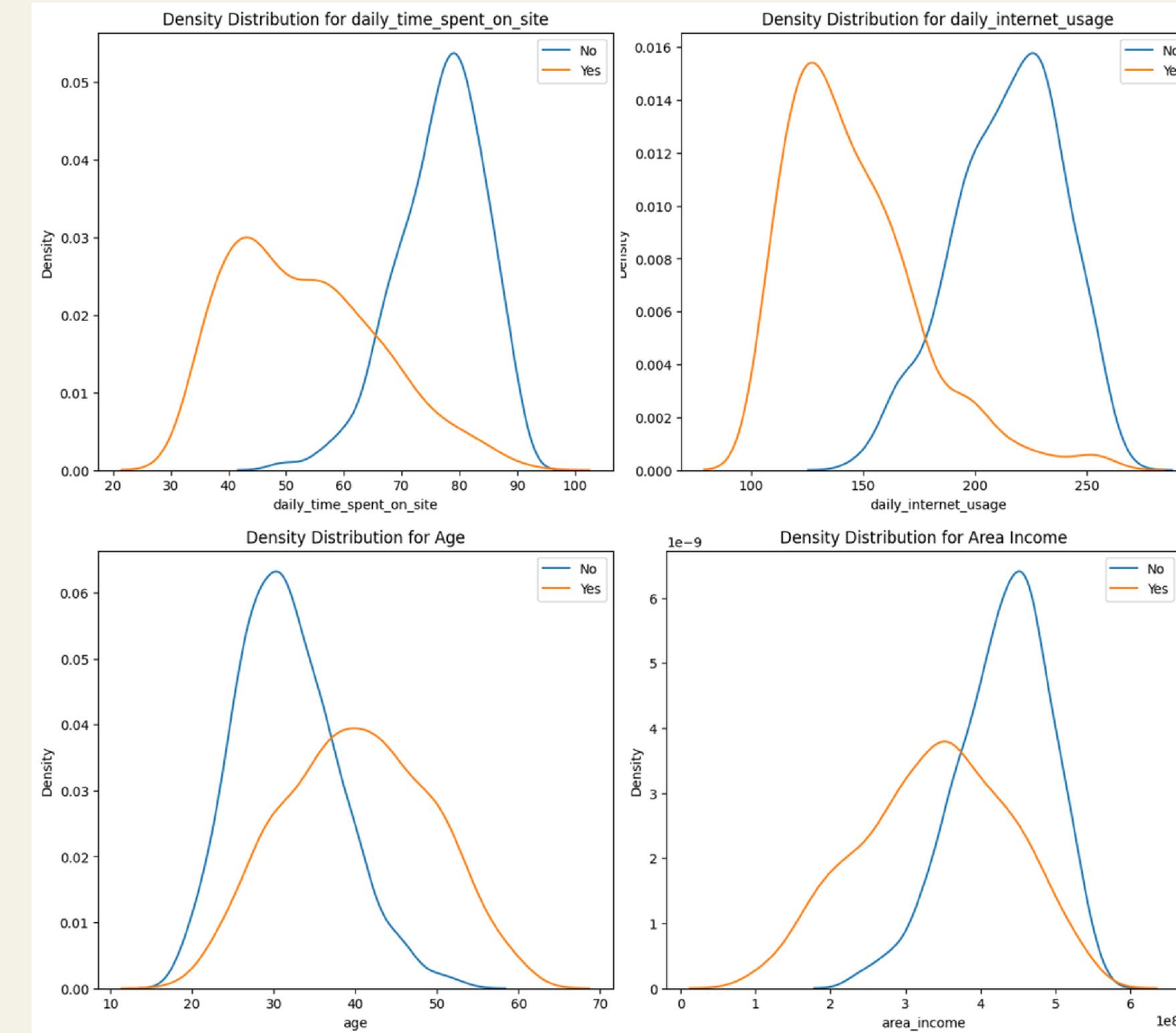
EDA

## Using PPScore for correlation

Diketahui bahwa feature yang berhubungan dengan target (clicked\_ad\_on) adalah : age, area\_income, daily\_usage\_internet, daily\_time\_spent\_on\_site dan province.

# Insight >>

Pelanggan yang cenderung meresponse adalah pelanggan dengan kriteria daily internet usage dan daily time spent pada website kecil, dengan usia yang relative lebih tua dan pendapatan yang relative lebih rendah. Sedangkan untuk pelanggan yang cenderung tidak mereponse adalah yang sebaliknya.



# Data Preprocessing



## Missing values, Data type, and Duplicated data

- Delete missing Values 1.3 %
- Change data type columns = 'timestamp' >> datetime
- Duplicated Data 0.0%



## Feature Encoding

- Feature encoding untuk columns yang bersifat binary Yes/No.
- Feature one hot encoding untuk columns yang non-binary

Feature One Hot Encoding



	male	timestamp	clicked_on_ad	city	province	category	is_Bali	is_Banten	is_Daerah	Khusus	is_Jawa	Barat	is_Jawa	Tengah
0	0	2016-05-30 02:34:00	0	Semarang	Jawa Tengah	Furniture	0	0	0	Khusus	0	0	0	0
1	1	2016-03-19 08:00:00	0	Medan	Sumatra Utara	Electronic	0	0	0	Ibukota	0	0	0	1
0	0	2016-04-29 07:49:00	0	Semarang	Jawa Tengah	Travel	0	0	0	Jakarta	0	0	0	0

Feature Encoding

## Business Metrics for Model Evaluation

Recall, dipertimbangkan untuk mengurangi False Negative yaitu keadaan dimana Machine Learning memprediksi pelanggan tidak akan meresponse namun ternyata meresponse >> Lost Cost

Precision, dengan tujuan untuk mengurangi False Positive yaitu keadaan dimana Machine Learning memprediksi pelanggan tidak meresponse namun ternyata meresponse >> Lost Revenue

## Revenue vs. Cost



# Modelling



## WITHOUT FEATURE SCALING

	model_name	model	accuracy	recall	precision	duration
0	K-Nearest Neighbor	KNeighborsClassifier()	0.671280	0.630137	0.691729	0.009359
1	Logistic Regression	LogisticRegression()	0.494810	0.000000	0.000000	0.024757
2	Decision Tree	DecisionTreeClassifier()	0.930796	0.952055	0.914474	0.020319
3	Random Forest	(DecisionTreeClassifier(max_features='sqrt', r...	0.948097	0.945205	0.951724	0.549195
4	Gradient Boosting	((DecisionTreeRegressor(criterion='friedman_ms...	0.944637	0.958904	0.933333	0.793807

## WITH FEATURE SCALING

	model_name	model	accuracy	recall	precision	duration
0	K-Nearest Neighbor	KNeighborsClassifier()	0.750865	0.732877	0.764286	0.008675
1	Logistic Regression	LogisticRegression()	0.948097	0.924658	0.971223	0.023625
2	Decision Tree	DecisionTreeClassifier()	0.910035	0.924658	0.900000	0.007551
3	Random Forest	(DecisionTreeClassifier(max_features='sqrt', r...	0.951557	0.952055	0.952055	0.582976
4	Gradient Boosting	((DecisionTreeRegressor(criterion='friedman_ms...	0.937716	0.952055	0.926667	0.942897

## Modelling With and Without Feature Scalling

- Feature scaling mampu meningkatkan performa khususnya pada model Logistic Regression
- Model dengan performa terbaik adalah Random Forest dengan Recall 0.95 dan Precision 0.95

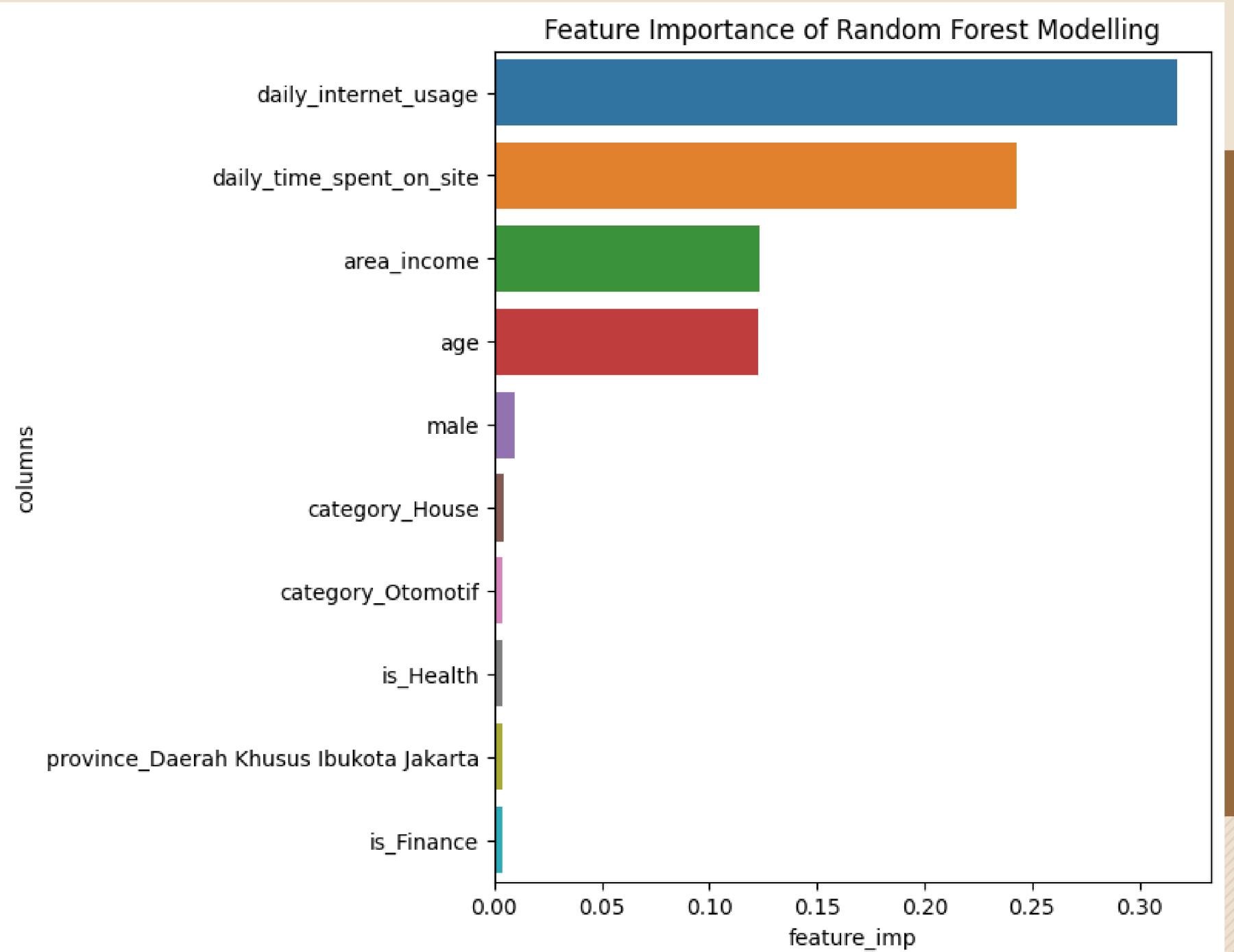
# Feature Importance

L<sub>I</sub>

Berikut ini merupakan feature Importance dari penggunaan model Random Forest yaitu :

1. Daily internet usage
2. Daily time spent on site
3. Income
4. Age

Feature Importance sejalan dengan hubungan correlation yang telah disebutkan sebelumnya.



# Business Simulation

Untuk melakukan simulasi, diasumsikan bahwa biaya yang perusahaan keluarkan adalah \$3/iklan. Kemudian, revenue yang didapatkan perusahaan adalah \$11/click. Maka :

$$\text{cost} = \$3 \times (289 \text{ iklan}) = \$ 867$$

$$\text{revenue} = \$11 \times (146 \text{ iklan}) = \$ 1.606$$

$$\text{revenue} - \text{cost} = 1.606 - 867 = \$ 739$$

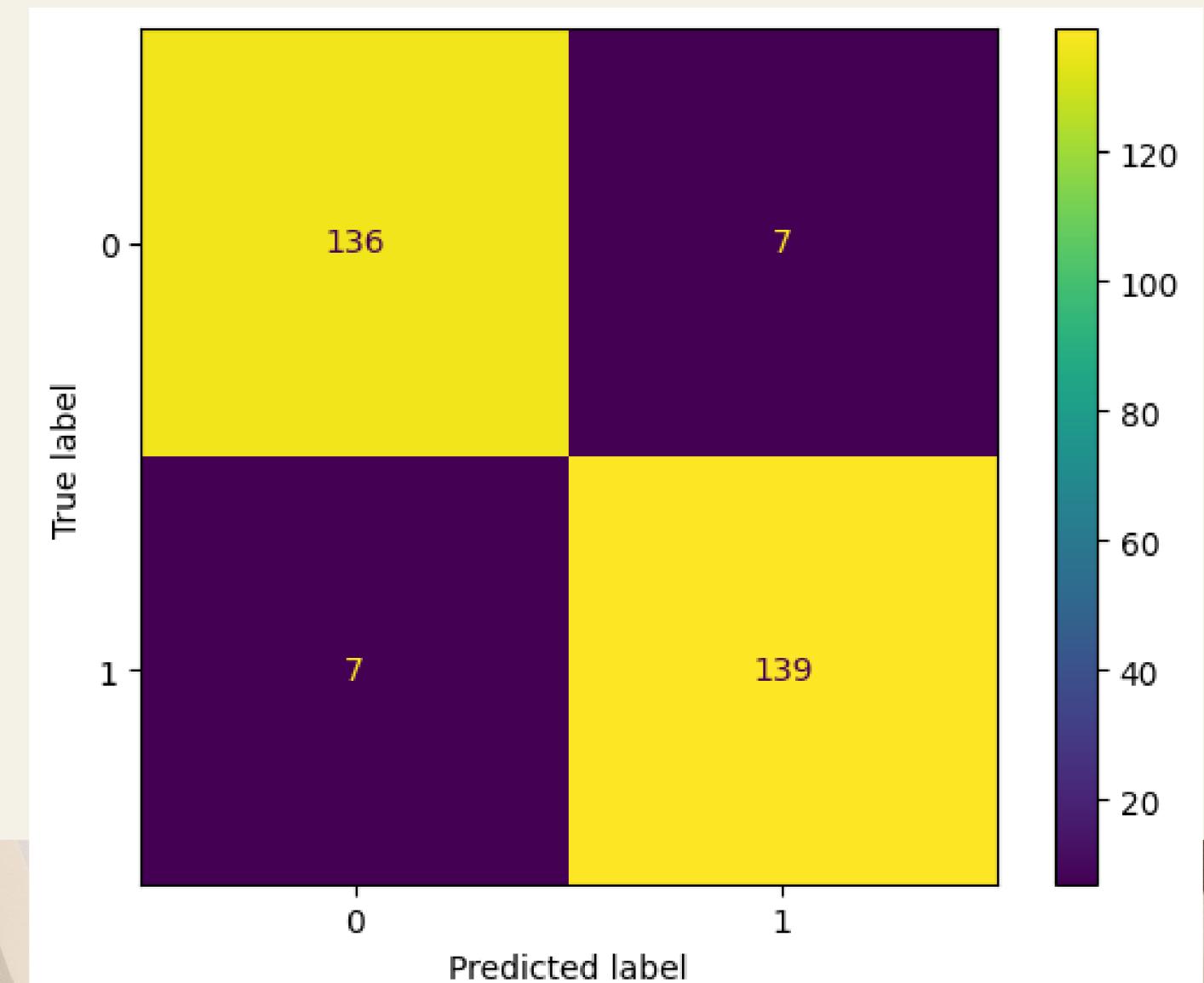
$$\text{cost} = \$3 \times (145 \text{ iklan}) = \$435$$

$$\text{revenue} = \$11 \times (139 \text{ iklan}) = \$1.529$$

$$\text{revenue} - \text{cost} = \$1.094$$

↑  
Sebelum (+ 185%) ↑  
Sesudah (+ 351%)

Dengan adanya model, perusahaan bisa lebih memfokuskan terhadap pelanggan yang diprediksi akan meresponse meskipun tidak mengurangi skala cost tsb.



# Customer Segmentation Result

Potential Customer	Target Customer
Kecenderungan tidak Meresponse	Kecenderungan Meresponse
Rentang usia relative lebih muda dan income yang tinggi	Rentang usia lebih tua dan income rendah
Daily internet usage dan time spent pada website yang cenderung tinggi	Daily internet usage dan time spent pada website yang cenderung rendah
• • •	• • •

## Soft Selling Strategy

Karena diasumsikan segment potential customer lebih melek digital sehingga tidak mudah tergiur iklan



## Topic and Content Strategy

Kemudahan untuk dipahami dan topik yang mampu menarik minat segment kelompok pelanggan

[BACK TO CONTENT PAGE](#)