# "CLICKED ON AD"

## USING PYTHON & MACHINE LEARNING

LIWU
DEC 21,2024

# AGENDA

# INTRODUCTION & BUSINESS PROBLEM

**"How can we predict if a user will click on an advertisement?"**

- **Boosts CTR (Click-Through Rate)**: Encourages more user interactions with ads.

- **Enhances Targeting**: Focuses on likely-to-engage users.

- **Actionable Insights**: Refines marketing strategies effectively.

- **Increases ROI (Return on Investment)** : Maximizes profit from ad spend.

- **Better User Experience**: Shows relevant ads to users.

# DATA DESCRIPTION

- 2016 **Jan-July**

- **3,000** observations

- Duplicates: **31**

- Missing Values:

| | Number of Missing | Percentage(%) |
|---|---|---|
| **City** | 33 | 1.10 |
| **Area Income** | 31 | 1.03 |
| **Country** | 28 | 0.93 |
| **Age** | 24 | 0.80 |
| **Gender** | 23 | 0.77 |

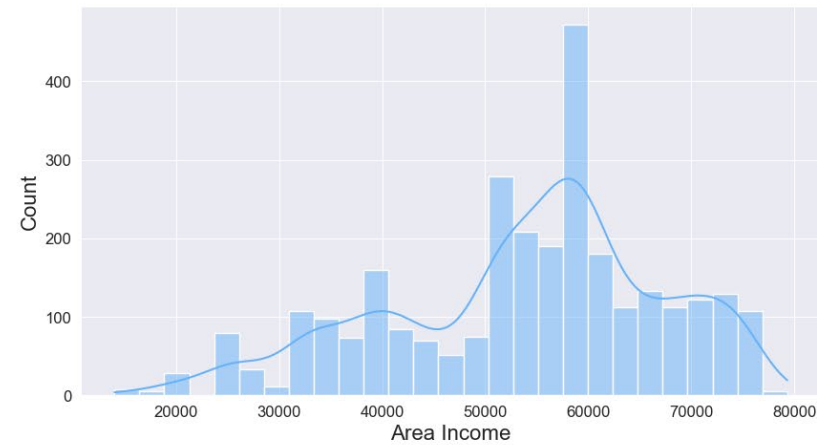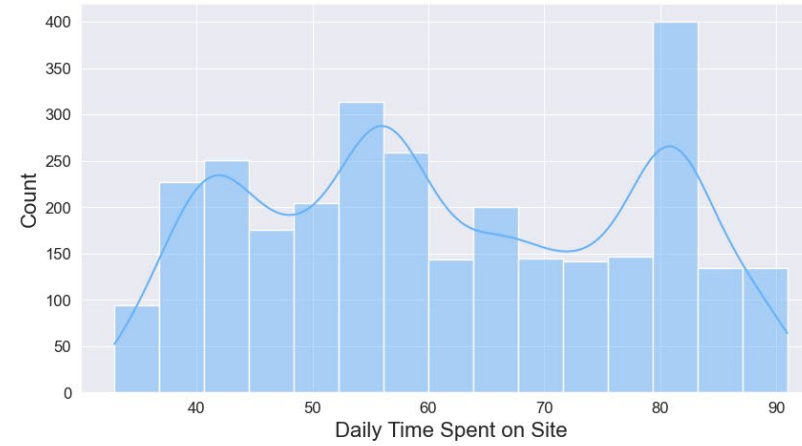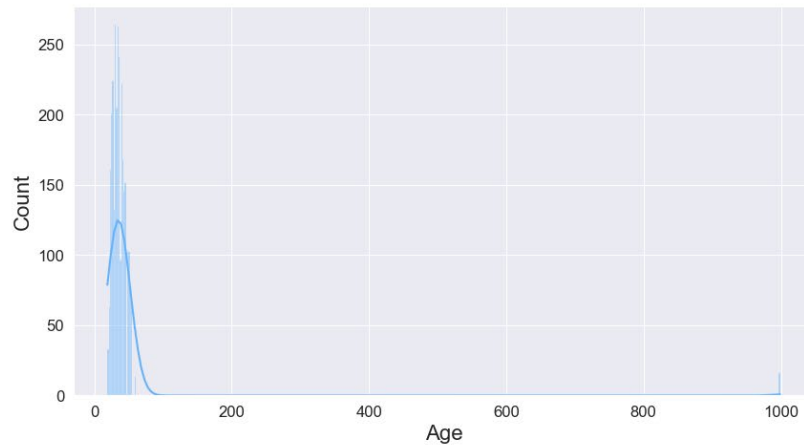| Target Variable | Clicked on Ad |
|---|---|
| **Numerical Variables** | • Timestamp (Datetime)<br>• Daily Time Spend on Site<br>• Daily Internet Usage<br>• Age<br>• Area Income |
| **Categorical Variables** | • Gender<br>• Ad Topic Line<br>• City<br>• Country |

# EXPLORATORY DATA ANALYSIS
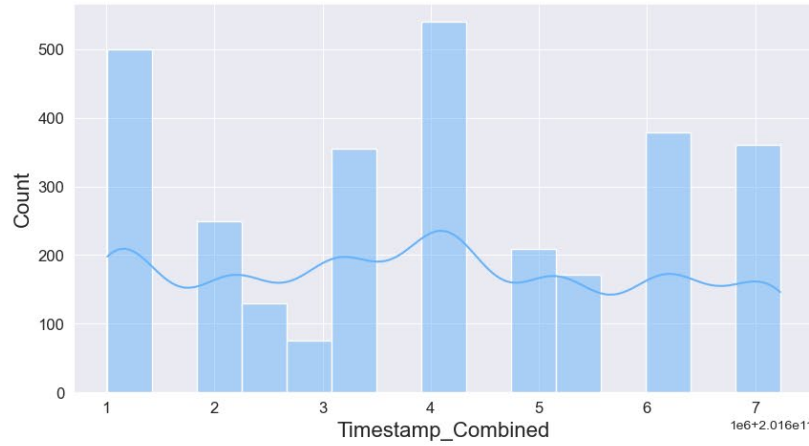
# "CLICKED ON AD"
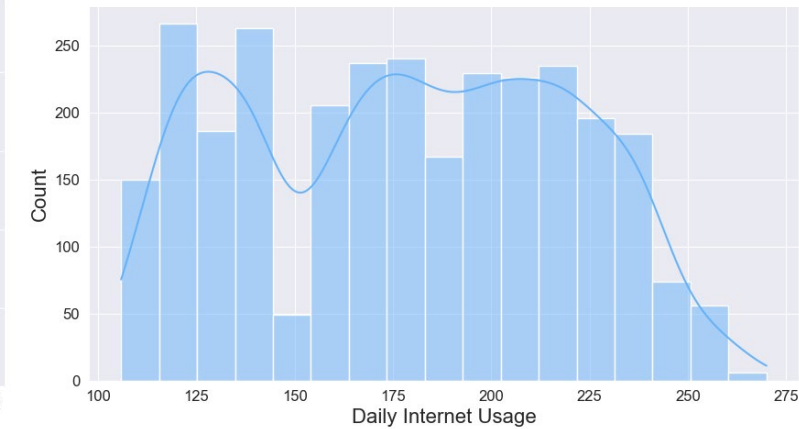


Not Clicked

50.6%

49.4%

Clicked
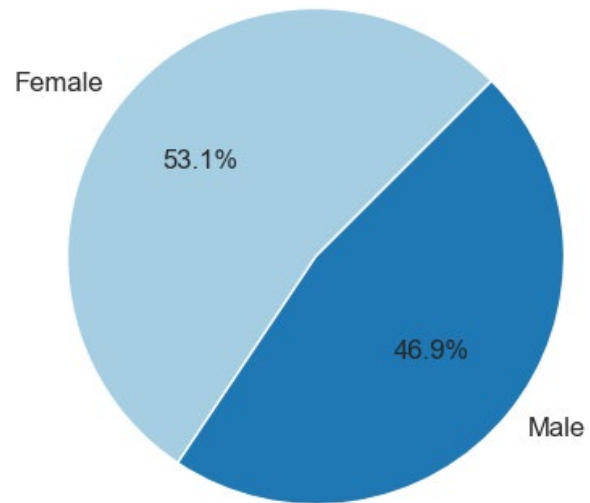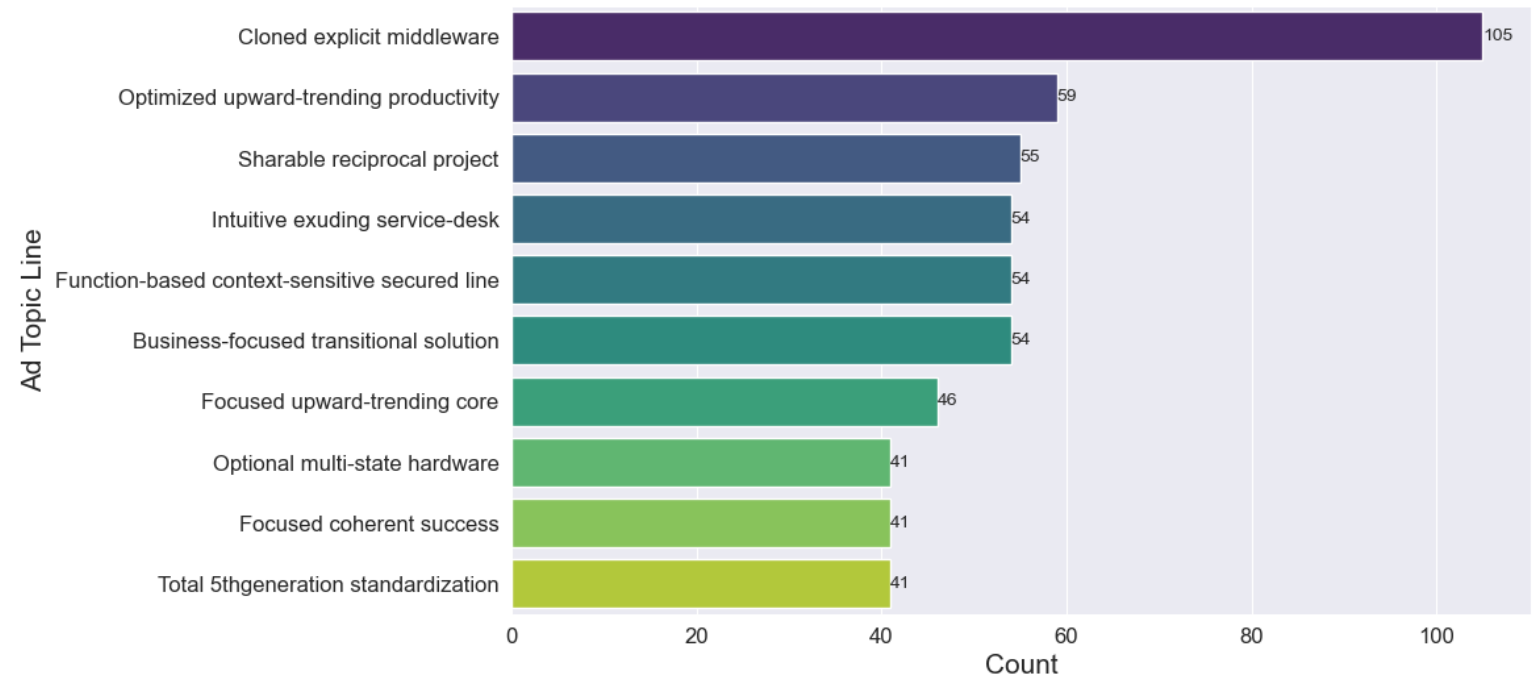
# NUMERICAL
# VARIABLES



## SKEWNESS

- Timestamp_Combined: **0.06**
- Daily Time Spend on Site: **0.09**
- Age: **13.17**
- Area Income: **-0.51**
- Daily Internet Usage: **-0.01**
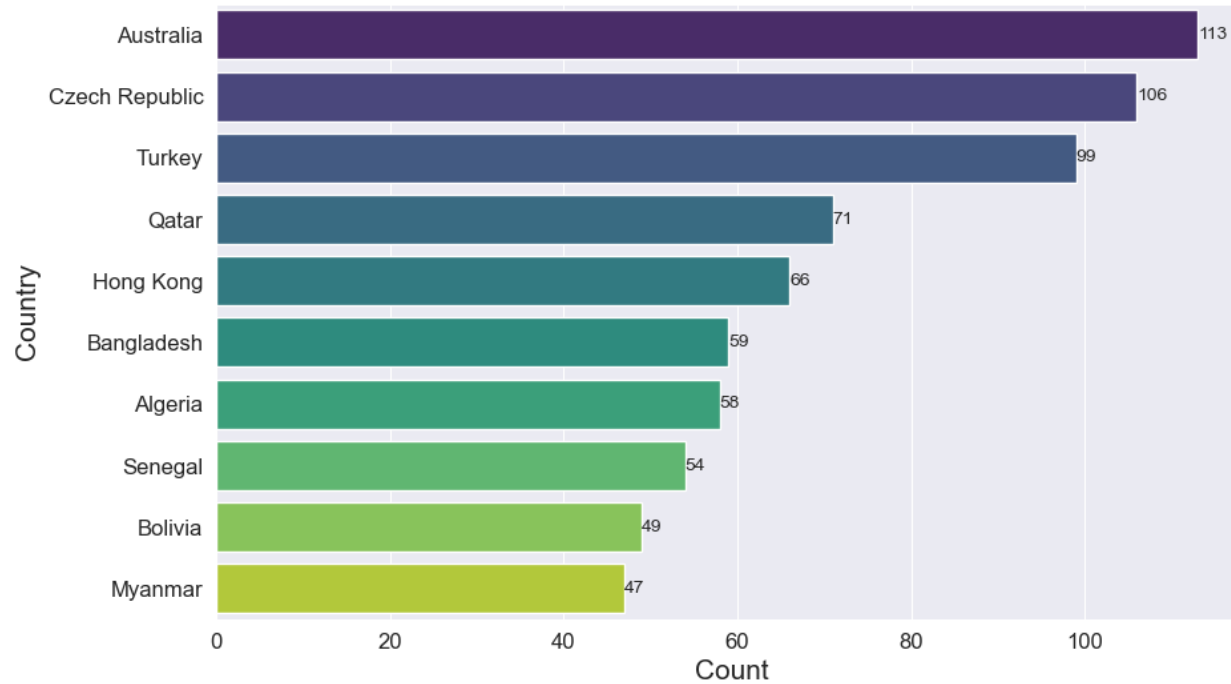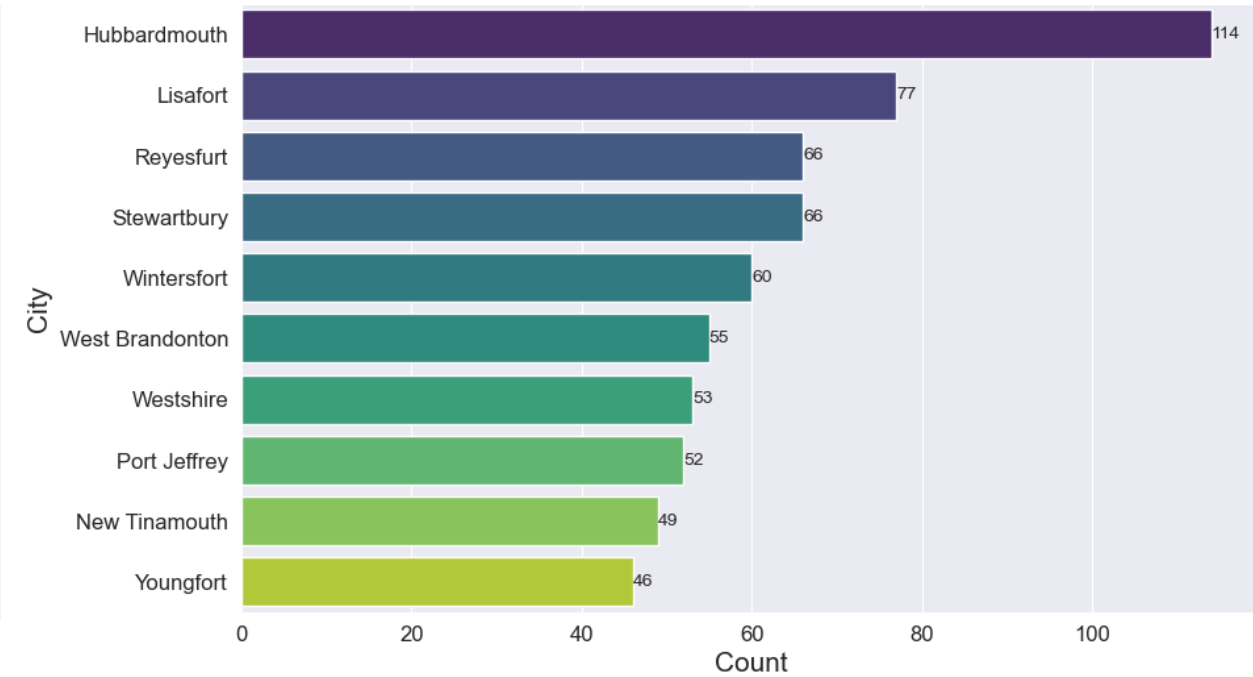
# CATEGORICAL
# VARIABLES



**Gender**

**Ad Topic Line** (417 levels)

# CATEGORICAL
# VARIABLES



**Country** (186 levels)

**City** (407 levels)

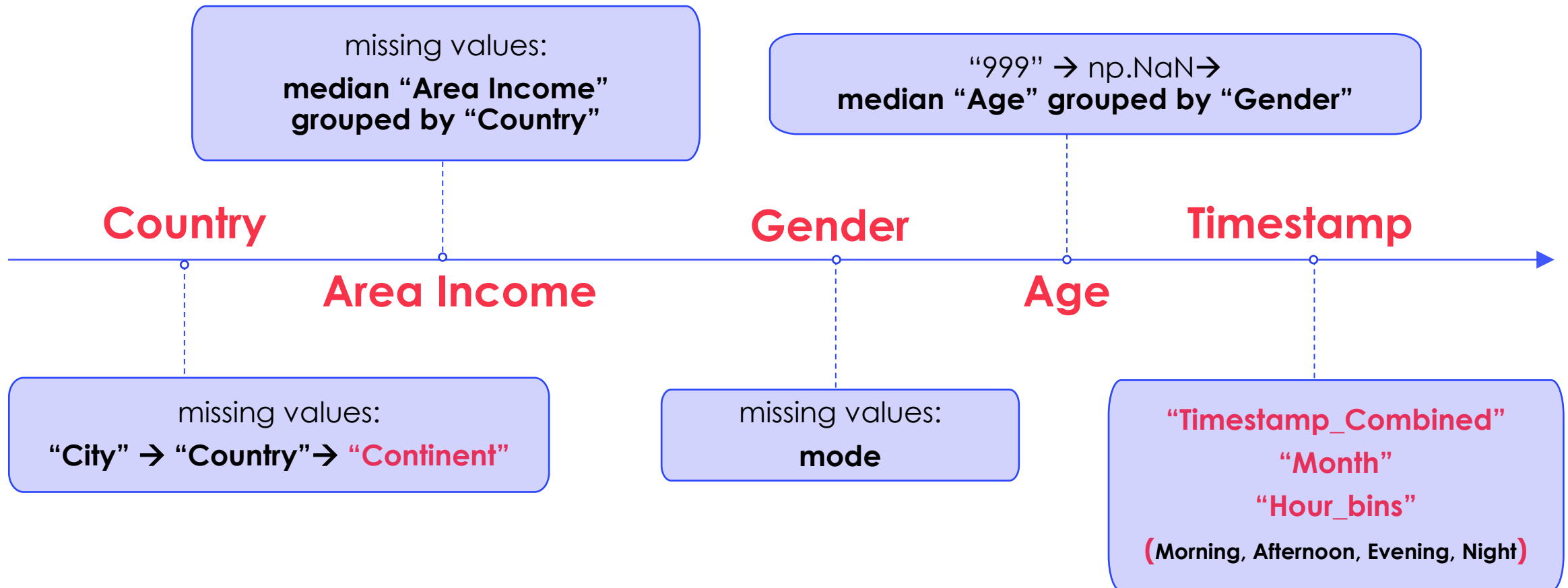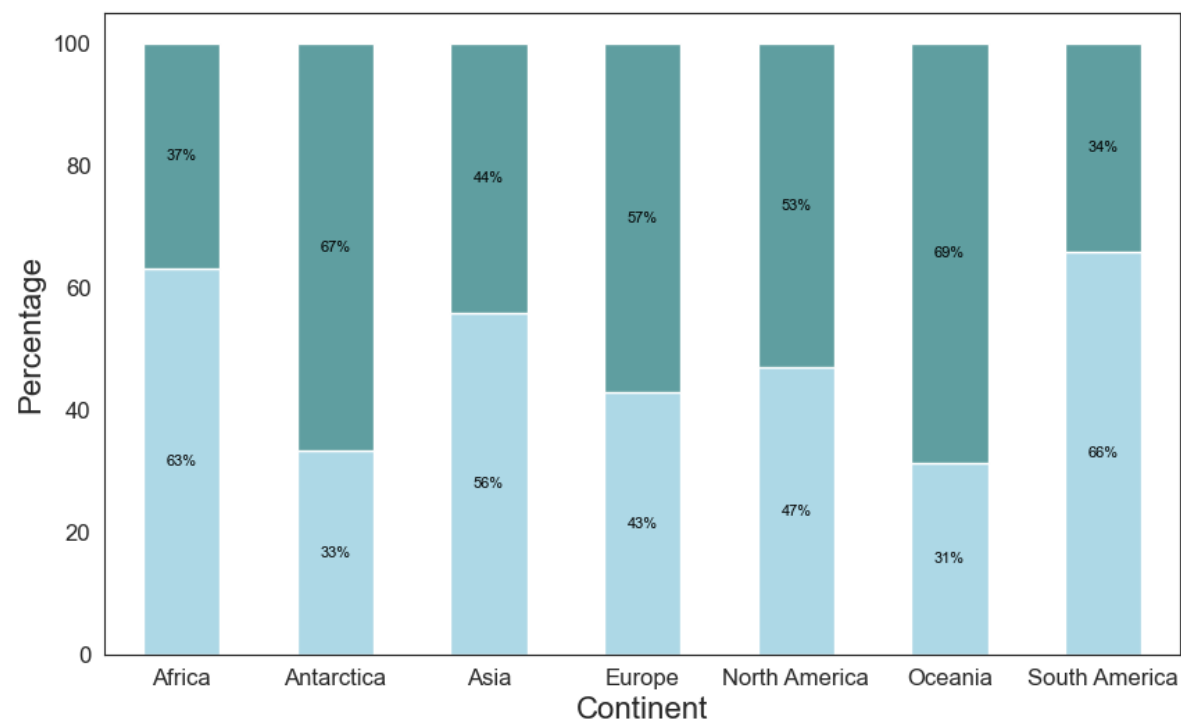# FEATURE ENGINEERING

# KEY STEPS

- Handle missing values

- Handle Outliers

- **Feature Creation** (Datetime: Month, Hour)

- **Feature Transformation** (Segmentation)

- **Feature Selection**
  - Statistical Tests (Chi-Square, T-test etc.)
  - Correlation, Multicollinearity (VIF)

- **One-Hot Encoding** for nominal categories (unordered)

- **Feature Scaling and Normalization** (Linear Models)

# FEATURE ENGINEERING

**missing values:**
**median "Area Income"**
**grouped by "Country"**

"999" → np.NaN→
**median "Age" grouped by "Gender"**

**Country**

**Gender**

**Timestamp**

**Area Income**

**Age**

missing values:
**"City" → "Country"→ "Continent"**

missing values:
**mode**

**"Timestamp_Combined"**
**"Month"**
**"Hour_bins"**
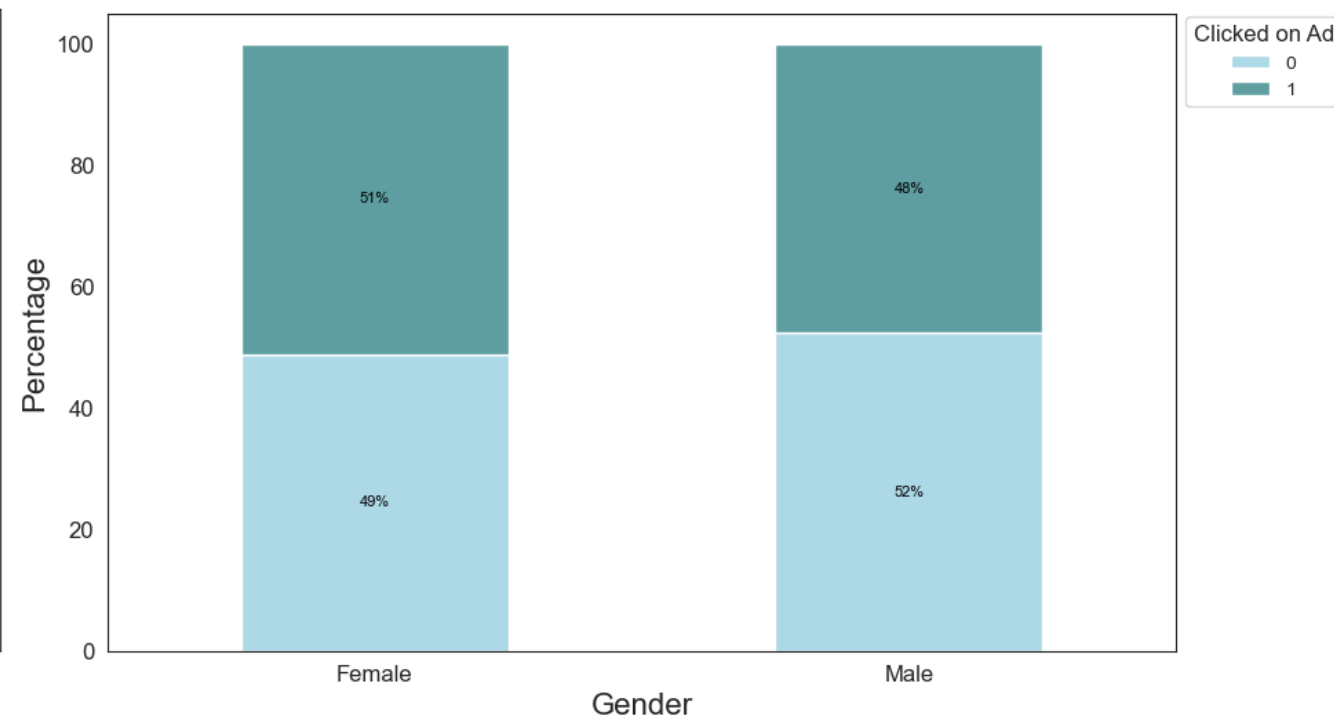**(Morning, Afternoon, Evening, Night)**

# "CLICKED ON AD" VS CATEGORICAL VARIABLES
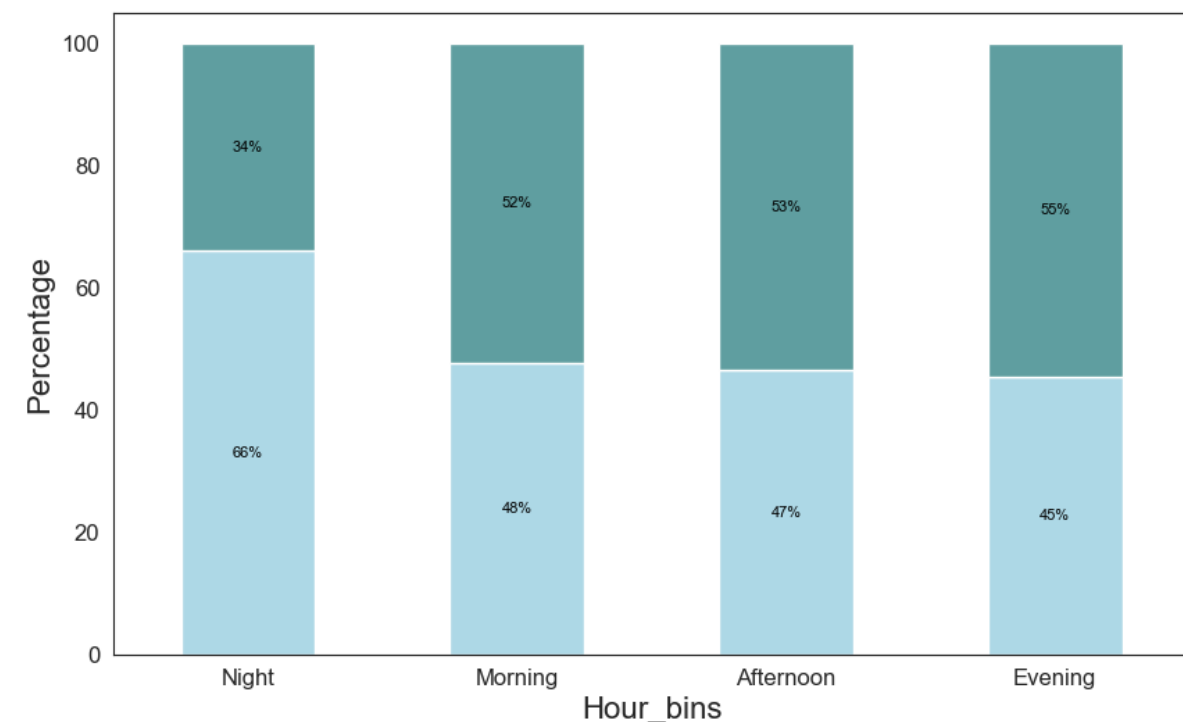
## Chi-Square test
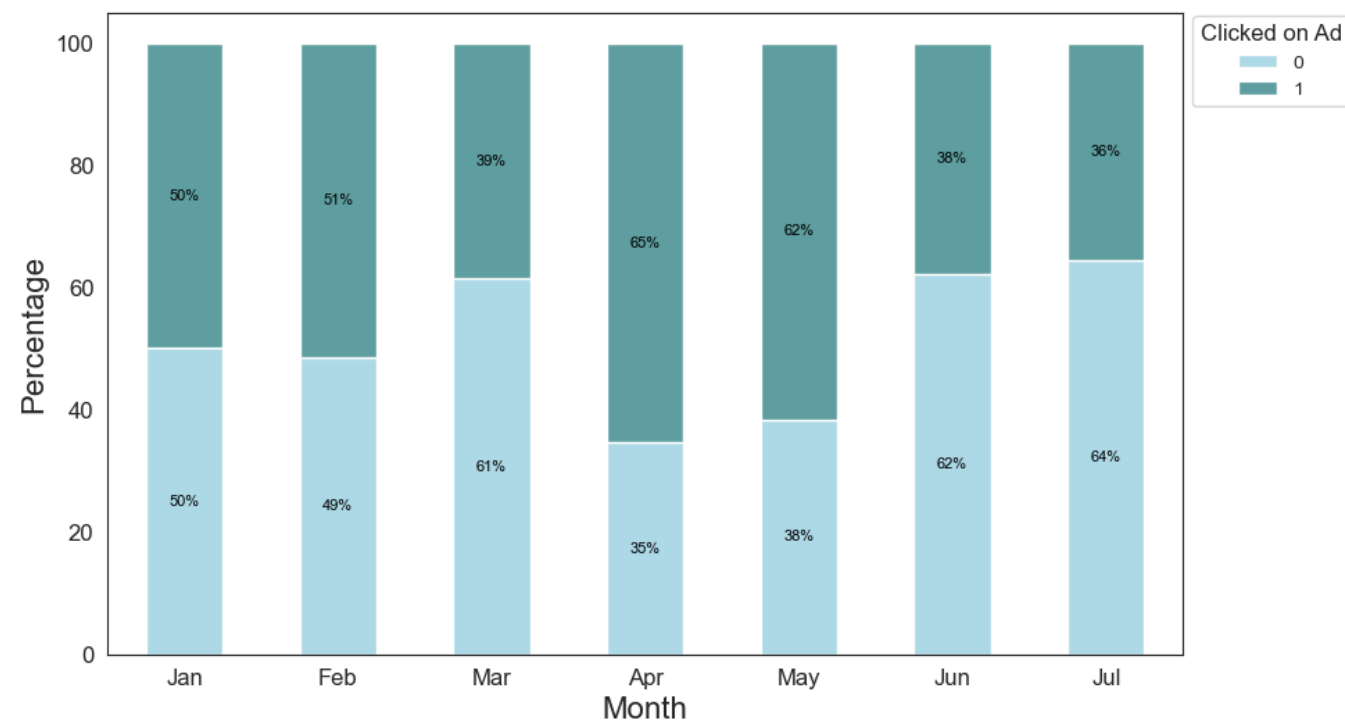


**Continent - Association**

**Gender - No Association**

# "CLICKED ON AD" VS CATEGORICAL VARIABLES
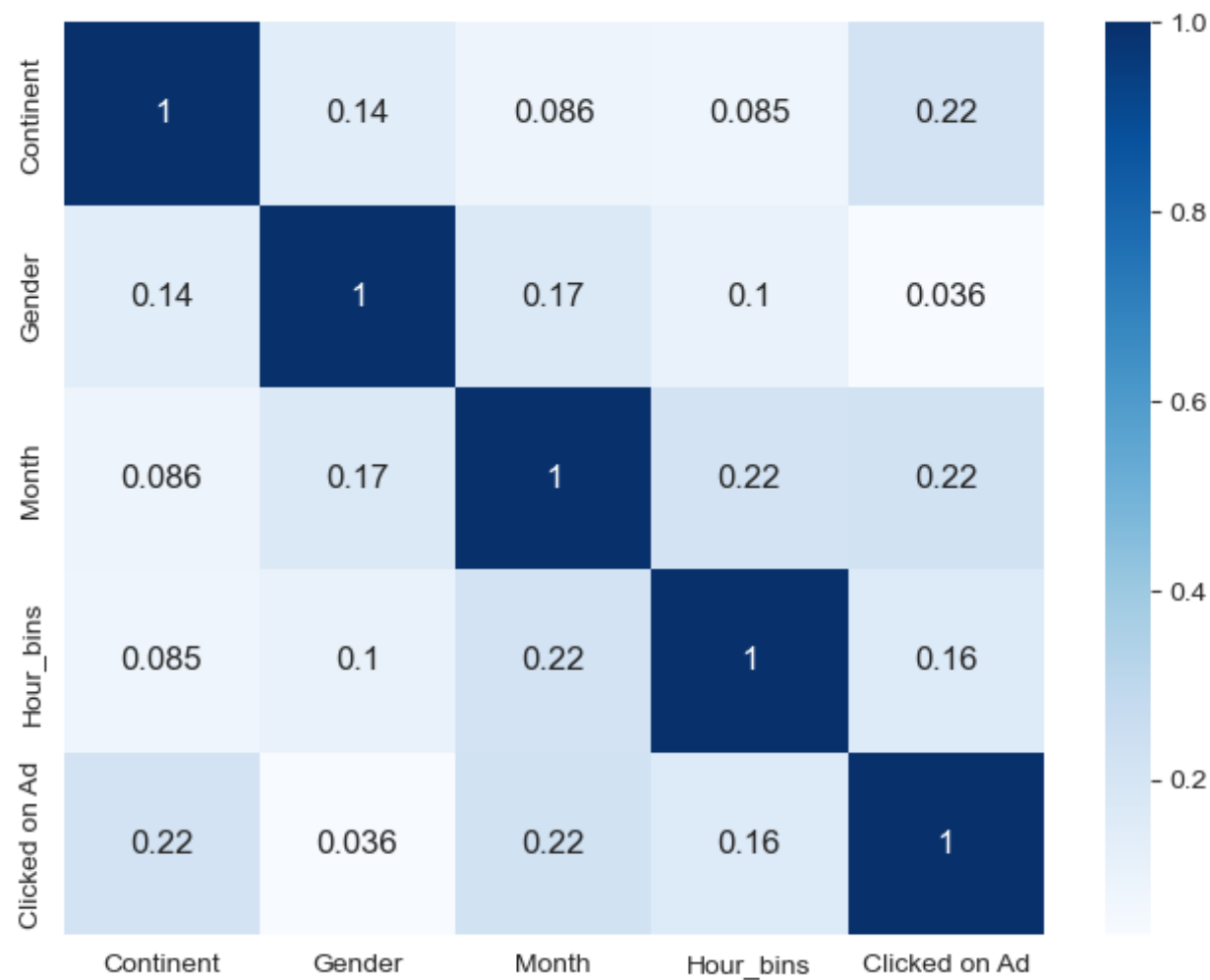
## Chi-Square test



**Hour_bins - Association**

**Month - Association**
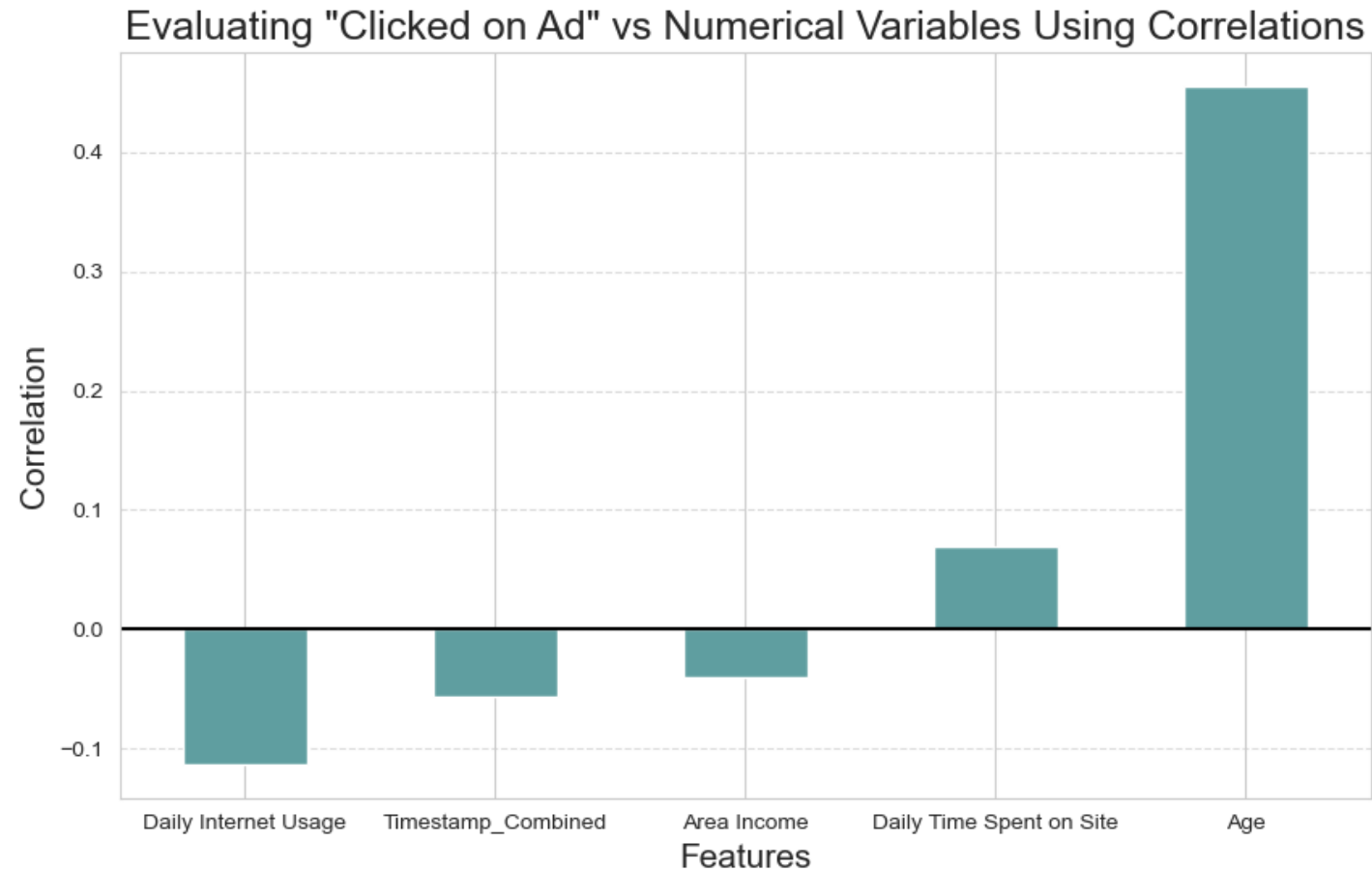
# "CLICKED ON AD" VS CATEGORICAL VARIABLES **Cramer's V**

# "CLICKED ON AD" VS
# NUMERICAL VARIABLES

## T-test



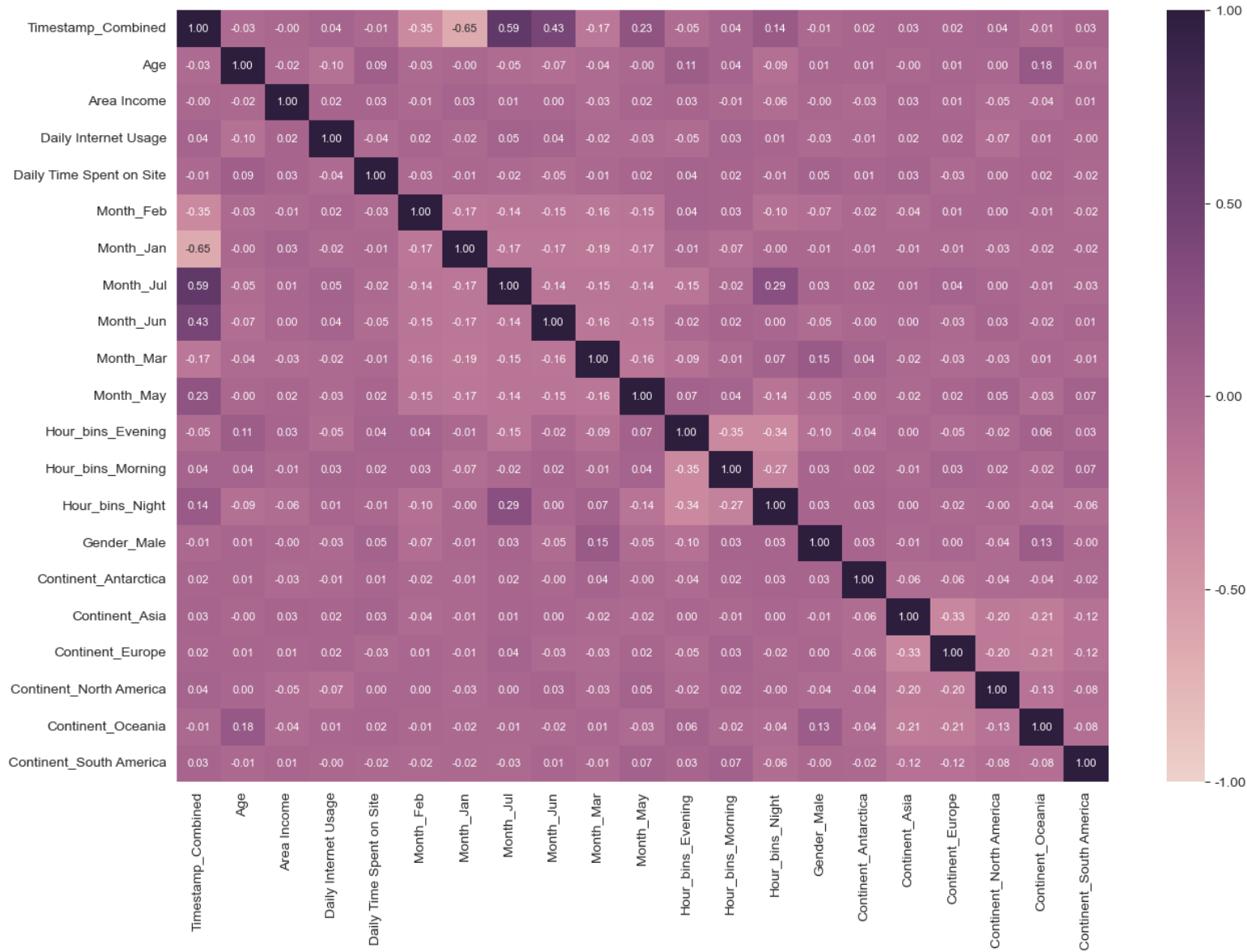Evaluating "Clicked on Ad" vs Numerical Variables Using Correlations

**"Age" :** Strongest Positive Correlation

# COLLINEARITY

# MULTICOLLINEARITY

## Linear Model
- Remove Features (VIF >10)
  - Unstable Coefficients
  - Reduced Interpretability
- Scale the data

## Tree-Based Models
- Keep All Features
- Scale not necessary

**Timestamp_Combined**
(VIF 85.83)

↓

**Age**
(VIF 15.48)

↓

**Daily Internet Usage**
(VIF 14.42)

↓

**Area Income**
(VIF 11.35)

| Feature | VIF (<=10) |
|---|---|
| Daily Time Spent on Site | 8.04 |
| Month_Feb | 1.49 |
| Month_Jan | 1.66 |
| Month_Jul | 1.65 |
| Month_Jun | 1.49 |
| Month_Mar | 1.64 |
| Month_May | 1.54 |
| Hour_bins_Evening | 2.05 |
| Hour_bins_Morning | 1.77 |
| Hour_bins_Night | 1.86 |
| Gender_Male | 1.96 |
| Continent_Antarctica | 1.05 |
| Continent_Asia | 1.95 |
| Continent_Europe | 1.90 |
| Continent_North America | 1.41 |
| Continent_Oceania | 1.50 |
| Continent_South America | 1.19 |

# MODEL DEVELOPMENT

**Logistic Regression**

**Decision Tree**
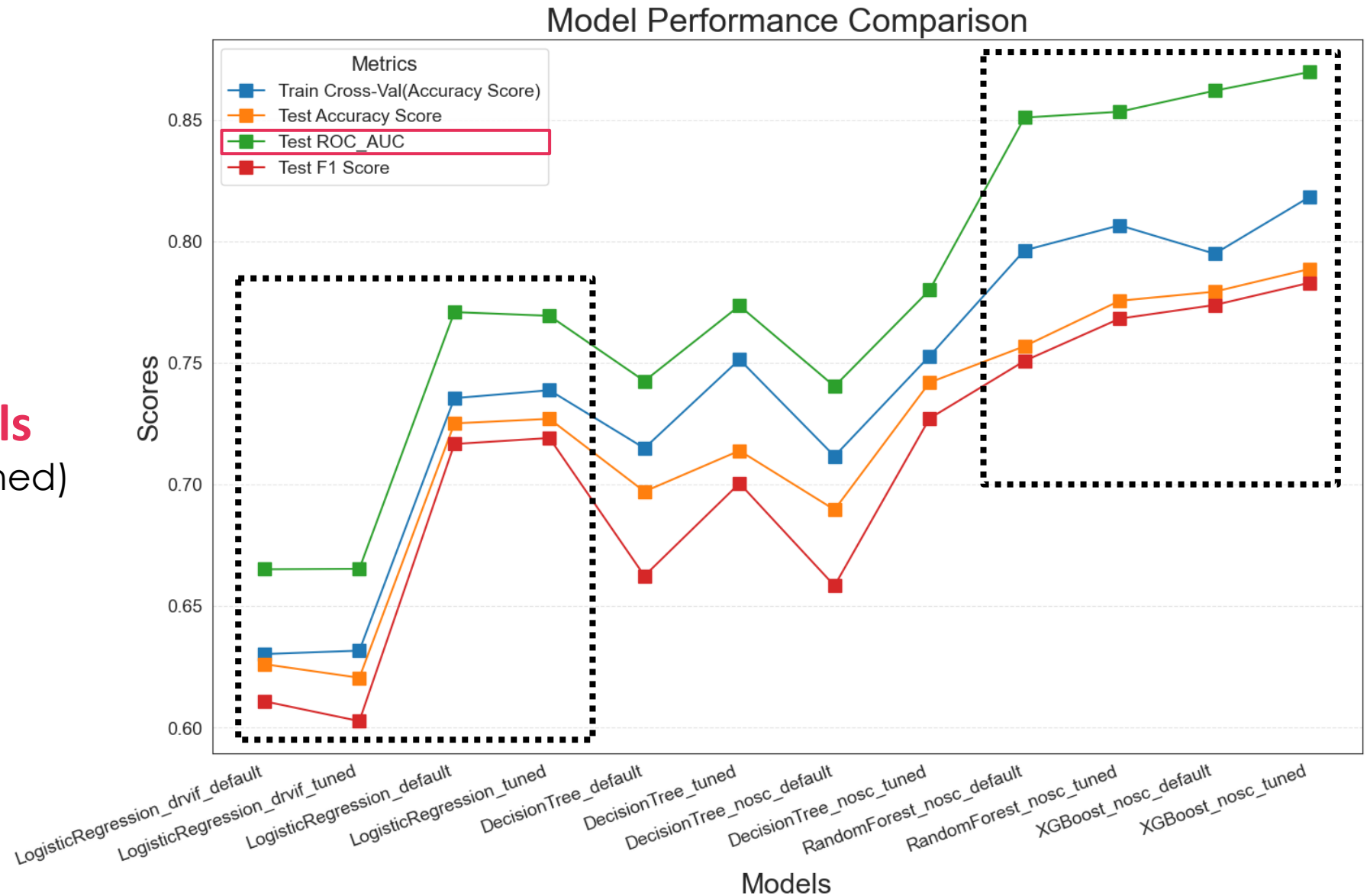
**Random Forest**

**XGBoost**

# ALL MODELS
Model Selection

## Linear Model

▪ Logistic Regression (Baseline)

## Tree-Based Models
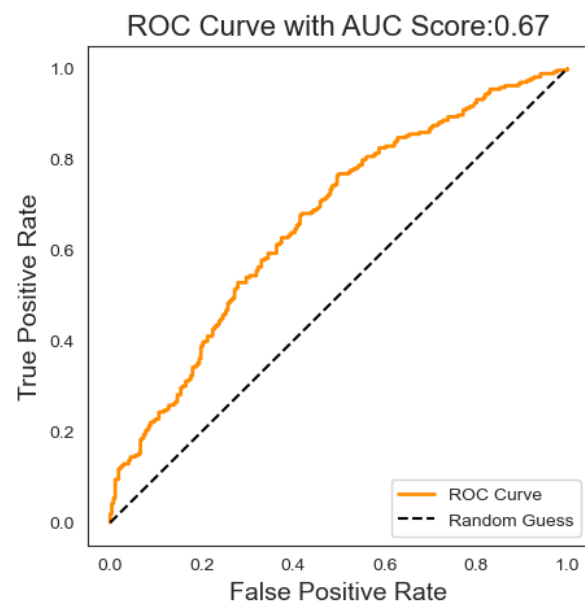
▪ Random Forest (Tuned)
▪ XGBoost (Tuned)



Model Performance Comparison

# ALL MODELS

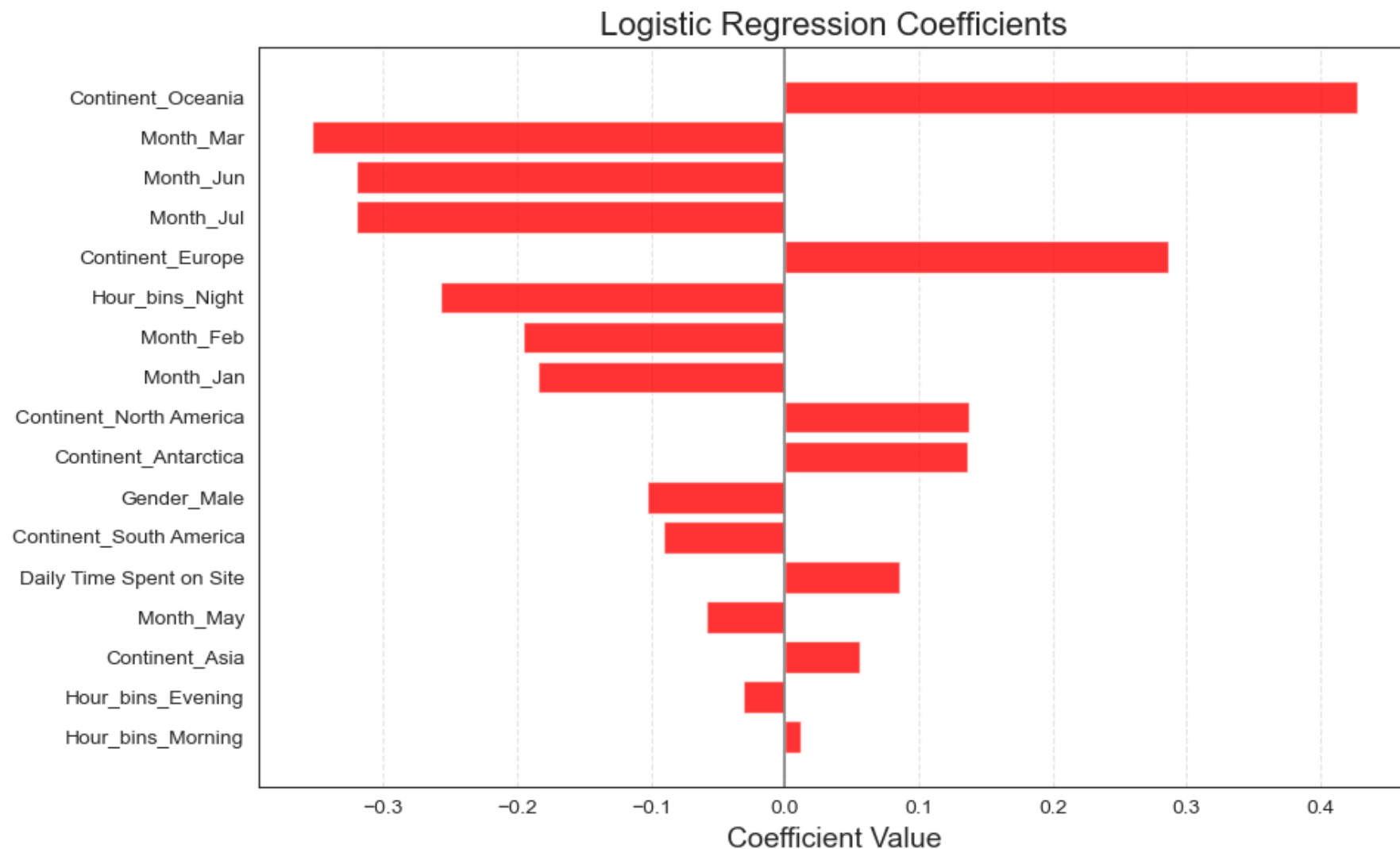| | Model | Train Cross-Val(Accuracy Score) | Test Accuracy Score | Train Cross-Val(ROC_AUC) | Test ROC_AUC | Train Cross-Val(F1 Score) | Test F1 Score |
|---|---|---|---|---|---|---|---|
| 0 | LogisticRegression_drvif_default | 0.630325 | 0.626168 | 0.673368 | 0.665218 | 0.615050 | 0.610895 |
| 1 | LogisticRegression_drvif_tuned | 0.631729 | 0.620561 | 0.673509 | 0.665399 | 0.617407 | 0.602740 |
| 2 | LogisticRegression_default | 0.735616 | 0.725234 | 0.801229 | 0.771038 | 0.723807 | 0.716763 |
| 3 | LogisticRegression_tuned | 0.738892 | 0.727103 | 0.801315 | 0.769515 | 0.727031 | 0.719231 |
| 4 | DecisionTree_default | 0.715027 | 0.697196 | 0.775662 | 0.742668 | 0.688912 | 0.662500 |
| 5 | DecisionTree_tuned | 0.751532 | 0.714019 | 0.800933 | 0.773631 | 0.734099 | 0.700587 |
| 6 | DecisionTree_nosc_default | 0.711751 | 0.689720 | 0.772482 | 0.740487 | 0.683956 | 0.658436 |
| 7 | DecisionTree_nosc_tuned | 0.752937 | 0.742056 | 0.800080 | 0.780097 | 0.740189 | 0.727273 |
| 8 | RandomForest_nosc_default | 0.796453 | 0.757009 | 0.875372 | 0.851061 | 0.789449 | 0.750958 |
| 9 | RandomForest_nosc_tuned | 0.806745 | 0.775701 | 0.877335 | 0.853479 | 0.800553 | 0.768340 |
| 10 | XGBoost_nosc_default | 0.795037 | 0.779439 | 0.881843 | 0.862125 | 0.792159 | 0.773946 |
| 11 | XGBoost_nosc_tuned | 0.818448 | 0.788785 | 0.892335 | 0.869842 | 0.812669 | 0.783109 |

# LOGISTIC REGRESSION

## Features VIF<=10

- Train Cross-Val Accuracy: **0.63**
- Test Accuracy: **0.62**
- Test ROC_AUC: **0.66**
- Test F1 Score: **0.60**



ROC Curve with AUC Score:0.67

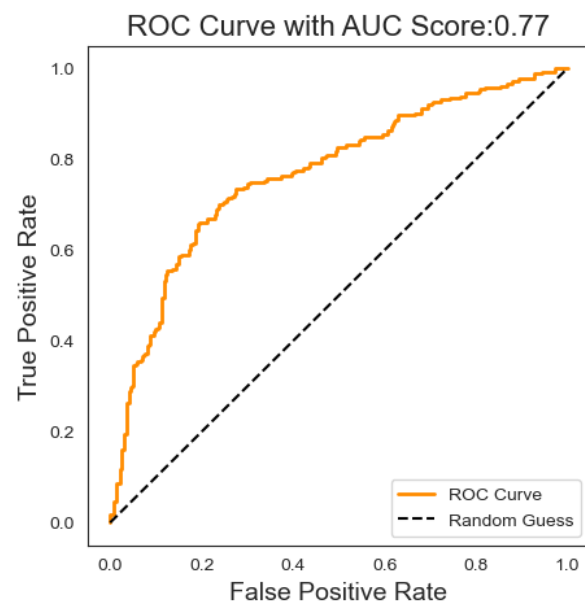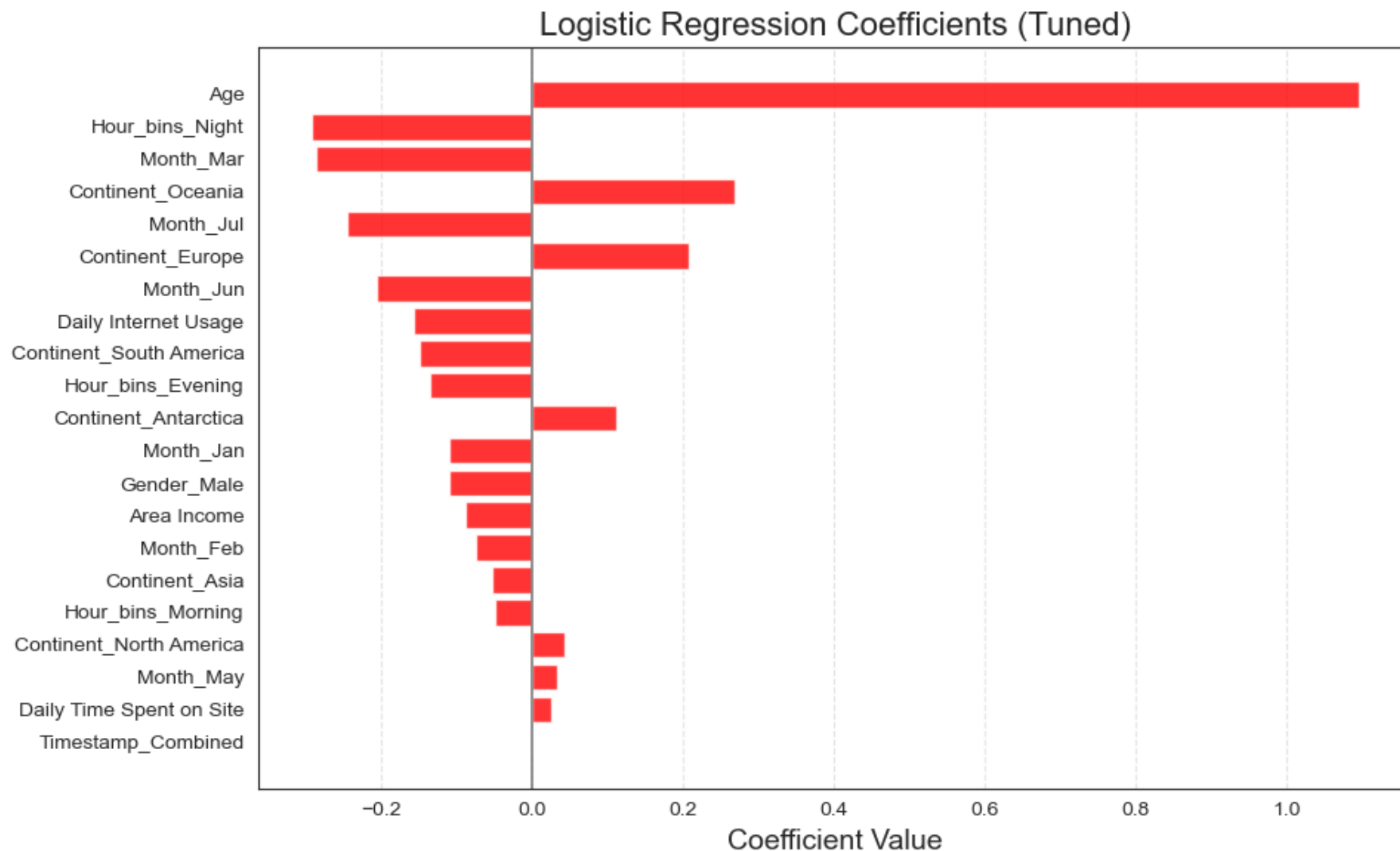

Logistic Regression Coefficients

# LOGISTIC REGRESSION

## All Columns

- Train Cross-Val Accuracy: **0.73**
- Test Accuracy: **0.72**
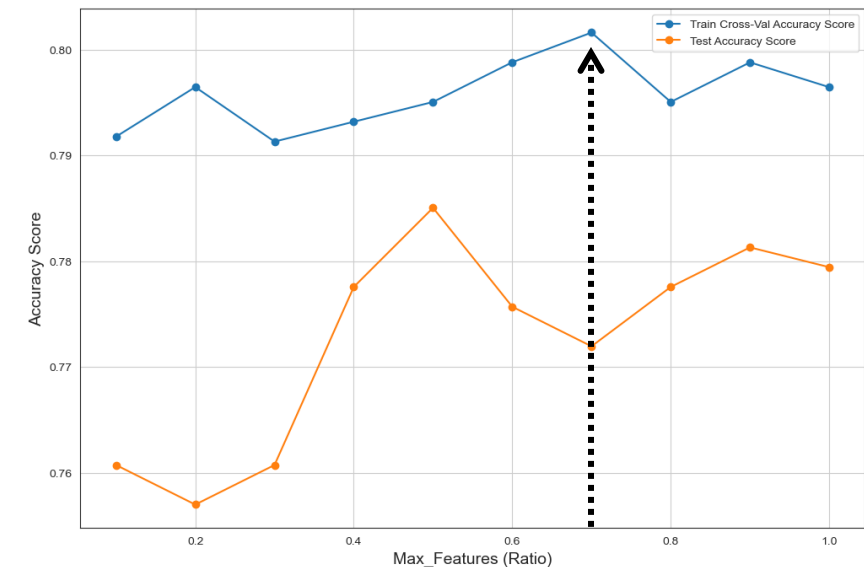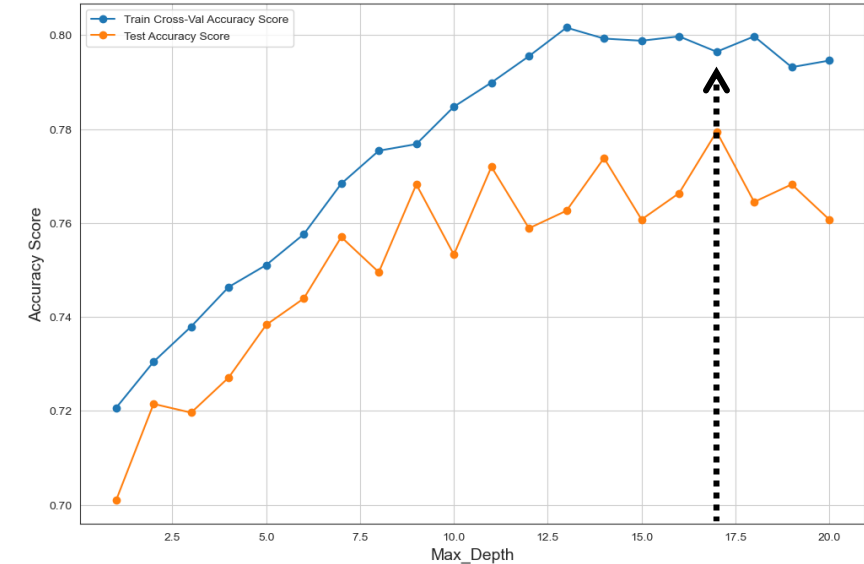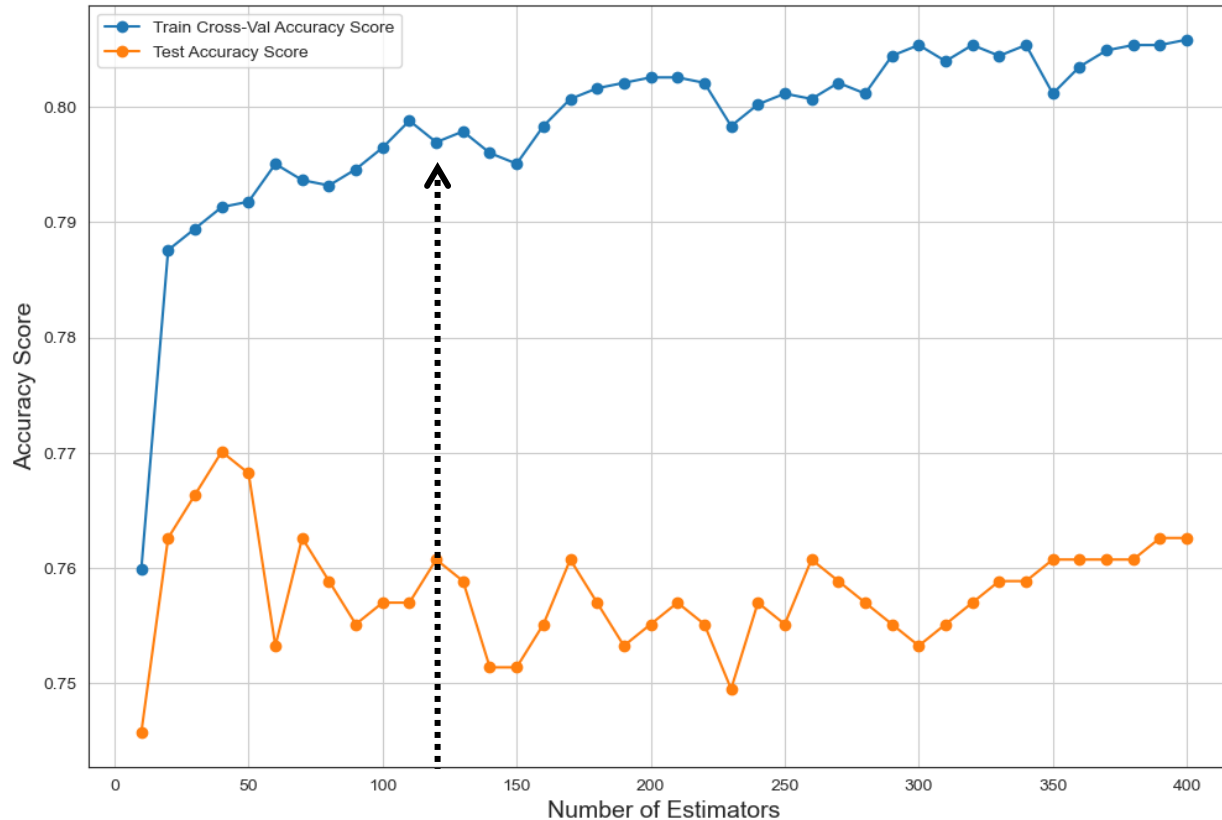- Test ROC_AUC: **0.76**
- Test F1 Score: **0.71**



ROC Curve with AUC Score:0.77



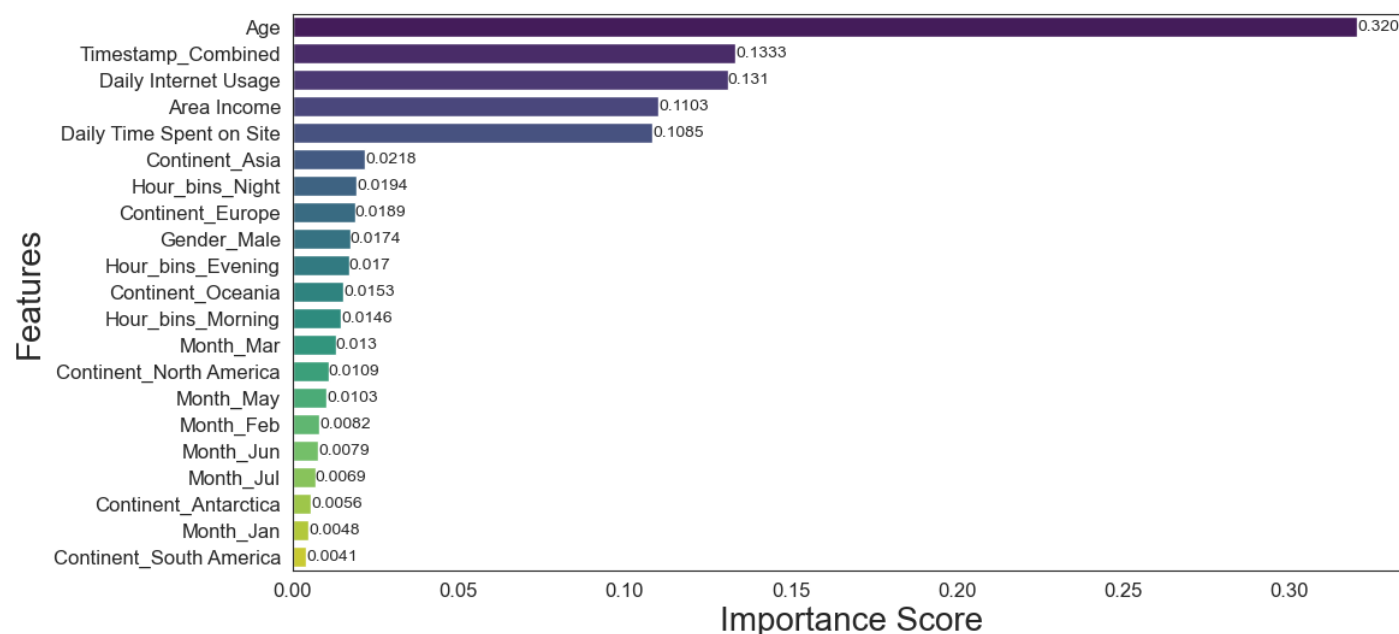Logistic Regression Coefficients (Tuned)
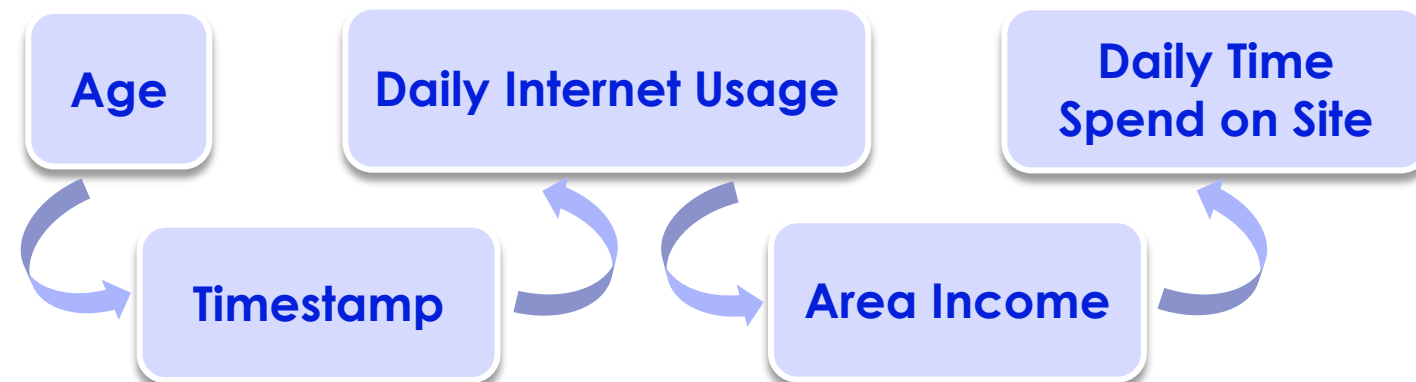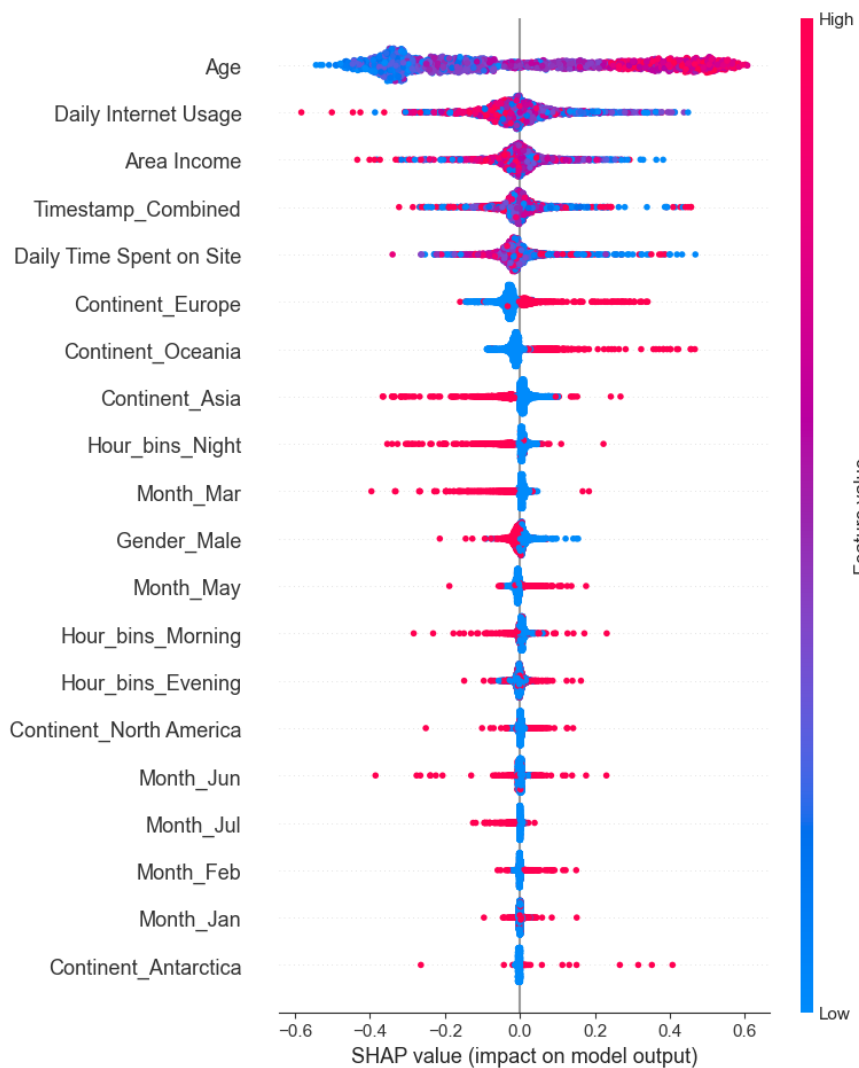
# TREE-BASED MODELS

## Random Forest

- Train Cross-Val Accuracy: **0.80**
- Test Accuracy: **0.77**
- Test ROC_AUC: **0.85**
- Test F1 Score: **0.76**

- n_estimators: 120
- max_depth: 17
- max_features: 0.7
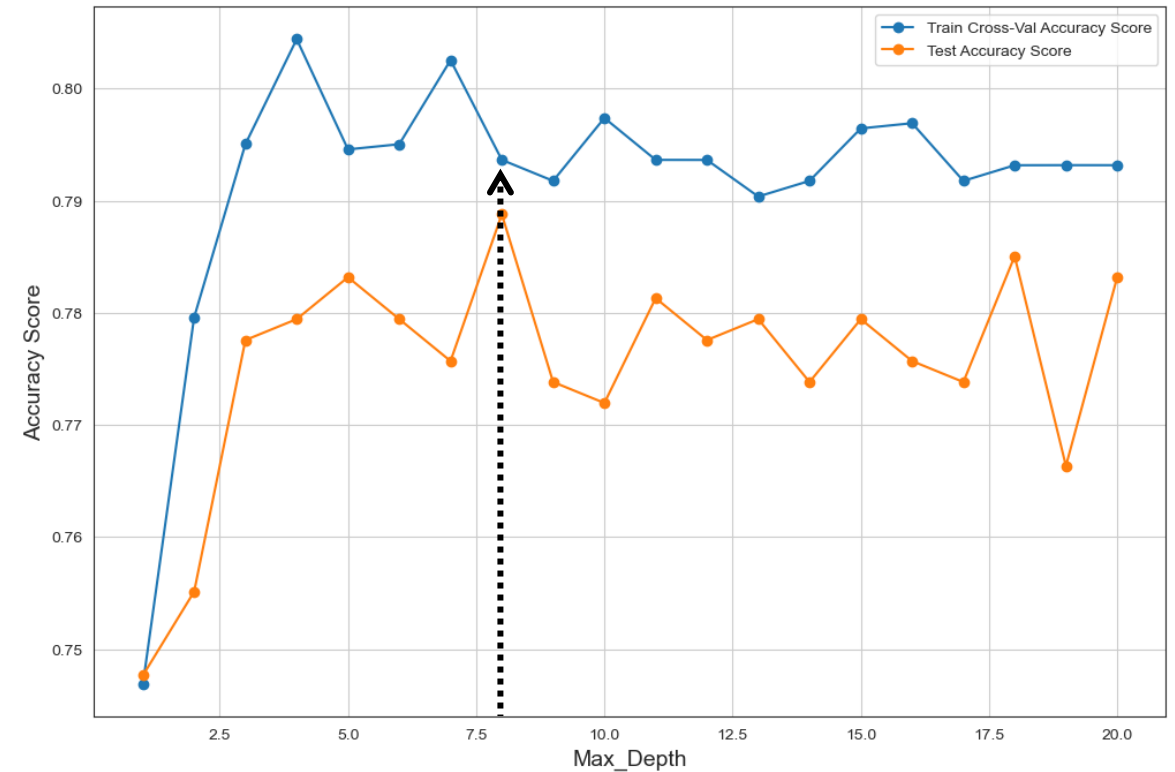
# TREE-BASED MODELS

**Random Forest**





Age → Timestamp → Daily Internet Usage → Area Income → Daily Time Spend on Site

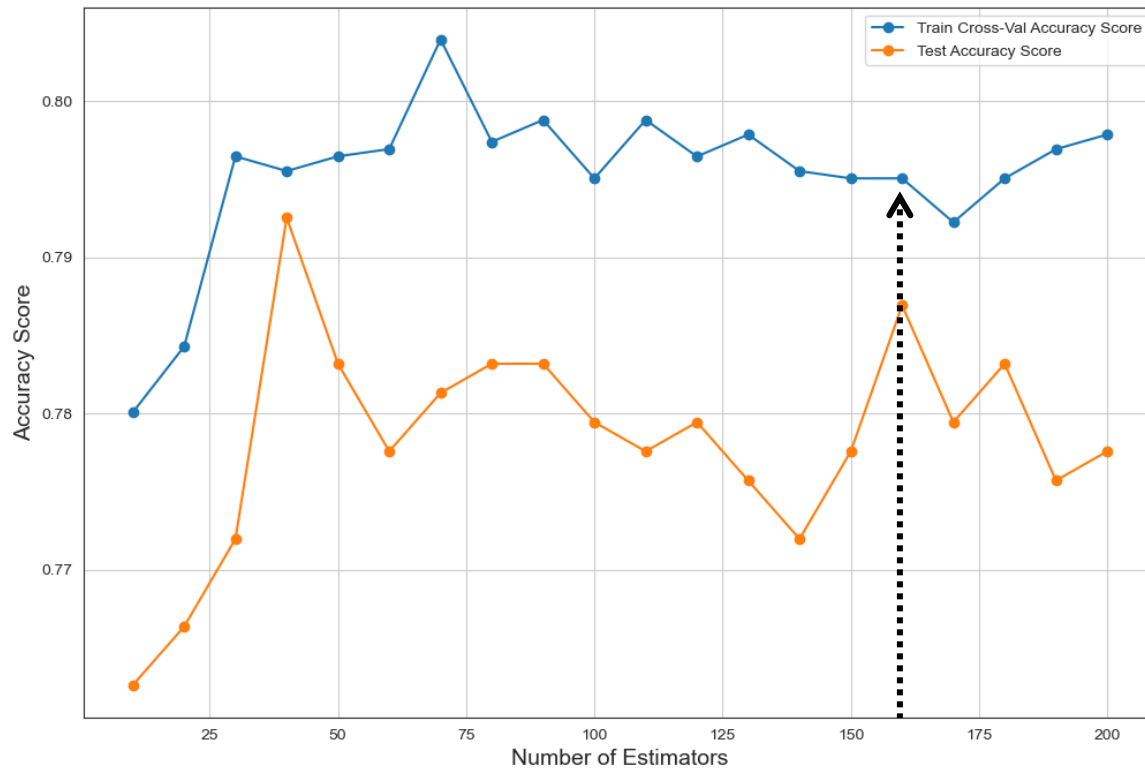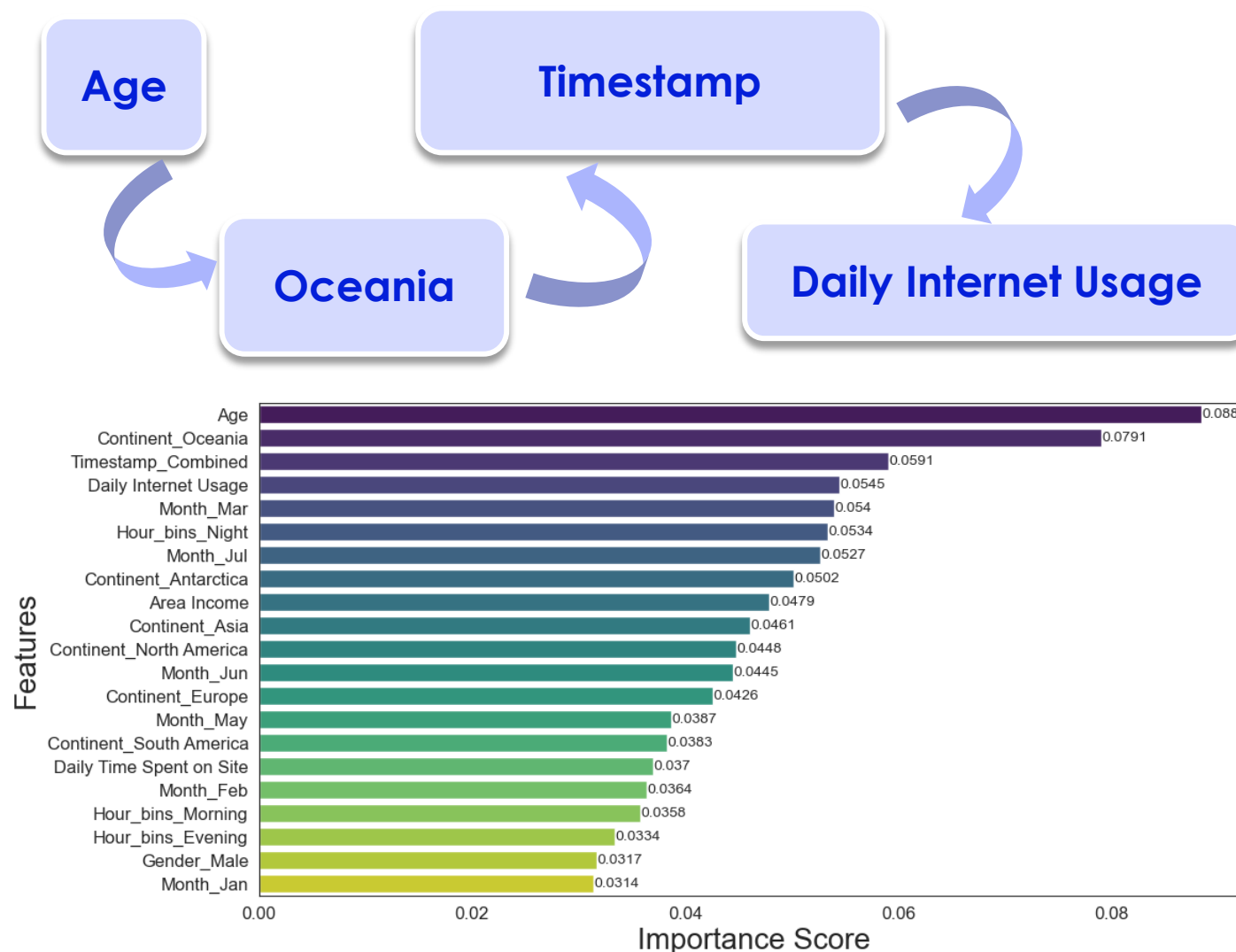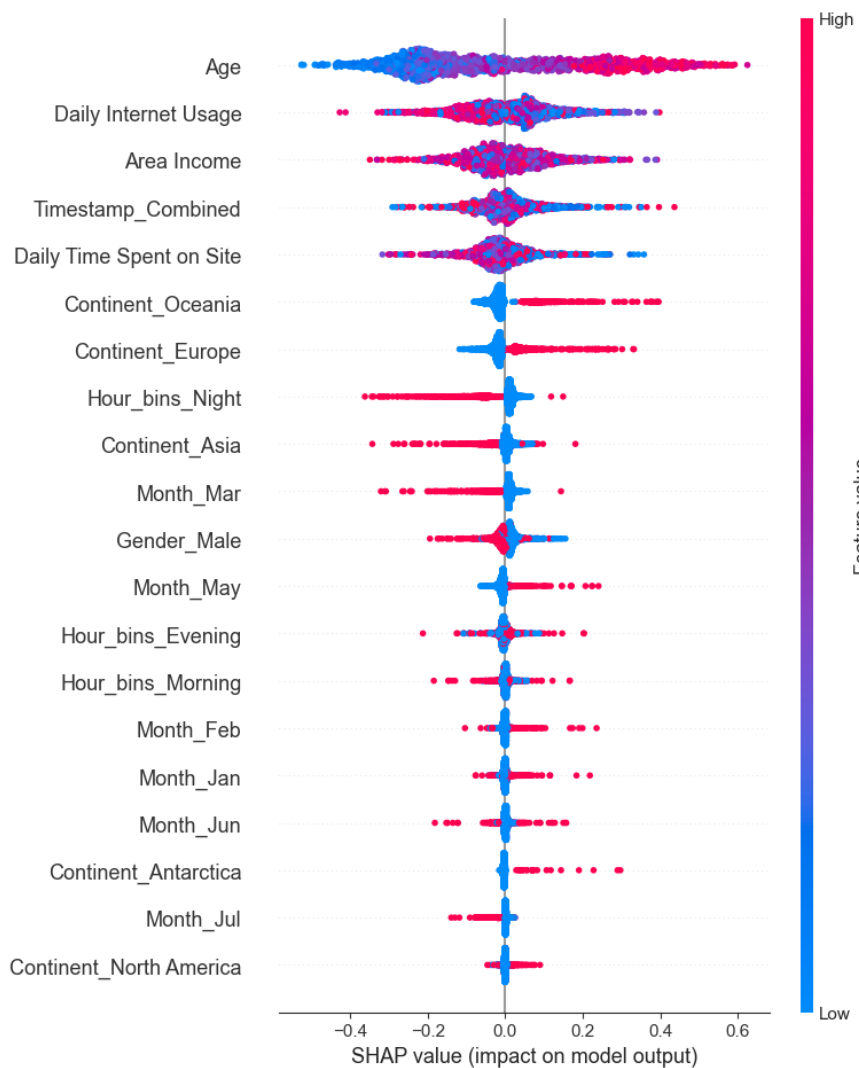| Feature | Importance Score |
| --- | --- |
| Age | 0.3206 |
| Timestamp_Combined | 0.1333 |
| Daily Internet Usage | 0.131 |
| Area Income | 0.1103 |
| Daily Time Spent on Site | 0.1085 |
| Continent_Asia | 0.0218 |
| Hour_bins_Night | 0.0194 |
| Continent_Europe | 0.0189 |
| Gender_Male | 0.0174 |
| Hour_bins_Evening | 0.017 |
| Continent_Oceania | 0.0153 |
| Hour_bins_Morning | 0.0146 |
| Month_Mar | 0.013 |
| Continent_North America | 0.0109 |
| Month_May | 0.0103 |
| Month_Feb | 0.0082 |
| Month_Jun | 0.0079 |
| Month_Jul | 0.0069 |
| Continent_Antarctica | 0.0056 |
| Month_Jan | 0.0048 |
| Continent_South America | 0.0041 |

# TREE-BASED MODELS

## XGBoost

- Train Cross-Val Accuracy: **0.81**
- Test Accuracy: **0.78**
- Test ROC_AUC: **0.86**
- Test F1 Score: **0.78**

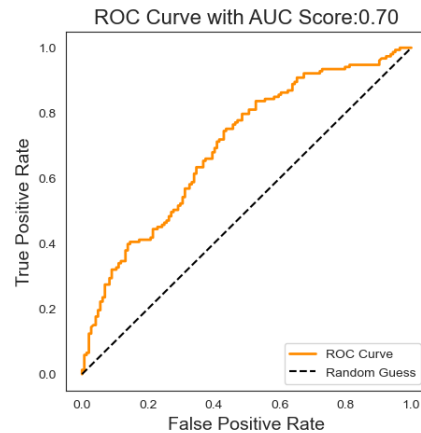- n_estimators: 160
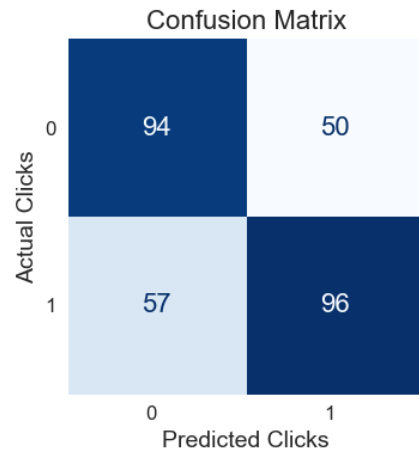- max_depth: 8

# TREE-BASED MODELS

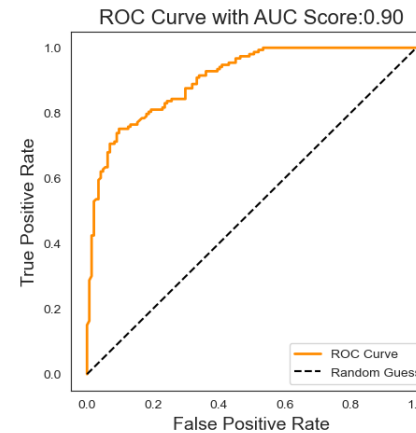**XGBoost**
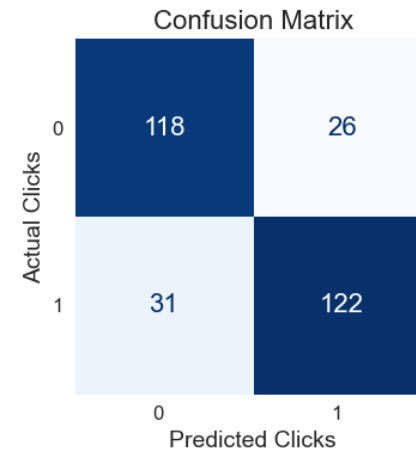
# MODEL EVALUATION

(UNSEEN DATA)

**Logistic Regression (Baseline)**

Accuracy: **0.63**
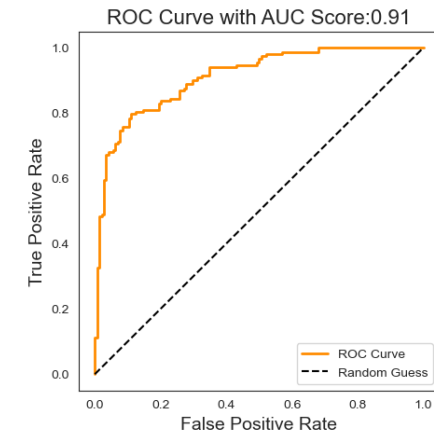**ROC_AUC: 0.70**
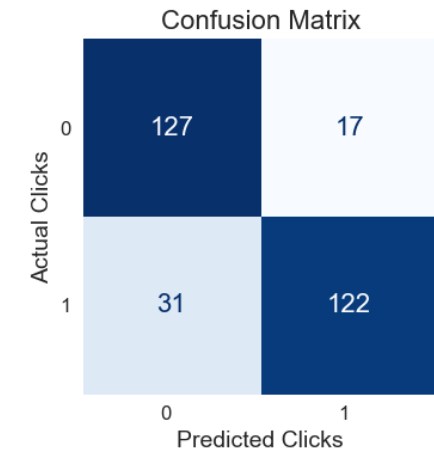F1 Score: **0.64**



**Random Forest (Tuned)**

Accuracy: **0.80**
**ROC_AUC: 0.90**
F1 Score: **0.81**



**XGBoost (Tuned)**

Accuracy: **0.83**
**ROC_AUC: 0.91**
F1 Score: **0.83**

# BUSINESS RECOMMENDATIONS

# MARKETING STRATEGIES

- **Target Older Audiences**
  Develop campaigns and messaging specifically targeted toward individuals in older age groups.

- **Focus on High Daily Internet Usage**
  Target these users through behavior-based segmentation and frequently visited platforms.

- **Optimize Ad Timing**
  Deliver ads during afternoon and evening hours when users are more likely to engage based on activity trends.

- **Prioritize High-Income Regions**
  Higher-income users show greater ad engagement. Target premium regions and promote high-value or luxury products.

- **Customize Campaigns by Regions**
  Analyze and tailor ad strategies for different regions (e.g., Oceania & Europe").

# THANK YOU

Li Wu
Instructor: Reza Moosavi