

## # Clickbait-Peak (CTR Prediction for ConnectSphere Digital)

Predicting which users are likely to click ads using real web analytics data and logistic regression.

See 'documentation.md' for full modeling process, code snippets and EDA

### ## Quickstart

1. Clone repo and install dependencies:  
pip install -r requirements.txt
2. Run EDA and modeling scripts.
3. Full results, plots, and performance details in 'documentation.md'.

### ## License

MIT  
For questions: Nitish Sehgal / scholasticnit@gmail.com

### # EdGlobe Minor Project Submission

#### # Predictive Modeling for Click-Through Rate Optimization at ConnectSphere Digital

##### # Problem Statement

ConnectSphere Digital manages many online ad campaigns but struggles with low ROI as ads are often shown to users unlikely to click. The goal is to create a data-driven model to identify and prioritize users most likely to engage, improving campaign effectiveness and client satisfaction.

##### # Project Goal

Develop and evaluate a logistic regression model using historical user data—Age, Area Income, Daily Time Spent on Site, Daily Internet Usage, and demographics like Male—to classify users as 'likely to click' or 'unlikely to click'. Evaluate performance with accuracy, precision, and recall.

##### # Business Objective

Enable ConnectSphere Digital to allocate ad budget more efficiently by targeting high-propensity users, thereby increasing overall Click-Through Rate (CTR) and boosting competitive advantage.

##### # Link for Dataset

[https://drive.google.com/file/d/1J5yqeeKVRpBZXrNxxwKN\\_uGG4pUybLView](https://drive.google.com/file/d/1J5yqeeKVRpBZXrNxxwKN_uGG4pUybLView)

### # Dataset Overview

- 1000 samples with 10 features: usage patterns, demographics, ad topic, location, timestamp, and clicks (target).
- No missing values detected.
- Balanced target variable: 50% clicked, 50% not clicked.

### # Exploratory Data Analysis (EDA) Outputs and Inferences

- \*\*Top 10 Cities and Countries by Frequency\*\* show a diverse user base; no dominant location bias.
- \*\*Unique Ad Topics\*\* count = 1000, meaning each row has a unique topic, thus dropped for modeling to avoid noise.
- \*\*Class balance\*\* is perfect (50%-50%), aiding unbiased model training.

### # Model Performance

- \*\*Accuracy: 95%\*\* — Model predicts clicks correctly 95% of the time.
- \*\*Precision: 95.5%\*\* — When predicted "likely to click," correct 95.5% of the time, minimizing wasted ads.
- \*\*Recall: 95.5%\*\* — Identifies 95.5% of all actual clickers, reducing missed opportunities.

### ### Confusion Matrix

`[[ 84 6] # True Negatives, False Positives`

`[ 5 106]] # False Negatives, True Positives`

text

### # Libraries Used and Their Purposes

- \*\*pandas\*\*: Efficient data loading and manipulation.
- \*\*matplotlib & seaborn\*\*: Visualization for EDA (distributions, counts, boxplots).
- \*\*scikit-learn\*\*: Machine learning modeling, including data splitting, scaling, logistic regression, and evaluation metrics.

### # Code Snippets Generating Outputs

- Missing values check:  
`print(df.isnull().sum())`

text

- Class distribution check:  
`print(df['Clicked on Ad'].value_counts(normalize=True))`

text

- Top locations frequency:  
`print(df['City'].value_counts().head(10))`  
`print(df['Country'].value_counts().head(10))`

text

- Model training and evaluation:  
from sklearn.linear\_model import LogisticRegression  
from sklearn.model\_selection import train\_test\_split  
from sklearn.preprocessing import StandardScaler  
from sklearn.metrics import accuracy\_score, precision\_score, recall\_score, confusion\_matrix  
  
X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)  
scaler = StandardScaler()  
X\_train[num\_features] = scaler.fit\_transform(X\_train[num\_features])  
X\_test[num\_features] = scaler.transform(X\_test[num\_features])  
  
model = LogisticRegression()  
model.fit(X\_train, y\_train)  
y\_pred = model.predict(X\_test)  
  
print('Accuracy:', accuracy\_score(y\_test, y\_pred))  
print('Precision:', precision\_score(y\_test, y\_pred))  
print('Recall:', recall\_score(y\_test, y\_pred))  
print('Confusion Matrix:\n', confusion\_matrix(y\_test, y\_pred))

text

### ## Final Thoughts

The model delivers robust predictions and can significantly improve targeted advertising at ConnectSphere Digital. Next steps might include error analysis on misclassifications, feature refinement, or testing advanced models. This repository contains full code, visualization, and evaluation to support deployment-ready insights.

---

\*This documentation was prepared for the EdGlobe Minor Project submission, showcasing a real-world data science solution with a pinch of fun and full technical rigor.\*

\*\*Screenshots:\*\*

