



**FATİH  
SULTAN  
MEHMET**  
VAKIF ÜNİVERSİTESİ

**T.C.  
FATİH SULTAN MEHMET VAKIF ÜNİVERSİTESİ**

**Engineering Faculty  
Computer Science**

## **Data Mining Project**

Name & Surname	Student ID
Muhittin AKIN	1821221006
Dilara ÇELİK	2021221042
Metin KAĞIT	1821221033

# Introduction

In this project, we will predict visit reason which is order or guest. We analyze relationships between features such as entry time and date, the block and apartment visited, and the individual visitor. To compare classification methods, we train six different models—Logistic Regression, K-Nearest Neighbors, SVM, Decision Trees, Random Forest, and XGBoost, Naive Bayes and MLP (Sklearn). For each model, we perform hyperparameter tuning with five-fold cross-validation (5-CV) to find the best parameter settings. Finally, we evaluate their performance on a held-out test set using accuracy, precision, recall, F1-score, and AUC.

## Data Collection and Preprocessing

We collect the data from security checkpoint of residential complex. And anonamize it to privacy. The dataset has 40367 rows and 7 columns. The columns name are:

- NO: Id. TARİH: Visit date.
- ZİYARET EDİLEN: Visited person block and apartmen no.
- TEYİT BİLGİSİ: Emty Column
- ZİYARET SEBEBİ: Visit reason
- GELİŞ SAATİ: Visit Time
- ÇIKIŞ SAATİ: Empty Colum

### 1. Cleaning empty cells

After loading the raw data from Excel with pandas, all empty cells (NaN) were removed.

### 2. Renaming columns

The original column headers were renamed to be more meaningful and consistent (e.g. ZİYARET SEBEBİ → visit\_reason).

### 3. Removing rows with missing values

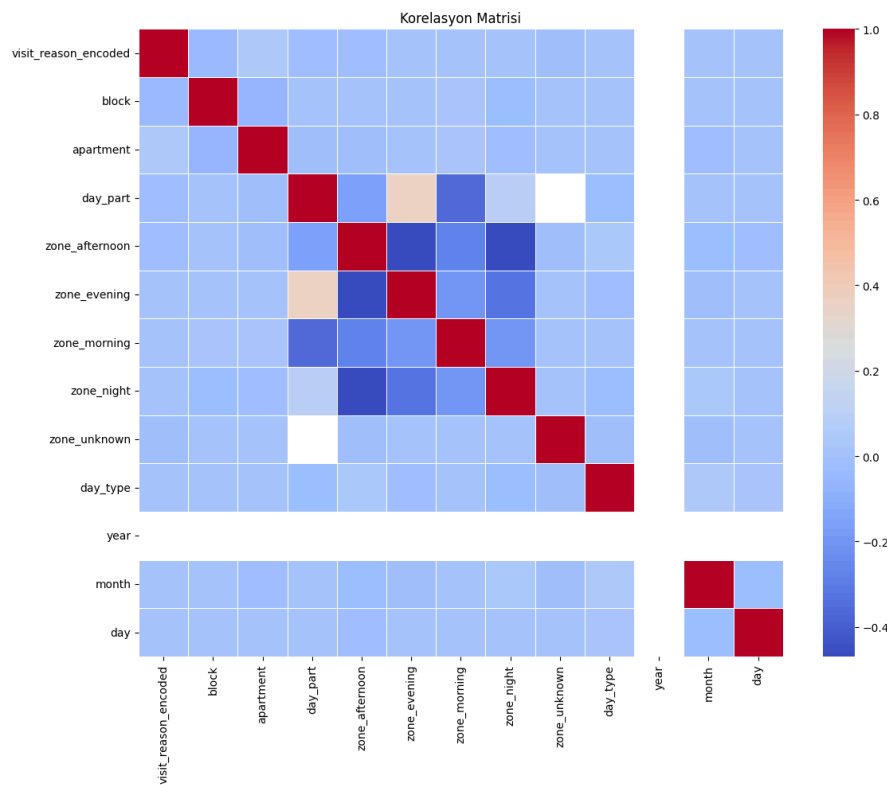
Any row containing a missing value in any column was dropped, because these gaps could not be logically or numerically imputed and might harm model performance.

## 4. Feature Engineering

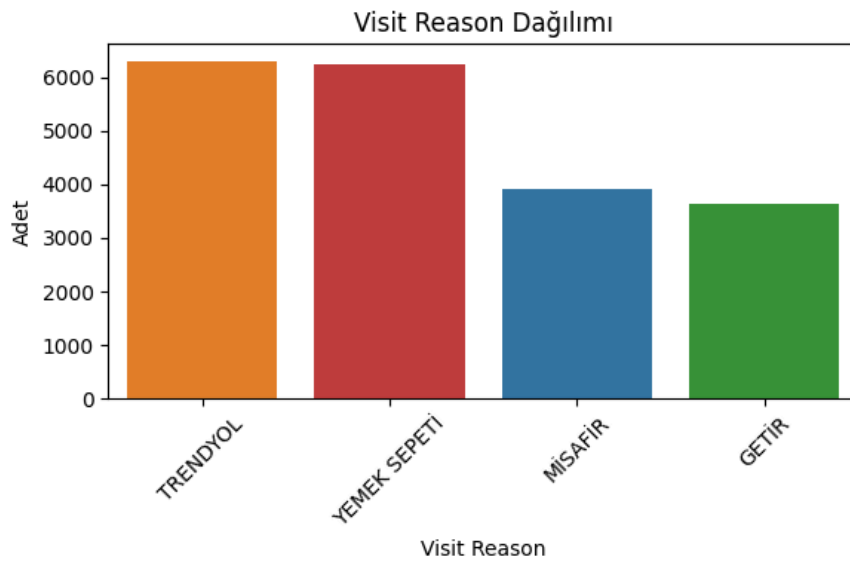
Since the original dataset had relatively few features, we created new variables:

- **visit\_reason\_encoded:** Encoded version of visit\_reason that include most repeated 4 unique value.
- **visit\_reason\_Binary:** Encoded version of visit\_reason that include 2 value is\_order 1 or not 0.
- **block & apartment:** Separate the visited column into two column, one of them just encoded block name and the other one is apartment number.
- **day\_part:** The arrival\_time column was split into five time-of-day categories—unknown, morning, afternoon, evening, and night. For each record, the column corresponding to its arrival slot is set to 1, and the others are set to 0.
- **day\_type:** A new column day\_type was created to classify each record as either weekday or weekend.
- **day&month&year:** date column separate 3 part which are day, month and year.

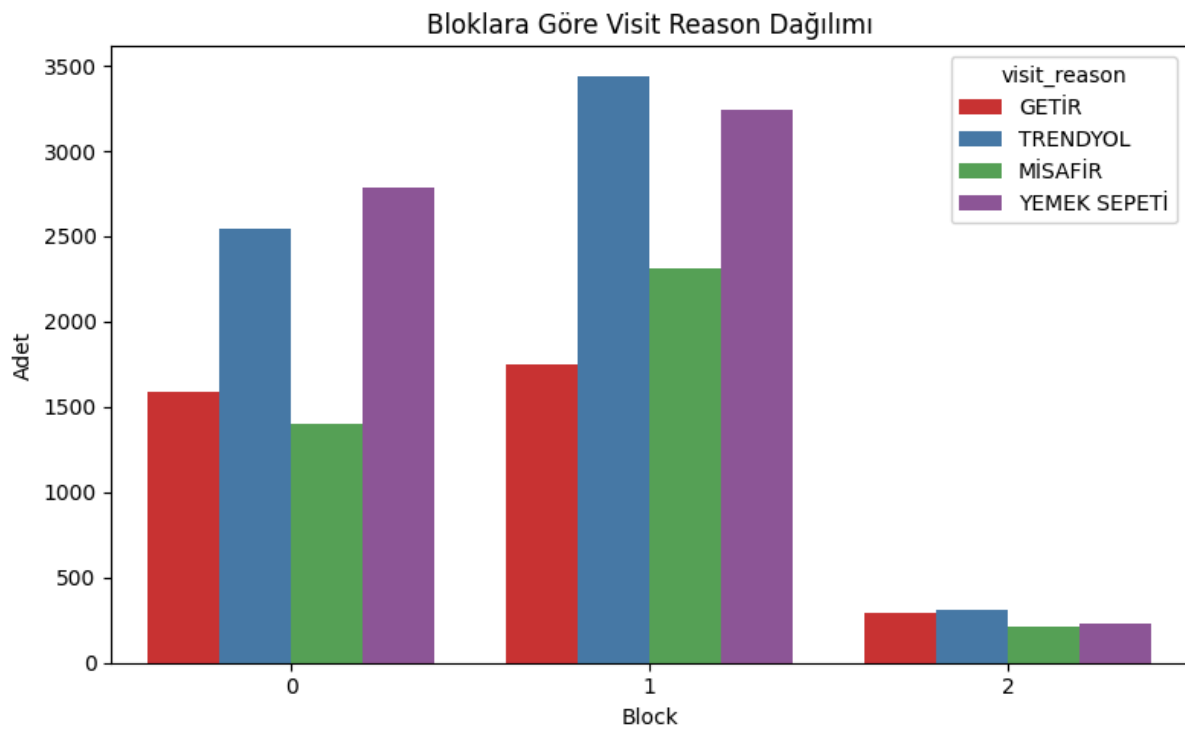
## Data Visualization



Since there is no linear correlation, the chosen models should be those that handle non-linear relationships well, such as Random Forest.

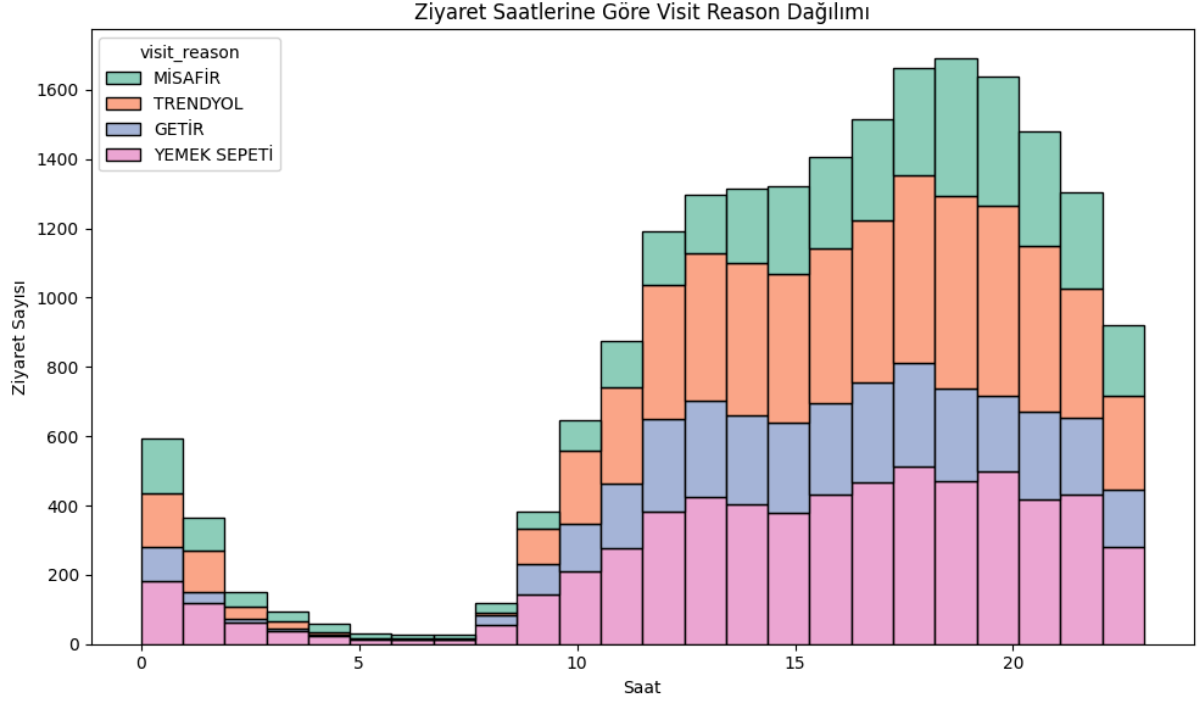


The graph show the amount of visit from most bigger 4 value.



0:A, 1:B, 2:C

The highest number of flats are B, A and C, respectively. It is observed that this situation also affects the number of visitors.



As we seen, there is a higher number of visitors in the afternoon and evening hours.

## Model Construction

In the model construction phase, we implemented six classification algorithms—Logistic Regression, K-Nearest Neighbors, Support Vector Machine, Decision Tree, Random Forest, XGBoost, Naive Bayes and MLP (Sklearn) — within scikit-learn pipelines and optimized their hyperparameters via GridSearchCV with cross-validation to identify the best performer for predicting visit reasons.

Firstly we trained our model without cross validation and tuning and the result can see in the table below.

Model	Set	Accuracy	Error	Precision	Recall	F1-score	AUC
Logistic Regression	Train	0.5296570898980537	0.4703429101019463	0.5310679611650485	0.5069508804448564	0.5187292555713608	0.5410894785148054
Logistic Regression	Test	0.5207556549838429	0.4792443450161571	0.8265204386839482	0.5120444718962323	0.6323417238749046	0.544248435182526
KNN	Train	0.8604031510658017	0.13959684893419833	0.8716845878136201	0.845227062094532	0.8582519703564286	0.9426201604259294
KNN	Test	0.6716380810340542	0.3283619189659458	0.8189683860232945	0.7600370599135269	0.7884030113727375	0.5628387815078113
SVM	Train	0.7197636700648748	0.28023632993512515	0.6667643439020102	0.8786685202347853	0.758188664156476	0.8036465567123416

<b>SVM</b>	Test	0.733532189 9080288	0.266467810091971 2	0.81354950781702 37	0.86781964175416 93	0.83980872683801 56	0.57317660897857 06
<b>Decision Tree</b>	Train	0.996022551 7454433	0.003977448254556 704	0.99953332814809 05	0.99250849552054 38	0.99600852547955 82	0.99996819577924 36
<b>Decision Tree</b>	Test	0.735520755 6549839	0.264479244345016 13	0.84139447236180 91	0.82736256948733 79	0.83431952662721 9	0.59294976453972 14
<b>Random Forest</b>	Train	0.996022551 7454433	0.003977448254556 704	0.99445684810224 03	0.99760580784677 17	0.99602883911015 15	0.99994114549020 92
<b>Random Forest</b>	Test	0.783494904 3002735	0.216505095699726 54	0.82720486591097 6	0.92402717726991 97	0.87293946024799 42	0.66481590035525 58
<b>XGBoost</b>	Train	0.868782823 6021007	0.131217176397899 32	0.95571673983584 65	0.77340129749768 3	0.85494749423717 24	0.97433955899165 2
<b>XGBoost</b>	Test	0.710166542 3813075	0.289833457618692 5	0.89572192513368 99	0.72421247683755 41	0.80088797814207 65	0.74807186161151 62
<b>Gaussian NB</b>	Train	0.596848934 1983317	0.403151065801668 26	0.59098824553765 78	0.62905468025949 95	0.60942760942760 94	0.62894535772423 6
<b>Gaussian NB</b>	Test	0.597315436 2416108	0.402684563758389 24	0.82753036437246 96	0.63125386040765 91	0.71618780658724 6	0.55714465562213 04
<b>MLP (Sklearn)</b>	Train	0.749459375 9654	0.2505406240346	0.70515752032520 33	0.85742971887550 2	0.77387425066220 55	0.83073368424387 88
<b>MLP (Sklearn)</b>	Test	0.721600795 4262987	0.278399204573701 26	0.81555423122765 19	0.84527486102532 42	0.83014861995753 72	0.57317011759244 33

Logistic Regression and Gaussian NB both show low performance, with training and test accuracies around 52–60% and AUC scores of about 0.55. KNN overfits, reaching 86% accuracy on training but dropping to 67% on test. SVM (72–73% accuracy, F1 ~84%, AUC ~0.57) and MLP (75–72% accuracy, F1 ~83%, AUC ~0.57) stay balanced on both training and test, with almost no overfitting.

Decision Tree (99% training, 73% test) and Random Forest (99% training, 78% test) also overfit, but Random Forest gives the best test results with 78% accuracy and F1 ~87%. XGBoost (87% training, 71% test) controls overfitting better and, with a test AUC of ~0.75, offers the most balanced class separation.

After that, we train the models with hyper parameter optimization and cross validation to improve model performance and stability.

Model	Accuracy	Precision	Recall	F1-score	AUC
<b>Logistic Regression</b>	<b>0.523743</b>	<b>0.824246</b>	<b>0.518966</b>	<b>0.636596</b>	<b>0.537356</b>
<b>KNN</b>	<b>0.687833</b>	<b>0.824358</b>	<b>0.777895</b>	<b>0.800435</b>	<b>0.57098</b>
<b>SVM</b>	<b>0.75317</b>	<b>0.816198</b>	<b>0.894848</b>	<b>0.853701</b>	<b>0.5752</b>
<b>Decision Tree</b>	<b>0.737358</b>	<b>0.847033</b>	<b>0.822131</b>	<b>0.834384</b>	<b>0.605229</b>
<b>Random Forest</b>	<b>0.79459</b>	<b>0.832988</b>	<b>0.931546</b>	<b>0.879513</b>	<b>0.684349</b>
<b>XGBoost</b>	<b>0.706678</b>	<b>0.89694</b>	<b>0.718091</b>	<b>0.797573</b>	<b>0.755434</b>
<b>Gaussian NB</b>	<b>0.592462</b>	<b>0.827852</b>	<b>0.623256</b>	<b>0.711097</b>	<b>0.561607</b>
<b>MLP (Sklearn)</b>	<b>0.731441</b>	<b>0.819947</b>	<b>0.853825</b>	<b>0.836481</b>	<b>0.582883</b>

According to these results, the best model is Random Forest, but XGBoost has the highest AUC score of all. Therefore, if you prioritize overall accuracy and recall, choose Random Forest; if you care more about ranking and discrimination ability, choose XGBoost.

## Conclusion

Our comprehensive evaluation compared various classification models using both a single train/test split and 5-fold cross-validation metrics, and we also tuning parameters by hyper parameter optimization and analyzed each model's AUC from the last fold's ROC curves. As a result, The less success model is Logistic Regression because Logistic Regression good to learned from linear relations our data is not contain linear relations. Also, XGBoost emerges as the most balanced model overall, while Random Forest ranks second for having the highest recall and test accuracy. If your goal is to maximize overall discrimination power and AUC, choose XGBoost; if you need to avoid missing any positive cases, Random Forest is the better choice.

