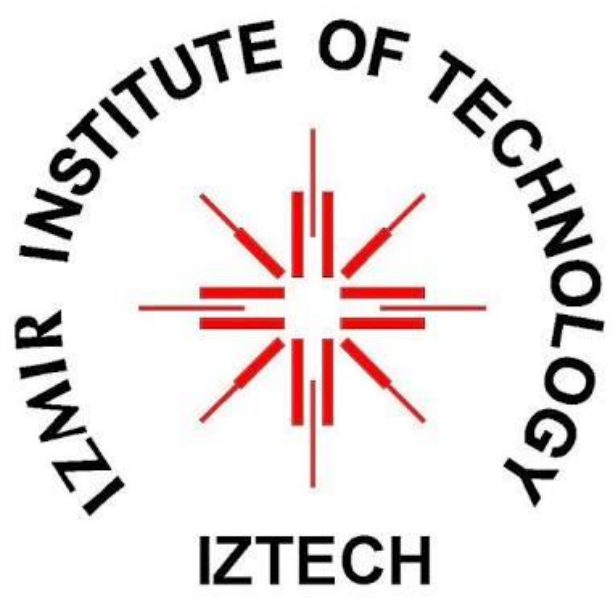# University Entrance Exam Data Science Studies

Berke Eren-Korkut Emre Arslantürk

250206008-250206039

Assist.Prof.Dr. Mehmet Serkan Apaydın

Izmir Institute of Technology

Electronics and Communication Engineering

## ABSTRACT

- Placement results data were collected from OSYM and Hacettepe University preference robot. [1], [2]
- Data were subjected to certain filtering processes and studied only on computer engineering data.
- There were 4 different categories of prediction:
  Entry threshold score, ranking, quota problem, and number of admitted student.
- The results obtained using different parameters and models were compared.

## INTRODUCTION

- Every year, both students and universities plan for the next year according to the placement results announced by OSYM.
- Our motivation is to estimate the desired values using extracted features.
- To get the best prediction performance, 4 different methods have been performed.

## DATA ANALYSIS

- The university entrance exam placement results were filtered for Computer Engineering.

### Dataset



- Then, universities with data between 2016 and 2021 were selected and each year compared with previous year.
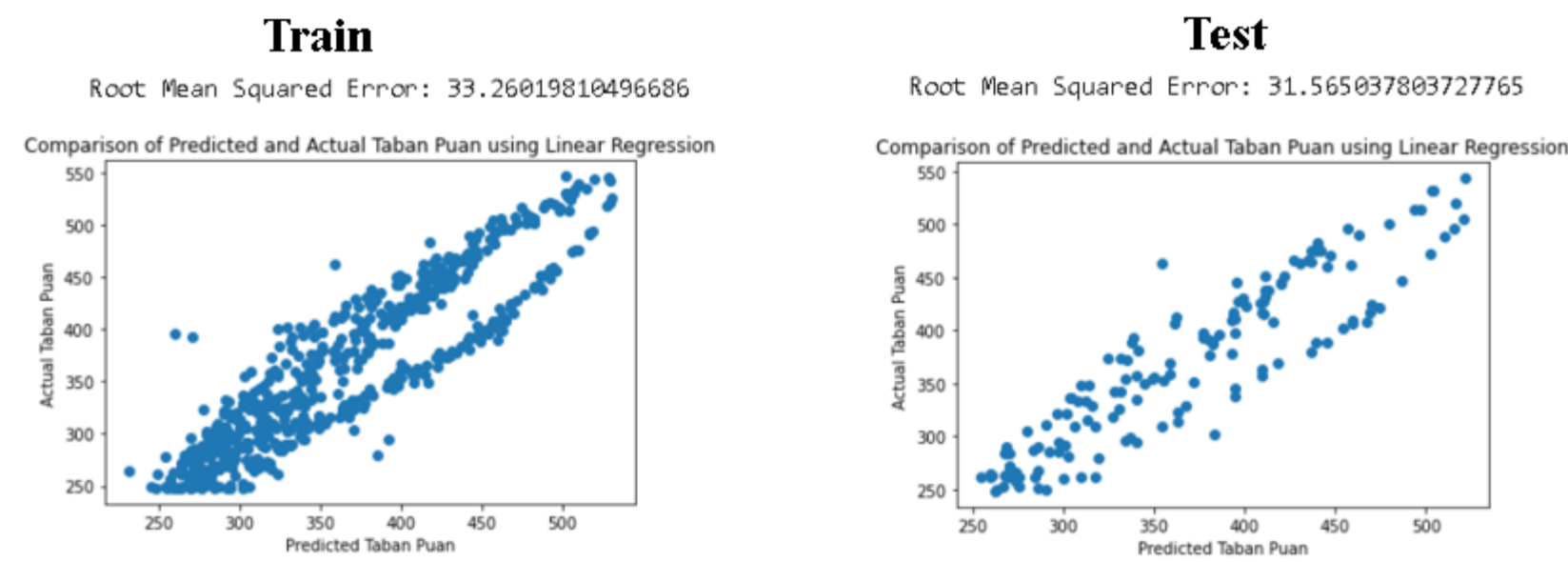- The dataset consists of 753 samples and 13 features.

## TABAN PUAN PREDICTION



- The features in figure are given to the model as input.
- 80% of the entire dataset is allocated in a random way for training and 20% for the test set.
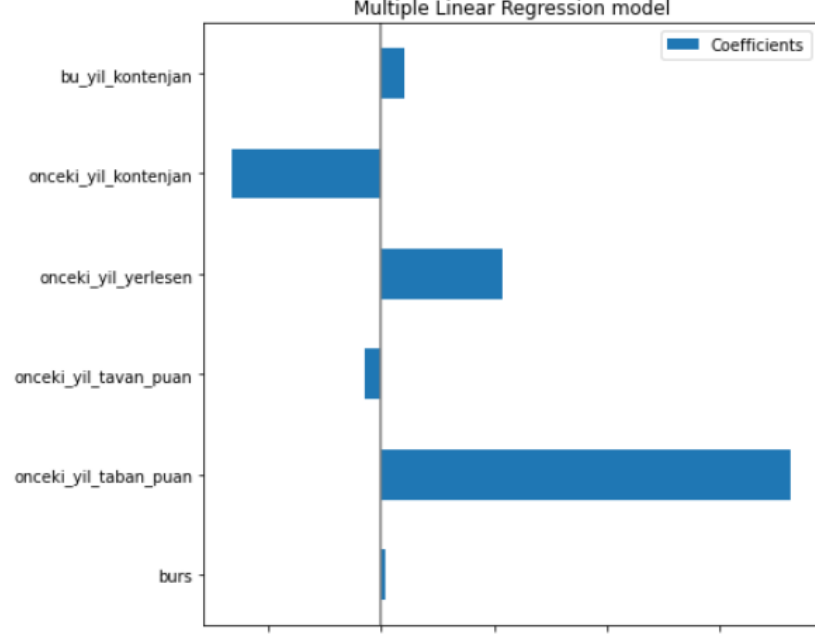- Linear regression utilized is a linear model, a model in which the input variables and the single output variable have a linear relationship.
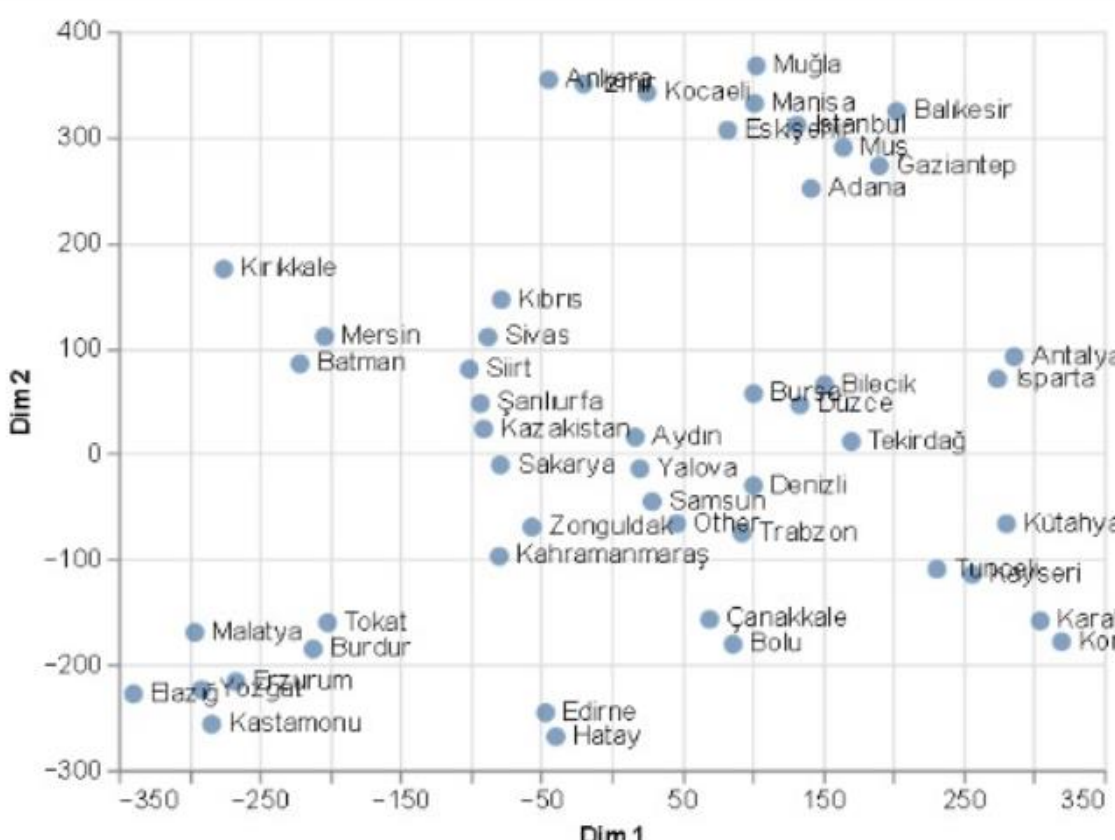


- The root mean squared error values obtained for both the train and the test set were calculated in terms of points.
- Comparison of predicted and actual threshold entry score using linear regression can be seen.

## Feature Importance



- Techniques that generate a score for all input characteristics for a particular model are referred to as feature importance.
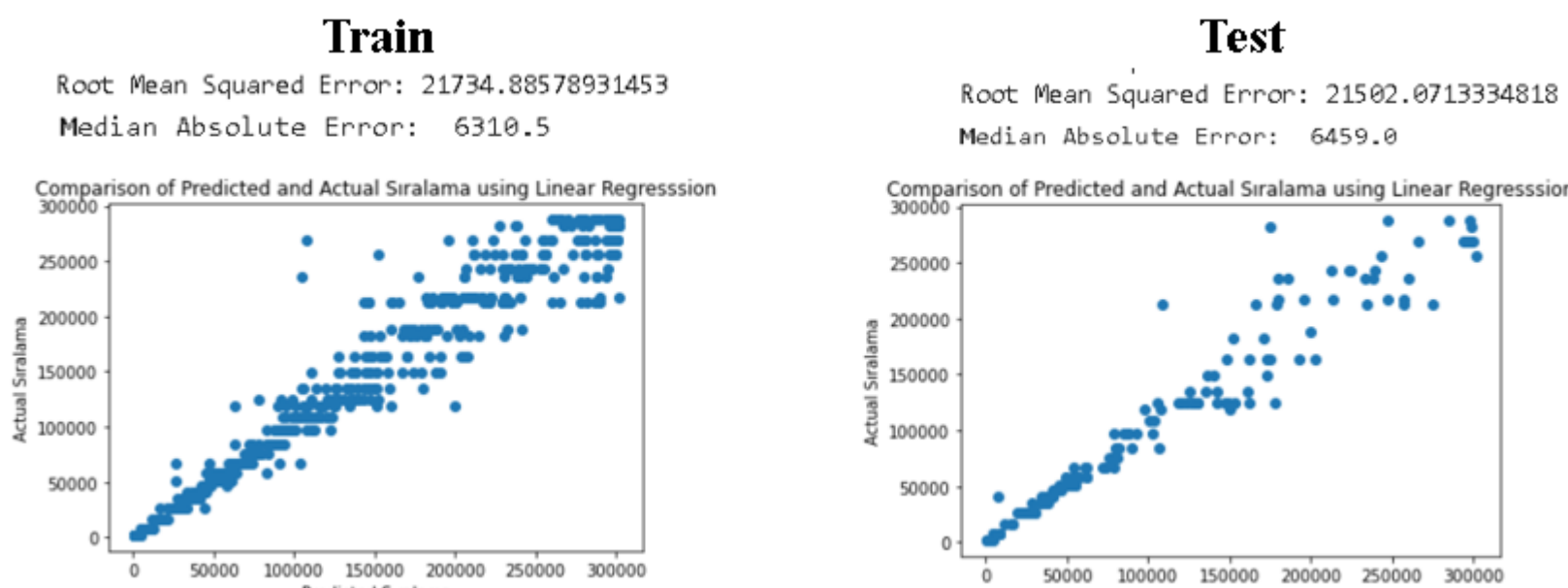
### Categorical Embedding



- To use these embedding layers, we first encoded the categorical variable with integer values.
- Each of these integers will then correspond to a vector representation of corresponding category.

## SIRALAMA PREDICTION



- The ranking is a parameter regardless of the difficulty of the exam. So, this prediction may give more reliable results.
- "City" and "university name" categorical features were added to the model by label encoding.



## THE QUOTA PROBLEM



- Class 0 indicates that the university cannot fill its quota, while class 1 indicates that it will fill it.
- Unbalanced data set has been balanced by applying oversampling to the data set.



- The results of the score estimation were also given to the model as input.



## YERLESEN PREDICTION



- A model was developed on calculating whether a university can fill its quota by estimating the number of people who have admitted.



## RESULTS

| Prediction | Categorical Features | Numerical Features | Model | Test RMSE | In terms of |
|---|---|---|---|---|---|
| Taban Puan | - | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan | Linear Regression | 31.56 | Points |
| Taban Puan | - | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan | Random Forest | 32.05 | Points |
| Taban Puan | sehir | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan | Decision Tree | 32.65 | Points |
| Taban Puan | sehir, okul | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan | Decision Tree | 34.53 | Points |
| Taban Puan | sehir, okul | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan | Neural Network | 31.59 | Points |
| Siralama | - | Common Features,onceki_yil_siralama | Linear Regression | 25735 | People |
| Siralama | sehir, okul | Common Features,onceki_yil_siralama | Linear Regression | 21502 | People |
| Yerlesen | sehir, okul | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan,tahmini_taban_puan | Decision Tree | 5.46 | People |
| Yerlesen | sehir, okul | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan,tahmini_taban_puan | Decision Tree | 4.46 | People |

| Prediction | Categorical Features | Numerical Features | Oversampling | F-1 Score(for 0) |
|---|---|---|---|---|
| Kontenjan | - | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan,onceki_yil_siralama | - | 0.65 |
| Kontenjan | - | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan,onceki_yil_siralama | + | 0.95 |
| Kontenjan | + | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan,onceki_yil_siralama | + | 0.94 |
| Kontenjan | + | Common Features,onceki_yil_taban_puan, onceki_yil_tavan_puan,onceki_yil_siralama,tahmini_taban_puan | + | 0.96 |

Common Features: burs, onceki_yil_yerlesen, onceki_yil_kontenjan, bu_yil_kontenjan

- For estimation the entry threshold score, the most successful result belongs to linear regression without "city" and "school" categorical features.
- In the prediction of yerlesen, it was observed that the estimated entry threshold score added to the model improved the result.
- It was observed that the best result in the quota was obtained in the model where categorical values were added, oversampling was applied to the dataset and the predicted entry threshold score was added.

## CONCLUSION

- It is decided that a trending ranking can be added as a feature by researching at the search numbers of universities and departments in Google search engine.
- It was decided to enlarge the dataset by adding it to the other departments as well as computer engineering.
- When creating a dataset, instead of making a comparison with just a year ago, it will be observed that as a trend.

## REFERENCES

[1] Yükseköğretim Kurumları Sınavı. 2021. (n.d.). Retrieved June 5, 2022, from https://www.osym.gov.tr/TR,21232/2021.html

[2] Osys Sonu Analalizor. (n.d.). Retrieved June 5, 2022, from http://yks.ee.hacettepe.edu.tr/

[3] Géron Aurélien. (2022). Hands-on machine learning with scikit-learn, Keras, and tensorflow concepts, tools, and techniques to build Intelligent Systems. O'Reilly.

[4] Making neural nets uncool again. (n.d.). Retrieved June 5, 2022, from https://www.fast.ai/