SCIENCE
PASSION
TECHNOLOGY

# Prediction Using Regression on Universities Data Set

**Group 14**:
David Mihola,
Ronald Infanger,
Thomas Sterner

18. 06. 2023

# Outline

# Data Set Analysis 1

- Universities from the United Kingdom,

- 21 columns and 145 row (131 unique rows).

| Idx | Column | Idx | Column |
|-----|--------|-----|--------|
| 1 | University_name | 12 | Student_satisfaction |
| 2 | Region | 13 | Student_enrollment |
| 3 | Founded_year | 14 | Academic_staff |
| 4 | Motto | 15 | Control_type |
| 5 | UK_rank | 16 | Academic_Calender |
| 6 | World_rank | 17 | Campus_setting |
| 7 | CWUR_score | 18 | Estimated_cost_of_living_per_year_(in_pounds) |
| 8 | Minimum_IELTS_score | 19 | Latitude |
| 9 | **UG_average_fees_(in_pounds)** | 20 | Longitude |
| 10 | **PG_average_fees_(in_pounds)** | 21 | Website |
| 11 | International_students | | |

**Group 14**: David Mihola, Ronald Infanger, Thomas Sterner
18. 06. 2023

# Data Set Analysis 2



Figure: Correlation heat map
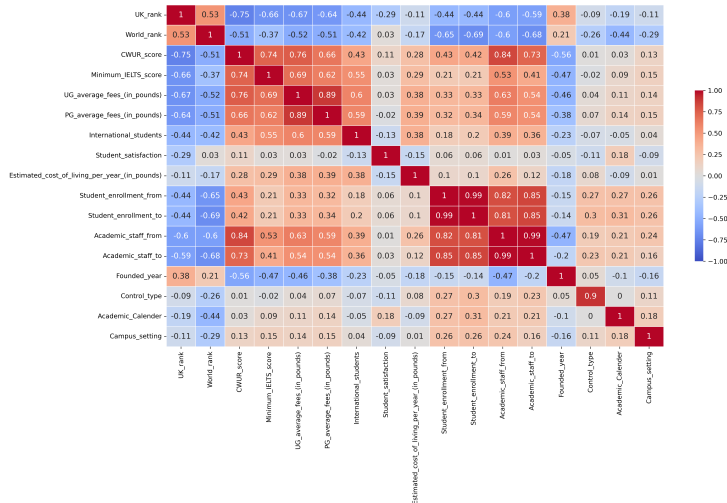
**Group 14**: David Mihola, Ronald Infanger, Thomas Sterner
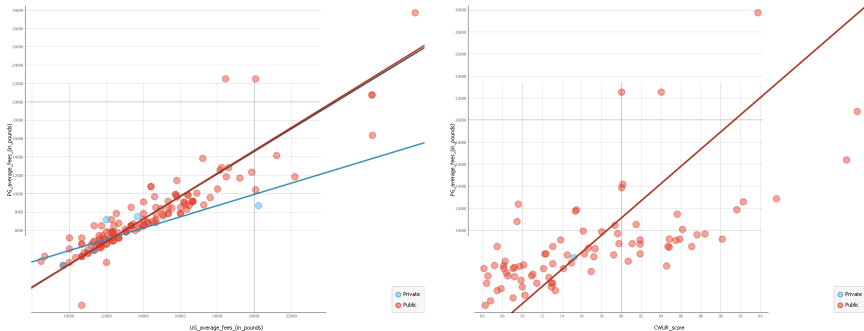18. 06. 2023

# Data Set Analysis 3



Figure: Linear dependency of PG fees, UG fees and CWUR

# Data Set Analysis 4



Figure: Linear dependency of World-rank, UK-rank and CWUR

# Data Set Analysis 5



Figure: Bias

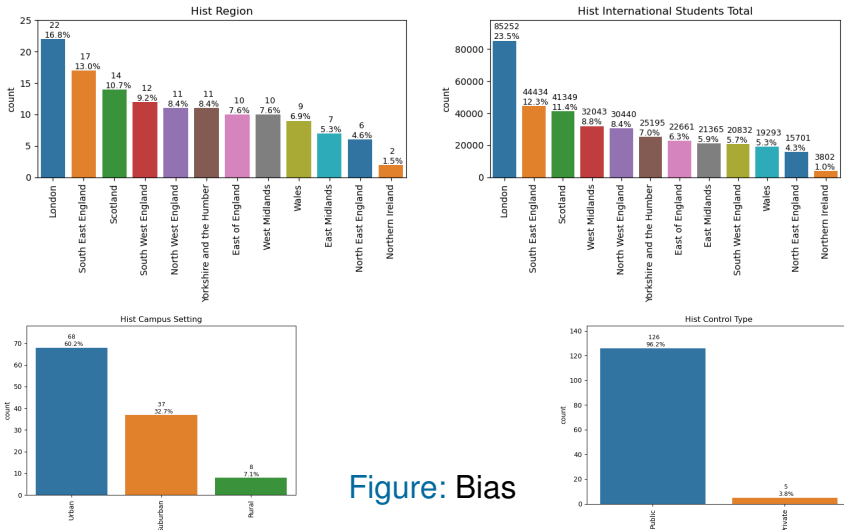**Group 14**: David Mihola, Ronald Infanger, Thomas Sterner
18. 06. 2023

# Cleaning and Pre-processing

1. Identification of missing values, split of compound columns, deduplication,

2. data set split,

3. missing value imputation,

4. normalization,

5. one-hot encoding,

6. removal of non-numeric columns.

# Imputation (mean, median, mixed)

- ## Missing values (6 out of 21 columns):

| Idx | Column | NaN | mean | median | mixed |
|-----|--------|-----|------|--------|-------|
| 1 | University_name | 14 | dropped | dropped | dropped |
| 3 | Founded_year | 14 | mean | median | researched online |
| 4 | Motto | 17 | dropped | dropped | dropped |
| 7 | CWUR_score | 47 | mean | median | linear regression on "UK_rank" |
| 16 | Academic_Calender | 26 | mode | mode | mode |
| 17 | Campus_setting | 18 | mean | median | KNN on "Latitude" and "Longitude" |

- ## Suspicious values (3 out of 21 columns):

| Idx | Column | Suspicious Value | Count | Imputation Approach |
|-----|--------|------------------|-------|---------------------|
| 3 | Founded_year | 9999 | 14 | researched online |
| 12 | Student_satisfaction | 0 | 6 | median |
| 14 | Academic_staff | over | 6 | 10000 |

**Group 14**: David Mihola, Ronald Infanger, Thomas Sterner
18. 06. 2023

# Training and Evaluation

- 5-fold cross validation across seeds $[40, 49]$,

- 3 subsets of the data set:

  - all continuous and categorical columns,

  - only continuous columns excluding `Latitude` and `Longitude`

  - columns with absolute value of correlation higher than 0.5 with the target variables,

- performance evaluation metrics:

  - MSE, MAE, RMSE, R2 score

- average performance across seeds $[40, 49]$.

**Group 14**: David Mihola, Ronald Infanger, Thomas Sterner
18. 06. 2023

# Models

Linear Regression (LR)

- baseline model to assess performance against,

- cross validation only for column subsets.

Fully Connected Neural Network (FCNN)

- 3 architectures with:

  - increasing number of hidden layers,

  - ReLU hidden activation, linear output activation.

# Models

## Random Forest (RF)

- grid search parameters:
  - max features $[1, 17]$
  - N estimators $[80, 100]$.

## Support Vector Regression (SVR)

- grid search parameters:
  - C (0.005, 0.01, 0.05, 0.1, 0.5, 1, 3, 5)
  - kernel (linear, rbf)
  - epsilon (0.0001, 0.0005, 0.005, 0.01, 0.25, 0.5, 1, 5)
  - gamma (0.0001, 0.001, 0.01, 0.1, 1)

# Prediction Performance, Ensemble

- no clear best imputation approach,

- no clear best subset of columns.

|  | MSE | MAE | RMSE | R2 score | Columns | Imputation |
|---|---|---|---|---|---|---|
| LR | 3708162.9 | 1408.8 | 1909.6 | 0.3075 | continuous | mixed |
| FCNN | 3389145.4 | 1304.6 | 1826.1 | 0.3659 | selected | median |
| RF | 3200786.3 | 1192.1 | 1761.7 | 0.4061 | continuous | median |
| SVR | 3643681.6 | 1219.7 | 1814.7 | 0.4625 | selected | mixed |
| Ensemble | 2363667.7 | 1108.6 | 1537.3 | 0.5178 | —— | —— |

Table: Average prediction performance across random seeds [40, 49]