# DAIMLER

Product Funding Prediction using Python Machine Learning
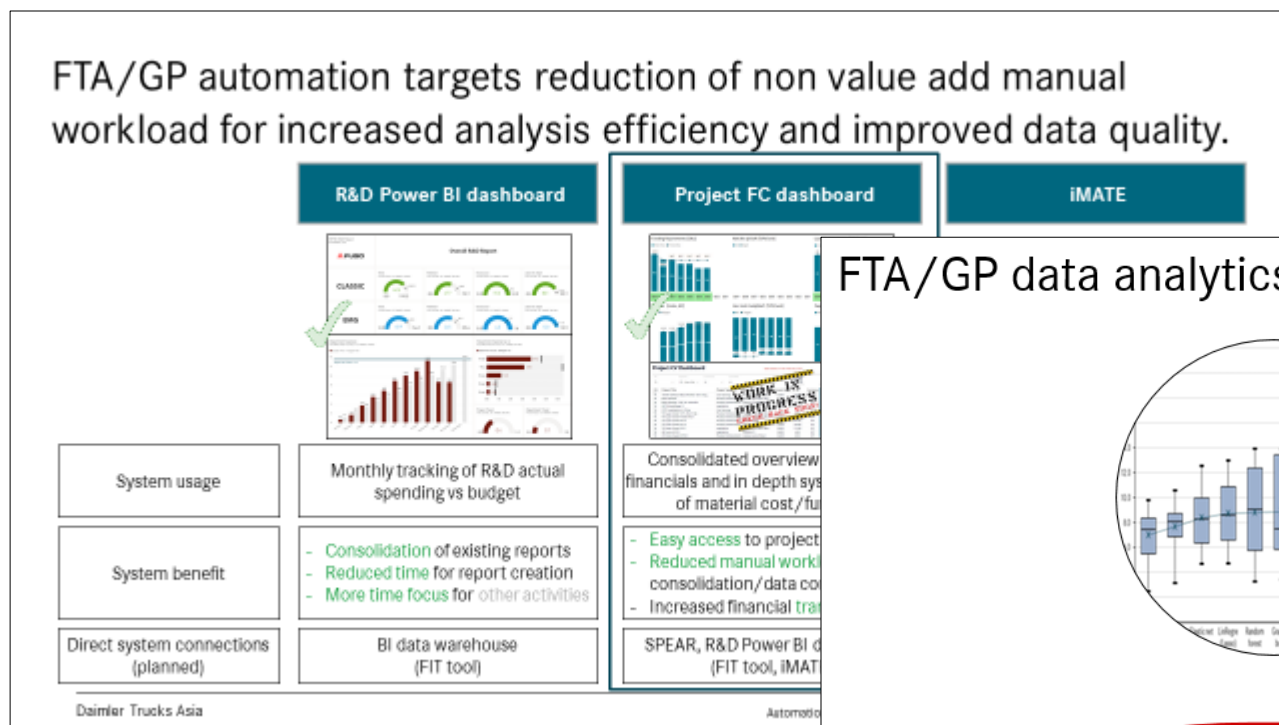13 April 2021
FTA/GP
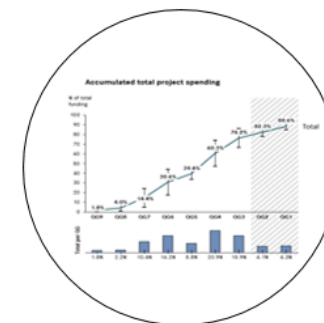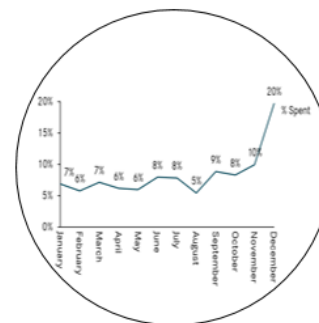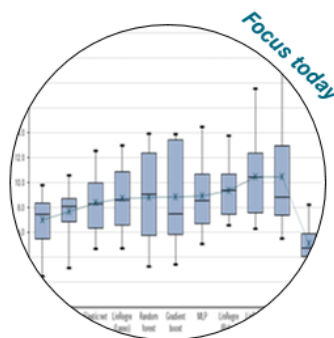
## Daimler Trucks Asia

FUSO

FUSO

BHARATBENZ

# In addition to report automation and iMate tool, FTA/GP focuses on data prediction using machine learning algorithms

FTA/GP automation targets reduction of non value add manual workload for increased analysis efficiency and improved data quality.

| R&D Power BI dashboard | Project FC dashboard | iMATE |
|---|---|---|

| System usage | Monthly tracking of R&D actual spending vs budget | Consolidated overview financials and in depth sys of material cost/fu |
| System benefit | - Consolidation of existing reports<br>- Reduced time for report creation<br>- More time focus for other activities | - Easy access to project<br>- Reduced manual worki<br>  consolidation/data co<br>- Increased financial tra |
| Direct system connections (planned) | BI data warehouse (FIT tool) | SPEAR, R&D Power BI (FIT tool, iMAT |

Daimler Trucks Asia

FTA/GP data analytics activities as the next step in finance digitization



Focus today

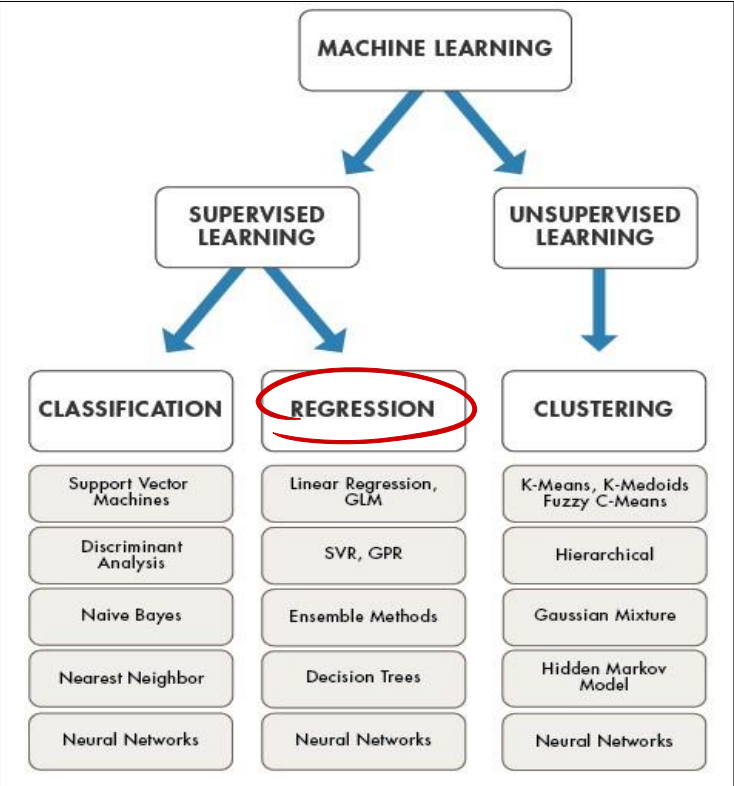| | **Project funding prediction** | **Time series-based periodic forecasting** | **Milestone-based funding cycle prediction** |
|---|---|---|---|
| Usage | Predictors-based funding estimation tool based on regression models. | Seasonality-adjusted auto-regression model. Based on historical data & rolling commitment | Project funding spending prediction. Based on historical data & project milestones. |
| Benefits | Early funding estimation with low-effort parametric input | Improved monthly/YE R&D forecast | Improved accuracy of the yearly budget planning for projects |

# Machine learning provides multiple finance applications. Regression models best suited as base for automated prediction tools

" Machine learning: "

Utilization of computer systems that are able to learn and adapt by recognizing patterns in data without following explicit instructions



**Tested models (10):** Linear regression, Lasso regression, Ridge regression, Elastic net, Random forest, Gradient boost, XGBoost, MLP, SVM, Stacking
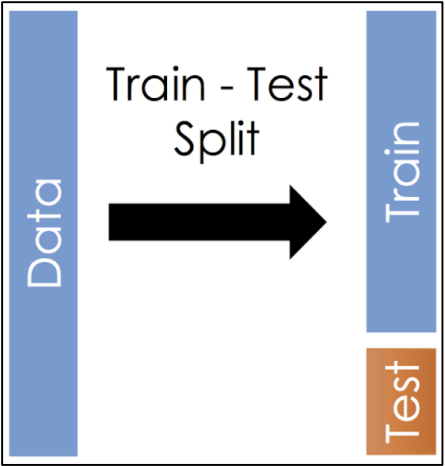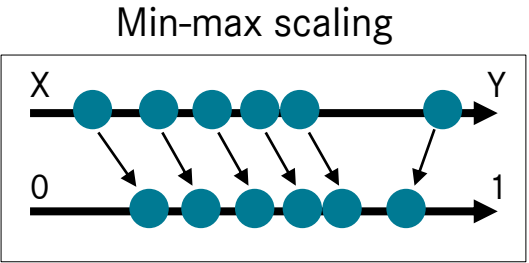
**Used tools:**

# Funding prediction activity follows standard machine learning process

## 1. Data preparation

**Dataset:**
**# of projects: 31**
**Selected data**



## 2. Pre-processing

Min-max scaling



Train - Test Split



## 3. Modeling



- Multivariate Linear Regression
- MV Lin Regression with Lasso Reg.
- MV Lin Regression with Ridge Reg.
- MV Lin Regression with Elastic Net Reg.
- Random Forest
- Gradient Boosting
- XGBOOST
- Multi Layer Perceptron
- Stacking (Grad. Boosting, LR Ridge, SV , RF)
- SVM

## 4. Validation



Model abs. mean deviation (random state n=10)
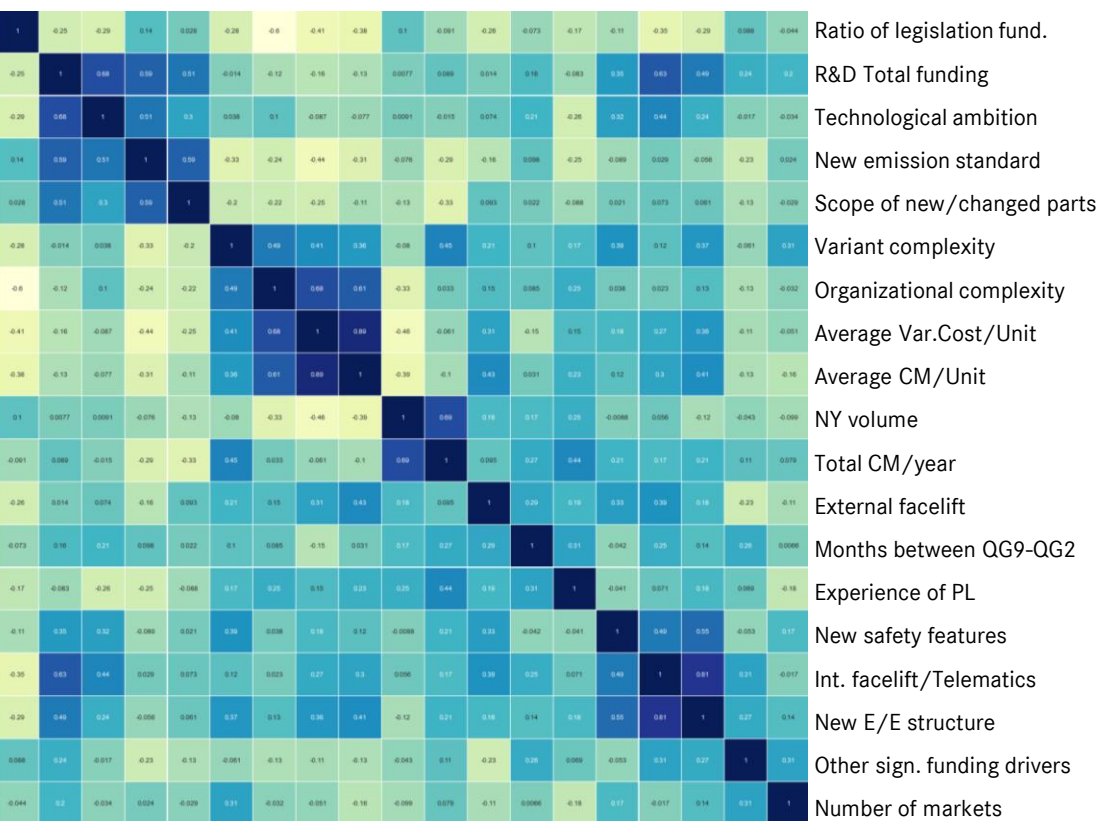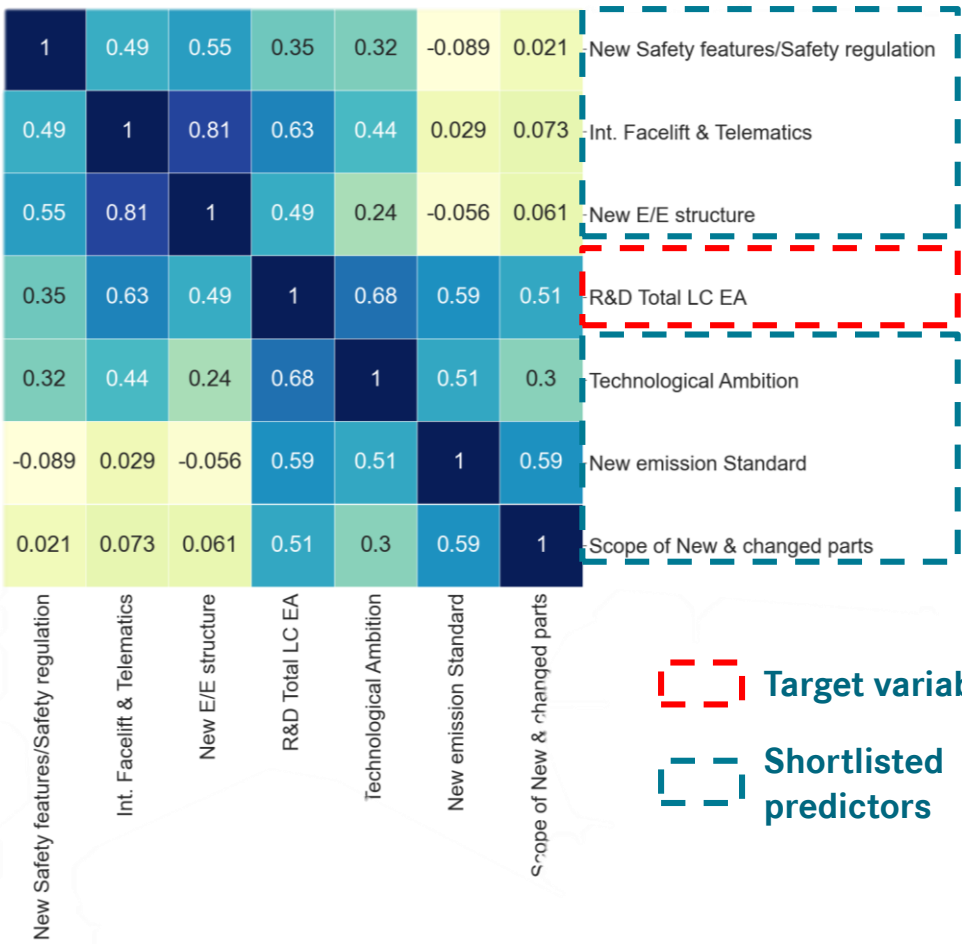
# During data exploration 6 most promising predictors out of 18 have been shortlisted to be used in prediction models

## Correlation heatmaps (Pearson)

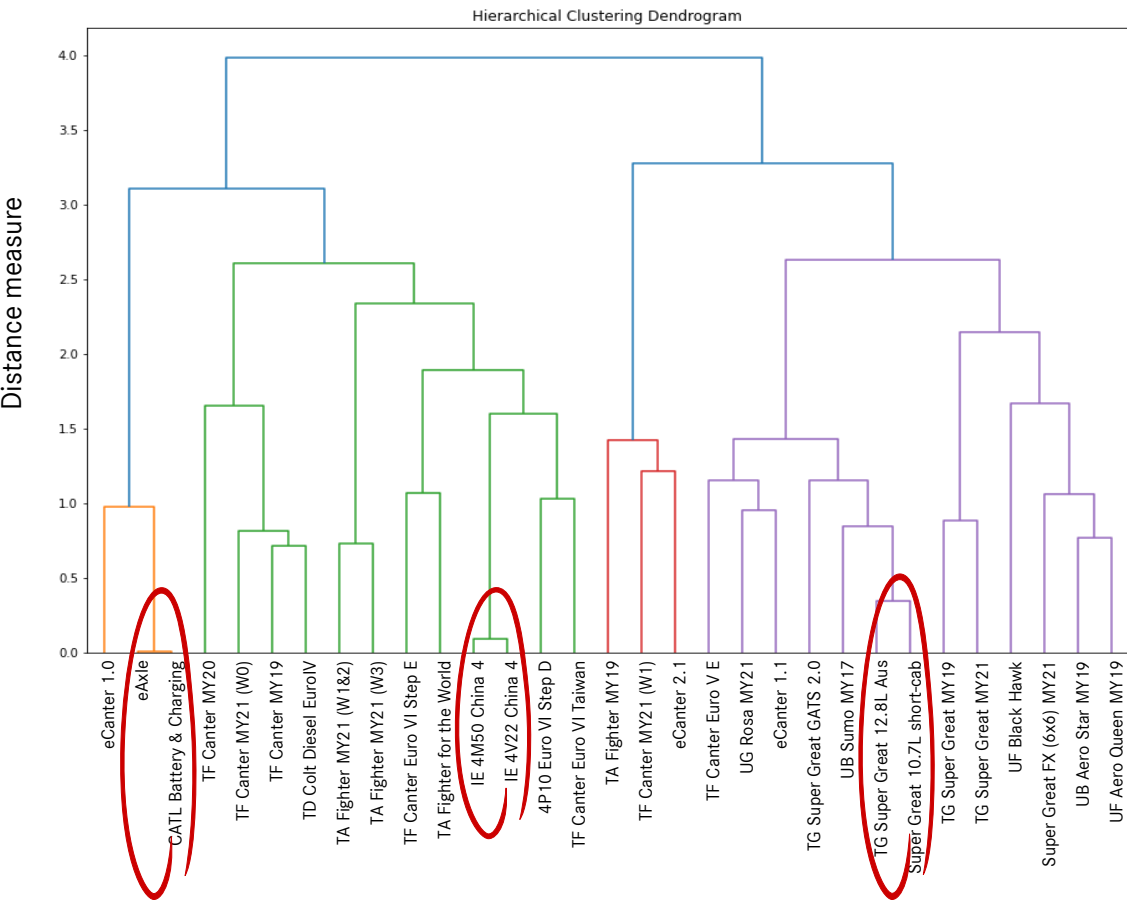**18 predictors that can be objectively attributed to a project**

**6 predictors with high levels of correlation to R&D funding**



Correlation over r=**0.35** considered moderate/strong

**Target variable**

**Shortlisted predictors**

# Heterogenic nature of dataset, limited sample size and input scarcity are main challenges that have to be addressed to achieve prediction quality

## Hierarchical agglomerative clustering (ward)



Hierarchical Clustering Dendrogram

**Cluster analysis shows high connection point for majority of projects** → **Highly heterogenic dataset**

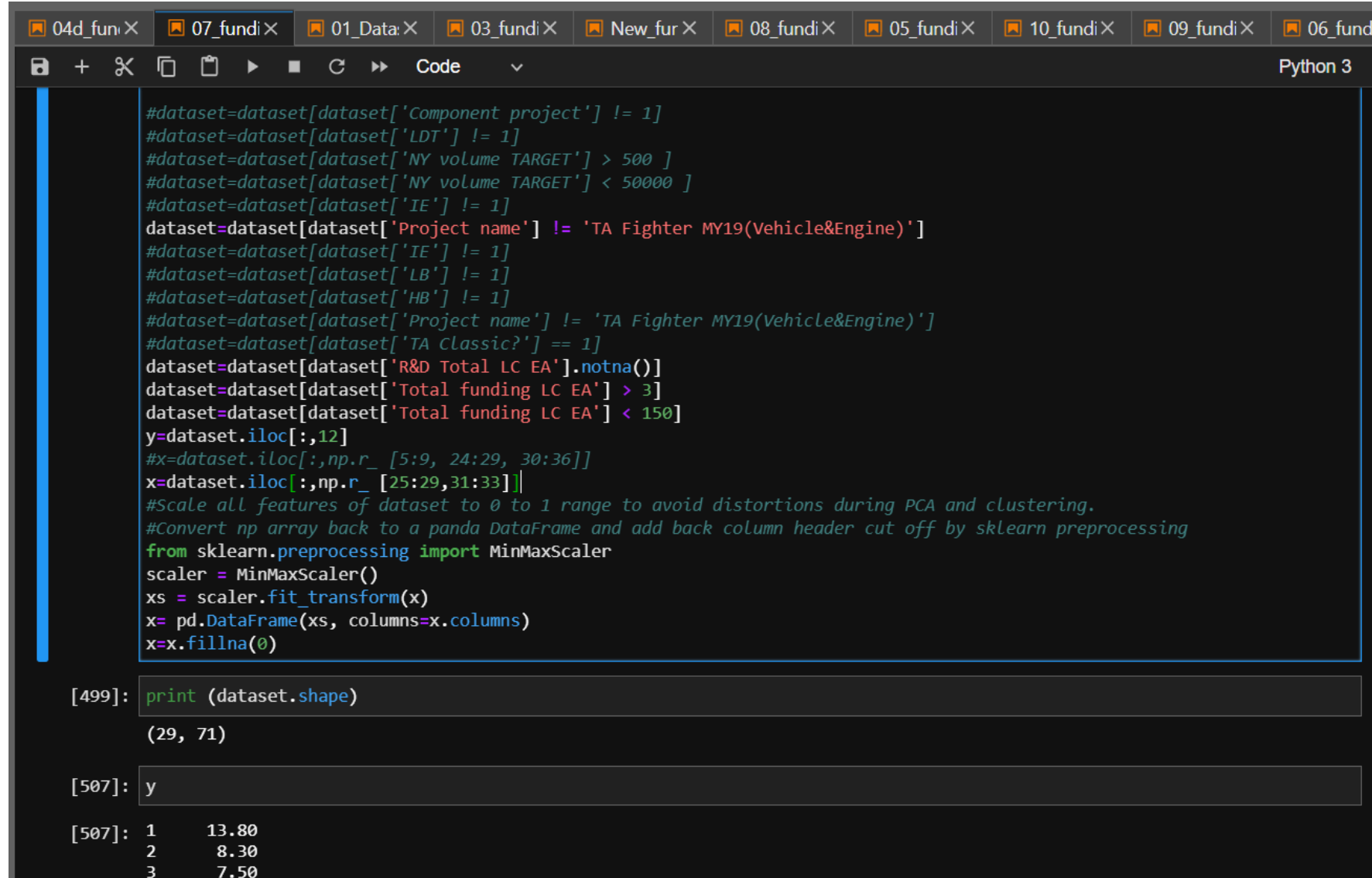## Excerpt of the dataset

**31 projects**　　　　**6 shortlisted predictors**

| Project name | Ext. Facelift | Int. Facelift & Telematics | New E/E structure | New Safety features/Safety regulation | New emission Standard | Optional: other expected, significant funding drivers |
|---|---|---|---|---|---|---|
| TA Fighter MY19(Vehicle&Engine) | 0.00 | 0.50 | 0.50 | 0.50 | 0.50 | 0.25 |
| 4P10 Euro VI Step D | 0.00 | 0.00 | 0.25 | 0.25 | 0.50 | 0.00 |
| UB Aero Star MY19 | 0.00 | 0.25 | 0.25 | 0.25 | 0.00 | 0.00 |
| UF Aero Queen MY19 | 0.25 | 0.50 | 0.25 | 0.50 | 0.00 | 0.00 |
| TF Canter Euro VI Taiwan | 0.00 | 0.00 | 0.00 | 0.25 | 0.25 | 0.00 |
| TF Canter MY19 | 0.00 | 0.00 | 0.00 | 0.25 | 0.25 | 0.00 |
| TF Canter MY20 | 1.00 | 0.50 | 0.00 | 0.25 | 0.00 | 0.00 |
| TG Supaer Great GATS 2.0 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.50 |
| TG Super Great 12.8L to AUS/NZ | 0.25 | 0.00 | 0.00 | 0.00 | 0.00 | 0.25 |
| TG Super Great MY19 | 0.00 | 0.25 | 0.25 | 0.50 | 0.00 | 0.00 |
| TG Super Great FX (6x6) MY21 | 0.25 | 0.25 | 0.25 | 0.00 | 0.25 | 0.50 |
| TD Colt Diesel EuroIV | 0.00 | 0.00 | 0.00 | 0.00 | 0.25 | 0.25 |
| TF Canter Euro V Emerging Markets | 0.00 | 0.00 | 0.00 | 0.25 | 0.25 | 0.25 |
| TF Canter Euro VI Step E | 0.50 | 0.25 | 0.00 | 0.75 | 0.75 | 0.00 |
| TA Fighter for the World (Vehicle&Engine) | 0.00 | 0.00 | 0.00 | 0.50 | 0.25 | 0.50 |
| UG Rosa MY21 | 0.00 | 0.25 | 0.25 | 0.50 | 0.25 | 0.50 |
| TF Canter MY21 (W0) | 0.00 | 0.00 | 0.00 | 0.00 | 0.25 | 0.25 |
| TF Canter MY21 (W1) | 0.00 | 1.00 | 0.50 | 0.25 | 0.00 | 1.00 |
| TG Super Great MY21 | 0.25 | 0.00 | 0.25 | 0.50 | 0.00 | 0.00 |
| TA Fighter MY21 (W1&2) | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 |
| TA Fighter MY21 (W3) | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 |
| TG Super Great 10.7L short-cab | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.25 |
| IE 4M50 China 4 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 |
| IE 4V22 China 4 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 |
| UF Black Hawk | 1.00 | 0.75 | 0.50 | 0.50 | 0.25 | 0.00 |
| UB Sumo MY17 | 0.00 | 0.00 | 0.00 | 0.00 | 0.25 | 0.00 |
| eCanter 1.0 | 0.00 | 0.25 | 0.00 | 0.00 | 1.00 | 0.00 |
| eCanter 1.1 | 0.00 | 0.00 | 0.00 | 0.25 | 0.00 | 0.00 |
| eCanter 2.1 | 0.00 | 1.00 | 0.50 | 0.50 | 1.00 | 1.00 |
| eAxle | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 |
| CATL Battery & Charging | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 |

**Small sample size (n=31) and heterogenic dataset** → **Careful selection of modeling approach required**

# Python Demo
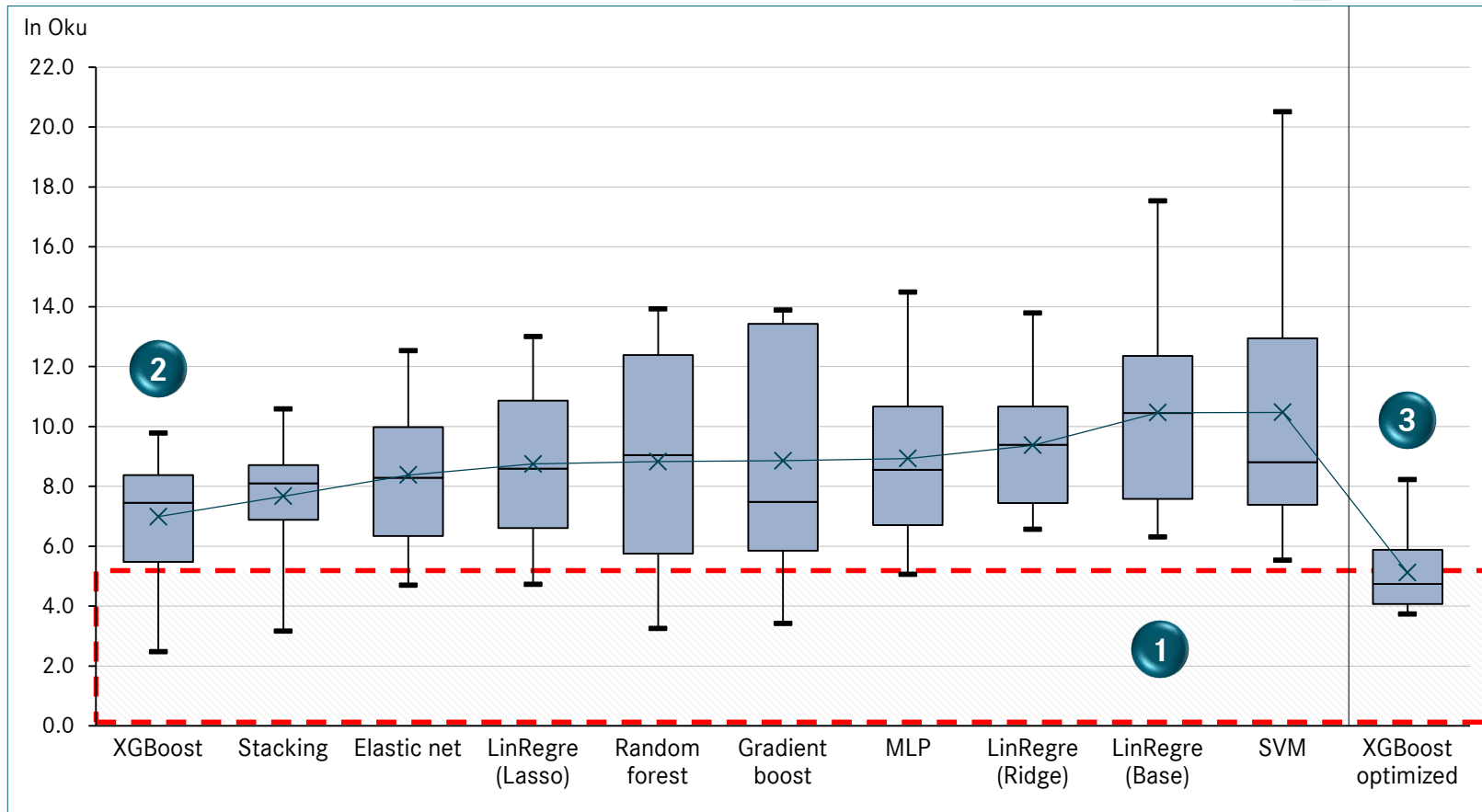
# XGBoost is the best model with average prediction accuracy close to ±5 Oku. Further optimization required to improve avg. accuracy and spread

**Absolute mean deviation by model (random samples n=10)**

✕ : Average ⬛ : 50% of values



**Testing and validation approach**

- R&D funding chosen as testing variable
- Model has been trained by 80% of random projects and tested on 20% of random projects
- Validation repeated x10 with random sets of projects
- Results represent average deviation of 6 random test projects (20%) from actual funding target

**Key takeaways**

1. Prediction corridor aspiration within 5 Oku of actual
2. XGBoost performs best among all tested projects even before optimization, but doesn't achieve the target prediction accuracy
3. XGBoost optimization with exclusion of just one outlier improves the prediction accuracy by 2 Oku and just reaches the target prediction range
- [Not shown] Prediction accuracy for total funding comparable with R&D funding

**Promising results. Further optimization of dataset & model to achieve target range**