



**TECHNISCHE
UNIVERSITÄT
DRESDEN**



**Big Data Analytics
in Transportation
@TU Dresden**

"Friedrich List" Faculty of Transport and Traffic Sciences Institute of Transport and Economics

Chair of Big Data Analytics in Transportation

Seminar Paper

Prediction of Freight Rail Transport in Germany

Thu Thuy Nguyen

27th August 2024

Supervising professor
Prof. Dr. rer. pol. Pascal Kerschke

Contents

1	Introduction	1
2	Project scope and data collection	2
3	Exploratory Data Analysis	3
3.1	Identify most important goods	3
3.2	Descriptive statistics of the dataset	3
3.3	Correlation Analysis	4
3.4	Inspecting temporal patterns in transport performance	5
3.5	Inspecting spatial patterns in transport performance	6
3.6	Inspecting combined temporal and spatial patterns	7
4	Developing machine learning model	8
4.1	Data preprocessing before modeling	8
4.2	Preliminary model selection and validation	8
4.3	Nested resampling and tuning hyperparameters	9
4.4	Nested resampling and feature selection	10
4.5	Model inspection and interpretation	11
4.6	Overall model assessment and limitations	13
5	Conclusion	14
	Bibliography	I

List of Figures

3.1	List of selected goods and their percentage in total transport performance . . .	3
3.2	Summary statistics of the dataset	4
3.3	Box-plots for transport volume and performance	4
3.4	Box-plots for all other numerical variables	4
3.5	Correlation matrix of all columns in the dataset	5
3.6	Change in transport performance by goods over 13 years	5
3.7	Transport performance in basic iron and steel and ferro-alloys etc. in 2023	6
3.8	Difference in transportation performance in basic iron and steel and ferro-alloys etc. between 2022 and 2023	6
3.9	Heatmap plot of routes with consistent operation across all year for specific types of goods	7
4.1	Scatter plots of predicted and target values for both models	11
4.2	Feature importance plot of RF	11
4.3	Top 5 most and least important features of XGB	12
4.4	Feature effect plots of RF and XGB	12

List of Tables

2.1	Explanation and sources of variables used in the data set	2
4.1	Performance metrics for different learners	9
4.2	Range of tuning and optimal configuration of 2 learners	9
4.3	Performance metrics before and after tuning hyperparameter	10
4.4	Performance metrics before and after feature selection	10

List of Abbreviations

NUTS	Nomenclature of Territorial Units for Statistics
GIS	Geo-graphic Information System
GDP	Gross Domestic Product
LR	Linear Regression
DT	Support Vector Machine
RF	Random Forest
SVM	Support Vector Machine
NN	Neural Network
XGB	Extreme Gradient Boosting
RMSE	Root Mean Squared Error
MAE	Mean Absolute Error
MAPE	Mean Absolute Percentage Error
R-sqrt	Coefficient of Determination
PDP	Partial Dependence Plot
ICE	Individual Conditional Expectation

1 Introduction

Railway freight transport plays a pivotal role in the national economy, serving as the backbone of supply chains and facilitating the efficient movement of goods across vast geographical areas. Its significance is underscored by the ability to handle large volumes of goods efficiently and sustainably. Predicting the performance of rail freight transport, in terms of volume and kilometers traveled, is essential for optimizing transport operations, enhancing capacity planning, and improving service reliability. Accurate forecasts enable stakeholders to make informed decisions regarding resource allocation, infrastructure investment, and logistical arrangements.

This paper inspects the activities of rail freight transport in Germany from 2011 to 2023 by analyzing historical data and conducting literature research on the relevant factors that could serve as prediction features for transport performance. Besides, temporal and spatial patterns and trends in rail transport activities will be revealed by tools of visualization. Additionally, machine learning models for prediction will be developed, optimized, and inspected. These models are expected not only to predict transport performance but also to provide insights into the causal relationships within the data, offering valuable perspectives for stakeholders to make informed decisions in the railway industry.

2 Project scope and data collection

This study will focus on rail freight transport performance, which is defined as the product of freight volume and transport distance. This indicator not only reflects the operational efficiency of the railway system but also provides data support for optimal resource allocation. The time scope for analysis is 13 years, from 2011 to 2023. The geographical scope is transport inside Germany, at NUTS level 2, or 40 government districts regions. This study focus on explore temporal and spatial patterns of transport activities by tools of visualization and develop simple machine learning models to predict rail freight transport performance.

The selected predictors and relevant datasets are detailed in table 2.1. Due to the lack of existing and public data on Germany regional GDP at the NUTS-2 level (government districts), the data on regional GDP at the NUTS-1 level (federal states) will be used. The study uses the production index of the energy-intensive and mining industries as features, recognizing their importance as key commodities for rail transport.

Column name	Explanation of the variable	Source
tpt_perfm	Yearly rail freight transport performance (tkm)	Destatis (2024g)
total_vol_ton	Yearly rail freight volume (tons)	Destatis (2024a)
gdp_origin	GDP of loading region at NUTS-1 level (billion Euro)	Destatis (2024e)
gdp_destination	GDP of unloading region at NUTS-1 level (billion Euro)	Destatis (2024e)
pro_ind_energy	Production index of extensive-energy industry	Destatis (2024d)
pro_ind_mining	Production index of mining n quarrying industry	Destatis (2024d)
year_rail_vol	Yearly total freight volume by rail (tons)	Destatis (2024b)
year_road_vol	Yearly total freight volume by road (tons)	Destatis (2024f)
year_water_vol	Yearly total inland waterway freight volume (tons)	Destatis (2024c)

Table 2.1: Explanation and sources of variables used in the data set

3 Exploratory Data Analysis

3.1 Identify most important goods

The project focuses on a few industries that are considered most important in rail transport. Two criteria were employed to identify the relevant industries. The first criterion is that the industry handles heavy, bulky goods, which predominantly rely on rail transport due to its ease and cost-efficiency. Additionally, the goods categories must be identifiable and represent a significant portion of the total rail volume. Based on these criteria, a list of 10 goods was selected. These goods demonstrate the highest transport performance among the identified categories and collectively account for approximately 74.69% of the total rail transport performance. Figure 3.1 illustrates these goods along with their respective percentage contributions.

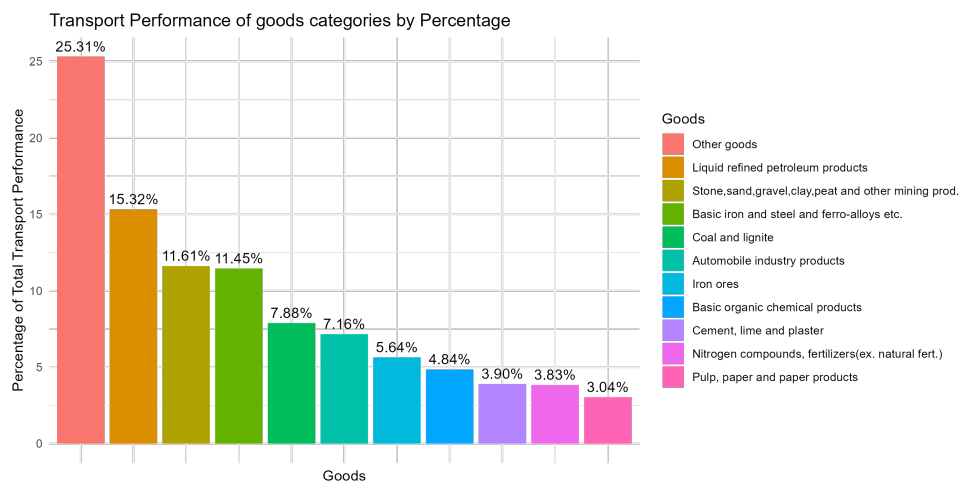


Figure 3.1: List of selected goods and their percentage in total transport performance

3.2 Descriptive statistics of the dataset

Descriptive statistics are essential for data analysis as they provide an understanding of a dataset's underlying characteristics and trends. Figure 3.2 shows the number of observations, mean, standard deviation, minimum, and maximum of the original dataset. The target variable, transport performance, has a minimum value of 22 and a maximum value of 1,607,624,680. The substantial range in these values is because transport performance is calculated as the product of the weight of transported goods and the transport distance.

Statistic	N	Mean	St. Dev.	Min	Max
Tpt_perfm	33,327	11,746,914.000	48,350,177.000	22	1,607,624,680
total_vol_ton	43,615	60,993.540	360,239.700	1	11,982,556
gdp_origin	234,000	271,528.600	195,219.000	27,245	577,123
gdp_destination	234,000	271,528.600	195,219.000	27,245	577,123
pro_ind_energy	239,330	102.104	6.579	83.508	108.100
pro_ind_mining	239,330	122.917	25.449	86.542	171.475
year_rail_volume	239,330	232,398.400	12,603.700	208,478.800	256,492.000
year_road_volume	239,330	2,945,875.000	109,369.500	2,764,804	3,104,055
year_water_volume	239,330	51,795.340	4,333.965	42,463.990	55,621.440

Figure 3.2: Summary statistics of the dataset

To inspect if outliers exist in these 2 variables, box plots were used to display the distribution of all numerical variables displayed in figure 3.3. While there are a few high values, it is difficult to conclude it's a typo error in the data or because of sudden increases in transport volume and performance. Therefore, these values will be kept in the dataset. For other numerical variables, their distribution was illustrated in figure 3.4. It can be seen that there are no outliers in those variables either.

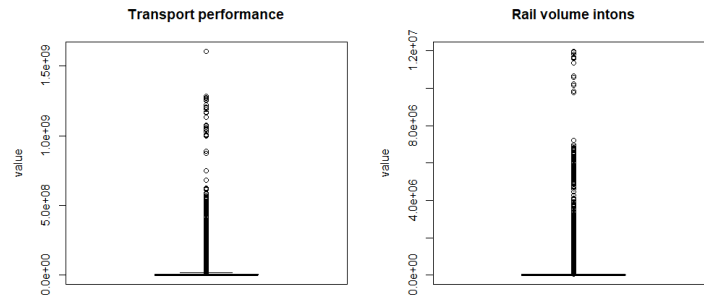


Figure 3.3: Box-plots for transport volume and performance

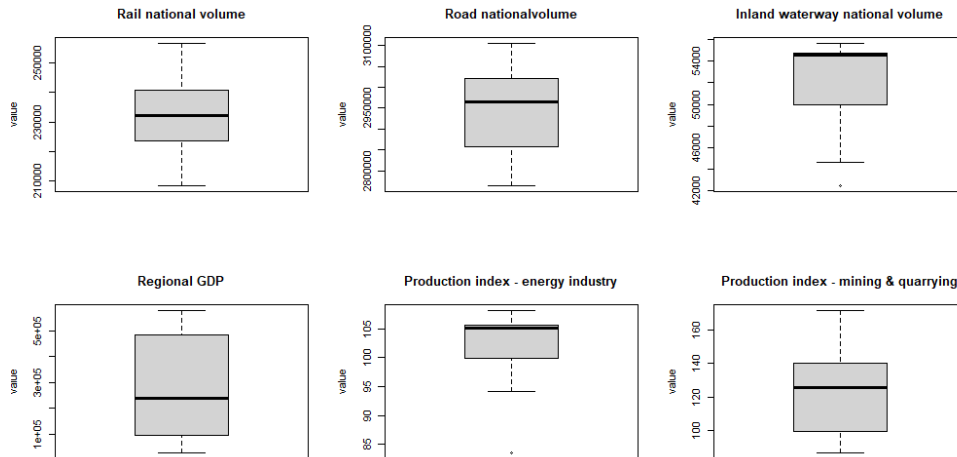


Figure 3.4: Box-plots for all other numerical variables

3.3 Correlation Analysis

The correlation matrix was conducted to understand relationships and potential collinearity between all features in the dataset. Figure 3.5 shows some interesting insights, where the target value (transport performance) highly correlates with the rail volume in tons, and has

small correlations with regional GDP values. Besides, the yearly water volume and production index of the energy industry have a highly positive correlation.

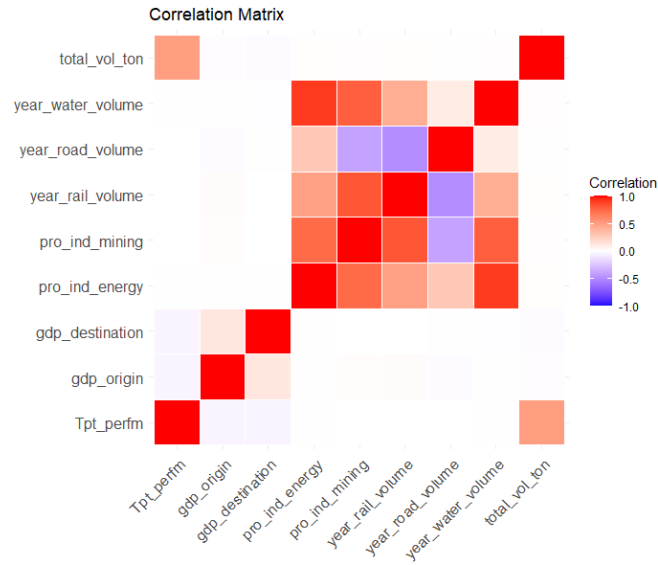


Figure 3.5: Correlation matrix of all columns in the dataset

3.4 Inspecting temporal patterns in transport performance

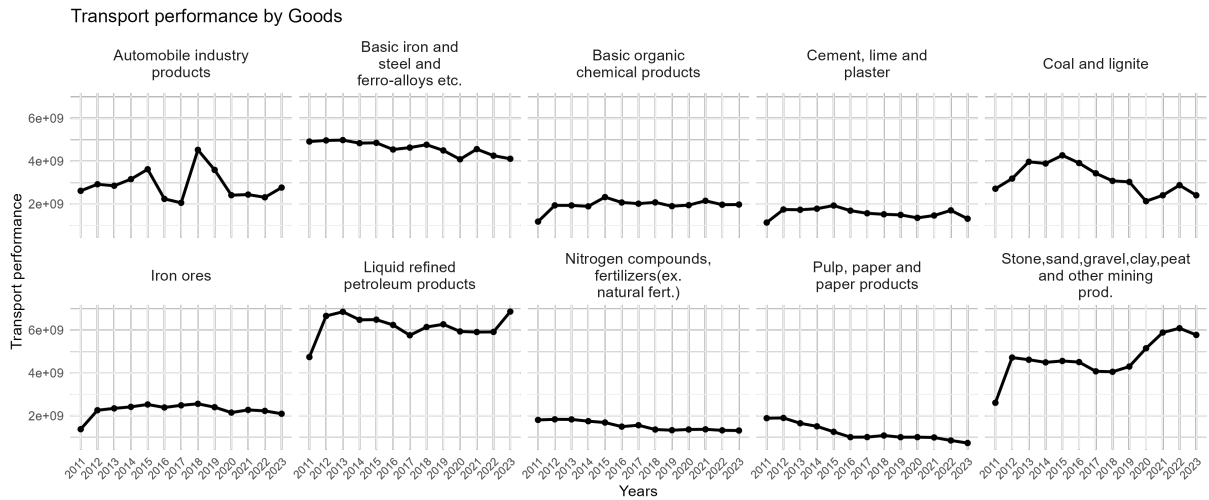


Figure 3.6: Change in transport performance by goods over 13 years

Figure 3.6 outlines the year-on-year changes in transport performance for each of the 10 selected goods. For example, when examining the trend for automobile industry products, the transport performance peaked in 2018 before experiencing a sharp decline, bringing it back to nearly the same level as in 2011 by 2023. It is important to note that the sudden increase observed in 2018 could be attributed to irregular factors. Specifically, the maximum value for transport performance, 1,607,624,680, was recorded for this product in 2018. This spike might represent an outlier or be influenced by external factors such as market demand fluctuations, and economic conditions.

3.5 Inspecting spatial patterns in transport performance

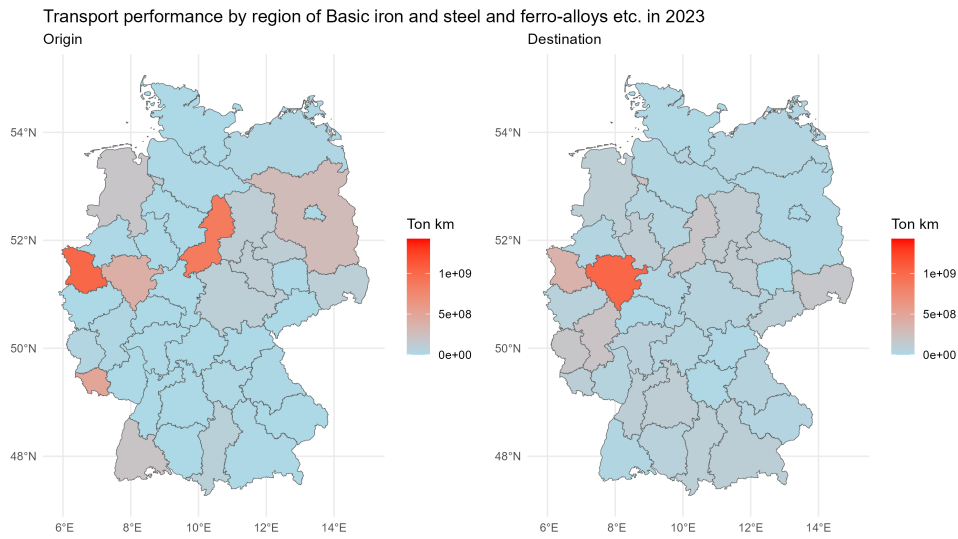


Figure 3.7: Transport performance in basic iron and steel and ferro-alloys etc. in 2023

Geographic information for Germany was sourced from GitHub, as documented by Schwarz et al. (2022). Transport performance values were visualized by distinguishing between origin and destination points on geographic maps. Figure 3.7 presents the transport performance for basic iron and steel and ferro-alloys etc. in 2023. High transport performance values are depicted in red, while low values are shown in blue. As an origin, Düsseldorf exhibits the highest transport performance, followed by Braunschweig. For destinations, Arnsberg records the highest transport performance, with products being delivered from various cities across Germany.

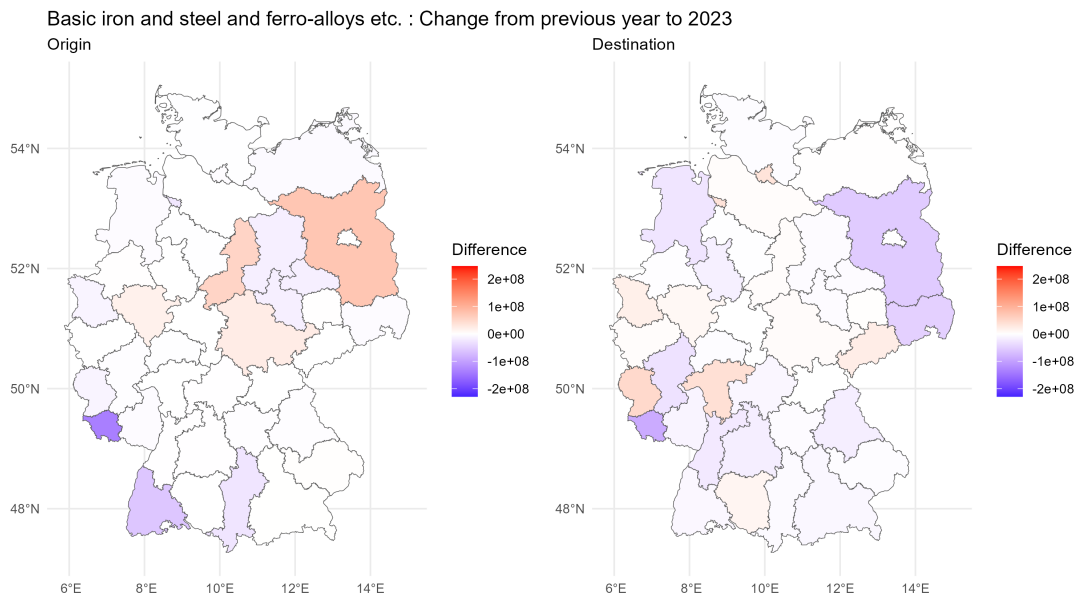


Figure 3.8: Difference in transportation performance in basic iron and steel and ferro-alloys etc. between 2022 and 2023

Figure 3.8 illustrates the changes in transport performance between 2022 and 2023. Regions with minimal change are depicted in white, while those with increased performance are shown in red and those with decreased performance in blue. As an origin, Saarland experienced the most significant decrease, whereas Brandenburg exhibited the highest increase. For destinations, Saarland again showed the greatest decrease, followed by Brandenburg.

3.6 Inspecting combined temporal and spatial patterns

One combined temporal and spatial aspect to investigate is to identify routes with consistent operation over the years, where transport activities occur annually for a specific set of origin-destination and goods. To examine this, a heatmap was created to count the number of years each good was actively transported along specific routes. Figure 3.9 is an example of an active transport route for coal and lignite across 13 years. The heatmap shows that Düsseldorf consistently exports to most regions over the 13-year period, while only a few importing regions, such as Sachsen-Anhalt, show continuous receiving activities across all 13 years. This visualization reveals key areas of goods flow, which can be areas with high levels of transport performance prediction accuracy.

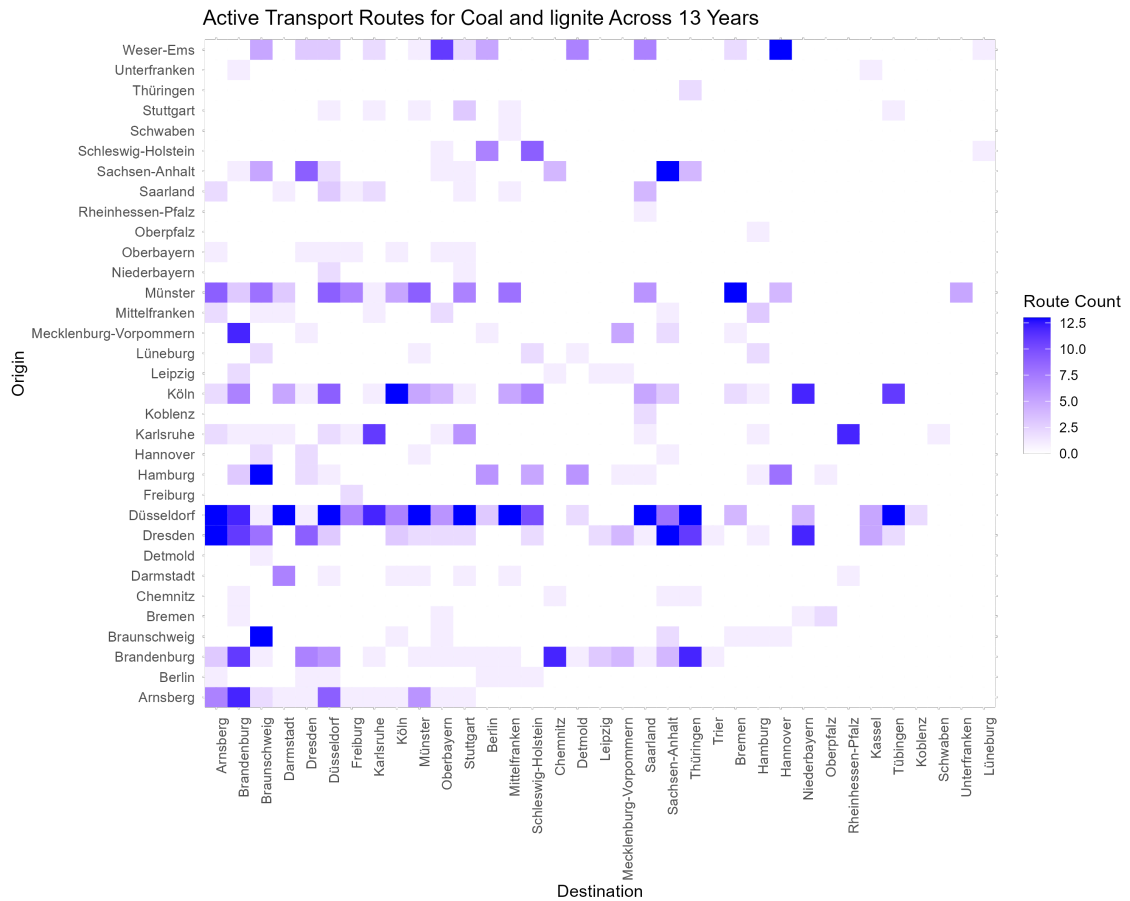


Figure 3.9: Heatmap plot of routes with consistent operation across all year for specific types of goods

4 Developing machine learning model

4.1 Data preprocessing before modeling

Following the insights from Chapter 2, the data was thoroughly preprocessed before modeling. The initial steps involved filtering out origin-destination sets for specific goods that did not operate consistently across all years, ensuring the reliability of the data for prediction. Only routes that have transport for all 13 years (from 2011 to 2023) were retained for the model. Additionally, all missing data points were removed to enhance the robustness of the subsequent analyses. The production index for the energy industry was excluded due to its high correlation with another existing feature, to prevent multicollinearity.

The total rail volume in tons is positively correlated with transport performance, which can be understood as transport performance is the product of volume in tons and distance in kilometers. This variable is the only one available at detailed temporal and spatial levels as the target variable. The remaining features are expected to explain the spatial aspects, allowing the machine learning model to accurately predict the target values.

After these preprocessing steps, the final dataset comprises 11,726 observations and 11 features. The data from 2011 to 2021 will be used for training and validating model performance during the feature selection and hyperparameter tuning phases, while the data for 2022 and 2023 are reserved for final model testing.

4.2 Preliminary model selection and validation

The task of predicting transportation performance is approached as a regression problem. This study explores various popular regression algorithms, with Linear Regression (LR) serving as the baseline model for comparative purposes. Other models, including Support Vector Machines (SVM), Neural Networks (NN), Decision Trees (DT), Random Forests (RF), and Extreme Gradient Boosting (XGB), were also evaluated. These advanced models are expected to outperform the baseline model as they are capable of capturing non-linear relationships, which were suggested by the correlation matrix, where many features showed non-linear relationships with the target values. For models such as SVM, NN, and XGB, which require numerical inputs, one-hot encoding was applied to all categorical variables.

Model performance was assessed using a set of metrics: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), and Coefficient of Determination (R-sqrt). These metrics were chosen for their complementary perspectives: RMSE

penalizes large errors more heavily, MAE provides a straightforward measure of average error magnitude, and R-sqrt quantifies the proportion of variance in the dependent variable explained by the model, and MAPE offers an accuracy measure in percentage terms, which is particularly useful for benchmarking against industry standards. Typically, a MAPE of around 50% is considered acceptable in industry settings, while in volatile sectors like transport, a MAPE around 20% can still indicate a strong model (Statology 2020), (Roberts 2023). The performance of different models is documented in Table 4.1.

Table 4.1: Performance metrics for different learners

Learner	RMSE (tkm)	MAE (tkm)	MAPE (%)	R-Squared
LR	4.525680×10^7	2.063403×10^7	274.3498	0.3815254
DT	2.552973×10^7	1.156166×10^7	78.18046	0.8031903
RF	1.136766×10^7	0.437979×10^7	24.18576	0.9609792
SVM	2.999978×10^7	0.995818×10^7	56.69612	0.7282371
NN	4.610986×10^7	1.978308×10^7	189.3553	0.3579900
XGB	0.957649×10^7	0.293481×10^7	5.132887	0.9723072

The results show that all models, except for Neural Networks, outperform the baseline model. Among the evaluated models, Gradient Boosting and Random Forests emerged as the top performers, particularly in terms of MAPE metrics. These models demonstrated acceptable levels of accuracy and will be subjected to further tuning and inspection in the subsequent sections.

4.3 Nested resampling and tuning hyperparameters

This section focuses on enhancing the performance of the Random Forest (RF) and Extreme Gradient Boosting (XGB) models through nested resampling and hyperparameter tuning. Both the inner and outer loops of the nested resampling process utilize a cross-validation strategy with a few number of folds, a choice driven by the size of the training dataset. For hyperparameter tuning, a random search strategy is employed. Given the current satisfactory MAPE and R-sqrt metrics, MAE is selected as the optimization metric to minimize errors. A comprehensive literature review informed the selection of hyperparameters for each model. The tuning range was carefully chosen to strike a balance between model accuracy and the risk of overfitting.

Table 4.2: Range of tuning and optimal configuration of 2 learners

Learner	Hyperparameter	Tuning range	Tuning result
RF	num.trees	[100, 1000]	140
	mtry	[3, 10]	8
	min.node.size	[1, 20]	2
	max.depth	[1, 30]	22
XGB	nrounds	[100, 1000]	319
	eta	[0.01, 0.3]	0.116596
	max_depth	[3, 10]	9
	subsample	[0.5, 1]	0.8763066
	colsample_bytree	[0.3, 1]	0.778358

For the Random Forest model, the `mtry` parameter, which controls the number of features considered for splitting at each node, was tuned within a range from \sqrt{p} to p , where p is the total number of features. This range is recommended to mitigate the risk of overfitting

(Breiman 2001). Additionally, the number of trees (`ntree`) and the maximum depth of the trees (`max_depth`) were tuned according to recommendation provided by Probst, Wright, and Boulesteix (2019). The minimum size of terminal nodes (`min_node_size`) was set to range from 1 to 20, which is a common practice in regression tasks (James et al. 2013).

For the XGB model, the tuning range for key hyperparameters such as the number of boosting iterations (`nrounds`), maximum tree depth (`max_depth`), subsample rate (`subsample`), and the proportion of features used (`colsample_bytree`) were chosen based on recommendations by Chen and Guestrin (2016) to balance model complexity and generalization. Additionally, the learning rate (`eta`) was tuned within a range of 0.01 to 0.3 to facilitate gradual learning and reduce the risk of overfitting (Hastie, Tibshirani, and Friedman 2009). Table 4.2 presents the selected hyperparameter ranges and the optimal configurations obtained after tuning.

Table 4.3: Performance metrics before and after tuning hyperparameter

Learner	RMSE (tkm)	MAE (tkm)	MAPE (%)	R-Squared
Before				
RF	1.136766×10^7	0.437979×10^7	24.18576	0.9609792
XGB	0.957649×10^7	0.293481×10^7	5.132887	0.9723072
After				
RF	1.018253×10^7	0.3355376×10^7	1.032573	0.9686913
XGB	0.791833×10^7	0.2638919×10^7	6.427276	0.9810669

According to table 4.3, the performance of the RF model improved across all metrics, especially in terms of MAPE after hyperparameter tuning, indicating enhanced predictive accuracy and consistency. Besides, while the XGB model also showed improvements across most metrics, there was a slight increase in MAPE. This suggests a potential trade-off, where the model reduced larger errors but may have slightly increased smaller ones, affecting overall prediction consistency.

4.4 Nested resampling and feature selection

Feature selection for the RF model poses challenges due to the relatively small number of features and constrained learner configuration. Conversely, the XGB model has a larger feature set due to the one-hot encoding of categorical variables. To explore whether the performance of the tuned XGB model could be further enhanced, nested resampling with feature selection was conducted. Both the inner and outer loops of the nested resampling used a cross-validation strategy with a few number of folds. Feature selection was carried out using a wrapper strategy with forward selection, which aims to identify an optimal subset of features for the model, ensuring unbiased results. MAE was used as the metric for improvement.

Table 4.4: Performance metrics before and after feature selection

	RMSE (tkm)	MAE (tkm)	MAPE (%)	R-Squared
Before	0.791833×10^7	0.263892×10^7	6.427276	0.9810669
After	6.134832×10^7	1.319488×10^7	3.956558	-0.1364743

The optimal subset of features identified for the XGB model included only the Origin from Saarland and the freight volume in tons. Table 4.4 details the model's performance before and after feature selection. The results indicate that after feature selection, the XGB model's performance deteriorated, with higher errors and a significant drop in its ability to explain the variance in the data. This outcome suggests that the selected feature subset may have been

overfitted to the training data, resulting in an overly simplistic model that fails to generalize well to unseen data.

4.5 Model inspection and interpretation

Best performing models, including random forest and extreme gradient boosting models after hyperparameters tuning, are inspected for further insights and understanding. Figure 4.1 shows the scatter plots of predicted and target values for both models. It can be seen that in RF models, there are some larger errors than in XGB models, depicted by points lying far away from the red line. However, overall, both models explained the data rather well.

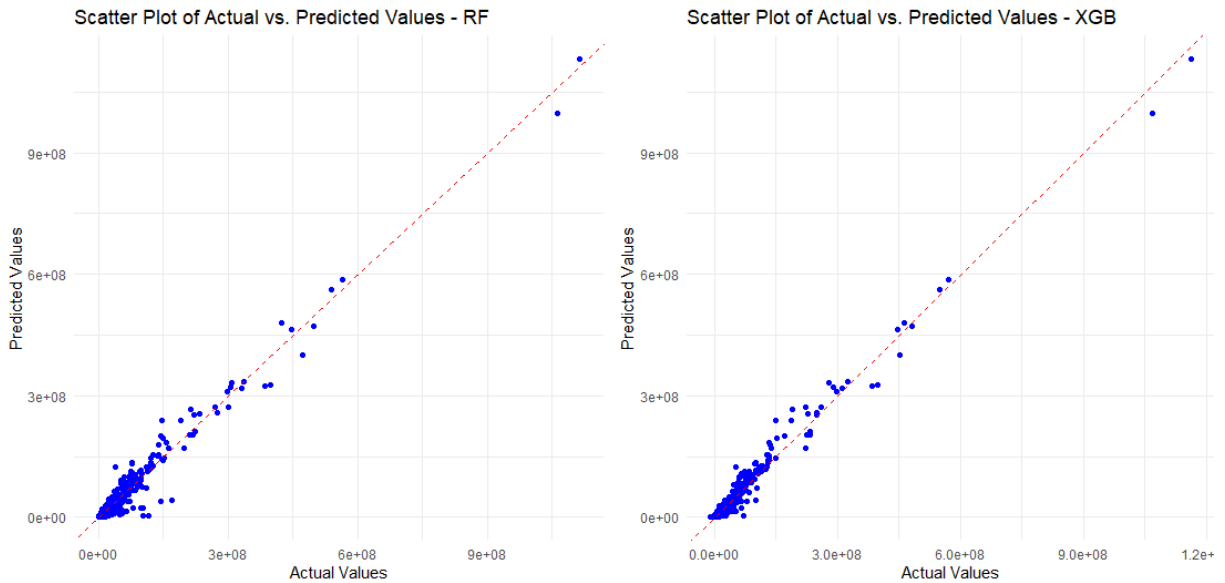


Figure 4.1: Scatter plots of predicted and target values for both models

To interpret how these models function, methods such as feature importance and feature effect plots were employed. Figure 4.2 presents the ranking of feature importance in the tuned RF model. In this model, rail freight volume in tons emerges as the most important feature, while the production index and yearly total volumes of rail, road, and inland water transport are the least important features.

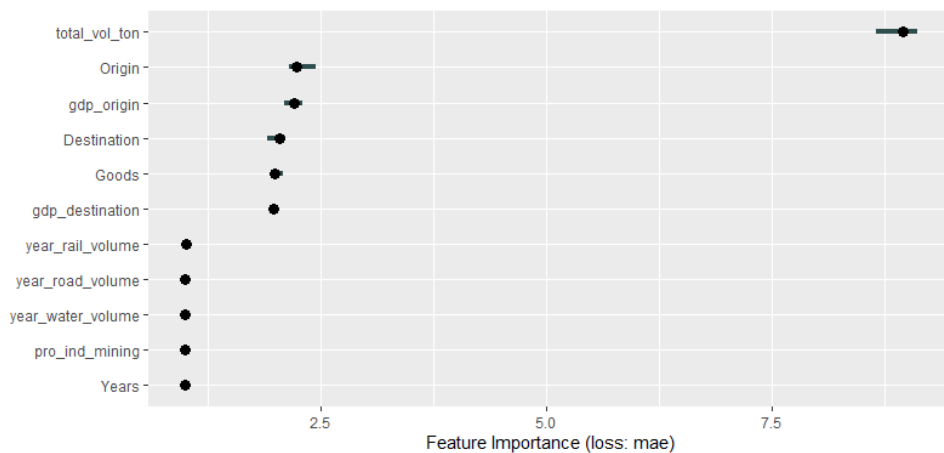


Figure 4.2: Feature importance plot of RF

For the XGB model, the number of feature columns is much larger due to the one-hot encoding of categorical variables. To facilitate visualization, the top 5 most and least important features are highlighted in 4.3. Similar to the RF model, total rail freight volume and regional GDP are the most predictive features in the XGB model, whereas the production index and yearly volumes of road and inland water transport are among the least important. Additionally, certain types of goods and specific regions for loading and unloading are more predictive than others.

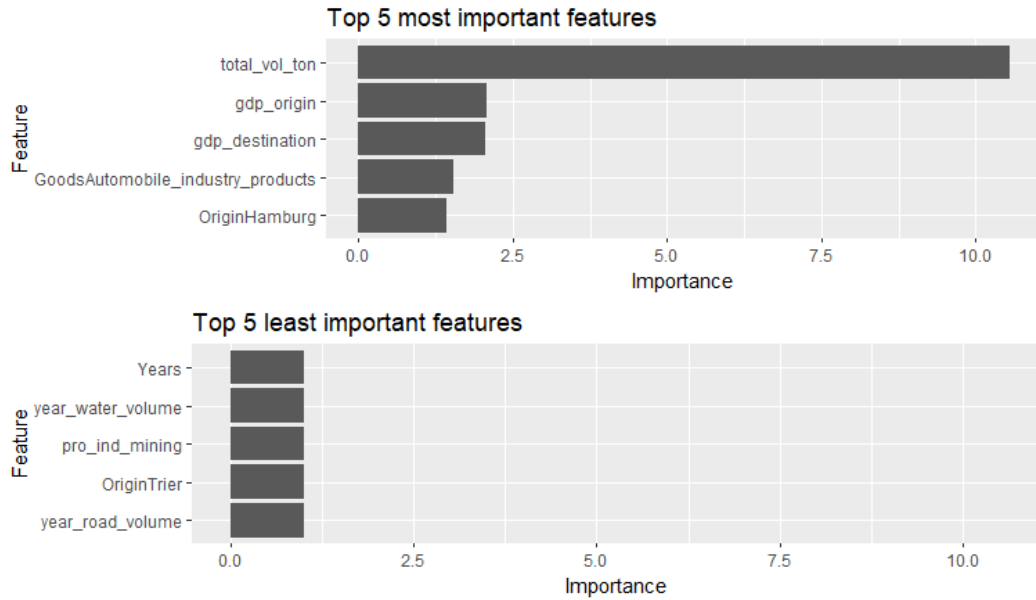


Figure 4.3: Top 5 most and least important features of XGB

The Individual Conditional Expectation (ICE) and Partial Dependence Plot (PDP) in figure 4.4 offers several key insights into the behavior of the model with respect to the rail freight volume in tons feature.

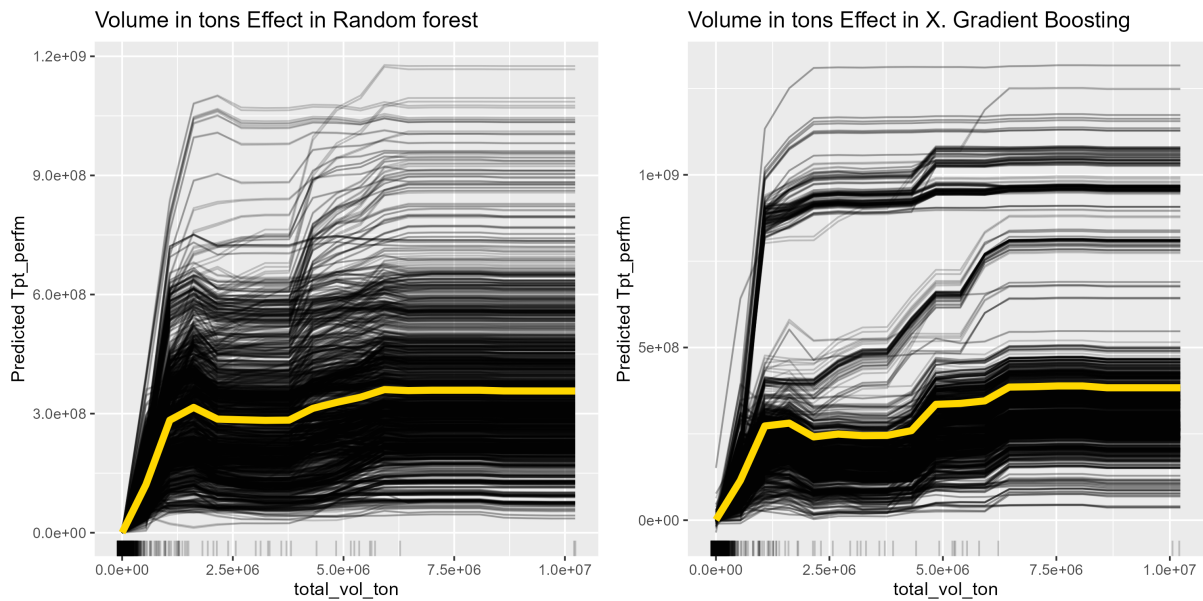


Figure 4.4: Feature effect plots of RF and XGB

The yellow PDP line indicates a predominantly positive relationship between the feature and

the target variable within the main distribution range of target values (under 3×10^8 tkm). This suggests that as the feature increases, the predicted target value also increases within this range, reflecting a general trend of positive association. However, beyond this threshold, the PDP line exhibits a slight decline followed by a plateau, indicating that the effect of the feature diminishes at higher levels.

The ICE plots further complement this analysis by revealing individual-level variability in the feature's impact on the target variable. While the majority of the ICE lines follow a pattern consistent with the PDP, there are certain ICE lines, represented in black, deviate from the general trend by exhibiting significantly higher target values and a different trajectory compared to the majority of the instances. These separate lines indicate that for certain data points, the feature's effect on the target variable is markedly different, potentially due to interactions with other unobserved features or specific conditions unique to these instances. This variability points to potential heterogeneity in the dataset, where certain data points are more sensitive to changes in Rail_vol_ton, indicating the need for further investigation into these interactions.

4.6 Overall model assessment and limitations

The models, particularly those utilizing Extreme Gradient Boosting (XGB) and Random Forest (RF) algorithms, exhibit strong performance in predicting transportation outcomes, as evidenced by metrics such as MAPE and R-sqrt. Hyperparameter tuning has enhanced the models' ability to capture overall trends, and feature importance analysis has provided insights into the relationships between various predictors and target values. However, the high MAE and RMSE values indicate that, despite their general effectiveness, these models may struggle with finer precision, likely due to the complex interactions within the data that are not fully captured.

The ICE plots further highlight the variability in how freight volume (measured in tons) influences predictions, suggesting that the models might miss some of the nuanced relationships within the data. Additionally, there are limitations to the developed models, including potential selection bias from excluding certain data before modeling and the possible loss of predictive power due to insufficient data granularity. The variability observed in the ICE plots also indicates that these models may not fully generalize across all scenarios. These factors emphasize the need for cautious interpretation of the results and suggest that future analyses should consider incorporating additional data or more complex modeling techniques to better capture the intricacies of the relationships involved.

5 Conclusion

This study provides a comprehensive analysis of yearly rail freight in Germany, encompassing the selection of relevant predictors, data visualization, modeling, and the application of basic machine learning methodologies. The dataset, compiled from a general literature review and public sources, has broad relevance to similar transportation challenges. The visualization techniques employed are particularly suited for analyzing transport data with complex spatial and temporal patterns. Additionally, the study includes a benchmark comparison of various machine learning algorithms, offering insights into their performance on this specific problem. It also explores foundational methods for enhancing and interpreting these models within the context of rail freight transport.

However, the study has certain limitations and areas for improvement. The granularity of some features is limited due to data aggregation at yearly intervals or over larger regions, potentially obscuring more detailed insights. Additionally, certain relevant features could not be included due to the lack of long-term datasets. The machine learning models developed for prediction were based on basic methodologies and may not fully capture the complex interactions within the data. Furthermore, data preprocessing steps, such as filtering specific routes or goods, might obscure significant insights into rail transport activities.

Future work could improve upon this study by incorporating a greater number of features and datasets, particularly those specific to smaller regions and industries. With more granular data, time series methods could be applied to enhance predictive accuracy. Moreover, as this study primarily employs basic methodologies, there is substantial potential to improve model performance and derive deeper insights by applying more advanced techniques.

Bibliography

- Breiman, L. (2001). "Random Forests". In: *Machine Learning* 45.1, pp. 5–32. DOI: 10.1023/A:1010933404324.
- Chen, T. and Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System". In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, pp. 785–794. DOI: 10.1145/2939672.2939785.
- Destatis (2024a). "Goods carried (freight transport by rail): Germany, years, region of loading, region of unloading, classification of goods (groups)". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=46131-0013&bypass=true&levelindex=0&levelid=1724230688794#abreadcrumb>.
- Destatis (2024b). "Goods carried, transport performance (freight transport by rail): Germany, years, main traffic relations". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=46131-0003&bypass=true&levelindex=0&levelid=1724230958314#abreadcrumb>.
- Destatis (2024c). "Goods carried, transport performance (inland waterway transport): Germany, years, main traffic relations". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=46321-0001&bypass=true&levelindex=0&levelid=1724231098004#abreadcrumb>.
- Destatis (2024d). "Index of production in manufacturing: Germany, months, original and adjusted data, economic activities (main groups and aggregates)". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=42153-0001&bypass=true&levelindex=1&levelid=1722644947372#abreadcrumb>.
- Destatis (2024e). "National accounts of the states (production accounts) - Gross domestic product at market prices (nominal): states, years (data BV4.1 trend for seasonal, calendar and trend adjusted)". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=82111-0001&bypass=true&levelindex=0&levelid=1724231409642#abreadcrumb>.
- Destatis (2024f). "Quantity of goods carried, transport performance (freight transport by road): Germany, years, mode of transport, transport categories". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=46231-0003&bypass=true&levelindex=0&levelid=1724231053812#abreadcrumb>.
- Destatis (2024g). "Transport performance (freight transport by rail): Germany, years, region of loading, region of unloading, classification of goods (groups)". In: URL: <https://www-genesis.destatis.de/genesis/online?operation=table&code=46131-0015&bypass=true&levelindex=0&levelid=1724230688794#abreadcrumb>.
- Hastie, Trevor, Tibshirani, Robert, and Friedman, Jerome (2009). "Overview of Supervised Learning". In: *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York,

- NY: Springer New York, pp. 9–41. ISBN: "978-0-387-84858-7". DOI: 10.1007/978-0-387-84858-7_2.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer. DOI: 10.1007/978-1-4614-7138-7.
- Probst, P., Wright, M. N., and Boulesteix, A. L. (2019). "Hyperparameters and Tuning Strategies for Random Forest". In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*. DOI: doi:10.1002/widm.1301.
- Roberts, Amber (2023). "Mean Absolute Percentage Error (MAPE): What You Need To Know". In: URL: <https://arize.com/blog-course/mean-absolute-percentage-error-mape-what-you-need-to-know/>.
- Schwarz, Francesco, R, Jannis, Stanislav, and Mikalai (2022). "deutschlandGeoJSON". In: GitHub. URL: https://github.com/isellsoap/deutschlandGeoJSON/blob/main/3_regierungsbezirke/3_mittel.geo.json.
- Statology (2020). "What is Considered a Good Value for MAPE?" In: URL: <https://www.statology.org/what-is-a-good-mape/>.