

Data Mining Project Outline HWS24

# **Gradient Descent Gururs - Predictive Modelling and Clustering of Train Delays**

Anh Tu Duong Nguyen<sup>2115931</sup>

October 13, 2024

Submitted to  
Data and Web Science Group  
Prof. Dr. Hertling  
University of Mannheim

# 1 Introduction

Railway systems play a crucial role in public transportation, and maintaining on-time train schedules is essential for operational efficiency and passenger satisfaction. In recent years, delays have become a growing concern for Deutsche Bahn, leading to disruptions that affect commuters and long-distance travelers alike. To address this, we aim to create a predictive model that categorizes train delays based on various factors, offering more detailed insights into the nature and severity of delays. By focusing on delay categories, we can provide a more structured and practical prediction framework. Our project aims to address this issue by developing a predictive model that categorizes future train delays, providing insights into the extent of the delay. Using historical data from Deutsche Bahn, we will analyze variables such as location, train type (e.g., RE, RB, S-Bahn), and additional factors like weather or time of day. From an academic perspective, this project will contribute to the field of predictive analytics, particularly in transportation systems. The findings will not only advance data mining techniques but also offer a practical application of machine learning models to real-world problems. From a business perspective, the ability to predict delays could enable Deutsche Bahn to implement proactive measures to minimize disruptions. By understanding patterns in delays, Deutsche Bahn can allocate resources more effectively, optimize scheduling, and ultimately improve customer satisfaction. Furthermore, passengers may benefit from more accurate information about expected delays, helping them plan their journeys better. In this sense, our project offers both an academic challenge and the potential for a meaningful business impact.

## 2 Methodology

This chapter outlines the methodology that is pursued for our group project. Section 2.1 explains the overall process on compiling the data for our project. Section 2.2 provides an overview on how we want to tackle the problems we have stated in Chapter 1.

### 2.1 Data Collection Methodology

While there is an existing dataset on Kaggle<sup>1</sup>, we want to propose that we crawl additional data using the same APIs for the existing dataset: Timetables and StaDa API from the DB API marketplace. The reason for this is that the existing dataset is limited to only one week in July, ranging from 08.07.2024 to 14.07.2024 and we would like to enrich this dataset with more recent data to achieve better performance as more data can help making a much more accurate analysis. However, if the crawling process turns out to be more difficult than we initially assumed, then we will proceed our project using the existing dataset. The dataset will include the following data:

- Station
- Line
- Path
- Size of Station
- State
- City
- Zip Code
- Planned Arrival
- Planned Departure
- Actual Arrival
- Actual Departure
- Announcements, i.e. construction on the track

---

<sup>1</sup><https://www.kaggle.com/datasets/nokkyu/deutsche-bahn-db-delays>

## 2.2 Data Mining Methodology

Due to the fact that we intend to crawl our data using the DB APIs, it becomes evident that there are several pre-processing steps required before the main work can start. At first we need to aggregate all data in a single dataset such that it matches the existing dataset. Some additional potential steps are:

1. Transformation, i.e. transforming discrete classes into numerical values
2. Binning according to regions, stations etc. for following clustering
3. Potentially detecting and removing outliers

In this work, we intend to leverage the property that Regression algorithms can be used for classification problems such that we can provide a prediction with a specific certainty for passengers. To predict the delay, we intend to use the interpolating Regression Trees to give a reasonable prediction taken from the training dataset interval. However, to compare our results we also aim to make use of other conventional methods for classification like Decision Trees and Nearest Centroids. If sufficient computing power is present, then it will be of great interest if the problems can be solved using Neural Networks.

## 3 Evaluation

The delay prediction model will be evaluated by partitioning the dataset, with 20% reserved for testing. To avoid overfitting, 10-fold cross-validation will be applied to 80% of the data. Performance will be primarily measured using the F1 score, which balances precision and recall. A confusion matrix will offer insight into misclassifications between delay categories, such as between medium and short delays. Additionally, a cost analysis will be conducted to assess the impact of these misclassifications, with a focus on high-impact errors.

## 4 Expected Results

The project aims to develop a model that accurately predicts train delays by category and duration, focusing on regional and express trains. While perfect accuracy is not expected, the model is intended to provide reliable predictions that inform decision-making in train scheduling. The expected results include identifying patterns that link delays to specific factors, such as stations, times of day, or regions. Certain stations may be identified as bottlenecks with more frequent or prolonged delays, while rush hours or high-traffic regions may show time-based delay trends. External factors, such as holidays or local events, may also contribute to longer delays in certain areas.

# Ehrenwörtliche Erklärung

Ich versichere, dass ich die beiliegende Bachelor-, Master-, Seminar-, oder Projektarbeit ohne Hilfe Dritter und ohne Benutzung anderer als der angegebenen Quellen und in der untenstehenden Tabelle angegebenen Hilfsmittel angefertigt und die den benutzten Quellen wörtlich oder inhaltlich entnommenen Stellen als solche kenntlich gemacht habe. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen. Ich bin mir bewusst, dass eine falsche Erklärung rechtliche Folgen haben wird.

**Declaration of Used AI Tools**

Tool	Purpose	Where?	Useful?
ChatGPT	Rephrasing	Throughout	+
DeepL	Translation	Throughout	+
Dall-E	Image generation	Figs. 2, 3	++
GPT-4	Code generation	functions.py	+
ChatGPT	Related work hallucination	Most of bibliography	++

Unterschrift

Mannheim, den 13.Oktober.2024