# Ordina train delay prediction

This project aims to make a model that predicts if there is a delay or not and how long the delay is going to be. The problem it is trying to solve is that sometimes people get stuck on the train because of the delays and with our idea people would be able to make a choice to take the car by looking at the predictions before going on the train. How realistic this idea is, is something that we will find out along the way.

Created by: e.bijker
Created on: June 7, 2022 11:40 AM
Changed on: June 9, 2022 11:36 AM

Context of use: Education
Level of education: Bachelor

## Impact on society
What impact is expected from your technology?

### What is exactly the problem? Is it really a problem? Are you sure?
Ordina has employees that need to travel all over the Benelux, these employees are likely to be reliant on public transport. Ordina tasked us with finding a way to improve the efficiency of their employees' travel, using machine learning. Ordina linked us to a public archive of train data for a head start.
We decided to use this train data to try to predict train departure delays in the Netherlands, Ordina will receive a trained model but no implementation thereof.

### Are you sure that this technology is solving the RIGHT problem?
Train departure delays can be caused by any multitude of reasons; weather, malfunctioning systems, pilot negligence, a long boarding process during rush hours, etc.

We cannot address these causes, as we have no data detailing the cause of each delay.

In this way, we don't address the causes at all and leave it to the algorithm to find patterns that might be used for predicting when these delays might happen and how long they might be.

### How is this technology going to solve the problem?
To improve Ordina employee travel we decided to attempt a machine learning prediction of train departure delays within the Netherlands. Delays don't happen often so we wanted to predict the big and impactful delays.
We are not aware of how such a prediction might impact Ordina employees' travel decisions, but we could provide advice to Ordina employees based on these predictions

### What negative effects do you expect from this technology?
Providing this delay prediction model to Ordina is unlikely to have a significant impact, be it positive or negative. If it did, however:

In the best-case scenario, we see Ordina employees using our model to build an application that provides their employees with train travel advice so that they might avoid routes with frequent, impactful delays

In the worst-case scenario we see Ordina using our delay prediction model to improve their employees' travel efficiency so much that given enough time, they might put their competitors in the Benelux out of business.

While not necessarily a bad effect, this could very well impact the IT consultancy landscape. For example, a customer in the future might have very few options for an IT consultant that fits their needs better or be forced to work with Ordina despite previous conflicts.

**In what way is this technology contributing to a world you want to live in?**
In making a model that can predict delays in public transport, we help keep commuters stay informed of possible mishaps during their commute and even how they might avoid them.

Visualizing the trends we found with delays, could also help researchers or public transport workers to investigate how these delays are caused and might be prevented.

**Now that you have thought hard about the impact of this technology on society (by filling out the questions above), what improvements would you like to make to the technology? List them below.**
To make our model more useful and trustworthy.
We would've liked to take the cause of delays into account. We would've needed more up-to-date data that is representative of the current state of the train infrastructure.
We would've also needed data from other countries' train transport.
We would then also need to find a solution for the natural imbalance in the data, delays that don't impact travel occur much more often than big impactful delays. It might be that the algorithm in the end simply cannot predict big unexpected delays because it couldn't find enough patterns with them to make a prediction

## Hateful and criminal actors
What can bad actors do with your technology?

*This category is not applicable for this technology.*

## Privacy

Are you considering the privacy & personal data of the users of your technology?

*This category is not applicable for this technology.*

## Human values

How does the technology affect your human values?

*This category is not applicable for this technology.*

## Stakeholders
Have you considered all stakeholders?

*This category is only partial filled.*

**Who are the main users/targetgroups/stakeholders for this technology? Think about the intended context by answering these questions.**

**Name of the stakeholder**
Ordina employees

**How is this stakeholder affected?**
Ordina employees are going to be using the potential application with the data from the Jupiter Notebook. The user of the final product can use it for travelling to work in a more reliable way. Given that the models are optimized at a satisfactory level they can be used to give predictions to the user of whether there will be a train delay and how long will be the duration of that delay. The added value for them would be are effectiveness and efficiency in their commute.

**Did you consult the stakeholder?**
Yes

**Are you going to take this stakeholder into account?**
Yes

**Name of the stakeholder**
NS

**How is this stakeholder affected?**
Data from travel with their train company is the majority of the data. Maybe they can contribute to optimizing the product by providing regular updates on the data for train travelling.

**Did you consult the stakeholder?**
No

**Are you going to take this stakeholder into account?**
No

**Did you consider all stakeholders, even the ones that might not be a user or target group, but still might be of interest?**
-

**Now that you have thought hard about all stakeholders, what improvements would you like to make? List them below.**
We would include a user manual of what data the potential app can use. In the beginning if the Jupiter notebook is the only thing used a manual can be done for that as well so Ordina employees know how to operate with it and what information they can retrieve from it.

## Data
Is data in your technology properly used?

### Are you familiar with the fundamental shortcomings and pitfalls of data and do you take this sufficiently into account in the technology?

We faced a few issues when we were working with the data. The biggest one was that the data was not balanced. There were much more train departures without delays than with delays. Also, the data was from 2016, taking into account the events that happened only in the last 2 years, we faced a pandemic, war and similar events that may disturb the train data with missing values or unpredictable behaviours. Because of that, it is difficult to predict delays accurately in real life, especially with old data. We also think that the data is just not enough to predict accurately if there is a delay or not since there are not enough features to get a score higher than 80%

### How does the technology organize continuous improvement  when it comes to the use of data?

The data is extracted from large XML files, using element names, such as train delay, departure station, etc. If the data changes in the future, if the names of the element in the XML are the same it will still get the correct data. In the beginning, we extracted the data using indexes, which may be a problem if changes in the data appear, since if one new column is added we need to change the indexes manually. To extract the data we needed a powerful computer and time, so it could be stored in a dataset. We can not just use the link from the data provider directly in our notebook. That means if new data is added, it should be manually converted and stored to a readable dataset and then used in the notebook. In conclusion, the notebook does not handle new data automatically.

### How will the technology keep the insights that it identifies with data sustainable over time?

The data we were provided from the company is public and can be accessible to everyone. It is stored in archives since it is from 2016. There is a chance that the data may be not accessible in the future. For example, the website that provides it may stop its hosting, because of cost issues and if the data is not transferred/saved it may be lost. Another scenario is that NS for example decides to make its data private and since they are the biggest train company in the Netherlands, a lot of data will be lost again. In case of the data is lost/deleted from the website, we have a copy of the data stored in our local drives, which is also used for the algorithms. Therefore we will always have the data we need for the model prediction, even if it got deleted from the website and the notebook will run without problems.

# Technology Impact Cycle Tool

**Ordina train delay prediction**

---

**In what way do you consider the fact that data is collected from the users?**

The data does not contain any user's personal information. It contains only the train departures information, including stations, time of a journey, delay, etc

**Now that you have thought hard about the impact of data on this technology, what improvements would you like to make? List them below.**

One of the important things we may do to improve our work is looking for other datasets, that can enrich the one we use. Now we don't have confident model results and we concluded that the dataset features are just not enough. Also, we can look at other data sources, since the data we currently have is highly unbalanced and sampling methods do not work that well. That way we can add more rows and balance the departures with and without delays. Another thing to improve is the handling of the data. We could make the extraction of data easier and more automated, instead of manually extracting every new dataset.

## Inclusivity
Is your technology fair for everyone?

### Will everyone have access to the technology?
We are trying to make a model that is able to predict train delays for Ordina. If they want to implement it, it would first be available for only their employees. So the short answer is no. But if the model is very successful in their company there is a chance that they will either sell it or create an application for the public.

### Does this technology have a built-in bias?
Our model or does not have a bias in the sense of being racist or anything but the data is unbalanced in the sense that there are overall more no delays than delays, there are more ns trains in there etc. but those biases are there in real life. the fact is that there are more no delays than delays (luckily otherwise it would kind of mean our train system isn't doing a good job). it is nothing that the people who collect changed.

### Does this technology make automatic decisions and how do you account for them?
Our AI model does automatically make decisions based on the max number of stops, time group (before rush hour, morning rush hour, between rush hours, afternoon rush hour, after rush hour), day of the week and departure station. To make sure that the decision is not biased we undersampled the data so the delay and no delay group is balanced. For now, we can verify our predictions with the actual data if there was a delay or not but if it was implemented the only way to verify it is to go into the train and see if the prediction is correct.

### Is everyone benefitting from the technology or only a a small group? Do you see this as a problem? Why/why not?
Since this model would first be only available for Ordina employees, there will indeed only be a small group benefitting from this model. But if the app would be very successful there is a chance that Ordina will either make it public or sell it for public use to NS for instance. Overall I don't think this would be a problem in the beginning (that only Ordina employees have access), but after they find out that it is working it would be best to make it public. Making such a thing public will help a lot of people, but after such an application is public, there will probably be a lot more trains nearly empty. So if you look at it that way there is a chance that if such a thing gets public train companies would suffer from it. But other forms of transportation will benefit from it. If people then decide to take the car instead there is a chance that it would have a negative impact on the climate/nature.

### Does the team that creates the technology represent the diversity of our society?

The people/team that made this project started with: 1 Dutch woman (white), 1 Dutch man (white), 1 Dutch man (mixed), 1 Portuguese-Dutch man (white), 1 Bulgarian woman (white), 1 Bulgarian man (white). while doing the project 2 people dropped the project (the 2 Dutch men (white & mixed)).
In conclusion, the team has a mix of different cultures but they are all western, the majority is white and in the beginning it was male dominant but that changed to 50/50 male/female.

The target group are Dutch male and females working for Ordina. Because all of the people developing the model are students, we don't have a big connection with the target group age-wise. Also, nationality wise, most of the people at the end are not (fully) dutch or from the Benelux, so that is probably also different from the target group (Assuming that most of the employees are dutch or from the Benelux). So we don't really represent the target group in those areas, but I do think we will be able to make an inclusive model. Also since the model now only works with train data it can not be racist.

### Now that you have thought hard about the inclusivity of the technology, what improvements would you like to make? List them below.

technology-wise there is nothing that could affect the inclusivity, but it would be better if we changed the target group to people using the trains overall instead of only the ones that work for Ordina. It would also be better if we had a team that was more representative of the target group and if we had some more people of colour and with different experiences, backgrounds, and cultures.

## Transparency
Are you transparent about how your technology works?

### Is it explained to the users/stakeholders how the technology works and how the business model works?
The technology we will deliver to the stakeholders is well explained and it gives a clear idea of what the project is about. We wrote a project proposal, including all important headers, which makes the understanding of it very easy for anyone. In addition, we put descriptions, conclusions and comments on every code we put in the notebooks so that everyone who reads it will understand what the code/graph is about. We are preparing a short presentation before each meeting with the stakeholders, where we put the most important information and progress we did in clear, easy to read slides.

### If the technology makes an (algorithmic) decision, is it explained to the users/stakeholders how the decision was reached?
We used 5 different machine learning models in our project. Each model is well explained with a reason why we chose it, a description of what the model does and a conclusion of the results. We did not explain much about why we used certain features for the models, but we chose the ones that give higher accuracy and are logical. Also, we use algorithms such as nearest neighbours, which are easy to explain and conclude, but also algorithms like clustering or random forest classifier, which are quite hard to explain how the results are formed. Therefore we described the results in the most understanding way and it remains clear to the stakeholder why and how the certain model is used.

### Is it possible to file a complaint or ask questions/get answers about this technology?
It is not possible to file a complaint or ask questions about the technology at this point. The users do not need to enter personal data and they will only use the application, which lowers the chances of users looking for support. In the future, a support line may be created for technology complaints.

### Is the technology (company) clear about possible negative consequences or shortcomings of the technology?
The negative consequence of the technology is that it may say that there is no delay, but in reality, there is an actual delay. In our case, the models return an accuracy of 77%, which is not high enough to make proper predictions. If we take into account real life, that score may be way lower, because of much more factors that we did not include in our prediction. This may be a problem since the employees of Ordina need to be on time for their clients and they need to trust the technology. However, it is not a huge problem, since what

will happen is that the employee may be a few minutes late.

**Now that you have thought hard about the transparency of this technology, what improvements would you like to make? List them below.**

One of the important improvements is the transparency regarding the models. Even though we put a lot of explanations to the models, we could spend more time on looking for models that fit better with the data or dive more into the algorithm work.

## Sustainability

Is your technology environmentally sustainable?

*This category is not applicable for this technology.*

## Future
Did you consider future impact?

### What could possibly happen with this technology in the future?

People could become over-reliant on our model and trust advice based on it so much that they don't even think of alternatives.
The model cannot predict any kind of accident or sudden fault that results in huge delays. Where any kind of accident happens at some point users might hold us or the implementers of our model responsible for not being able to predict it.

### Sketch a or some future scenario (s) (20-50 years up front) regarding the technology with the help of storytelling. Start with at least one utopian scenario.

In the future, this model could be made so much more accurate and reliable so that its predictions help NS repair or improve certain inefficiencies or causes of delay.

Train companies like NS would be held to a higher standard, ultimately causing a huge improvement in timely train travel as well as preventing financial losses due to delay fines, lawsuits or just saving money from efficient train planning.

Even later the same model could be improved by collecting data measured during the improvement and be implemented in neighbouring countries and improve travel there. Eventually spreading so far that all connected train infrastructure is completely managed by the same system.

It could end up displacing all other kinds of ground-based public transport in favour of the new completely unified and efficient train system.

### Sketch a or some future scenario (s) (20-50 years up front) regarding the technology with the help of storytelling. Start with at least one dystopian scenario.

It would be very unlikely but it could come to pass that after taking ownership of our models, Ordina would use additional data to improve its accuracy and reliability.

We don't know what this data would be but we can imagine it might include things such as:
Weather data
Causes of delay
Region data.

# Technology Impact Cycle Tool

**Ordina train delay prediction**

---

Eventually, Ordina might find some of the causes of delay are very prevalent and cause most delays. Following that they might potentially find these causes might be things such as conductor negligence and suicide attempts.

If Ordina then wishes to further improve the predictions they might find a way to predict these causes and use those predictions to predict the probability of train delays.

In a very far fetched scenario, they might succeed in doing so, with some consequences.

In the best-case scenario: They would breach personal privacy by using these individuals' data.

In the worst-case scenario: They would use more and more data eventually predicting any aspect of a person's actions just to predict delays.

### Would you like to live in one of this scenario's? Why? Why not?

From a group discussion, we roughly agreed the utopian scenario would be a beneficial scenario for everyone assuming the entity that owned the international train system followed regulations from each of the different countries the system encompassed. The regulations would prevent this entity from making abuse of its monopoly on train travel without getting in the way of this very efficient and unified system.

### What happens if the technology (which you have thought of as ethically well-considered) is bought or taken over by another party?

Assuming Ordina or a company that buys our model from them gets ownership of our technology and decides to use for ethically questionable means, there is very little to nothing we have done to stop it.
The delay data could maybe be expanded to include the train conductors id to help with more accurately predicting delays, using the conductors trends in the track record.
This could breach a conductors privacy, and harm their job security, in many situations this might actually be illegal or unions might prevent access to this information in the first place.

### Impact Improvement: Now that you have thought hard about the future impact of the technology, what improvements would you like to make? List them below.

Seeing as these negative scenarios are always the case of an external entity taking ownership of our built model and modifying/improving it to fit their needs, we cannot modify or improve our technology to prevent that.

---

What we could do is license our model under a license that would prevent modification or prevent consequences of modification to be our responsibility, however seeing as this technology will be provided to Ordina no questions asked, this is not doable in this case.