

Data Science Practical

---

## Data quality metrics for text data

---

Elite Master Program Data Science  
Ludwig-Maximilians-Universität München  
in collaboration with  
BMW Group

Vladana Dakovic  
Valari Pai  
Ekaterina Shmaneva

Munich, June 7<sup>th</sup>, 2022



Supervised by Dr Maka Karalashvili and Dr Matthias Schubert

## Abstract

When (or after) the vehicle is produced, different defects occur. These defects are recorded and stored in the database that summarizes similar defects and assigns them to the prebuild defect cluster. All these defects are composed by humans and contain a high amount of text, thus, analysis of such data is complicated and time-consuming. Current project was proposed to ease the investigation of such defects.

The goal of the project is to build a model that will preprocess the human-created text data, summarize it, classify it according to the given labels, based on the "sense" of the generated summary and evaluate its quality of it.

To explain the used methods, applied to solve the main tasks, the authors present the current protocol. It shows how the application of the Google T5 model for summarization and DistilBERT for classification can help achieve the project's aim and presents the results of the implemented models.

# Contents

<b>1</b>	<b>Motivation</b>	<b>1</b>
<b>2</b>	<b>Related work</b>	<b>2</b>
<b>3</b>	<b>Theoretical aspects</b>	<b>3</b>
3.1	The task of summarization . . . . .	3
3.2	The task of classification . . . . .	3
3.3	Existing methods and approaches . . . . .	3
3.4	Model choice . . . . .	6
3.4.1	Summarization . . . . .	6
3.4.2	Classification . . . . .	7
<b>4</b>	<b>Background &amp; prerequisites</b>	<b>8</b>
4.1	Project goals . . . . .	8
4.2	Data overview . . . . .	8
4.3	Data processing and baseline models . . . . .	9
<b>5</b>	<b>Implementation storyline</b>	<b>11</b>
5.1	Project storyline . . . . .	11
5.2	Implementation aspects and set up . . . . .	13
5.3	Performance analysis . . . . .	19
<b>6</b>	<b>Future work</b>	<b>23</b>
<b>7</b>	<b>Conclusion</b>	<b>24</b>
<b>8</b>	<b>References</b>	<b>25</b>
<b>A</b>	<b>Appendix</b>	<b>V</b>
A.1	The list of the Python packages, used in the project . . . . .	V
A.2	Scores outputs in graphical form . . . . .	VII

# 1 Motivation

Before starting a vehicle model series production, the vehicle concept and prototype engineering are tested. Above all, the company collects test result data to track quality defects that require some rework. This data is referred to as the Prototype. Afterwards, a vehicle model goes into a series production in a plant, and, during the production phase, similar data is collected similarly for a similar purpose. The latter defects get reported as a ticket that needs to be resolved by the end of vehicle production. This data is referred to as the Production.

All quality defect recordings that were created during either of the mentioned phases contain a human, free text description. For better maintenance of these defects, a more manageable, superordinate data source is built to summarize similar quality defects in Production and Prototype. Specifically, each recording in these data should describe a prebuilt, known defect cluster. Besides similar defects, this data source should summarize similar steps conducted to fix those defects.

However, the analysis of such an amount of the human-written text data requires a high amount of humane workers and their working hours. Thus, there arose a necessity to create a model which is able to not only preprocess such data but also analyse the text of the mentioned defects, summarize them for better and easier readability and classify it.

## 2 Related work

Natural Language Processing (NLP) is a machine learning field that allows computers to analyze, control, and possibly create human dialect. Over the last years, it is getting more automated, allowing people to perform all kinds of tasks: from text recognition to text generation. Many new libraries and methods were created to help data scientists to proceed in this area.

Thus, Devlin et al. [1] introduced the BERT - Bidirectional Encoder Representations from Transformers - a language representation model that can perform eleven different NLP tasks (e.g. next sentence prediction, question answering, etc.). Different BERT models were trained with various approaches, such as diverting the number of layers, hidden units, and attention heads using the same hyperparameters and training procedure. These tests have shown that BERT is a competitive model that is effective for both fine-tuning and feature-based approaches.

However, BERT is a technology helping to generate “contextualized” word embeddings/vectors which makes it very compute-intensive at inference time. As an answer to such limitations, many new models for optimizing the BERT were generated, for example, RoBERTa (a robustly optimized BERT approach) ([2, 3]). The main difference between them is that the latter was trained on the bigger amount of data and also on the longer sequences. Moreover, tokenization in RoBERTa is performed with a byte-level Byte-Pair Encoding (BPE) encoding scheme, and its library contains 50K subword units, whereas BERT’s character-level BPE only has a 30K vocabulary [4]. Being pretrained on the raw texts only, without human-created labels, RoBERTa is perfectly suited for such NLP tasks as Summarization. Singh [5] suggests using it to create an abstractive summary of the reviews that Amazon users wrote for the products they purchased.

Another example of the models capable of doing a summary are the OpenAI’s GPT engines ([6, 7, 8]). Generative Pre-trained Transformer can perform different NLP tasks as well: question answering, text Summarization etc. without any supervised training. Also, they require exceptionally few to no examples to get the tasks and perform equivalent or better than the state-of-the-art supervised trained models [9].

Lately, transfer learning has proven to be a powerful technique in NLP. The idea behind it is to pre-train the model on a data-rich task first before fine-tuning it on a downstream task. To set a new state of the art in the field, the Google research team offered an approach to transfer learning in NLP: Text-to-Text Transfer Transformer (T5) ([10, 11]). T5 model has 11 billion parameters and showed a great performance on 17 NLP tasks, for example, text classification [12], data to text generation [13] and Summarization [14]. As a matter of fact, the text-to-text architecture of the T5 made it easy to feed structured data into the model. However, in the case of bigger data, it tends to ignore some of the information in the data input. As for text classification tasks, T5 performs even better, reaching over 90% accuracy threshold, despite considering only 512 tokens of the input.

Nevertheless, the main goal of the current project is solving the task of the text Summarization, which can be performed by all of the above-mentioned models and following the classification of the generated summaries. The next sections of this report elaborate in detail on the theoretical aspects of both tasks and state the application of the Google T5 model to it.

## 3 Theoretical aspects

### 3.1 The task of summarization

According to [15], a **summary** can be defined as a text that is produced from one or more texts that contain a significant portion of the information in the original text(s), and that is no longer than half of the original text(s). Hence, training a computer to produce such a summary is called the task of **automatic summarization**.

Summarization aims to condense some text data into a shorter version while preserving most of its meaning. This ultimately saves storage and time resources that processing the long text requires. Summarization also helps to discard irrelevant information and focus on the central ideas of the text.

Generally, machine (automatic) summarization is split into two types [16]:

1. **Extractive:** Here, important text or sentences are extracted as they appear in the original document and grouped to form a concise summary. Most extractive summarization techniques focus on finding and extracting Keywords from the parent text. It can be compared to highlighting the most crucial parts of the text with a marker.
2. **Abstractive:** This approach focuses on generating summaries using the important ideas or facts that the document contains without repeating them verbatim. It is similar to the summary that a human would write after reading the text.

### 3.2 The task of classification

The task of **classification** can be defined as categorizing open-ended text into two or more predefined classes based on some rules or similarities between these texts. It provides valuable insights about unstructured text data as it divides them into classes.

There are three main approaches to machine-based classification tasks: [17]:

1. **Rule-based systems:** In this approach, the text is classified by using a set of linguistic rules that can be defined by the user. Usually, the rule is based on some keywords that indicate the text belonging to a particular group.
2. **Machine learning-based systems:** A machine learning algorithm learns to make classifications based on past observations. Training data with labelled examples is vital for this approach.
3. **Hybrid systems:** These are a combination of both of the above-mentioned approaches. They are useful to build classifiers for a unique task for greater precision.

### 3.3 Existing methods and approaches

The most common approaches are reviewed in terms of their usability for classification and Summarization in this section. All models are separated into three groups, depending on the tasks they can be performed on.

## 1. Models, only used for summarization tasks

- *Sumy* is a library that provides a variety of algorithms for text summarization. Some of these algorithms are LexRank, Luhn, Latent Semantic Analysis (LSA), and KL-Sum. All of them are based on different concepts, which are suitable for different tasks. Sumy is also easy to use, as the algorithm can be imported without much coding or fine-tuning. However, most of the algorithms in Sumy are supposed to be used for extractive summarization [18].

## 2. Models, only used for classification tasks

- *Naive Bayes* algorithm provides a probabilistic classifier that is based on the Bayes' Theorem. The classification is implemented by calculating the probability of each 'tag' or 'class' for the given text and then determining the label with the highest probability.
- *Support Vector Machines (SVM)* calculate a divisionary line between two or more classes. Such a line is known as the decision boundary and determines the best result between vectors that belong to the classes and also the ones that do not. However, the main drawback of SVMs is that they perform well only when there is a limited amount of data [17].

## 3. Models, used for both summarization and classification tasks

- *Gensim* is a python library specifically engineered for Natural Language Processing (further: NLP) tasks.
  - **Summarization:** Gensim performs extractive text summarization using the TextRank algorithm. TextRank algorithm deems the sentences that contain words that occur most frequently as significant and assigns them a 'Rank'. The sentences with the highest rank are extracted to form a summary [18].
  - **Classification:** The Gensim library provides the Doc2vec algorithm that is strong enough to perform Multi-class text classification. Doc2vec is similar to word2vec but uses a Distributed Bag of Words (DBOW) instead of Continuous Bag of Words (CBOW) or Skip-gram [19].
- *Deep learning models (CNN and RNN):* Deep learning is a very important field of machine learning which represents multiple layered Neural networks that are designed to mimic the human brain [20]. For NLP, the most widely used Deep learning algorithms are Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN). CNNs are traditionally used in computer vision tasks, however, recent research has shown that they are as well very effective on NLP tasks. RNNs are specifically designed to process sequential information.
  - **Summarization:** CNNs have mostly been implemented to only perform extractive text summarization. However, these models have a complex architecture, are highly computationally expensive, and are hard to interpret. It is also hard to implement deep bidirectionality using CNNs [21].

Recurrent Neural Networks (RNN) have also been used to perform extractive text summarization with state-of-the-art performance. GRU and LSTM models have an easy-to-interpret approach but can only deliver high performance in specific cases. The main drawback of RNN models is that they cannot handle long-term dependencies [22].

- **Classification:** CNNs can be used for classification by utilizing a feature that is applied to words or n-grams to extract high-level features [20].

RNNs are also effective in performing classification as they have the ability to memorize the previous output and use that information to base the next one [17].

- *BERT and BERT-based models:* BERT is a bidirectional transformer-based model which was implemented to overcome the drawbacks of the RNN models. BERT (Bidirectional Encoder Representations from Transformers) [1] is a pre-trained model which can be easily fine-tuned to perform multiple tasks. BERT was a revolutionary model which provided a strong architectural base for many other models. These models mostly focus on improving BERT's performance or making it more efficient. Some examples of such models are ALBERT (A smaller model with stronger performance) [23], RoBERTa (A larger model with more parameters aimed at making BERT more robust) [3], and DistilBERT (A distilled version of BERT that is faster, smaller and lighter) [24], etc.

Since the fine-tuning of BERT and Co to perform any NLP tasks is unchallenging, they can perform both classification and Summarization. BERT is also very effective in performing abstractive summarization.

- *T5:* Google's text-to-text transfer transformer model is trained end-to-end with a text string as the input and a modified text string as the output. This gives the T5 model an advantage over BERT-based models as the latter only returns a class label.

The T5 model is used to perform multiple NLP tasks with state-of-the-art performance including abstractive summarization. This is a pre-trained model which is trained on the unlabelled large text corpus called C4 (Colossal Clean Crawled Corpus) using deep learning [14].

There are five different versions of the pre-trained T5 model available on [HuggingFace: T5](#) depending on the size of the model. The smallest is the "T5-small" with 60 million parameters, whereas the largest, "T5-11B", has 11 billion parameters.

T5 is implemented using HuggingFace transformers and can be fine-tuned to the required NLP task. So, it can perform both Classification and Summarization tasks.

- *GPT models:* OpenAI's GPT (Generative Pre-trained Transformer) is one of the most well-known NLP models out there. The latest version, GPT-3, has 175 billion parameters that give the model a tremendous amount of power. GPT-3 can be used for all sorts of NLP tasks and outperforms many state-of-the-art models [8]. However, GPT-3 is not open-sourced and, hence, can only be used via an API after registration.



- *XLNet*: The XLNet model can be interpreted as a modification of the BERT model. It is a bidirectional transformer-based model which is pre-trained in a regressive manner, similar to the GPT family of models. It comes in two versions, which differentiate in size: XLNet-base-cased and XLNet-large-cased. Because of its size, XLNet is very expensive to evaluate the SotA (State of the Art) results of the XLNet-large model. However, it generally gives very good results on downstream language tasks like question answering, sentiment analysis, etc [25]. Though, when it comes to summarization, it is outperformed by T5 [6].

### 3.4 Model choice

Over the past few years, transfer learning has led to a new wave of state-of-the-art results in natural language processing. Transfer learning’s effectiveness comes from pre-training a model on abundantly available unlabeled text data with a self-supervised task, such as language modelling or filling in the missing words. After that, the model can be fine-tuned on smaller labelled datasets, often resulting in a better performance than training on the labelled data alone. The recent success of transfer learning was ignited in 2018 by GPT, ULMFiT, ELMo, and BERT, and 2019 saw the development of a huge diversity of new methods like XLNet, RoBERTa, ALBERT, Reformer, and MT-DNN. The rate of progress in the field has made it complicated to evaluate which improvements are most meaningful and how effective they are when combined. [14]

#### 3.4.1 Summarization

First research showed that the best model (from BERT-styled models) for summarization is RoBERTa. RoBERTa is an encoder model similar to BERT, but it uses dynamic MASKing. So, RoBERTa sees the same sequence masked differently, unlike BERT who sees the MASKed sequence only once. It also completely discards the NSP objective and uses a much larger corpus (160GB) during pre-training instead. This provides RoBERTa with much better results than BERT and XLNet model [3].

After implementation of RoBERTa authors discovered that this model was not the best suitable for the task. RoBERTa is just an encoder-based model and, thus, does not perform well on summarization tasks. Research showed that picking either an encoder-decoder based model or only a decoder based model will provide better results for summarization.

Wanting to explore the limits of Transfer Learning, researchers at Google wanted to create a unique model which could be applied to many NLP tasks such as summarization, translation, questions, and answers. The model was named Text-To-Text Transfer Transformer (T5). Unlike BERT, which had only encoder blocks, T5 uses both encoder and decoder blocks. Moreover, T5 does not output a label or a span of the input to the input sentence, and the output is a text string as well. This reason makes the T5 model more suitable for summarization tasks than any BERT-styled model. Due to the lack of computational resources, the authors decided to confine to the "T5-small" version, which was as well pre-trained on a multi-task mixture of unsupervised and supervised tasks, and performs not worse, than its extended variations [26].

### 3.4.2 Classification

Due to the benefits of transfer learning, the authors decided to implement and fine-tune a pre-trained model. Since classification and sentiment analysis is a task much simpler than summarization, BERT-style models are still a good choice. The developers of this project were limited in computational resources, so it was decided to implement the DistilBERT model. DistilBERT is a much smaller, faster and cheaper model compared to BERT and has provided SOTA results for classification tasks. DistilBERT is created by distilling the BERT base model. As a result, it has about 40% fewer parameters than BERT which gives it the ability to run 60% faster. Despite this, DistilBERT is capable of retaining around 95% of BERT's performance. The authors of DistilBERT have tried to minimize inductive biases which large models usually learn to manage during pretraining by using a triple loss approach which focuses on language modelling, distillation and cosine-distance losses. Further, this model can also be found on HuggingFace just like BERT giving it the same level of flexibility. Because of the above-stated reasons, DistilBERT is an obvious choice for the task of classification under this particular scenario. [24].

Implementation aspects and computational results are as well provided in further chapters of this work (c.f. Sec. 5.2 for the details on implementation and Sec. 5.3 for the performance analysis).

## 4 Background & prerequisites

### 4.1 Project goals

The main goal of the project is to derive reasonable quality metrics for text data in the data sources, analyze the free text description data, and generate and analyze a summary of it.

Normally, projects of such scope can have two target audiences: the first one is a *Data Steward*, who takes ownership of the data, works with the business to define the programme's objectives [27], and, thus, not necessarily equipped with a Data Science background. The second one is a *Data Scientist* - analytical data expert, who has technical skills to solve complex problems [28]. This implies that every task within the project should be easily understood by persons without a Data Science background as well.

One of the project tasks is the derivation of the suitable metrics from the data, the choice of which strongly depends on the available data. Hence, metric derivation and output of the expected results should be included in the list of goals as well.

Due to several data security issues, analyzing and handling the initial data was impossible. Consequently, the research for the open-source data, which has the most similar structure, was conducted. This resulted in the use of the Amazon Product Reviews dataset, which has different categories and a great number of reviews that could not be processed with existing resources. Thereby, the authors narrowed the data to utilizing the Quality Food reviews. This data set was suitable to obtain a good performance of the model and had a diversity, that was most similar to the original data, for which this project was designed to handle. Based on the new data, a set of new goals was defined, including data preprocessing, Summarization, and analyzing the goodness of the summary.

### 4.2 Data overview

Amazon Product Review dataset is a publicly available dataset [29], which contains 568.454 reviews. The data is stored within 10 columns: *Id*, *ProductId*, *UserId*, *ProfileName*, *HelpfulnessNumerator*, *HelpfulnessDenominator*, *Score*, *Time*, *Summary*, and *Text* (Fig. 1).

Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all	This is a confection that has been around a fe...
4	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient i...
5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy	Great taffy at a great price. There was a wid...

Figure 1: A preview of data

To proceed with the further analysis, first of all, all reviews, that do not have a value, are dropped. This cleans the dataset up to 568.427 reviews.

Moreover, not all columns are needed for Summarization tasks: only columns *Summary* and *Text* are kept, whereas other columns are not important and are discarded. For more precise summaries, stop words were removed using the additional filter on the data. Another filter helped to exclude all reviews that were too long (reviews longer than 512 tokens).

The distribution analysis of tokens in columns *Text* and *Summary* within the remaining data can help to get a better understanding of the data (Fig. 2).

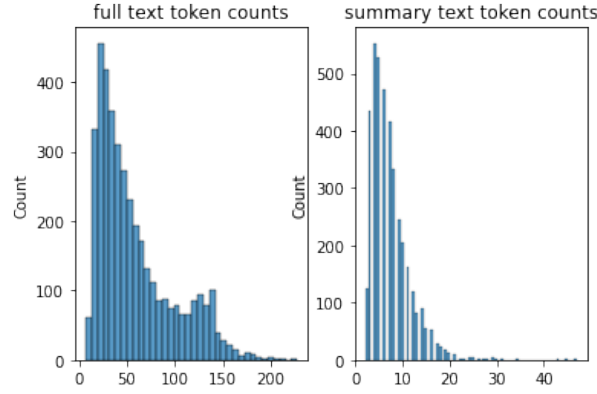


Figure 2: Example of the distribution of the number of tokens in *Text* and *Summaries*

As a baseline model, a function, the Git repository for which was linked in the Sec. 5.2, was used to untokenize the reviews back as the text. Additionally, due to the lack of computational resources which were at their disposal, the authors could not perform the training task on the whole dataset. Hence, a larger subsamples of the reviews of different size were used.

Similarly, for classification task, authors only kept columns *Score*, *Summary*, and *Text*. Column *Score* refers to the rating (on a scale of 1 to 5) provided by the reviewer for the variety of food products, that amazon offers. This *Score* was then used to calculate the *Sentiment*: a Boolean value (positive or negative) to indicate the sentiment of the review based on the rule: if the score is greater than or equal to three, positive, otherwise negative.

### 4.3 Data processing and baseline models

Summarization and classification tasks can be performed with various models, most of which are available in the Huggingface library in Python. [Huggingface library](#) is specially designed for NLP Transformers implementation, and it supports other widely used Python libraries. As a Summarization model, the authors used the Google T5 model, whereas for the classification task, the DistilBERT model was chosen. Both models can be imported from the Transformers package. Additionally, to enable the training of the model, data needs to be encoded in an appropriate way using a predefined Tokenizer.

In Python, the T5 model is implemented in several sizes: *t5-small*, *t5-base*, *t5-large*, *t5-3b* and *t5-11b*. The difference between models is illustrated in Fig. 3. As a consequence of available computation power, we have implemented a *t5-small* model for generating summaries of Reviews.

<i>Model size variants</i>						
Model	Parameters	# layers	$d_{\text{model}}$	$d_{\text{ff}}$	$d_{\text{kv}}$	# heads
Small	60M	6	512	2048	64	8
Base	220M	12	768	3072	64	12
Large	770M	24	1024	4096	64	16
3B	3B	24	1024	16384	128	32
11B	11B	24	1024	65536	128	128

Figure 3: T5 model size variants [11]

For the T5-small model, the authors have split the dataset into train, validation, and test datasets in the ratio of 80-10-10, which is one of the most used techniques [30]. Afterwards, each dataset was encoded using **T5TokenizerFast** from the Transformers package to suit the desired model input. Moreover, the maximum length of tokens was set to 512 tokens for *Text* and 128 for *Summaries*. All larger *Text* and *Summary* were cut after respecting the maximum token length.

Similar processing was utilized while implementing the Classification model. Here, instead of the T5, the DistilBERT model was used, which is as well importable via the Transformers package. The *Text* field is used as input with *Sentiment* as a label to train the model. Different from the split of the dataset for the T5 Summarization model, here authors used a stratified split using the *train\_test\_split* function from [sklearn library](#). To train and fit the model, the text input was first encoded using the pre-trained DistilBert Tokenizer. This is done to remove punctuation splitting and word pieces. The authors then proceed to fit a pre-trained DistilBert Classification model on the prepared dataset.

## 5 Implementation storyline

### 5.1 Project storyline

The project started on the 18<sup>th</sup> of October 2021 with the Kick-off from the BMW side. The developers were made familiar with the task and dataset to use and suggested starting the research on the possible methods. However, as stated in the 4.1, due to various issues with the access to the data and data security, there arose the necessity to conduct additional research on the plausible for the task of open-source data. That affected the timeline and the scope of the project (Fig. 4).

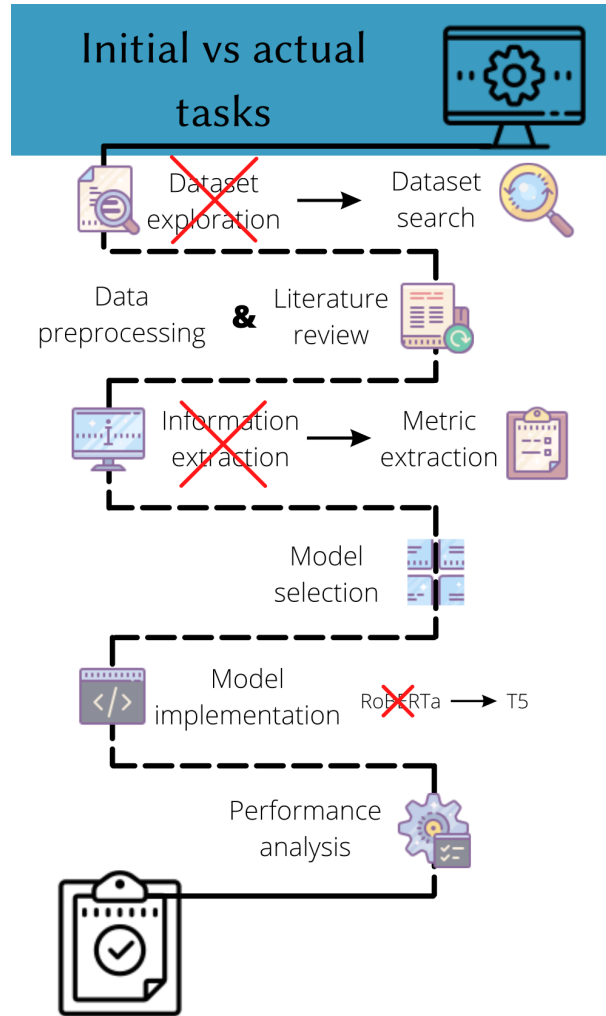


Figure 4: The scope of the project, including tasks, changed due to unexpected issues

All further tasks, displayed within the scope are described in more detail consecutively in this section.

**Metric selection and extraction** As a result of the research, the Amazon Product Reviews datasets [29] were found and examined. The first focus fell on the Musical

Instruments reviews and product metadata. Being compressed in the JSON format, the data contained the following dimensions: *reviewerID* (ID of the reviewer), *asin* (ID of the product), *reviewerName* (name of the reviewer), *helpful* (helpfulness rating of the review), *reviewText* (text of the review), *overall* (rating of the product), *summary* (summary of the review, written by user), *unixReviewTime* (time of the review (unix time)), *reviewTime* (time of the review (raw)).

Developers proposed using three metrics, provided in the dataset:

1. **helpfulness** analysis as a first layer rule-based evaluation. If the helpfulness exceeded some preset threshold, **reviewText** analysis is performed.
2. **reviewText** analysis should be performed according to some characteristics of the text, such as length, most occurring words, repetitions, the mood of the review, and product names within it.
3. **overall** score, given by the user, should be consistent with the mood of the review itself.

However, during the approach discussion, it was decided to reduce the analysis to the **reviewText** dimension, as it is the most important metric for the initial task. In parallel, keyword search and Summarization techniques were researched and evaluated.

**Model selection** The model choice is described more precisely in the Sec. 3.4. Moreover, together with the model selection, a bigger and more applicable dataset was found for a better analysis. The new dataset was as well containing the Amazon Product reviews, yet still, a different preprocessing was needed (c.f. Sec. 4.2).

At the same time, literature research and review were conducted (c.f. Sec. 2). After the above-mentioned steps project entered its practical phase, including the tasks of metrics analysis for the new dataset, evaluating existing Summarization and classification techniques and corresponding Python packages.

**Model implementation** Implementation aspects are discussed in Sec. 5.2 of the current report. The first model that was selected for implementation was RoBERTa [4]. As mentioned in Sec. 3.4.1, it is an encoder model that suffers from some issues with weight assignment to the decoder. Thus, although the model was correctly implemented and adapted to chosen data, its training was performed incorrectly. Hence, the T5 model was chosen as it is more robust and easier in issue handling.

**Performance analysis** All outputs and their interpretation are thoroughly discussed in the Sec. 5.3.

Furthermore, the project is tied to time, which made time planning crucial. Though, as discussed in the paragraphs above and the Sec. 3.4, due to unexpected issues with the RoBERTa model, the authors were forced to change the model used for Summarization. Consequently, the deadlines of the project had to be moved to a later period, resulting in the new timeline (Fig. 5).

Needless to say, every complicated model requires a lot of work and bug fixing. Hence, the authors of this report came across several issues during the project that as well

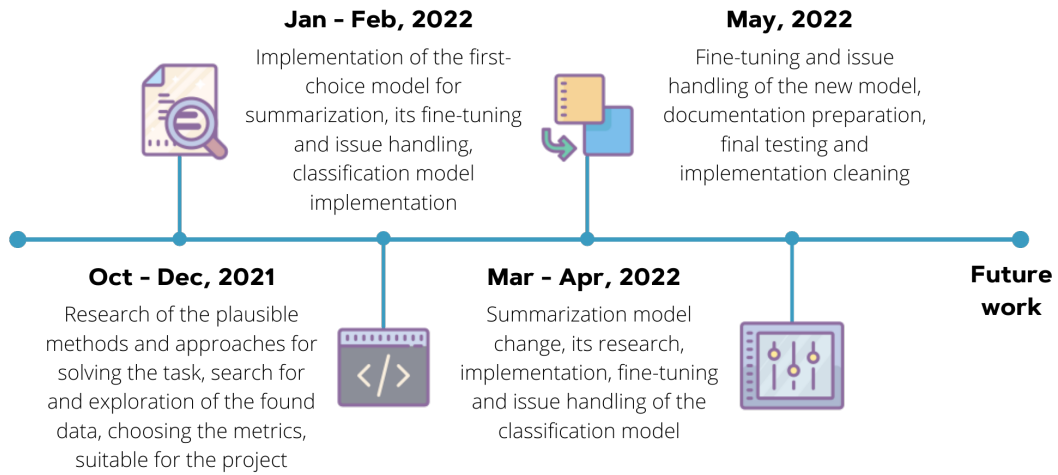


Figure 5: The updated timeline of the finalized project

affected the deadlines and resulted in the prolongation of the project. The most crucial and relevant of them are shown in Fig. 6.

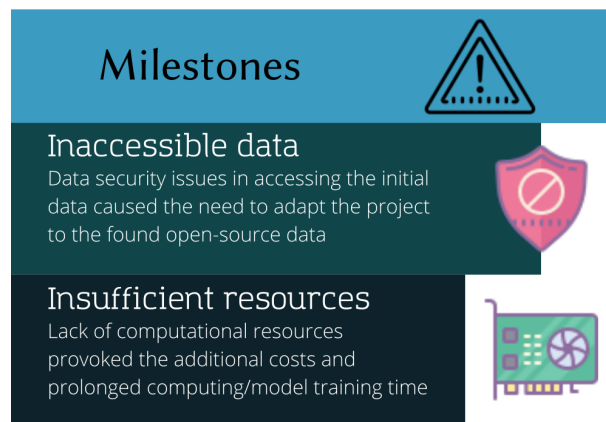


Figure 6: The most troublesome aspects of the project, that resulted in the necessity to expand the deadlines

## 5.2 Implementation aspects and set up

The implementation and documentation regarding the project can be found in the corresponding [GIT Repository](#)

**Project set up** To compute the models, the following set up was used to proceed with the described models:

1. **Environment:** Python ver. 3.8.5 and above



2. **Computational system:** Google Colab (Pro)
3. **Processing machine:** Google GPU Accelerator

**Summarization implementation** Section 3.1 of the current document refers to the different types of Summarizations. Initially, the standard BERT Library [1] was chosen. Although, for easier and faster computation, its distilled modification (DistilBERT [24]) was used. Additionally, the robustly optimized BERT approach, RoBERTa [3], was applied for creating Summarizations.

The vast majority of the NLP models process the data on the token level. Thus, tokenization is a crucial step while conducting such tasks as Summarization. However, during the training, RoBERTa showed some shortcomings, namely, the weakness of its tokenizer: some weights of the model checkpoint were not used when initializing the model. Hence, it could not learn anything during training. The authors performed further research and after the trial-test phase, shifted their attention to the application of the Text-to-Text Transfer Transformer (T5). Despite being bulkier than its predecessor, the T5 model has proven to be more accurate in performing Summarization tasks. Additionally, its tokenizer is more robust. Thus, the in-built *T5TokenizerFast* module was used in the model.

To encode the initial dataset to the correct input of the model, the new class (*SummaryDataModule*) was defined. Within this class, authors defined different methods, which will adapt the train, validation and test data.

---

```

1  class SummaryDataModule(pl.LightningDataModule):
2      def __init__(
3          self,
4          train_df: pd.DataFrame,
5          test_df: pd.DataFrame,
6          val_df: pd.DataFrame,
7          tokenizer: T5TokenizerFast,
8          batch_size: int = 8,
9          text_max_token_len: int = 512,
10         summary_max_token_len: int = 128
11     ):
12
13         super().__init__()
14         self.train_df = train_df
15         self.test_df = test_df
16         self.val_df = val_df
17         self.batch_size = batch_size
18         self.tokenizer = tokenizer
19         self.text_max_token_len = text_max_token_len
20         self.summary_max_token_len = summary_max_token_len
21
22

```

```
23 def setup(self, stage=None) :
24     self.train_dataset = SummaryDataset(
25         self.train_df,
26         self.tokenizer,
27         self.text_max_token_len,
28         self.summary_max_token_len
29     )
30
31     self.test_dataset = SummaryDataset(
32         self.test_df,
33         self.tokenizer,
34         self.text_max_token_len,
35         self.summary_max_token_len
36     )
37
38     self.val_dataset = SummaryDataset(
39         self.val_df,
40         self.tokenizer,
41         self.text_max_token_len,
42         self.summary_max_token_len
43     )
44
45 def train_dataloader(self):
46     return DataLoader(
47         self.train_dataset,
48         batch_size=self.batch_size,
49         shuffle=True,
50         num_workers=2
51     )
52
53 def val_dataloader(self):
54     return DataLoader(
55         self.val_dataset,
56         batch_size=self.batch_size,
57         shuffle=False,
58         num_workers=2
59     )
60
61 def test_dataloader(self):
62     return DataLoader(
63         self.test_dataset,
64         batch_size=self.batch_size,
65         shuffle=False,
66         num_workers=2
67     )
```

---

After having data tokenized as input of the T5 model, the authors have created a new class, describing the model, where the training and validation steps were specified. Furthermore, *AdamW* was used as an optimizer, which is normally applied as a default optimizer for *T5TokenizerFast* (e.g. [31]).

---

```

1  class SummaryModel(pl.LightningModule):
2
3      def __init__(self):
4          super().__init__()
5          #initializing model
6          self.model = T5ForConditionalGeneration.from_pretrained(modelName, return_dict=True)
7
8
9      # Defining forward function and it output
10     def forward(self, input_ids, attention_mask, decoder_attention_mask, labels=None):
11         output = self.model(
12             input_ids,
13             attention_mask=attention_mask,
14             labels=labels,
15             decoder_attention_mask=decoder_attention_mask
16         )
17         return output.loss, output.logits
18
19
20     def training_step(self, batch, batch_idx):
21         input_ids = batch[ "text_input_ids"]
22         attention_mask = batch["text_attention_mask"]
23         labels = batch["labels"]
24         x = batch[ "text_input_ids"]
25         labels_attention_mask = batch["labels_attention_mask"]
26
27         loss, outputs = self(
28             input_ids=input_ids,
29             attention_mask=attention_mask,
30             decoder_attention_mask=labels_attention_mask,
31             labels=labels
32         )
33
34         batch_dictionary={ "loss": loss, "labels": labels}
35
36         self.log("Loss/Train (Batch)", loss, prog_bar=True, logger=True)
37         self.logger.experiment.add_scalar("Loss/Train (Epoch)", loss, self.current_epoch)
38         #return loss
39         return batch_dictionary
40

```

---

```

41 def validation_step(self, batch, batch_idx):
42     input_ids = batch[ "text_input_ids"]
43     attention_mask = batch["text_attention_mask"]
44     labels = batch["labels"]
45     labels_attention_mask = batch["labels_attention_mask"]
46
47     loss, outputs = self(
48         input_ids=input_ids,
49         attention_mask=attention_mask,
50         decoder_attention_mask=labels_attention_mask,
51         labels=labels
52     )
53
54     self.logger.experiment.add_scalar("Loss/Val (epoch)",loss,self.current_epoch)
55     self.log("Loss/Val (Batch)", loss, prog_bar=True,logger=True)
56     epoch_dictionary={'loss': loss}
57     return epoch_dictionary
58
59
60 def test_step(self, batch, batch_idx):
61     input_ids = batch[ "text_input_ids"]
62     attention_mask = batch["text_attention_mask"]
63     labels = batch["labels"]
64     labels_attention_mask = batch["labels_attention_mask"]
65
66     loss, outputs = self(
67         attention_mask=attention_mask,
68         decoder_attention_mask=labels_attention_mask,
69         labels=labels
70     )
71     self.logger.experiment.add_scalar("Loss/Test",loss,self.current_epoch)
72     self.log("test_loss", loss, prog_bar=True,logger=True)
73     return {'loss': loss}
74
75 # Configuring optimizer as most used one AdamW
76 def configure_optimizers(self):
77     return AdamW(self.parameters(), lr=0.0001)

```

---

**Classification implementation** Another goal of the project was to classify generated summaries after they were created. Another model was created for that part, however, in this case, remaining with the use of BERT (precisely, DistilBERT) and applying NLTK library additionally proved itself suitable.

In this case, DistilBERT Tokenizer was used to encode the data. Its implementation within the model is shown in the code below:

---

```

1 tokenizer([training_sentences[0]], truncation=True, padding=True, max_length=128)
2     # Tokenizing the data
3 train_encodings = tokenizer(training_sentences, truncation=True, padding=True)
4 val_encodings = tokenizer(validation_sentences, truncation=True, padding=True)
5 # Slicing the dataset
6 train_dataset = tf.data.Dataset.from_tensor_slices((dict(train_encodings), training_labels))
7 val_dataset = tf.data.Dataset.from_tensor_slices((dict(val_encodings), validation_labels))
8
9 # Loading the DistilBert model from transformers
10 model = TFDistilBertForSequenceClassification.from_pretrained
11     ('distilbert-base-uncased', num_labels=2)
12
13 # Defining and fitting the model on the training data
14 optimizer = tf.keras.optimizers.Adam(learning_rate=5e-5, epsilon=1e-08)
15 callbacks=tf.keras.callbacks.EarlyStopping(monitor='accuracy',
16                                           min_delta=0.0001,
17                                           patience=3,
18                                           mode='auto',
19                                           verbose=2,
20                                           baseline=None)
21
22 model.compile(optimizer=optimizer, loss=model.compute_loss, metrics=['accuracy'])
23 model.fit(train_dataset.shuffle(100).batch(16),
24         epochs=5,
25         batch_size=16,
26         validation_data=val_dataset.shuffle(100).batch(16), callbacks=callbacks)

```

---

For newly generated summaries, two classes were defined: *Positive* and *Negative*, meaning the level of customers' satisfaction with the reviewed product. The analysis was done for both

1. initial text of the review:

---

```

1 # Classifying the text into sentiment classes
2 for i in range(0, len(df)):
3     predict_input_text = tokenizer.encode(df['Text'][i],
4                                           truncation=True,
5                                           padding=True,
6                                           return_tensors="tf")
7     tf_output_text = loaded_model.predict(predict_input_text)[0]
8     tf_prediction_text = tf.nn.softmax(tf_output_text, axis=1)
9     labels = ['Negative', 'Positive']
10    label_text = tf.argmax(tf_prediction_text, axis=1)
11    label_text = label_text.numpy()
12    df["Sentiment_text"][i] = (labels[label_text[0]])

```

- 
2. and a newly generated summary:

---

```

1  # Classifying the generated summaries into sentiment classes
2  for i in range(0, len(df)):
3      predict_input_sum = tokenizer.encode(df['Generated_summary'][i],
4                                          truncation=True,
5                                          padding=True,
6                                          return_tensors="tf")
7      tf_output_sum = loaded_model.predict(predict_input_sum)[0]
8      tf_prediction_sum = tf.nn.softmax(tf_output_sum, axis=1)
9      labels = ['Negative', 'Positive']
10     label_sum = tf.argmax(tf_prediction_sum, axis=1)
11     label_sum = label_sum.numpy()
12     df["Sentiment_summary"][i] = (labels[label_sum[0]])

```

---

During the training, the performance of the model was analysed and resulted in over 90% accuracy. Moreover, the results between classified texts and summaries were compared and the error in classes for generated summary and text input did not exceed 10%. Next Section presents the detailed report on the achieved results and their interpretation.

### 5.3 Performance analysis

To analyze the overall performance of the models and to understand how far the authors managed to achieve the goal of the project, several metrics were chosen and calculated during the project for each task:

1. **Running walltime** shows the amount of time that was required to perform the training (only). It depends on the amount of the data sample and number of epochs that were set up for the training. It is calculated as the sum of the amount of time that was needed for every epoch.
2. **Accuracy** is the number of correctly predicted data points out of all the data points.
3. **Validation** and **Training loss** that describe the performance of the model, indicating how well it is fitting the training and the new data correspondingly.
4. **Rouge score** stands for Recall-Oriented Understudy for Gisting Evaluation. It works by comparing an automatically produced summary or translation against a set of reference summaries (typically human-produced). There are several calculation techniques, but authors use the *ROUGE-1* (overlap of unigrams between the system summary and reference summary), *ROUGE-2* (overlap of bigrams between the system and reference summaries) and *ROUGE-L* (longest matching sequence of words using LCS) [32].

*Recall* refers to how much of the reference summary the system summary is recovering or capturing.

*Precision* measures, how much of the system summary was in fact relevant or needed.

*F-measure* provides a single score that balances both the concerns of precision and recall in one number.

5. **Classification error** is calculated as a percentage of the generated summaries, that were classified differently, than corresponding initial review text (i.e. if the summary is classified as negative whereas the review text was marked as positive, the counter increases).

**Summarization** To analyse the performance of the Summarization model, the authors have chosen the Training and Validation loss and Running time as the standard performance metrics. As for the analysis of the generated summaries, the Rouge score which, as stated above, is the metric that would best measure its goodness. Only the Rouge-1 and Rouge-L scores were used.

With available resources, the authors have performed T5 Summarization for two data samples: 10 000 reviews and 100.000 reviews and aimed to compare the model performance in both cases. Due to the unavailability of the BMW Cluster and time limitations, the full training was not performed. The results that were achieved by the end of the training are presented in the table 1:

Metric name	Result when run on the smaller sample (10.000 data sentences)	Result when run on the bigger sample (100.000 data sentences)
<b>Running time</b>	approx. 30 minutes	approx. 3 hours
<b>Train loss</b>	2.481	2.849
<b>Validation loss</b>	3.532	3.859
<b>Rouge scores</b>		
<i>Rouge-1</i>		
recall	0.266	0.143
precision	1.0	1.0
f-measure	0.421	0.250
<i>Rouge-L</i>		
recall	0.266	0.143
precision	1.0	1.0
f-measure	0.421	0.250

Table 1: The results achieved after training the proposed summarization model for two data samples

**Classification** Similarly, to analyse the performance of the classification model, authors again applied commonly used metrics: Running time, Train and Validation loss. However, differently from the summarization model, Accuracy and Classification error were analyzed. These two metrics can describe the performance of the model, and give

valid information about it. Though, using them for the summarization model would not provide a good understanding of the attained results about summaries and goodness of them as the Rouge score would. The results that were achieved for the classification model by the end of the training are shown in the table 2:

Metric name	Result when run on the smaller sample (10.000 data sentences)	Result when run on the bigger sample (100.000 data sentences)
Running time	approx. 36 minutes	approx. 3.5 hours
Train loss	0.137	0.094
Accuracy	0.953	0.963
Validation loss	0.274	0.194
Validation accuracy	0.901	0.935
Classification error	0.1	0.1

Table 2: The results achieved after training the proposed classification model for two data samples

**Interpretation of the results** The results of training of both models on the smaller data sample of 10.000 entries showed that the running time of the T5 Summarization was around 30 minutes and around 36 minutes for the DistilBERT classification model. In the case of the bigger sample (100.000 entries), it was 3 and 3.5 hours correspondingly.

As for the training loss, it can be seen that it is rather high for training on both dataset samples (being approximated by 2.5 for the smaller sample and 2.8 for the bigger one). That means that the T5 summarization model provided by the authors still requires further parameter fine-tuning for a better fit to the used data. However, as the project was tied to time, this improvement was left for future work. Needless to say, this also affects the validation loss that shows the overall performance of the T5 model on the validation dataset. Though, the graphs representing the validation losses during the training (Fig. 7b and 9b in Sec. A.2) display that it improves with the numbers of epoch.

The loss graphs for the bigger sample training also displayed a similar tendency of the curves, which can be interpreted as a sign of a good model fitting. Hence, for future work, the authors suggest not only performing the training on the full dataset but also conducting it for the bigger number of epochs.

To measure the goodness of the created summaries, authors used Rouge-1 and Rouge-L scores and calculated recall, precision and f-measure for both of them. For the most examples in both samples, the Rouge-1 score metrics match the Rouge-L metrics. That means that the matching unigrams between the model-generated and user-created summaries composite the longest matching sequences.

It is important to elaborate that the authors suggest evaluating the rouge scores of the model against the main text of the review and not the review summary of the user. The reason for it is that most the users tend to express their attitudes via emojis, jokes or sarcasm which can be falsely interpreted by both the summarizer and classifier in the later stages of the analysis. This fact affects the scope of the Rouge scores: plenty of the recall metrics are pretty low (i.e. not many of the n-grams in the reference (Text of the



review) are also present in the generated summary), whereas the precision for the most of the generated summaries lies within 1.0 (i.e. all of the n-grams in the generated summary are also present in the reference (Text of the review)).

The overall performance of the classification model can be assessed as very good. It can be seen that the train and validation accuracy are above 90% for both samples which means that the classifier works well on both training and test data. Furthermore, a rather low (0.137 for the smaller and 0.094 for a bigger sample) train loss proves the close to the perfect fit of the model to the used data. The validation loss shows the goodness of the model as well.

Indeed, the 2-label classification is conducted well and shows a classification error rate of 10%. That means that only 10% of generated summaries were classified differently from the reviewer-written summary (a.k.a. the name of the review, given by a user). This result is true for both data samples which on top of that demonstrates that the generated summaries convey the mood of the review in a very precise manner, despite sometimes differentiating from the user-written one.

## 6 Future work

Due to the limited time, many different adaptations, tests and experiments are left for the future concerning the more profound analysis of specific mechanisms and methods. Moreover, there is a number of additional areas for future research.

First of all, as the project was meant for a specific task, hence, the models should be adapted to be able to utilize the BMW data. That will include some rework on the data loader and data processing.

The way summarization model is constructed can be changed as well. For instance, its optimization: although AdamW is the mostly used optimizer when adapting the T5 model to the Summarization task, different approaches can be tested as well. Namely, gradient descent and its variations are also commonly used when training the NLP models. Further fine-tuning possibilities could also help to make Summarizations more precise and suitable for defect analysis of the produced vehicles.

As for the classification model, based on the utilized data, the authors only used 2-label classification. However, depending on the desired outcome, the number of classes can be expanded according to the BMW production team's expectations. Moreover, in contrast to summarization model, here, the graphical output of the scores was not implemented. Hence, another output is possible as well.

Additionally, in accordance with the data and requirements from the BMW side, the classification task can be updated to the clusterisation one. As the defects, during the prototype and production phase will have different meanings and should be separated in a different way clusterisation seems to be more suitable for solving this issue.

## 7 Conclusion

During the project, the authors aimed to build a model that will solve two tasks:

1. takes the text input and generates the summary of it,
2. takes the generated summaries as the input and classifies it according to the mood of the review.

To sum up, the stated goals were accomplished. First of all, the authors conducted thorough research on automated summarization and classification and presented the definitions and theoretical base of it. Additionally, two models were developed to solve the above-mentioned tasks.

For the summarization generation, Google's Text-to-Text Transfer Transformer model was successfully applied to the Amazon Products Reviews (category: quality food). For ease of computation, the authors used the "t5-small" variation of the model. Moreover, the samples of 10.000 and 100.000 were taken to train and analyse its performance. Results showed that further parameter fine-tuning is needed, although the overall performance of the model tends to improve with a bigger sample of the data.

As for the second task, the authors decided on the classification based on the Distil-BERT model with 2 labels, corresponding to the level of satisfaction of the customer with the product ('Negative' label was used in case of dissatisfaction and 'Positive' otherwise). The same setup and sampling were used here as well. In the end, the suggested model performed with an accuracy of over 90% which shows that the authors managed to solve the task successfully.

What is more, despite being bound to use the external data, the authors developed a model that can be easily adapted to the data of the beneficiary of the project (BMW Group) or even further expanded according to the further needs of the BMW.

## 8 References

- [1] Jacob Devlin **and others**. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. arXiv:1810.04805. type: article. arXiv, 24 **may** 2019. DOI: [10.48550/arXiv.1810.04805](https://doi.org/10.48550/arXiv.1810.04805). arXiv: [1810.04805\[cs\]](https://arxiv.org/abs/1810.04805). URL: <http://arxiv.org/abs/1810.04805> (**urlseen** 22/05/2022).
- [2] Aastha Singh. *Evolving with BERT: Introduction to RoBERTa*. Analytics Vidhya. 9 **july** 2021. URL: <https://medium.com/analytics-vidhya/evolving-with-bert-introduction-to-roberta-5174ec0e7c82> (**urlseen** 22/05/2022).
- [3] Yinhan Liu **and others**. *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. arXiv:1907.11692. type: article. arXiv, 26 **july** 2019. arXiv: [1907.11692\[cs\]](https://arxiv.org/abs/1907.11692). URL: <http://arxiv.org/abs/1907.11692> (**urlseen** 22/05/2022).
- [4] Rohan Jagtap. *RoBERTa: Robustly Optimized BERT-Pretraining Approach*. DataSeries. 19 **august** 2020. URL: <https://medium.com/dataseries/roberta-robustly-optimized-bert-pretraining-approach-d033464bd946> (**urlseen** 22/05/2022).
- [5] Anubhav. *Step by Step Guide: Abstractive Text Summarization Using RoBERTa*. Medium. 20 **december** 2020. URL: <https://anubhav20057.medium.com/step-by-step-guide-abstractive-text-summarization-using-roberta-e93978234a90> (**urlseen** 22/05/2022).
- [6] Manmohan Singh. *Summarize Reddit Comments using T5, BART, GPT-2, XLNet Models*. Medium. 4 **january** 2021. URL: <https://towardsdatascience.com/summarize-reddit-comments-using-t5-bart-gpt-2-xlnet-models-a3e78a5ab944> (**urlseen** 22/05/2022).
- [7] Sukanya Bag. *Text Summarization using BERT, GPT2, XLNet*. Analytics Vidhya. 26 **april** 2021. URL: <https://medium.com/analytics-vidhya/text-summarization-using-bert-gpt2-xlnet-5ee80608e961> (**urlseen** 22/05/2022).
- [8] kdnuggets. *Summarization with GPT-3*. KDnuggets. Section: Products and Services. URL: <https://www.kdnuggets.com/summarization-with-gpt-3.html/> (**urlseen** 22/05/2022).
- [9] Priya Shree. *The Journey of Open AI GPT models*. Walmart Global Tech Blog. 10 **november** 2020. URL: <https://medium.com/walmartglobaltech/the-journey-of-open-ai-gpt-models-32d95b7b7fb2> (**urlseen** 25/05/2022).
- [10] Maria Yao. *10 Leading Language Models For NLP In 2021*. TOPBOTS. 11 **may** 2021. URL: <https://www.topbots.com/leading-nlp-language-models-2020/> (**urlseen** 22/05/2022).
- [11] Qiurui Chen. *T5: a detailed explanation*. Analytics Vidhya. 11 **june** 2020. URL: <https://medium.com/analytics-vidhya/t5-a-detailed-explanation-a0ac9bc53e51> (**urlseen** 22/05/2022).
- [12] pedrormarques. *Fine tuning a T5 text-classification model on colab*. pedrormarques. 21 **october** 2021. URL: <https://pedrormarques.wordpress.com/2021/10/21/fine-tuning-a-t5-text-classification-model-on-colab/> (**urlseen** 24/05/2022).

- [13] Mathew Alexander. *Data to Text generation with T5; Building a simple yet advanced NLG model*. Medium. 10 **april** 2021. URL: <https://towardsdatascience.com/data-to-text-generation-with-t5-building-a-simple-yet-advanced-nlg-model-b5cce5a6df45> (urlseen 24/05/2022).
- [14] turbolab. *Abstractive Summarization Using Google's T5*. Turbolab Technologies. Section: Technology. 4 **october** 2021. URL: <https://turbolab.in/abstractive-summarization-using-googles-t5/> (urlseen 22/05/2022).
- [15] Eduard Hovy. *Text Summarization*. The Oxford Handbook of Computational Linguistics. ISBN: 9780199276349. 13 **january** 2005. DOI: [10.1093/oxfordhb/9780199276349.013.0032](https://doi.org/10.1093/oxfordhb/9780199276349.013.0032). URL: <https://www.oxfordhandbooks.com/view/10.1093/oxfordhb/9780199276349.001.0001/oxfordhb-9780199276349-e-32> (urlseen 22/05/2022).
- [16] Praveen Dubey. *Understand Text Summarization and create your own summarizer in python*. Medium. 23 **december** 2018. URL: <https://towardsdatascience.com/understand-text-summarization-and-create-your-own-summarizer-in-python-b26a9f09fc70> (urlseen 29/05/2022).
- [17] monkeylearn. *Text Classification: What it is And Why it Matters*. MonkeyLearn. URL: <https://monkeylearn.com/text-classification/> (urlseen 29/05/2022).
- [18] Shrivar Sheni. *Text Summarization Approaches for NLP - Practical Guide with Generative Examples*. Machine Learning Plus. 24 **october** 2020. URL: <https://www.machinelearningplus.com/nlp/text-summarization-approaches-nlp-example/> (urlseen 29/05/2022).
- [19] Susan Li. *Multi-Class Text Classification with Doc2Vec & Logistic Regression*. Medium. 4 **december** 2018. URL: <https://towardsdatascience.com/multi-class-text-classification-with-doc2vec-logistic-regression-9da9947b43f4> (urlseen 22/05/2022).
- [20] Inés Roldós. *Go-to Guide for Text Classification with Machine Learning*. MonkeyLearn Blog. Section: Text Classification. 2 **march** 2020. URL: <https://monkeylearn.com/blog/text-classification-machine-learning/> (urlseen 22/05/2022).
- [21] Leo Laugier. *Extractive Document Summarization Using Convolutional Neural Networks-Reimplementation*. 2018. URL: <https://www.semanticscholar.org/paper/Extractive-Document-Summarization-Using-Neural-Laugier/ed0f189bbbbbccceefb41ccb36> (urlseen 29/05/2022).
- [22] Christian Heumann **and others**. *LMU Course: Deep Learning for Natural Language Processing*. URL: <https://moodle.lmu.de/course/view.php?id=17645> (urlseen 15/06/2022).
- [23] Zhenzhong Lan **and others**. *ALBERT: A Lite BERT for Self-supervised Learning of Language Representations*. arXiv:1909.11942. type: article. arXiv, 8 **february** 2020. DOI: [10.48550/arXiv.1909.11942](https://doi.org/10.48550/arXiv.1909.11942). arXiv: [1909.11942\[cs\]](https://arxiv.org/abs/1909.11942). URL: <http://arxiv.org/abs/1909.11942> (urlseen 22/05/2022).

- [24] Victor Sanh **and others**. *DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter*. arXiv:1910.01108. type: article. arXiv, 29 **february** 2020. DOI: [10.48550/arXiv.1910.01108](https://doi.org/10.48550/arXiv.1910.01108). arXiv: [1910.01108\[cs\]](https://arxiv.org/abs/1910.01108). URL: <http://arxiv.org/abs/1910.01108> (**urlseen** 22/05/2022).
- [25] Zihang Dai. *XLNet: Generalized Autoregressive Pretraining for Language Understanding*. original-date: 2019-06-19T08:16:46Z. 22 **may** 2022. URL: <https://github.com/zihangdai/xlnet> (**urlseen** 22/05/2022).
- [26] Zhanpeng Wang. “Smart Auto-completion in Live Chat Utilizing the Power of T5”. phdthesis.
- [27] experian. *What is a Data Steward? — Experian Business*. URL: <https://www.experian.co.uk/business/glossary/data-steward/> (**urlseen** 29/05/2022).
- [28] SAS-institute. *What is a data scientist?* URL: [https://www.sas.com/en\\_us/insights/analytics/what-is-a-data-scientist.html](https://www.sas.com/en_us/insights/analytics/what-is-a-data-scientist.html) (**urlseen** 29/05/2022).
- [29] Amazon. *Amazon review data*. URL: <https://jmcauley.ucsd.edu/data/amazon/> (**urlseen** 22/05/2022).
- [30] Rachel Draelos. *Best Use of Train/Val/Test Splits, with Tips for Medical Data*. Glass Box. 15 **september** 2019. URL: <https://glassboxmedicine.com/2019/09/15/best-use-of-train-val-test-splits-with-tips-for-medical-data/> (**urlseen** 29/05/2022).
- [31] Andrew Fogarty. *Summarization: T5*. URL: [http://seekinginference.com/applied\\_nlp/T5.html#rouge-metrics](http://seekinginference.com/applied_nlp/T5.html#rouge-metrics) (**urlseen** 03/06/2022).
- [32] Kavita Ganesan. *An intro to ROUGE, and how to use it to evaluate summaries*. freeCodeCamp.org. 26 **january** 2017. URL: <https://www.freecodecamp.org/news/what-is-rouge-and-how-it-works-for-evaluation-of-summaries-e059fb8ac840/> (**urlseen** 22/05/2022).

## A Appendix

### A.1 The list of the Python packages, used in the project

This section contains information on all the libraries, modules and other packages, that were used for implementation of the tasks of the project. Corresponding documentation can be found in the linked references.

Name of the library	Required for the task of	Reference	Commentar (version)
AdamW	summarization	<a href="#">Pytorch: AdamW</a>	
dataclass	summarization	<a href="#">Dataclasses: dataclass</a>	
DataLoader	summarization	<a href="#">Pytorch: DataLoader</a>	
Dataset	summarization	<a href="#">Pytorch: Dataset</a>	
DistilBertTokenizerFast	classification	<a href="#">HuggingFace: DistilBert</a>	
field	summarization	<a href="#">Dataclasses.field</a>	
logging	summarization	<a href="#">logging</a>	
matplotlib.pyplot	summarization	<a href="#">Pyplot</a>	
ModelCheckpoint	summarization	<a href="#">tf: ModelCheckpoint</a>	
nltk	classification, summarization	<a href="#">NLTK</a>	
numpy	classification, summarization	<a href="#">NumPy</a>	1.19.5 (classification)
Optional	summarization	<a href="#">Typing.optional</a>	
pandas	classification, summarization	<a href="#">Pandas</a>	
pytorch_lightning	summarization	<a href="#">PyTorch: lightning</a>	
re	classification, summarization	<a href="#">Regular Expressions</a>	
rouge_score	summarization	<a href="#">Rouge score</a>	
sacremoses	classification	<a href="#">PyPi: sacremoses</a>	0.0.45
stopwords	summarization	<a href="#">nltk: stopwords</a>	
StratifiedShuffleSplit	classification	<a href="#">sklearn: Str. Shuffle Split</a>	
T5ForConditionalGeneration	summarization	<a href="#">HuggingFace: T5</a>	
T5TokenizerFast	summarization	<a href="#">HuggingFace: T5 Tokenizer</a>	
tensorboard	summarization	<a href="#">tf: TensorBoard</a>	
TensorBoardLogger	summarization	<a href="#">Pytorch: tf Board logger</a>	
tensorflow	classification, summarization	<a href="#">Tensorflow</a>	2.7.0 (classification)
tensorflow_datasets	classification	<a href="#">tf: datasets</a>	
TFDistilBertForSequenceClassification	classification	<a href="#">HuggingFace: DistilBert</a>	
torch	summarization	<a href="#">Pytorch</a>	

train_test_split	classification	<a href="#">sklearn: train<sub>test</sub>split</a>	
Trainer	summarization	<a href="#">Transformers: Trainer</a>	
TrainingArguments	summarization	<a href="#">Transformers: TrainingArgs</a>	
transformers	classification, summarization	<a href="#">HuggingFace: Transformers</a>	4.5 (summarization), 4.7.0 (classification)
viewitems	classification	<a href="#">six: viewitems dictionary</a>	
word_tokenize	summarization	<a href="#">nltk: word.tokenize</a>	



## A.2 Scores outputs in graphical form

Sec. 5.3 shows the performance analysis metrics and their interpretation. However, for the T5 summarization model additional graphical output of the scores was implemented. Current section shows the stated outputs for the summarization model performed on two samples: 10.000 entries (Fig. 7 and 8) and 100.000 entries (Fig. 9 and 10).

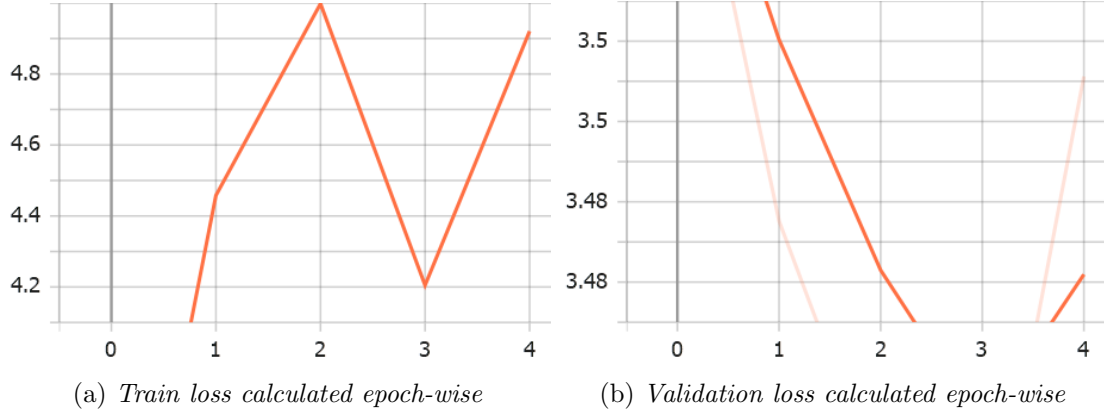


Figure 7: Train and validation losses calculated every epoch for the training on the smaller sample of 10.000 entries

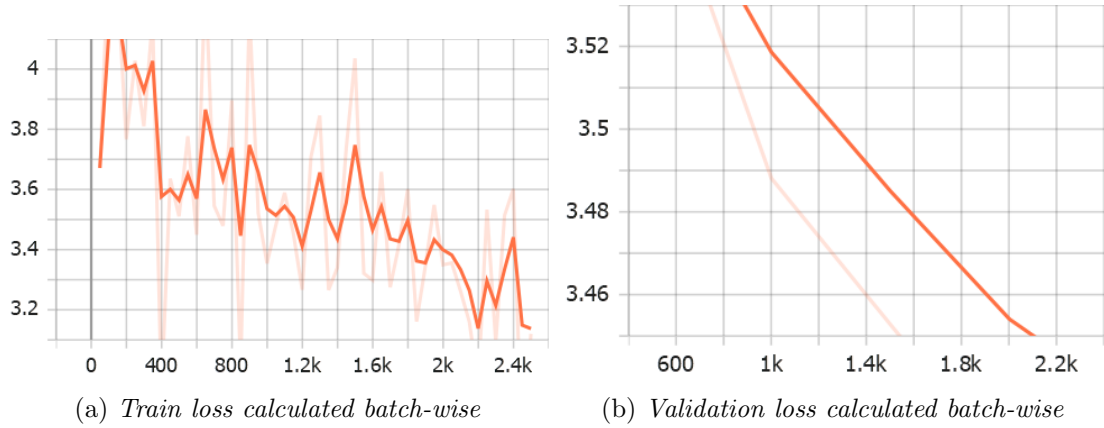


Figure 8: Train and validation losses calculated per batch for the training on the smaller sample of 10.000 entries

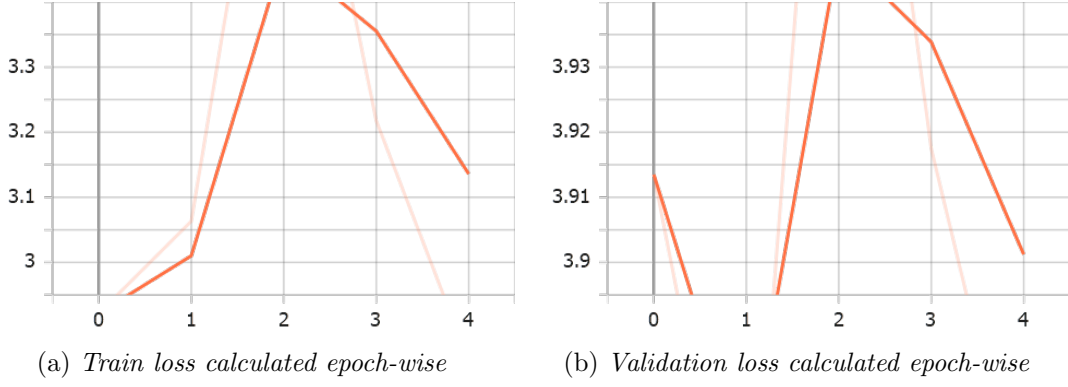


Figure 9: Train and validation losses calculated every epoch for the training on the smaller sample of 100.000 entries

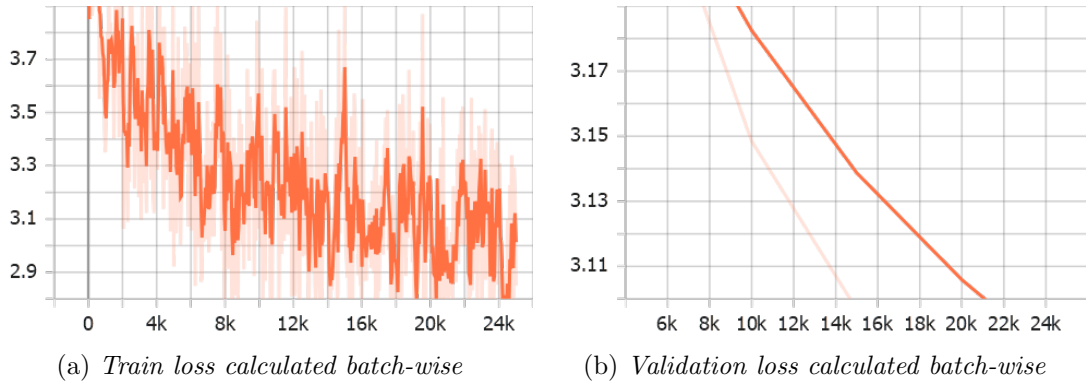


Figure 10: Train and validation losses calculated per batch for the training on the smaller sample of 100.000 entries