

# **ETIM Klassifikation mit ML-Approach**

**Vorhersage der ETIM-Klasse basierend auf  
Daten aus der Geberit Datenbank**

**Studienarbeit FS 2021**

**Version: 1.0**

Autoren: Etienne Baumgartner, Nathanael Gall  
Betreuer: Prof. Oliver Augenstein  
Industriepartner: Geberit AG, Jona  
Experte: Nico Schmid  
Themengebiet: Machine Learning  
Studiengang: Informatik

---

## Abstract

---

<b>Situation</b>	Das ETIM Klassifikationsmodell setzt sich auf dem internationalen Markt immer mehr durch. Die ETIM Klassifikation vereinfacht einerseits den Datenaustausch zwischen Händler und Hersteller und andererseits die Klassifizierung der Produkte. Die Geberit AG nutzt diesen Standard für ihren Produktkatalog und die Klassifizierung wird manuell durch einen Sachbearbeiter vorgenommen.
<b>Ziel</b>	Das Ziel der Arbeit ist die Entwicklung und der Vergleich von zwei Verfahren zur automatischen Bestimmung der ETIM Klasse auf der Basis von maschinellem Lernen
<b>Methode / Vorgehen</b>	<p>Der Erfolg der Arbeit beruht auf der richtigen Verwertung der von Geberit zur Verfügung gestellten Produktdaten. Auf Grund der Analyse wurden zwei Ansätze entwickelt, um die ETIM Klassifizierung zu erlernen.</p> <p>One-Hot-Encoding Die Daten werden für das Trainieren auf einem Deep Neural Network One-Hot-Encoded. Um die Anzahl Dimensionen des Netzwerks möglichst klein zu halten, wird anhand des Bayes Errors die optimale Spaltenkombination ermittelt. Auf dieser Auswahl wird ausserdem ein Lookup Table erstellt, um eindeutige Datensätze direkt zu klassifizieren.</p> <p>Textembedding Viele der Spalten weisen einen grossen Anteil an deutschem Text auf. Es wird mit der fastText Library ein Textembedding trainiert, welches die Artikel den jeweiligen ETIM Klassen zuordnet.</p>
<b>Wesentliche Ergebnisse</b>	<p>Nach der Optimierung dieser zwei Ansätze stellt sich heraus, dass der textbasierte fastText Algorithmus genauere Resultate liefert. Mit dem fastText Model wird auf dem vorab abgetrennten finalen Testset eine Top-1-Accuracy von 0.961 und eine Top-3-Accuracy von 0.982 erreicht.</p> <p>Die höhere Accuracy ist auf die bessere Verwertung der Spaltentexte zurückzuführen, weil das Textembedding die Nähe zweier verschiedener Werte abbilden kann, während beim One-Hot-Encoding alle unterschiedlichen Werte äquidistant sind.</p>
<b>Empfehlungen</b>	Die Empfehlung ist, die drei wahrscheinlichsten Klassen des fastText Top-3-Algorithmus zu präsentieren und die finale Entscheidung einem Benutzer zu überlassen
<b>Schlüsselwörter</b>	Machine Learning, ETIM, Classification, fastText, Neural Network

## Inhaltsverzeichnis

<b>Abstract</b> .....	<b>2</b>
<b>1 Einleitung</b> .....	<b>5</b>
<b>2 Definitionen</b> .....	<b>6</b>
2.1 ETIM .....	6
<b>3 Technologien und Umgebung</b> .....	<b>8</b>
<b>4 Datenanalyse</b> .....	<b>9</b>
4.1 Ausgangslage .....	9
4.2 Beobachtungen .....	10
4.3 Bayes Error Analyse.....	15
4.4 Schlussfolgerungen.....	18
<b>5 Datenaufbereitung</b> .....	<b>20</b>
5.1 Encoding .....	20
5.1.1 Dimensionen .....	20
5.2 Übersicht der Aktionen .....	21
5.2.1 Löschen.....	21
5.2.2 Ersetzen.....	21
5.2.3 Auswählen.....	22
5.3 Unterteilung Test-, Validierungs- und Trainingsset .....	22
<b>6 Algorithmen</b> .....	<b>26</b>
6.1 Definitionen .....	26
6.2 Übersicht.....	26
6.3 Benchmark.....	27
6.4 Lookup Table .....	28
6.4.1 Algorithmus .....	28
6.4.2 Analyse .....	29
6.4.3 Erweiterung der Spaltenauswahl .....	31
6.5 Neuronales Netzwerk .....	32
6.5.1 Algorithmus .....	32
6.5.2 Verbesserungen .....	34
6.6 Text Embedding mit fastText.....	37
6.6.1 Algorithmus .....	37
6.6.2 Verbesserungen .....	38
6.7 Kombinationen .....	39
6.7.1 Top 3.....	39
6.7.2 Lookup und DNN.....	40
<b>7 Vergleich der Resultate</b> .....	<b>41</b>
7.1 Analyse der Resultate .....	41
7.2 Rücksprache mit Geberit.....	42
<b>8 Lösung</b> .....	<b>43</b>
8.1 Klassifizierung .....	43
8.2 Auswertung .....	43

---

8.3	GUI.....	44
8.3.1	Funktionalität.....	44
<b>9</b>	<b>Bewertung, Methodenreflexion, Empfehlungen.....</b>	<b>47</b>
<b>10</b>	<b>Literatur und Quellenverzeichnis.....</b>	<b>49</b>
10.1	Quellenverzeichniss .....	49
10.2	Abbildungen .....	49
10.3	Tabellen .....	49
<b>Anhang</b>	<b>.....</b>	<b>50</b>
I	Auswertung der reservierten Daten für die Bewertung .....	50
II	Bedienungsanleitung GUI.....	51
III	Abgabeordner .....	52
IV	Management Summary .....	55

# 1 Einleitung

---

<b>Auftrag</b>	<p>Diese Semesterarbeit befasst sich mit einem Klassifizierungsproblem, das anhand maschinellen Lernens gelöst werden soll.</p> <p>Die Ausgangslage stellt eine Ansammlung von ungefähr 25'000 Datensätzen dar. Die Aufgabe besteht nun darin verschiedene Algorithmen aufzubauen und eine möglichst gute positive Klassifizierungsrate auf den Daten zu erhalten. Es werden keine Beschränkungen bezüglich der Algorithmen-Wahl gesetzt.</p>
<b>Motivation</b>	<p>Das Interesse von Geberit liegt in der maschinellen Klassifizierung ihrer Produkte. Genauer gesagt geht es um die Zuordnung ihrer Produktpalette zum weltweit genutzten ETIM Standard.</p> <p>Für diese Zuordnung wurde bisher eine manuelle Zuteilung anhand interner Produktdaten durchgeführt. Die Grundlage der manuellen Zuteilungen wurde in Form der 25'000 Datensätze zur Verfügung gestellt.</p> <p>Geberit möchte ihre Möglichkeiten für eine optimierte, automatische Klassifizierung prüfen und übergibt dieses Anliegen in Form einer Semesterarbeit an die Hochschule OST. Der Auftrag setzt sich aus der Analyse der Daten, Aufbereitung der nutzbaren Informationen und dem Vergleich der möglichen Ansätze zusammen.</p>
<b>Ziel</b>	<p>Das Ziel der Arbeit ist die konzeptuelle Untersuchung der Möglichkeiten, den bisher manuellen Klassifizierungsprozess von Produkten zur respektiven ETIM Standard Klasse mit Hilfe von maschinellem Lernen zu erleichtern.</p> <p>Hierbei geht es vor allem darum aufzuzeigen, dass ein maschineller Ansatz zur Klassifizierung einen Mehrwert gegenüber der manuellen Bearbeitung bieten kann. Damit lässt sich eine Zuordnung von neuen Produkten schneller und genauer durchführen. Dabei sollen verschiedene Ansätze untersucht und verglichen werden, um eine passende Lösung zu präsentieren.</p> <p>Optimalerweise wird eine potenzielle Grundlage für eine spätere Software-Lösung in Form eines Prototyps bereitgestellt.</p>
<b>Anwendung der Resultate</b>	<p>Die Voraussetzung dieser Arbeit ist nicht die perfekte Zuordnung der Daten. Ein guter Algorithmus soll eine Einschränkung der möglichen Zuweisungen erlauben und unter anderem potenzielle ETIM Klassifizierungen vorschlagen. Der empfohlene Algorithmus soll einen menschlichen Sachbearbeiter in der Zuordnung der ETIM Klassen weitgehend durch diese vorgeschlagenen Klassifizierungen unterstützen.</p>
<b>Vorgehen</b>	<p>Zu Beginn werden die Daten von Geberit untersucht. Durch diese Analyse werden wichtige Informationen und grundlegende Charakteristiken der Daten aufgedeckt. Auf der Basis dieser Eigenschaften wird die Datenaufbereitung lanciert. Dabei werden die Daten so weit präpariert, dass ein optimales Lernen möglich ist und keine unerwünschten Nebeneffekte, zum Beispiel durch Ausreisser, auftreten. Die aufgearbeiteten Daten werden den jeweiligen Algorithmen eingespeist und die Resultate anschliessend verglichen. Auf Grund der Resultate wird zum Abschluss eine Empfehlung abgegeben und falls zeitlich möglich ein Prototyp in Form einer Anwendung bereitgestellt.</p> <p>Zusammenfassend stellt sich diese Arbeit somit wie folgt zusammen:</p> <ol style="list-style-type: none"><li>1. Datenanalyse</li><li>2. Datenaufbereitung</li><li>4. Algorithmen-Implementation</li><li>5. Resultats Vergleich</li><li>6. Fazit und Empfehlung</li></ol>

## 2 Definitionen

**Allgemein** Bevor mit der eigentlichen Bearbeitung des Themas begonnen werden kann, wird anbei der ETIM Standard erklärt.

### 2.1 ETIM

**Definition** «ETIM ist ein zweistufiges Klassifikationsmodell, das Produktdaten aus dem Fachbereich Elektrotechnik strukturiert erfasst und einen standardisierten elektronischen Datenaustausch von Produktdaten zwischen Herstellern und Handel ermöglicht.» [1]

**Zusammensetzung** Einem Produkt wird eine eindeutige Artikelklasse zugeordnet. Diese Artikelklasse ist wiederum genau einer Artikelgruppe unterstellt. Jede Artikelklasse wird durch ein eindeutiges Set an Merkmalen wie *Name*, *Typ*, *Einheit*, usw., identifiziert.

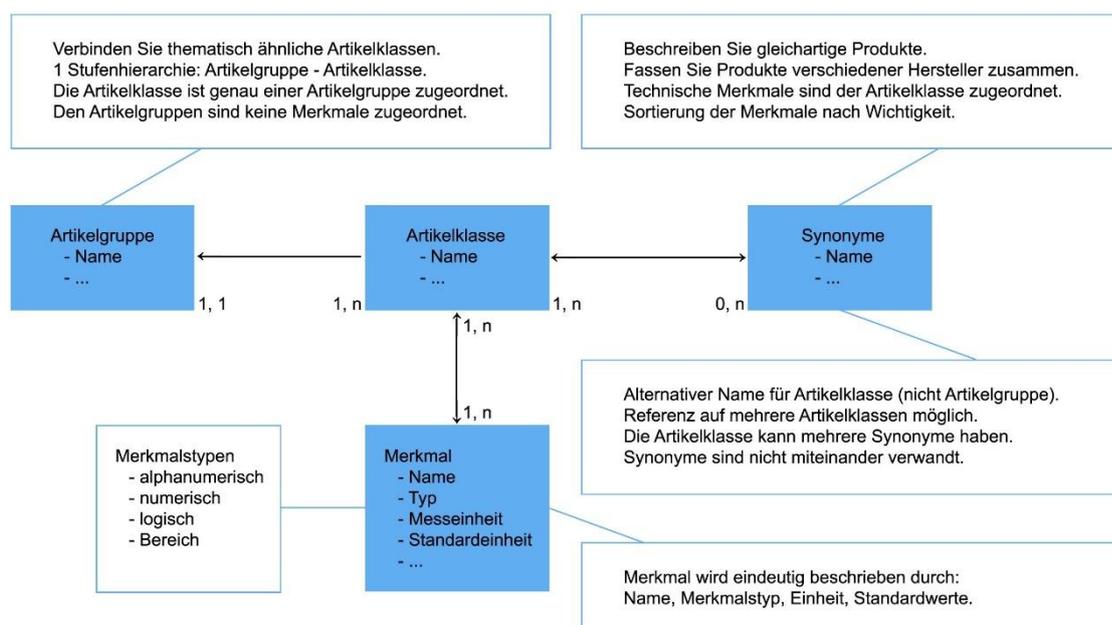


Abbildung 1: ETIM Zusammensetzung, Quelle <https://www.etim.ch/de/klassifizierung/modell-informationen>

**ETIM Klasse** In der Semesterarbeit ist nur die ETIM Klasse von Bedeutung. Die ETIM Klasse vereinigt somit ähnliche Produkte verschiedener Hersteller zu einer Gruppierung. Die Produkte einer Klasse weisen ähnliche Spezifikationen auf und können anhand dieser auch klassifiziert werden.

**ETIM Beschreibung** Jedes Produkt hat neben der ETIM Klasse eine Beschreibung.

**Format** Das Format der ETIM Klasse ist immer gleich: EC + 6 Zeichen  
Das Format der ETIM Beschreibung besteht aus maximal 80 Zeichen und ist textbasiert.

- Relevanz** Die Daten von Geberit enthalten jeweils die ETIM Klasse und die dazugehörige ETIM Beschreibung. Diese liegen in einem 1:1 Verhältnis und daher ist nur die ETIM Klasse relevant.
- Nutzen der Klassifizierung** Sobald einem Produkt eine ETIM Artikelklasse zugeordnet werden kann, können die Merkmale für den ETIM Standard automatisch aus der Produktdatenbank von Geberit extrahiert werden. Diese Zuordnung erfolgt momentan jedoch noch manuell und ist mühselig. Wenn sich ein maschineller Ansatz zur Klassifizierung als plausibel herausstellt, kann auch der Export des Produktkatalogs von Geberit viel effizienter durchgeführt werden.

### 3 Technologien und Umgebung

---

**Sprache** Als Programmiersprache wurde Python gewählt.

Python hat eine reiche Ansammlung an Frameworks und Libraries, die auf Datenanalysen und maschinelles Lernen ausgelegt sind.

Ausserdem gibt Tools, um einfach und schnell Desktopanwendungen zu erstellen.

**Python Pakete** Anbei ist eine Liste der Python Pakete aufgeteilt nach Bereich:

Bereich	Pakete
<b>Analyse</b>	pandas, matplotlib, tabulate, itertools
<b>Aufbereitung</b>	pandas, numpy, sklearn
<b>Algorithmen</b>	pandas, numpy, sklearn, tensorflow, pickle, joblib, fasttext
<b>Prototyp</b>	pandas, numpy, tkinter, pandastable, erstellte Algorithmen.

Genauere Angaben liegen in den jeweiligen Notebooks des Abgabeordners vor.

**Technologien** Der gesamte Code für die Analyse, die Aufbereitung und den Aufbau der Algorithmen wurde in Jupyter Notebooks verfasst.

Der Prototyp ist eine Python Anwendung.

## 4 Datenanalyse

### 4.1 Ausgangslage

**Allgemein**

Die Daten wurden als CSV-Datei zur Verfügung gestellt und bestehen aus Informationen bezüglich 25'826 Artikel in der Geberit Datenbank. Jeder Datensatz hat die folgenden Informationen:

**Datensätze**

Spaltenname	Form	Type	Beschreibung
<b>article_id</b>	ART_{Nummer}	Text/Integer	
<b>name</b>	Combination of numbers and text	Text/Integer	
<b>aricle_number</b>	Nummer vom <Name>	Integer(.)	
<b>etim_description</b>	Beschreibung des Nutzungsraums	Text	
<b>etim_id</b>	EC{Nummer} 6 Digits	Text/Integer	
<b>product_hierarchie</b>	xx..yyy..zzzz + Textbeschreibung	Integer(./)Text	
<b>short_text</b>	Text (+Dimension)	Text	
<b>product_name</b>	{Unternehmen} + Text	Text	
<b>path</b>	{Text}\{Text}\{Text}\...	Text(\)	
<b>system_id</b>	SYS_{Nummer} 6 Digits	Text/Integer	
<b>system_name</b>	Text	Text	2. Pfadtext
<b>category_id</b>	CAT_{Nummer} 6 Digits	Text/Integer	
<b>category_name</b>	Text	Text	3. Pfadtext
<b>family_id</b>	FAM_{Nummer} 6+ Digits	Text/Integer	
<b>family_name</b>	Text	Text	4. Pfadtext
<b>group_id</b>	GRP_{Nummer} 6+ Digits	Text/Integer	
<b>group_name</b>	Text	Text	5. Pfadtext
<b>series_id</b>	SER_{Nummer} 6+ Digits	Text/Integer	
<b>series_name</b>	Text	Text	6. Pfadtext
<b>type_id</b>	TYP_{Nummer} 6+ Digits	Text/Integer	
<b>type_name</b>	Text	Text	

Tabelle 1: Spaltenbeschreibung

Einfachheitshalber werden die Spaltennamen für die weitere Bearbeitung angepasst. Diese ist ein Schritt, der in der Datenaufbereitung durchgeführt wird, jedoch bereits jetzt angewendet wird.

## 4.2 Beobachtungen

**Struktur Geberit** Während der ETIM Standard<sup>1</sup> auf zweistufiger Klassifizierung beruht, scheint Geberit grob einer sechsstufigen Struktur zu folgen. Artikel werden der Reihe nach einem System, einer Kategorie, einer Familie, einer Gruppe, einer Serie und zuletzt einem Typ zugeordnet. Diese Zuordnung korreliert mit der ETIM Klassifizierung, es gibt jedoch kein direktes Mapping.

**ETIM Klassen** Die Produkte sind in 260 der 6'382 ETIM Klasseneingeteilt oder haben die Bezeichnung '\EC000000', welche keiner Klasse des ETIM Standards entspricht. Die Aufteilung in ETIM Klassen ist annähernd Pareto verteilt, fällt jedoch stärker ab.

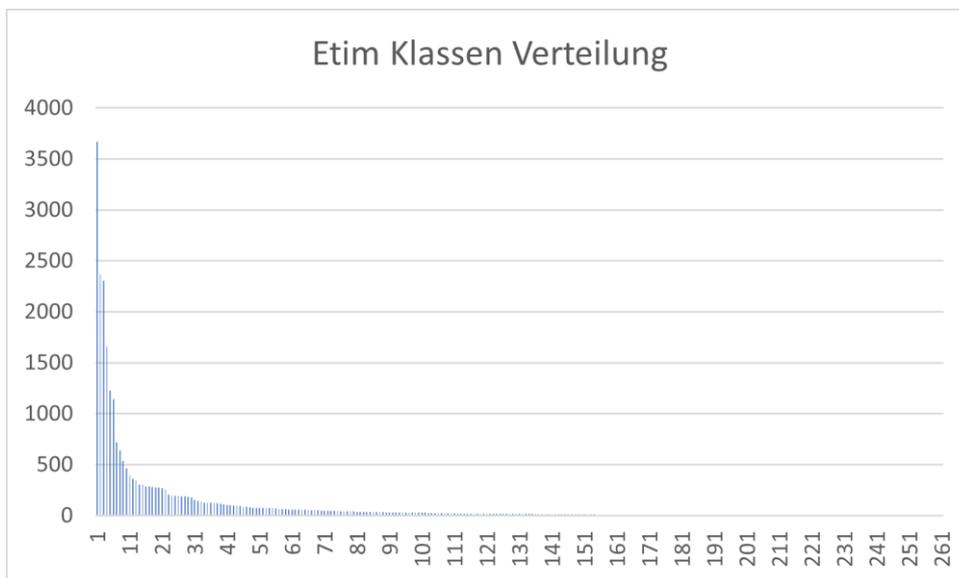


Abbildung 2: ETIM Verteilung

<sup>1</sup> <https://prod.etim-international.com>, Zugriff: 01.06.2021



In den obenstehenden Abbildungen werden die ETIM Klassen oder `etim_ids` und ihre numerischen Vorkommen in unterschiedlichen Darstellungen aufgezeigt. Daraus wird die Annahme gewonnen, dass eine Mindestanzahl an Vorkommen vorausgesetzt werden kann, bevor die Daten von Algorithmen verarbeitet werden.

**Unique Values**

<b>column_name</b>	<b>unique_values</b>	<b>column_name</b>	<b>unique_values</b>
<b>article_id</b>	25826	<b>category_id</b>	57
<b>name</b>	25826	<b>category_name</b>	57
<b>article_number</b>	25826	<b>family_id</b>	225
<b>etim_description</b>	260	<b>family_name</b>	132
<b>etim_id</b>	261	<b>group_id</b>	470
<b>product_hierarchie</b>	930	<b>group_name</b>	296
<b>short_text</b>	23535	<b>series_id</b>	1029
<b>product_name</b>	8712	<b>series_name</b>	519
<b>path</b>	8784	<b>type_id</b>	1774
<b>system_id</b>	7	<b>type_name</b>	590
<b>system_name</b>	7		

Tabelle 2 Unique Values

**article\_id, name und article\_number**

Die Spalten `article_id`, `name` und `article_number` enthalten ausschliesslich einzigartige Werte. Ein Algorithmus, welcher auf diesen Spalten trainiert, erreicht eine Genauigkeit von 1.0 und würde jede `etim_id` aus dem Datenset perfekt zuweisen. Sobald jedoch ein neues Produkt hinzukommt, ist der Algorithmus ahnungslos.

**product\_name**

Die Spalte `product_name` enthält ebenfalls ausschliesslich einzigartige Werte. Diese setzen sich jedoch zusammen aus einzelnen Worten wie dem Namen des Herstellers, Materialien und Grössenangaben des Artikels. Diese könnten theoretisch extrahiert werden.

**etim\_id etim\_description**

Überraschenderweise haben die `etim_id` und die `etim_description` unterschiedlich viele unique Values. Dies ist darauf zurückzuführen, dass die nichtexistierende `etim_id 'EC000000'` ebenfalls dazugezählt wird, jedoch keine `etim_description` dafür existiert.

Wie bereits angesprochen liegen `etim_id` und `etim_description` in einer 1:1 Beziehung zueinander.

**short\_text**

Der `short_text` besteht ebenfalls aus vielen einzigartige Werte.

Im Gegensatz zu `article_id`, `name` und `article_number` ist es jedoch zu voreilig diese Spalte abzuschreiben. Die Werte bestehen ausschliesslich aus Text, woraus potenziell viele Informationen gewonnen werden können.

**path** Aus der Tabelle 1 geht hervor, dass der `path` aus den Spalten `[]_ids` und `[]_names` besteht. Der `path` ist somit eine mögliche Information Akkumulation bestehend aus den `[]_name` Spalten.

Die Grundkonstellation ist die Folgende:

`Systeme\system_name\category_name?family_name\...`

Der Rest des `path` ist besteht aus abgeänderten oder umgeschriebenen Variationen der `group_name`, `series_name` und `type_name` Spalten.

Dies ist der Fall, da der `path` die Ordnerstruktur auf dem Geberit Server widerspiegelt.

**product\_hierarchie** Die erste Hälfte der `product_hierarchie` setzt sich zusammen aus drei Zahlen, die eine Systematische Hierarchie widerspiegelt, nach welcher Geberit ihre Produkte einordnet, aber weder mit ihrer Ordnerstruktur noch mit dem ETIM System übereinstimmt. Die zweite Hälfte ist eine Liste von Wörtern, welche aus den Technischen Daten des Artikels generiert wird.

**[]\_id, []\_name:** Es ist sinnvoll diese Spalten genauer zu untersuchen. Unter anderem weisen sie eine akzeptable Menge von einzigartigen Werten auf und sind die Bausteine aus denen der `path` zusammengesetzt ist. Dabei fällt auf, dass sich die Menge der einzigartigen Werte stetig erhöhen, je «tiefer» die jeweilige Position im `path`.

Ebenfalls auffällig ist die unterschiedlichen Anzahl der einzigartigen Werte für ein `[]_id` und `[]_name` Paar. Eine nähere Betrachtung der Korrelationen zwischen `[]_id` und `[]_name` führt zu folgendem Ergebnis:

1. Alle `[]_id` Werte können **eindeutig** dem jeweiligen `[]_name` zugeordnet werden.
2. `[]_name` Werte können **unterschiedlichen** `[]_id` zugeordnet werden.

Das heisst, dass sich ein Wert aus `[]_id` stets auf den selben `[]_name` Wert abbilden lässt. Umgekehrt gilt jedoch, dass ein `[]_name` Wert mehreren `[]_id` Werten zugeordnet werden kann, z.B.:

1. Die `type_id` 'TYP\_101063' wird immer auf den `type_name` 'Anschlusswinkel für Gaszähler' abgebildet.
2. Der `type_name` 'Anschlusswinkel für Gaszähler' wird aber auf die beiden `type_ids` 'TYP\_101063' und 'TYP\_100973' abgebildet.

Dies führt zu der Annahme, dass die `[]_id` potenziell mehr Informationen trägt als der `[]_name`.

Untenstehend ist eine Tabelle mit Anzahl der Fälle in denen der Wert aus `[]_namen` mehr als einem Wert aus `[]_id` zugeordnet ist.

<code>[]_name</code>	Anzahl <code>[]_id</code> Zuteilungen
<code>system_name</code>	0
<code>category_name</code>	0

family_name	16
group_name	54
series_name	185
type_name	250

Tabelle 3: Zuteilungen der Namen und Ids

NaN

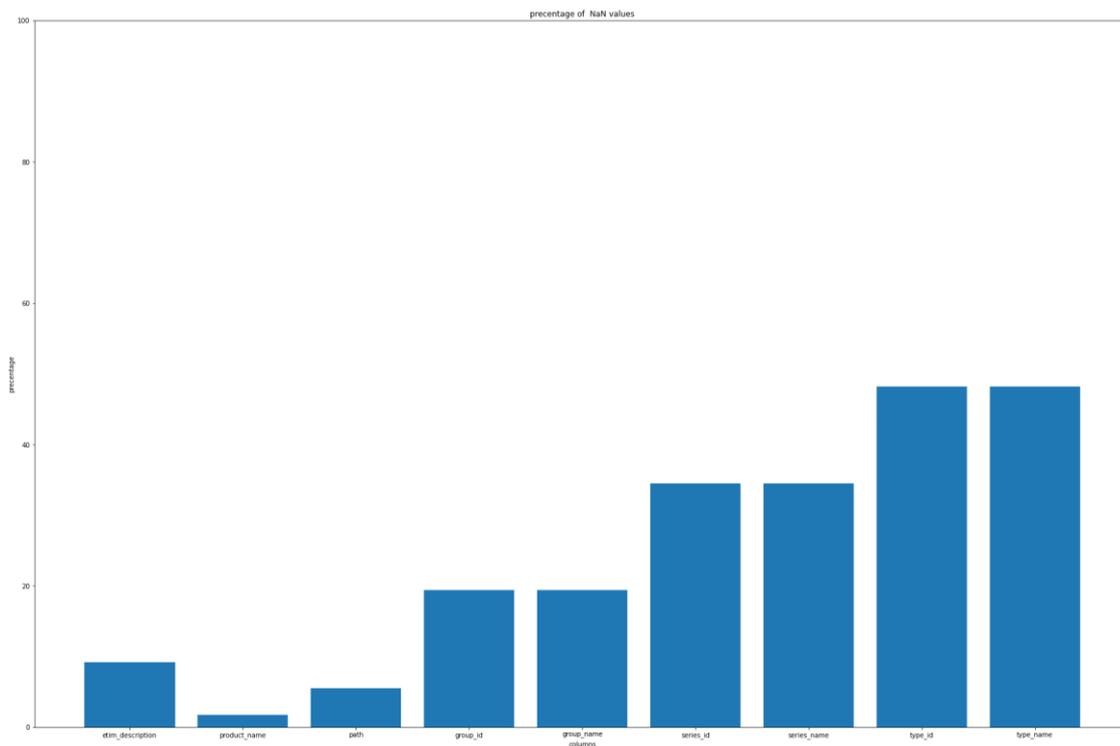


Abbildung 5: Prozentanteil der NaN Werte

Eine weitere wichtige Beobachtung beschreibt die Anzahl der NaN Werte in den jeweiligen Spalten.

Interessant ist, dass sich die NaN Werte mit wenigen Ausnahmen auf die `group_name`, `series_name` und `type_name` Spalten beschränken. Diese Entdeckung unterstützt die Annahme, dass der `path` im Wesentlichen aus den Informationen der Spalten `system_name`, `category_name` und `family_name` bestehen.

Eine weitere spannende Beobachtung ist, dass ungefähr 50% der `type_id` und `type_name` NaN Werte sind.

Der Prozentsatz der NaN Werte für die `etim_description` ist damit zu erklären, dass der `etim_id` 'EC00000' keine `etim_description` zugeordnet werden kann.

### 4.3 Bayes Error Analyse

- Motivation** Der Bayes Error soll per Definition die kleinstmögliche Error Rate berechnen. Aus den gewonnenen Informationen werden Entscheidungen für die weitere Entwicklung der Algorithmen hervorgehen.
- Fehlerursache** Mit den gegebenen Daten gibt es nur eine einzige mögliche Ursache für unvermeidbare Fehler. Wenn es in den Trainingsdaten mehrere Datensätze gibt, welche identischen Input aber unterschiedliche `etim_ids` haben, gibt es keine bessere Strategie, als die Häufigste dieser `etim_ids` auszugeben.
- Probleme** Die Geberit Daten haben ein Problem, was den Bayes Error betrifft. Es gibt einige Spalten, die einzigartig auf allen Datensätzen sind. Dies führt zu einem Bayes Error von Null. Deshalb wird der Bayes Error auf diesen Spalten nicht berechnet.
- Vorgehen** Für die Entwicklung eines Algorithmus braucht es Informationen bezüglich der Spalten auf denen gelernt werden soll. Für dieses Projekt wird ein Bayes Algorithmus aufgebaut, sodass einzelne irreduzible Error Analysen auf einer bestimmten Anzahl Spalten durchgeführt werden können. Dieser Algorithmus liefert Informationen für eine bestmögliche Auswahl an Spalten, auf denen trainiert werden sollte.
- Auswahlverfahren** Für die Auswahl der Spalten präsentieren sich einige Kriterien:
- Anzahl Kolonnen
  - Anzahl einzigartigen Werte
  - Bayes Error Rate
- Hinweis** Für den Bayes Algorithmus wurden die Spalten `article_id`, `name` und `article_number` bereits ausgenommen. Sie weisen eine einzigartige Zuordnung zu den `etim_ids` auf und sind daher nicht brauchbar.
- Algorithmus** Der Algorithmus nimmt eine Menge von Inputspalten und gruppiert alle Datensätze nach Kombinationen dieser Spaltenwerte. Jede dieser Gruppen wird nach der `etim_id` gruppiert. Für jede Kombination der Spaltenwerte des Inputs wird diese `etim_id` Gruppe gelöscht, welche die höchste Kardinalität aufweist. Die restlichen Datensätze werden von einem beliebigen Algorithmus, der die Inputdaten perfekt gelernt hat, nie korrekt klassifiziert. Teilt man die Anzahl Datensätze die übrig geblieben sind durch die gesamte Anzahl Datensätze, ergibt sich daraus der Bayes Error.

**Resultat**

Component 0	uniques	Bayes Error
system_id	7	0.749090064
category_id	57	0.557384032
family_name	132	0.428599086
family_id	225	0.353403547
type_name	590	0.31749944
group_name	296	0.28394469

product_hierarchie	930	0.25811198
group_id	470	0.220424429
series_name	519	0.173180394
series_id	1029	0.098385857
type_id	1774	0.045559788
product_name	8712	0.011586207
path	8784	0.007129687
short_text	23535	0.002671726

**Tabelle 4: Bayes Error mit einzelnen Komponenten**

Die Zusammenstellung der einzelnen Spalten und ihren Bayes Errors ist logisch nachvollziehbar. Je mehr einzigartige Werte in einer Spalte, desto einfacher ist eine Zuordnung der `etim_id`.

Somit wäre es auf den bestehenden Daten theoretisch möglich eine korrekte Zuordnung mit Wahrscheinlichkeit von 0.997328274 zu erreichen, falls man sich ausschliesslich auf den `short_text` fokussiert.

Jedoch wäre es keinem Algorithmus möglich eine Zuordnung durchzuführen, wenn ein neues Produkt mit neuem `short_text` hinzugefügt werden soll.

Für die ursprüngliche Auswahl kommen also nur diejenigen Spalten in Frage, welche nicht zu viele einzigartigen Werte aufweisen.

Deshalb wird ab hier für die Auswahl der Spalten die Anzahl der einzigartigen Werte auf 1000 beschränkt.

Component 0	unique	Component 1	unique	Bayes Error
series_name	519	type_name	590	0.058376156
group_id	470	type_name	590	0.063803005
product_hierarchie	930	type_name	590	0.066547166
group_name	296	type_name	590	0.072142239
product_hierarchie	930	series_name	519	0.089930822
group_id	470	series_name	519	0.09122191

**Tabelle 5: Bayes Error mit zwei Komponenten**

Diese Tabelle enthält alle 2-er Kombinationen, die unter einem Bayes Error von 0.1 liegen.

Es kann bereits eine leichte Tendenz herausgelesen werden. Die Kombinationen bestehen ausschliesslich aus den Spalten `product_hierarchie`, `group_id`, `series_name` und `type_name`. Es ist zu erwarten, dass sich dieser Trend auch mit 3-er Kombinationen fortsetzen wird.

Component 0	Component 1	Component 2	Bayes Error
-------------	-------------	-------------	-------------

product_hierarchie	series_name	type_name	0.029085303
group_id	series_name	type_name	0.035952637
group_name	series_name	type_name	0.040043057
product_hierarchie	group_id	type_name	0.042246873
product_hierarchie	group_name	type_name	0.044921721
family_id	series_name	type_name	0.047790339
family_name	series_name	type_name	0.047995889

**Tabelle 6: Bayes Error mit drei Komponenten**

Für obenstehende Auswahl der 3-er Kombinationen wurde der Bayes Error auf 0.05 gesenkt.

Die vorhergesagte Tendenz bezüglich der Spalten erfüllt sich mit einer Ausnahme. Eine Kombination mit `family_id` schafft es ebenfalls, die Kriterien zu erfüllen.

Ab diesem Punkt ist es möglich eine Voraussage bezüglich den 4-er Kombinationen zu machen. Dem bereits in den 2-er Kombinationen bemerkten Trend folgend, wird die beste Kombination aus `product_hierarchie`, `group_id`, `series_name` und `type_name` bestehen.

Component 0	Component 1	Component 2	Component 3	Bayes Error
product_hier.	group_id	series_name	type_name	0.02583423
product_hier.	group_name	series_name	type_name	0.026372443
product_hier.	family_name	series_name	type_name	0.028468654
product_hier.	family_id	series_name	type_name	0.028468654
product_hier.	system_id	series_name	type_name	0.029085303
product_hier.	category_id	series_name	type_name	0.029085303

**Tabelle 7: Bayes Error mit vier Komponenten**

Die Voraussage trifft zu und somit ergibt sich Folgendes als beste Kombination:

- Spalten: `product_hierarchie`, `group_id`, `series_name`, `type_name`
- Anzahl einzigartigen Werte: < 1000
- Bayes Error: 0.0258342303552207

Die gesamten Daten zur Bayes Analyse sind im Abgabeordner unter `analysis/bayes_error_all_data_kombinationen.xlsx` abgelegt.

**Schlussfolgerung** Anhand der Bayes Error Analyse können die vier Spalten `product_hierarchie`, `group_id`, `series_name`, `type_name` als aussagekräftigste Kombination identifiziert werden.

Es gibt auf den originalen Daten also einen theoretischen irreduziblen Fehler von 0.02583423, solange die ausgewählten Spalten als Feature Spalte genutzt werden.

Die Bayes Analyse hat einen weiteren Vorteil.

Es ist möglich den Bayes Error auf einem der Trainingsset zu berechnen, um eine obere Grenze für die Trainingsaccuracy von einem Algorithmus zu bestimmen.

## 4.4 Schlussfolgerungen

<b>EC000000</b>	<p>Alle Datensätze, die eine <code>etim_id</code> von <code>'EC000000'</code> aufweisen, werden aus den Daten entfernt.</p> <p>Geberit verwendet die ETIM Klasse <code>'EC000000'</code> für alle Artikel, welche aus unbestimmten Gründen nicht durch einen Sachbearbeiter manuell klassifiziert werden können. Dies geschieht einerseits durch Human Error oder ist darauf zurückzuführen, dass die passende ETIM Klasse im ETIM Standard noch nicht vorhanden ist.</p> <p>Es wäre interessant, ein Algorithmus zu trainieren, der vorhersagt, in welchen Fällen für einen Artikel keine ETIM Klasse existiert. Es gibt anhand der verfügbaren Daten jedoch keine Indizien dafür, ob dieser Fall vorliegt oder ob es sich schlichtweg um einen Human Error handelt. Basierend auf dieser Überlegung werden die <code>'EC000000'</code> für das Training ignoriert.</p>
<b>Einzigartige Werte</b>	<p>Die Spalten <code>article_id</code>, <code>name</code> und <code>article_number</code> weisen für jeden Datensatz einzigartige Werte auf und haben für die Problembehandlung keinen Nutzen.</p>
<b>system_id system_name</b>	<p>Diese beiden Spalten sind redundant. Sie weisen eine eindeutige Zuweisung, also eine 1:1 Beziehung zueinander auf.</p>
<b>category_id category_name</b>	<p>Diese beiden Spalten sind redundant. Sie weisen eine eindeutige Zuweisung, also eine 1:1 Beziehung zueinander auf.</p>
<b>Top etim_id</b>	<p>Die acht meist vorkommenden <code>etim_ids</code> machen ungefähr 50% der Zuteilungen aus. Viele <code>etim_ids</code> weisen ein Vorkommen von weniger als zwei auf. Diese Zeilen werden ignoriert.</p>
<b>Bayes Error</b>	<p>Aus der Bayes Analyse geht hervor, dass <code>product_hierarchie</code>, <code>group_id</code>, <code>series_name</code>, <code>type_name</code> voraussichtlich gut geeignete Spalten für das Trainieren bilden.</p>
<b>path</b>	<p>Der <code>path</code> besteht im Wesentlichen aus den Informationen der <code>system_name</code>, <code>category_name</code> und <code>family_name</code> Spalten. Diese Spalte weist daher eine akkumulierte Informationsansammlung auf.</p>
<b>Duplikate</b>	<p>Viele Datensätze sind in ihren Spaltenwerten identisch. Die Duplikationen werden nach einer optimalen Spaltenauswahl entfernt. Dies ist nötig, damit ein Algorithmus einen mehrfach vorkommenden Dateninput nicht unproportional gut lernt.</p>

**NaN Werte**

NaN Werte werden ohne Informationsverlust mit dem String: 'no information' ersetzt.

---

## 5 Datenaufbereitung

---

### 5.1 Encoding

---

**Allgemein** Da ausschliesslich mit kategorischen Daten gearbeitet wird, ist es unumgänglich mit Encodings zu arbeiten. Die Art des Encodings ist abhängig von der Wahl der Algorithmen. Für die Benutzung der meisten maschinell lernenden Algorithmen sind numerischen Werten vorzusetzen.

**One-Hot-Encoding** Dafür bietet sich das One-Hot-Encoding an. Es bildet jeden einzigartigen Wert auf einen einzigartigen n-dimensionalen Vektor aus Nullen und Einsen ab.

Damit können die kategorischen textbasierten Informationen auf ein numerisch basiertes System abgebildet werden, die jeweils denselben Abstand zueinander haben. Das ist vorteilhaft, da lernende Algorithmen dadurch keine Ordnung auf kategorischen Daten herauslesen können.

Dieses Encoding wird für den Aufbau und Training eines Neuralen Netzwerks benutzt.

**Embeddings** Eine zweite Möglichkeit ist der Gebrauch eines Textembeddings. Die Aufgabe von Textembeddings besteht darin, aus textbasierten Informationen Zusammenhänge zu erlernen und diese auf einen Vektor abzubilden.

Da ein Grossteil der Spaltenwerte deutschsprachiger Text ist, könnte ein schon trainiertes deutsches Embedding verwendet werden. Dieses wäre zum Beispiel auf der Gesamtheit der deutschen Wikipedia Seite trainiert worden.

Dies hat jedoch zwei Nachteile:

1. Das Training ist «unsupervised». Das heisst der Zusammenhang zwischen dem Resultierenden Vektor und der `etim_id` noch muss noch separat erlernt werden.
2. In den Daten sind die meisten Wörter aus dem Sanitärbereich. Im schlimmsten Fall sogar Wörter, welche es gar nicht gibt wie zum Beispiel '50x50'. Dies führt dazu, dass ein ungemein grosses Model gebraucht wird, damit in sich so ähnlichen Texten wesentliche Unterschiede erkannt werden.

Aus den obig genannten Gründen ergibt es Sinn ein eigenes Embedding mit der fast-Text Library von Facebook<sup>2</sup> zu erstellen.

#### 5.1.1 Dimensionen

---

**Allgemein** Während beim Textembedding die Anzahl Dimensionen willkürlich festgelegt werden können, ist dies beim One-Hot-Encoding nicht so einfach möglich. Durch die Datenanalyse wird aufgezeigt, wie viele einzigartige Werte in jeder Spalte vorhanden sind. Die Kodierung all dieser Werte mit dem oben beschriebenen One-Hot-Encodings führt zu einer riesigen Anzahl an Dimensionen.

---

<sup>2</sup> <https://fasttext.cc/>, Zugriff: 01.06.2021

---

<b>Dimensionalität</b>	Um diesen riesigen Dimensionen entgegenzuwirken, werden durch die Analyse des Bayes Errors diejenigen Spalten ausgewählt, welche mit möglichst geringer Anzahl verschiedener Werte, einen möglichst geringen Bayes Error erreichen. Dies Spalten weisen das beste Verhältnis von Information zur Grösse der Dimensionen auf.
<b>Resultat</b>	Mit Hilfe der Bayes Error Analyse, wurden <code>etim_id</code> , <code>product_hierarchie</code> , <code>group_id</code> , <code>series_name</code> , <code>type_name</code> und, wie noch aufgezeigt wird, <code>path</code> als beste Spaltenkombination identifiziert und somit werden die Dimensionen der One-Hot-Encodings drastisch reduziert.

## 5.2 Übersicht der Aktionen

---

<b>Übersicht</b>	In diesem Abschnitt werden allgemeine Tatsachen der Daten und die daraus resultierenden Aktionen durch die Aufbereitung beschrieben. Diese Informationen gehen direkt aus den Schlussfolgerungen der Datenanalyse hervor und werden hier übersichtlich dargestellt.  Die Aktionen finden in der nachfolgenden Sektion bei der Unterteilung der Datensets ihre Anwendung.
------------------	--

### 5.2.1 Löschen

---

<b>EC000000</b>	Dies stellt die einfachste Entscheidung für die Aufbereitung dar. Datensätze mit einer <code>etim_id</code> <code>'EC000000'</code> werden im maschinellen Lernen nicht gebraucht. Diese 2367 Zeilen werden daher aus dem Datensatz gelöscht.
<b>Vorkommen</b>	Alle <code>etim_ids</code> , die weniger als 2-mal im Datenset vorkommen, sind nicht aussagekräftig genug und werden deshalb entfernt.
<b>Duplikate</b>	Duplikate der Datensätze müssen gelöscht werden, da sie den Lernprozess der Algorithmen negativ beeinflussen können.

### 5.2.2 Ersetzen

---

<b>NaN</b>	Beim Import entstehen aus nicht ausgefüllten Zellen NaN Zellen. Diese Zellen werden durch den String <code>'no information'</code> ersetzt.
<b>Einmalig vorkommende Keys</b>	Die Werte, die einmalig pro Spalte vorkommen, haben entweder keinerlei oder zu viel Einfluss auf die Entscheidung des Labels. Diese Zellen werden mit dem String <code>'unique'</code> ersetzt.

### 5.2.3 Auswählen

- article\_id** Die `article_id` besteht nur aus einmalig vorkommenden Werten. Sie hat daher keine Relevanz für die Problemlösung.
- article\_number** Die `article_number` besteht nur aus einmalig vorkommenden Werten. Sie hat daher keine Relevanz für die Problemlösung.
- Hinweis** Aus der Datenanalyse geht hervor, dass neben `article_id` und `article_number` noch weitere Spalten als redundant oder als nicht relevant eingestuft wurden.
- Diese Spalten wurden hier bewusst nicht gelistet, da sie für das Textembedding durchaus einen Mehrwert haben können. Vor allem die Spalten `name` und `product_name` weisen einen hohen potenziellen textbasierten Informationsgehalt auf.

### 5.3 Unterteilung Test-, Validierungs- und Trainingsset

**Daten zur Beurteilung der Studienarbeit** Vor der Datenübergabe an die Studierenden wurden 1000 Datensätze entfernt, anhand derer das Resultat der Arbeit gemessen werden kann. Diese Daten werden den Studierenden vorenthalten. Im Rahmen dieser Arbeit wird davon ausgegangen, dass die übrig gebliebenen `geberit_data_original` alle Daten darstellen.

Die Auswertung mit diesen 1000 Datensätzen wird jedoch im Anhang **Fehler! Verweisquelle konnte nicht gefunden werden.** noch aufgeführt.

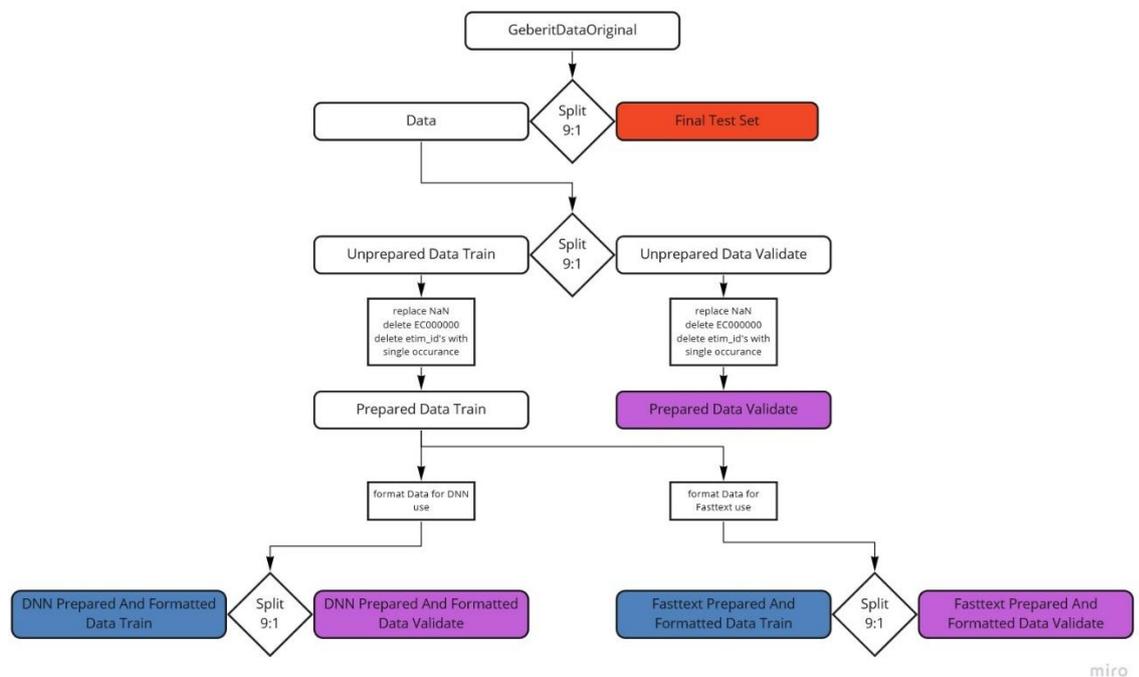


Abbildung 6: Graphische Darstellung der Datenaufbereitung erstellt mit <http://www.miro.com>

---

<b>Mehrstufige Parameter</b>	Die Entwicklung der Algorithmen spielt sich auf drei Ebenen ab: <ol style="list-style-type: none"><li>1. Das Training der Modelle.</li><li>2. Das Optimieren der Hyperparameter.</li><li>3. Das Entscheiden wie gut die einzelnen Modelle sind und welche davon Geberit empfohlen werden sollen auf <code>prepared_data_validate</code>.</li></ol>
<b>Testset</b>	Um sich ein Bild zu verschaffen, wie sich ein Algorithmus auf neuen Daten verhält, ist ein Test-, oder Validierungsset nötig mit Daten, mit welchen der Algorithmus bis zu diesem Punkt noch nicht interagiert hat.
<b>final_test_set</b>	Um am Ende der Arbeit eine gute Metrik zu haben, wie gut der Algorithmus ist, wird das <code>final_test_set</code> im Verhältnis 1:9 von <code>geberit_data_original</code> abgetrennt. Es umfasst 2583 unbehandelte Daten.
<b>Unprepared Data</b>	Die aus dem ersten Split verbleibenden 23'243 Datensätze werden nun weiter in <code>unprepared_data_validate</code> und <code>unprepared_data_train</code> aufgeteilt. Der zweite Split wird ebenfalls im Verhältnis 1:9 durchgeführt.
<b>Prepared Data</b>	Auf <code>unprepared_data_validate</code> und <code>unprepared_data_train</code> wird nun eine erste Datenaufbereitung durchgeführt. <ol style="list-style-type: none"><li>1. NaN Werte werden durch 'no information' ersetzt</li><li>2. Records mit <code>etim_id</code> 'EC000000' werden gelöscht</li><li>3. <code>etim_ids</code> mit weniger als 2 Vorkommen werden gelöscht</li></ol> <p>Die beiden bearbeiteten Sets werden nun als <code>prepared_data_validate</code> und <code>prepared_data_train</code> abgespeichert.</p> <p><code>prepared_data_validate</code> wird gebraucht, um die Modelle, welche aus den verschiedenen Ansätzen entstehen, zu vergleichen.</p> <p><code>prepared_data_train</code> stellt die Grundlage des Trainierens aller Modelle dar.</p> <p>Falls ein Algorithmus Hyperparameter optimiert, muss das Validierungsset dafür ebenfalls aus <code>prepared_data_train</code> kommen.</p>
<b>Lookup_prepared_and_formatted_data</b>	Da der Lookup Table keine Hyperparameter hat und direkt trainieren kann, müssen dafür keine weiteren Massnahmen getroffen werden.
<b>dnn_prepared_and_formatted_data</b>	Das DNN ist der erste Algorithmus, welcher zusätzliche Datenaufbereitung benötigt: <ol style="list-style-type: none"><li>1. Spaltenauswahl: <code>etim_id, product_hierarchie, group_id, path, series_name, type_name, path.</code></li><li>2. Löschen der Duplikate der ausgewählten Spalten.</li><li>3. Einzigartige Werte ersetzen durch 'unique'</li></ol> <p>Da die Daten sehr unregelmässig verteilt sind, tendiert das DNN dazu, seltene Inputs zu ignorieren und stattdessen die häufigsten Inputs zu lernen. Der 2. Schritt soll dem entgegenwirken. Er verhindert, dass derselbe Input mehrmals pro Durchlauf trainiert wird.</p>

Die im 3. Schritt als `'unique'` gekennzeichneten Werte werden in einem separatem Dictionary festgehalten, damit sie bei der Trainingsvalidierung ebenfalls ersetzt werden können.

Für das Training der Hyperparameter muss von dem formatierten Trainingsset ein Validierungsset abgetrennt werden. Um zu verhindern, dass im Validierungsset `etim_ids` vorkommen, welche im Trainingset nicht vorkommen, werden die Daten auf jeder `etim_id` separat im Verhältnis 9:1 auf `dnn_prepared_and_formatted_data_train` und `dnn_prepared_and_formatted_data_validate` aufgeteilt.

**fasttext\_prepared\_and\_formatted\_data**

Für fastText werden die Daten folgendermassen formatiert:

1. Konkatenieren des Inhalts der Spalten `name`, `short_text`, `product_name`, `path`, `system_name`, `category_name`, `group_name`, `series_name` und `type_name` sowie die zweite Hälfte der `product_hierarchie` zu einem String `text`.
2. Umwandeln zur Form:  
`'__label__' ++ etim_id ++ ' ' ++ text`
3. Zufälliges Aufteilen der Daten im Verhältnis 9:1 auf `fasttext_prepared_and_formatted_data_train.txt` und `fasttext_prepared_and_formatted_data_validate.txt`.

**Zusammenfassung** Anbei eine Zusammenfassung der Aktionen auf den jeweiligen Sets:

Daten Set	Aktionen
<b>original_data</b>	Split
<b>data</b>	Split
<b>final_test_set</b>	-
<b>unprepared_data_train</b>	-
<b>unprepared_data_validate</b>	-
<b>prepared_data_train</b>	Löschen & Ersetzen
<b>prepared_data_validate</b>	Löschen & Ersetzen
<b>dnn_prepared_and_formatted_data_train</b>	Auswahl & Löschen
<b>dnn_prepared_and_formatted_data_validate</b>	Auswahl & Löschen
<b>fasttext_prepared_and_formatted_data_train</b>	Auswahl
<b>fasttext_prepared_and_formatted_data_validate</b>	Auswahl

Tabelle 8: Datenbearbeitung Aktionen

**Übersicht und numerische Informationen** Die angrenzende Tabelle enthält numerische Informationen bezüglich der verschiedene Sets:

Daten Set	# Records	# etim_id
original_data	25826	261
data	23243	255
final_test_set	2583	176
unprepared_data_train	20918	254
unprepared_data_validate	2325	170
prepared_data_train	18982	218
prepared_data_validate	2109	169
dnn_prepared_and_formatted_data_train	8213	218
dnn_prepared_and_formatted_data_validate	1027	203
fasttext_prepared_and_formatted_data_train	17083	218
fasttext_prepared_and_formatted_data_validate	1899	172

Tabelle 9: Datenbearbeitung numerische Details

**Nutzung der Sets** Folgend wurde der Verwendungszweck all der Datensets aufgelistet:

Daten Set	Verwendung
data	Finales Trainieren der Modelle
final_test_set	Finale Validierung der Modelle
unprepared_data_train	Zwischenschritt
unprepared_data_validate	Zwischenschritt
prepared_data_train	Basis der Trainingsdaten, Lookup
prepared_data_validate	Validierung Trainingsmodelle, bestimmen, welche Kombination davon sinnvoll sind.
dnn_prepared_and_formatted_data_train	Training des DNN Modells
dnn_prepared_and_formatted_data_validate	Hyperparameter Optimierung des DNN Modells
fasttext_prepared_and_formatted_data_train	Training des fastText Modells
fasttext_prepared_and_formatted_data_validate	Hyperparameter Optimierung des fast-Text Modells

Tabelle 10: Datenbearbeitung Zweck

## 6 Algorithmen

### 6.1 Definitionen

**Top 1** Ein Top-1 Modell liefert jeweils die wahrscheinlichste EIMT Klasse für einen bestimmten Input.  
Bei jedem Top-1 Ergebnis entscheidet jeweils der Algorithmus, welche `etim_id` ausgewählt wird. Aus der Problemstellung geht jedoch heraus, dass dies bisher durch manuelle Bearbeitung gemacht wurde.

**Top 3** Bei den Top-3 Modellen ist ein Benutzer gezwungen, sich am Entscheidungsprozess zu beteiligen, indem ihm drei `etim_id` zur Auswahl präsentiert werden.  
Top-3 Algorithmen sind somit eine Kombination aus maschineller Vorarbeit und menschlicher Schlussentscheidung.

**Resultate und Metriken** Das Resultat eines Algorithmus setzt sich aus folgenden Metriken zusammen:

Resultat	Beschreibung
<b>Trainingsaccuracy</b>	Die Trainingsaccuracy beschreibt die Genauigkeit des Modells auf den Trainingsdaten
<b>Validationaccuracy</b>	Die Validationaccuracy beschreibt die Genauigkeit des Modells auf den Validierungsdaten
<b>Bayes/Trainings Differenz</b>	Die Bayes/Trainings Differenz beschreibt Differenz zwischen der Genauigkeit des Modells auf den Trainingsdaten und der theoretisch bestmöglichen Genauigkeit auf den Trainingsdaten.
<b>Top 1 Accuracy</b>	Ein Top-1 Resultat enthält jeweils nur die wahrscheinlichste ETIM Klasse und beschreibt die vollautomatische Klassifizierung der Produkte. Die Top-1-Accuracy beschreibt die Genauigkeit des Top-1 Resultats auf den Testdaten <code>prepared_data_validate</code> .
<b>Top 3 Accuracy</b>	Die Top-3-Accuracy beschreibt den Recall für $k=3$ den Testdaten <code>prepared_data_validate</code> . Oder in anderen Worten, die Wahrscheinlichkeit, dass wenn ein Modell die drei wahrscheinlichsten <code>etim_ids</code> ausgibt, die Richtige dabei ist.

Tabelle 11 Verwendete Metriken

### 6.2 Übersicht

**Benchmark** Nach der Analyse und Aufbereitung der Daten, wird zuerst ein Benchmark Algorithmus erstellt. Das Ziel dieses Algorithmus ist, eine grobe Einschätzung zu machen, wie gut die Daten lernbar sind.

<b>Lookup Table</b>	Ein Lookup Table auf den Daten soll alle eindeutig zuweisbaren Fälle übernehmen und kann in Kombination mit weiteren Modellen als Vorarbeiter agieren.
<b>DNN</b>	Ein Deep Neural Network oder DNN soll die Zuweisungen der <code>etim_ids</code> mit Hilfe von One-Hot-Encodings auf ausgewählten Spalten erlernen.
<b>fastText</b>	Der fastText Algorithmus soll die Zuweisung der <code>etim_ids</code> anhand des hohen Informationsgehaltes der Spaltenwerte durch Text Embeddings erlernen.
<b>Kombinationen</b>	Kombinationen und Varianten der oben genannten Algorithmen werden in die Analyse miteinbezogen.

## 6.3 Benchmark

<b>Data Set</b>	Für den Benchmark wird das Data Set <code>prepared_data_train</code> benutzt.										
<b>Data Selection</b>	Der Benchmark läuft ausschliesslich auf den Kolonnen <code>etim_id</code> und <code>product_hierarchie</code> .										
<b>Vorgehen</b>	Für die Benchmark wird die <code>product_hierarchie</code> One-Hot encodiert, um damit ein Neuronales Netzwerk zu trainieren.										
<b>Ausführung</b>	<p>Diese Daten werden dann im Verhältnis 3:1 in ein Trainings-, und Validierungsset aufgeteilt trainiert:</p> <table> <tr> <td><b>Epochen</b></td> <td>30</td> </tr> <tr> <td><b>Batchsize</b></td> <td>128</td> </tr> <tr> <td><b>Optimizer</b></td> <td>Adam</td> </tr> <tr> <td><b>Learning rate</b></td> <td>0.01</td> </tr> <tr> <td><b>Netzwerk</b></td> <td> <pre> model = Sequential() model.add(Dense(800, activation='relu', kernel_initializer='he_normal', input_shape=(n_features,))) model.add(Dense(700, activation='relu', kernel_initializer='he_normal')) model.add(Dense(300, activation='relu', kernel_initializer='he_normal')) model.add(Dense(260, activation='relu', kernel_initializer='he_normal')) model.add(Dense(260, activation='softmax')) </pre> </td> </tr> </table>	<b>Epochen</b>	30	<b>Batchsize</b>	128	<b>Optimizer</b>	Adam	<b>Learning rate</b>	0.01	<b>Netzwerk</b>	<pre> model = Sequential() model.add(Dense(800, activation='relu', kernel_initializer='he_normal', input_shape=(n_features,))) model.add(Dense(700, activation='relu', kernel_initializer='he_normal')) model.add(Dense(300, activation='relu', kernel_initializer='he_normal')) model.add(Dense(260, activation='relu', kernel_initializer='he_normal')) model.add(Dense(260, activation='softmax')) </pre>
<b>Epochen</b>	30										
<b>Batchsize</b>	128										
<b>Optimizer</b>	Adam										
<b>Learning rate</b>	0.01										
<b>Netzwerk</b>	<pre> model = Sequential() model.add(Dense(800, activation='relu', kernel_initializer='he_normal', input_shape=(n_features,))) model.add(Dense(700, activation='relu', kernel_initializer='he_normal')) model.add(Dense(300, activation='relu', kernel_initializer='he_normal')) model.add(Dense(260, activation='relu', kernel_initializer='he_normal')) model.add(Dense(260, activation='softmax')) </pre>										

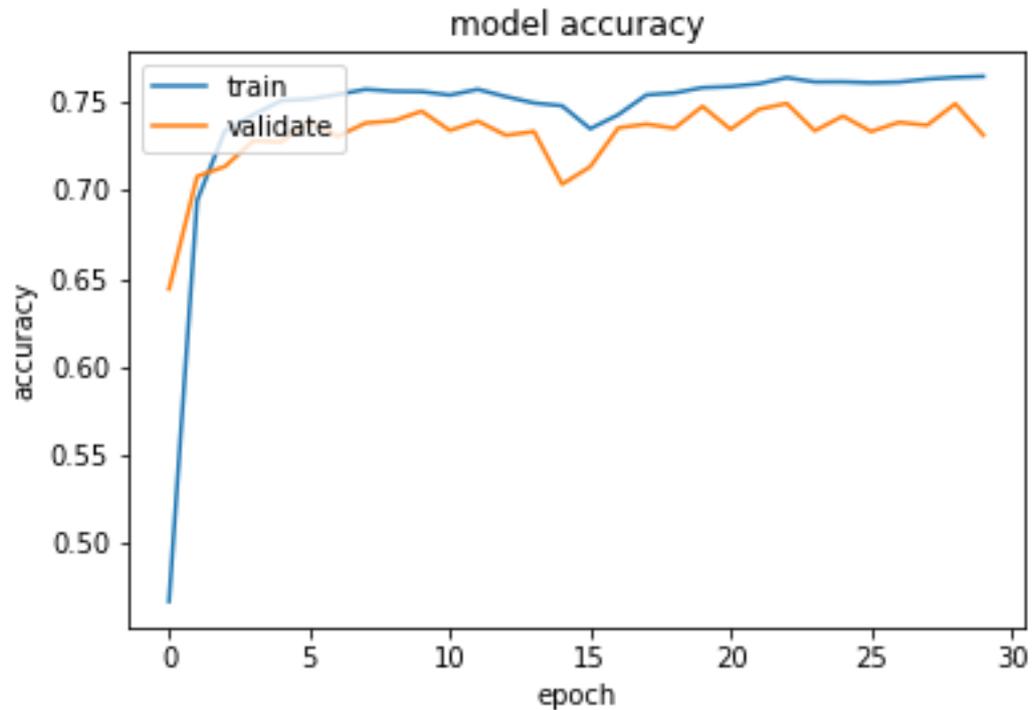


Abbildung 7: Benchmark Train/Validate Accuracy

**Train Accuracy** Der Benchmark erreicht eine Trainingaccuracy von 0.767.

**Validation Accuracy** Ohne weitere Optimierung kann der Benchmark Algorithmus Daten aus dem Validierungsset mit einer Accuracy von 0.749 richtig klassifizieren. Wenn also ein Algorithmus diese Baseline nicht klar übertrifft, wird er verworfen.

## 6.4 Lookup Table

### 6.4.1 Algorithmus

**Daten** Für den Aufbau des Lookup Table wird `prepared_data_train` benutzt. Diese Daten enthalten keine NaN Werte, keine EC000000 und keine `etim_ids` mit weniger als zwei Vorkommen. Diese Umstände bilden eine passende Grundlage für den Lookup Table.

**Daten Auswahl** Die Basis des Lookup bildet folgende Kombination von Spalten:

```
etim_id, product_hierarchie, group_id, series_name,
type_name
```

Die Auswahl ist auf die Analyse des Bayes Errors gestützt.

**Lookup Dictionary** Der Aufbau des Lookup Table verläuft wie folgt:

1. Daten auf die oben beschriebenen Spalten beschränken
2. Aussonderung der Duplikate
3. `etim_id` wird in einem Set abgelegt und anschliessend im Dictionary mit den Spaltenwerten als Key abgespeichert

**Resultat** Es gibt drei mögliche Resultate für einen Lookup:

1. Einen `etim_id`
2. Mehrere `etim_ids`
3. Key Error

**Eine `etim_id`** Die Kombination aus `product_hierarchie`, `group_id`, `series_name` und `type_name` verweist auf eine einzelne `etim_id` und ist klar zuweisbar.

**Mehrer `etim_ids`** Die Kombination aus `product_hierarchie`, `group_id`, `series_name` und `type_name` auf verweist auf mehrere `etim_id`'s und ist nicht klar zuweisbar.

**Key Error** Die Kombination aus `product_hierarchie`, `group_id`, `series_name` und `type_name` verweist auf keine `etim_id` und ist daher nicht zuweisbar. Die Kombination wird nicht im Datenset des Lookup Table gefunden und ist dem Algorithmus nicht bekannt.

**Nutzen** Die Aufgabe des Lookup Table besteht darin, alle eindeutig zugewiesenen `etim_ids` zurückzugeben und alle anderen Fälle an einen weiteren Algorithmus weiterzureichen.

## 6.4.2 Analyse

**Data Set** Da es sich beim Lookup Table ausschliesslich um Auswendiglernen handelt, wird eine Analyse direkt auf dem `prepared_data_validate` durchgeführt.

**Certainty Rate** Um die Performance des Lookup Table näher zu betrachten, wird die Certainty Rate eingeführt. Sie beschreibt wie sicher der Lookup bei seinen Entscheiden sein muss, bevor ein Resultat zurückgegeben wird.

Die Resultate mit mehreren `etim_ids` wird wie folgt behandelt:

Berechnen des prozentualen Vorkommens der einzelnen `etim_ids`. Falls die häufigste vorkommende `etim_id` prozentual über der Certainty Rate liegt wird diese gewählt, ansonsten bleibt sie nicht zugewiesen.

**Performance des Lookup** Die Analyse des Lookup Table führt zu folgenden übergreifenden Ergebnissen: Hierbei wurde in 0.05 Schritten von 0.5 bis 1.0 über die Certainty Rate iteriert.

Certainty Rate	Richtige Zuordnung	Falsche Zuordnung	Mehrere Werte	Key Error
0.0	0.8360	0.0621	0.0	0.1019
0.5	0.8222	0.0375	0.0384	0.1019
0.55	0.8137	0.0313	0.0531	0.1019
0.60	0.8070	0.0247	0.0664	0.1019
0.65	0.8028	0.0223	0.0730	0.1019
0.70	0.7961	0.0185	0.0835	0.1019
0.75	0.7914	0.0156	0.0910	0.1019
0.80	0.7838	0.0147	0.0996	0.1019
0.85	0.7776	0.0138	0.1067	0.1019
0.90	0.7643	0.0119	0.1219	0.1019
0.95	0.7435	0.0090	0.1456	0.1019
1.00	0.7421	0.0090	0.1470	0.1019

**Tabelle 12: Lookup Performance**

\* Gerundet auf 2 Nachkommastellen

**Richtige Zuordnung**

 Die *Richtige Zuordnungs* Rate bestimmt den Prozentsatz, der eindeutig richtig zugeteilten `etim_id` zu den jeweiligen Kolonnenwerte.

**Falsche Zuordnung** Die *Falsche Zuordnungs* Rate bestimmt den Prozentsatz der eindeutig falsch zugewiesenen `etim_id` auf die jeweiligen wirklichen `etim_id`. Der Wert der `etim_id` aus dem Lookup entspricht also nicht dem eigentlichen Wert.

Dieses Vorkommen kann auf eine potenzielle Falschzuordnung zurückzuführen sein oder aus mehreren Werten wurde zwar die häufigste ausgewählt, jedoch war diese Auswahl falsch.

**Key Error**

 Der *Key Error* entspricht einer Kombination aus Spaltenwerten, die in dieser Form nicht im Trainingset vorhanden war.

**Mehrere Werte**

 Die *Mehreren Werte* beschreiben jene Kombinationen, welche sich nicht eindeutig auf eine `etim_id` zuweisen lassen und die häufigste `etim_id` hat die Certainty Rate nicht erreicht.

**Schlussfolgerung**

 Der Lookup Table erreicht seine obere Grenze bei 0.836 richtigen Zuordnungen, solange ihm totale Freiheiten bei Entscheiden, also die jeweils wahrscheinlichste `etim_id`, gegeben werden. Jedoch ist seine Error Rate in diesem Falle mit 0.0621 zu hoch.

Die Aufgabe des Lookup soll darin bestehen, alle eindeutig zuweisbaren `etim_ids` in einem ersten Durchlauf zu klassifizieren. Den Rest soll an einen Algorithmus weitergegeben werden, der zuverlässigere Entscheidungen treffen kann.

Diese Aufgabe erfüllt der Lookup Table bei einer Certainty Rate von 1.0, wobei folgendes gilt:

<b>Nicht klassifiziert</b>	0.1470 (Mehrere Werte) + 0.1019 (Key Errors) = 0.2489
<b>Klassifiziert (Richtig und Falsch)</b>	1 - 0.2489 = 0.7511
<b>Precision</b>	0.7421 / (0.7421 + 0.0090) = 0.988

Tabelle 13: Lookup Schlussfolgerung

Das Resultat ist eine Accuracy von 98.8% auf denjenigen Daten, die eindeutig zugeordnet werden und 24.89% der Datensätze, die an einen zweiten Algorithmus zur Klassifizierung weitergegeben werden.

### 6.4.3 Erweiterung der Spaltenauswahl

**Ausgangslage** Aus dem obigen Abschnitt geht hervor, dass 14.70 % der Daten nicht eindeutig einer `etim_id` zugeordnet werden können. Auf dem `prepared_data_train` Set sind das 307 Datensätze.

Eine weitere Analyse auf den nicht zugewiesenen Daten ist sinnvoll, um zu verstehen, wie diese Zuordnungen verbessert werden könnten.

Dies spielt vor allem für den Algorithmus eine Rolle, an den diese nicht durch den Lookup zuordbaren Daten weitergegeben werden, sprich dem DNN.

**Vorgehen** Es werden zusätzlich alle Spalten untersucht, die nicht im Key des Lookup Table vorkommen. Daraus lässt sich auf potenziell informative Spalten schliessen, die bei einer Klassifizierung der Daten helfen.

Vorab müssen die Datensätze nach den Lookup Table Keys gruppiert werden. Anschliessend werden alle möglichen Werte für die Kolonnen `path`, `product_name`, `system_id`, `category_id`, `category_name`, `family_id`, `family_name`, `group_name`, `series_id` und `type_id` des jeweiligen Keys aufgelistet und abgespeichert.

Jeder Key muss untersucht werden, um zu überprüfen, ob in einer der Spalten mehr als ein Wert vorhanden ist. Dies wäre ein Indiz dafür, dass wenn diese Spalte zum Entscheidungsprozess hinzugezogen würde, bessere Resultate folgen.

**Resultat** Die untenstehende Tabelle zeigt auf, welche Kolonne wie oft mehrere Werte enthalten.

Kolonne	Anzahl Fälle
<code>path</code>	20

<b>product_name</b>	20
<b>system_id</b>	0
<b>category_id</b>	1
<b>category_name</b>	1
<b>family_id</b>	1
<b>family_name</b>	0
<b>group_name</b>	0
<b>series_id</b>	0
<b>type_id</b>	0

Tabelle 14: Lookup Analyse

Zum Beispiel werde für 20 verschiedene Einträge im Lookup Table mehrere Einträge in der Kolonne `path` gefunden.

**Schlussfolgerung** Durch diese Analyse werden die Kolonnen `path` und `product_name` als zusätzliche potenzielle Informationsquelle identifiziert.

In der Bayes Analyse wurden bereits Informationen über die Spalten gesammelt.

Column	Unique Values	Bayes Error
<b>product_name</b>	8712	0.011586207
<b>path</b>	8784	0.007129687

Tabelle 15: Bayes Error

Da `path` einen kleineren irreduziblen Fehler aufweist, ist es sinnvoll diese Spalte für den Neuralen Netzwerk Algorithmus im nächsten Abschnitt miteinzubinden. Ausserdem besteht der `path` im Wesentlichen aus den Informationen der `system_name`, `category_name` und `family_name` Spalten. Dementsprechend werden Informationen von Kombinationen dieser Spalten ebenso einfließen.

Das Hinzufügen wird sicherlich den Trainings Error des DNN verkleinern, hat jedoch auch den Nachteil, dass es eine starke Erhöhung der Dimensionen zur Folge hat. Das liegt an den 8784 einzigartige Werte in der Spalte `path`.

## 6.5 Neuronales Netzwerk

### 6.5.1 Algorithmus

**Daten** Als Trainingset wird `dnn_prepared_and_formatted_data_train` benutzt.  
 Als Validierungsset wird `dnn_prepared_and_formatted_data_validate` benutzt.  
 Als Testset wird `prepared_data_validate` benutzt.

**Daten Auswahl** Die Basis des Neural Netzwerk bilden folgende Spalten:

```
etim_id, product_hierarchie, group_id, series_name,
type_name und path
```

Die Auswahl ist auf die Analyse des Bayes Errors und des Lookup Table gestützt. Diese Spalten werden aus den `dnn_prepared_and_formatted_data_train` extrahiert und One-Hot codiert.

**Encoder** Als One-Hot-Encoder wird der `OneHotEncoder` von `sklearn` verwendet. Die Einstellung `unknown='ignore'` führt dazu, dass Werte, welche im Trainingsset nicht anzu-treffen sind, auf den Nullvektor abgebildet werden. Dieser `OneHotEncoder` wird als `Xencode.joblib` abgespeichert. Für die `etim_id's` wird der `LabelEncoder` von `sklearn` verwendet. Dieser wird unter `yencode.joblib` gespeichert.

**Netzwerk** Das Netzwerk wird durch manuelles Anpassen der Hyperparameter auf eine gewisse Netzwerk-Dimension eingestellt. Diese Einstellung erlaubt eine möglichst hohe Trainingsgenauigkeit. Das resultierende trainierte Modell wird unter `dnn_model.h5` abgespeichert.

**Output Layer** Das Output Layer ist ein Softmax Layer, welches zu Beginn des Trainings mit der relativen Häufigkeit der `etim_ids` im Trainingsset gewichtet wird, um eine schnellere Konvergenz des Netzwerks zu erreichen.

**Bayes Error** Die Analyse des Bayes Error auf `prepared_and_formatted_data_train` ergibt folgendes Resultat:

Comp. 0	Comp. 1	Comp. 2	Comp. 3	Comp. 4	Bayes Error
product_hier.	group_id	series_name	type_name	path	0.0657494

Der Algorithmus soll folglich eine Trainings-Accuracy von 0.934 anstreben.

**Ausführung**

<b>Epochen</b>	30
<b>Batchsize</b>	32
<b>Optimizer</b>	Adam
<b>Netzwerk</b>	<pre>model = Sequential() model.add(Dense(1900, activation='relu', kernel_initializer='he_normal', input_shape=(n_features,))) model.add(Dense(1300, activation='relu', kernel_initializer='he_normal')) model.add(Dense(900, activation='relu', kernel_initializer='he_normal')) model.add(Dense(450, activation='relu', kernel_initializer='he_normal')) model.add(Dense(250, activation='relu', kernel_initializer='he_normal')) model.add(Dense(148, activation='relu', kernel_initializer='he_normal')) model.add(Dense(180, activation='relu', kernel_initializer='he_normal')) model.add(Dense(y_LabelEncoder.classes_.shape[0], activation='softmax'))</pre>

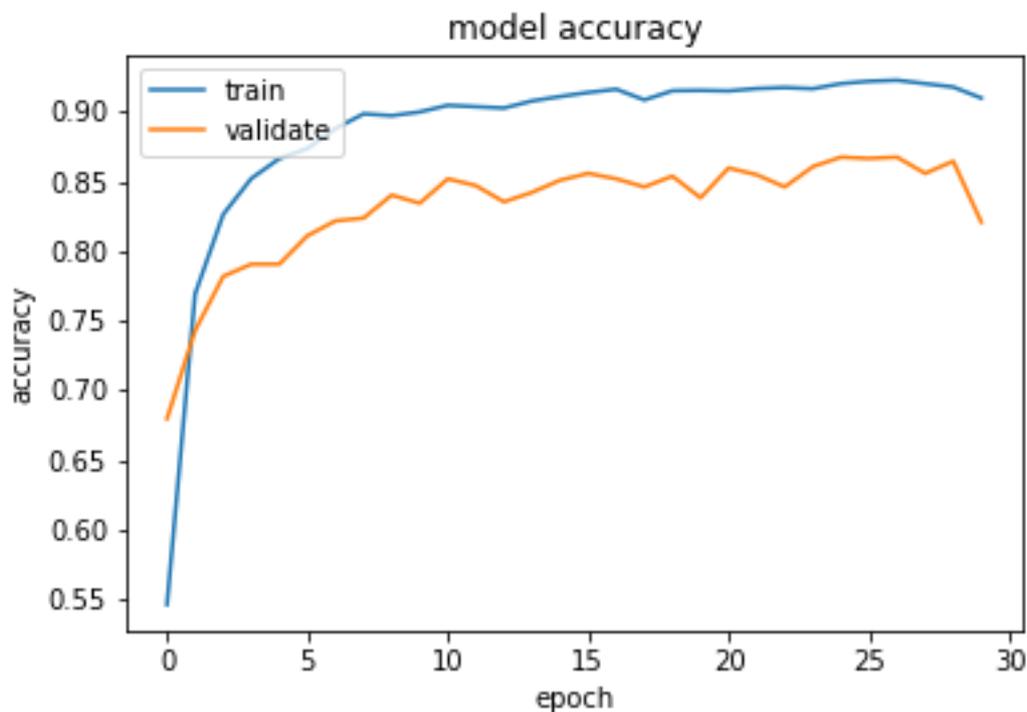


Abbildung 8: DNN Train/Validate Accuracy

Resultate	Training	Validation	Bayes	Bayes/Train Diff.	Top 1
	0.910	0.821	0.934	0.024	0.881

**Trainings & Validierungs Accuracy** - Diese Konfiguration ergibt ein Model mit einer Trainingsaccuracy von 0.910 und einer Validierungsaccuracy von 0.821.

**Bayes Error vs. Trainings Accuracy** Das Ziel des Trainings ist, sich dem Bayes Error von 0.934 so nahe wie möglich anzunähern. Die Differenz ist mit 0.024 jedoch zu hoch und muss verbessert werden.

**Test Accuracy** Auf `prepared_data_validate` hat dieses Model eine Top-1-Accuracy von 0.881.

### 6.5.2 Verbesserungen

**Metaparameter** Um das Netzwerk und den Lernprozess weiter zu verbessern, wurden diverse Metaparameter angepasst:

```
regularisation_parameter = 0.002
dropout_rate = 0.2
batch_size = 256
epochs = 72
```

Es wurde ein Dropout Layer zwischen allen Layern im DNN ausser dem Softmax Layer eingefügt, welches während dem Training alle Nodes des vorherigen Layers mit der Wahrscheinlichkeit `dropout_rate` ignoriert.

Für den `regularisation_parameter` wurde in allen Layern ausser dem Softmax Layer ein L1 Regularisierer mit dem Gewicht `regularisation_parameter` eingefügt.

**Lernprozess** Die Parameter, welche beeinflussen wie schnell und wie genau das DNN die Trainingsdaten lernt, sind die `epochs` und die `batch_size` Parameter. Diese wurden von Hand durch Ausprobieren festgelegt.

**Regularisierung** Die Parameter, welche beeinflussen wie gut das DNN Informationen aus den Trainingsdaten generalisiert, sind die `dropout_rate` und die `regularisation_parameter`.

Der Algorithmus optimiert diese Parameter, indem systematisch verschiedene Kombinationen ausprobiert werden. Ausgewählt werde diejenigen, welche die höchste Validierungsaccuracy aufweisen.

<b>Ausführung</b>	<b>Epochen</b>	72
	<b>Batchsize</b>	256
	<b>Regularisation</b>	0.002
	<b>Dropout</b>	0.2
	<b>Optimizer</b>	Adam
	<b>Netzwerk</b>	<pre> model = Sequential() model.add(Dense(1900, activation='relu', kernel_initializer='he_normal', input_shape=(n_features,) , activity_regularizer=l1(regularisation_parameter))) model.add(Dropout(dropout_rate)) model.add(Dense(1300, activation='relu', kernel_initializer='he_normal', activity_regularizer=l1(regularisation_parameter))) model.add(Dropout(dropout_rate)) model.add(Dense(900, activation='relu', kernel_initializer='he_normal', activity_regularizer=l1(regularisation_parameter))) model.add(Dropout(dropout_rate)) model.add(Dense(450, activation='relu', kernel_initializer='he_normal', activity_regularizer=l1(regularisation_parameter))) model.add(Dropout(dropout_rate)) model.add(Dense(250, activation='relu', kernel_initializer='he_normal', activity_regularizer=l1(regularisation_parameter))) model.add(Dropout(dropout_rate)) model.add(Dense(148, activation='relu', kernel_initializer='he_normal', activity_regularizer=l1(regularisation_parameter))) model.add(Dropout(dropout_rate)) model.add(Dense(180, activation='relu', kernel_initializer='he_normal', activity_regularizer=l1(regularisation_parameter))) model.add(Dense(y_LabelEncoder.classes_.shape[0], activation='softmax'))           </pre>

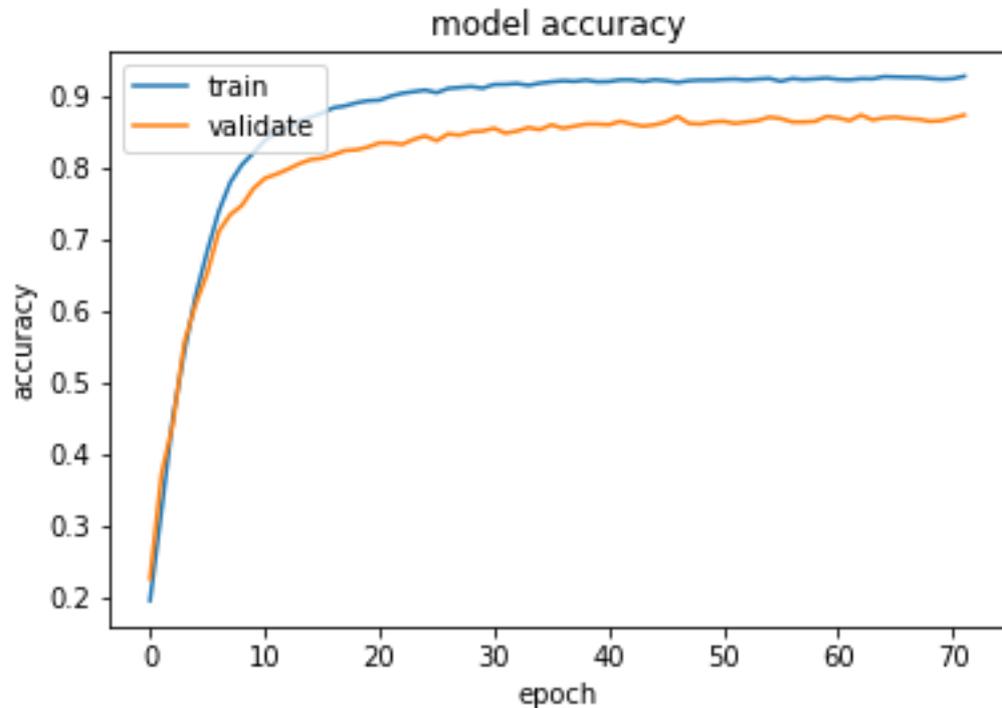


Abbildung 9: DNN Optimierte Train/Validate Accuracy

Wahrscheinlichkeit	Training	Validation	1-Bayes Err.	Bayes/Train Diff.	Top 1
	0.928	0.874	0.934	0.006	0.911

**Trainings & Validierungs Accuracy** Durch die optimierten Metaparametern erreicht das Model eine Trainings-Accuracy von 0.928 und eine Validierungs-Accuracy von 0.874.

**Bayes Error vs. Training Accuracy** Eine Annäherung der Trainings-Accuracy an den Bayes Error von 0.006 ist zufriedenstellend.

**Test Accuracy** Auf `prepared_data_validate` hat dieses Model eine Top-1-Accuracy von 0.911.

**Validierungs vs. Test Accuracy** Es ist auffallend, dass das trainierte Model auf dem Validierungsset bedeutend schlechter abschneidet als auf dem Testset. Das ist auf die unterschiedliche Häufigkeit der `etim_ids` im `dnn_prepared_and_formatted_data_validate` gegenüber dem `prepared_data_validate` zurückzuführen. Durch die Datenaufbereitung, spezifisch durch das Löschen der Duplikate und dem Train-Validate Split, wurden die Anzahl Vorkommen von häufigen ETIM Klassen deren der seltenen angepasst.

## 6.6 Text Embedding mit fastText

### 6.6.1 Algorithmus

**Idee** Der letzte Algorithmus braucht einen anderen Ansatz, um die Dimensionen der Trainingsdaten zu reduzieren. Anstelle des One-Hot-Encodings wird ein Textembedding benutzt. Da viele der Spalten in den ursprünglichen `prepared_data_train` Deutschen Text beinhalten, werden die Textinformationen extrahieren, um ein Textembedding darauf zu trainieren.

**Daten** Als Trainingsset wird `fasttext_prepared_and_formatted_data_train` benutzt. Als Validierungsset wird `fasttext_prepared_and_formatted_data_validate` benutzt. Als Testset wird `prepared_data_validate` benutzt.

**Daten Auswahl** Die Basis des fastText bilden folgende Spalten:

```
product_hierarchie, short_text, product_name, path, system_id, category_id, family_id, family_name, group_id, group_name, series_id, series_name, type_id, type_name und name
```

Die Auswahl wird in folgendes Format umgewandelt:

```
'__label__' ++ etim_id ++ ' ' ++ X
```

(Wobei `X` ein String mit allen textbasierten Informationen aus dem jeweiligen Datensatz ist)

Abgespeichert werden diese Strings in einem `.txt` File und später für das Training und die Validation hinzugezogen.

**Vorgehen** Um das Textembedding zu trainieren, wird die Open Source Library fastText benutzt. Mit der Funktion `train_supervised` wird ein Model trainiert, das für jedes `X` und jede `etim_id` die Wahrscheinlichkeit voraussagt, dass `X` als `etim_id` klassifiziert wird. Das Model wird mit Learning Rate `lr=1.0` und der Anzahl Epochen `epoch=25` erstellt und anschliessend unter `fasttext_model.bin` abgespeichert.

**Bayes Error** Eine Analyse anhand des Bayes Error macht im Fall des fastText Algorithmus keinen Sinn. Dies ist auf die Spaltenauswahl zurückzuführen. Da die Spalte `name` in der Auswahl enthalten ist erreicht der Bayes Error einen Wert von 0.0, da es sich in dieser Spalte ausschliesslich einzigartige Werte finden lassen.

Zusammenfassung	Training	Validation	Top 1
	1.0	0.966	0.956

**Trainings & Validierungs Accuracy** fastText erreicht ohne Optimierung eine Trainingsaccuracy von 1.000 und eine Validierungsaccuracy von 0.966

**Test Accuracy** Das fastText Model erreicht eine Top-1-Accuracy auf `prepared_data_validate` von 0.967.

## 6.6.2 Verbesserungen

**Metaparameter** FastText bietet die Funktion an, optimale Metaparameter zu bestimmen. Hierfür wird `fasttext_prepared_and_formatted_data_validate` verwendet. Dafür werden beim Trainieren und anschliessend beim Validieren gewisse Metriken wie der F1 Score automatisch abgeglichen, um die optimale Parameter Einstellung zu finden.

**Optimierte Parameter** Nach einer Optimierungszeit von fünf Stunden sehen die Metaparameter wie folgt aus:

```
lr=5.0
dim=167
ws=5
epoch=100
minCount=1
minn=3
maxn=6
neg=5
wordNgrams=1
loss='softmax'
bucket=72689
thread=3
lrUpdateRate=100
t=0.0001
label='__label__'
```

Zusammenfassung	Training	Validation	Top 1
	1.0	0.974	0.963

**Trainings & Validierungs Accuracy** Die Trainingsaccuracy ist 1.000. Dies ist jedoch nicht verwunderlich, da unter anderem die Spalte `name` benutzt wird, welche für jeden Artikel einzigartig ist. Die Spannende Metrik ist die Validierungsaccuracy von 0.974, welche misst, wie gut das Model generalisiert.

**Test Accuracy** Mit optimierten Metaparametern erreicht fastText auf `prepared_data_validate` eine Top-1-Accuracy von 0.9634.

## 6.7 Kombinationen

### 6.7.1 Top 3

**Problem** Bei jedem Top-1 Ergebnis entscheidet jeweils der Algorithmus, welche `etim_id` ausgewählt wird. Aus der Problemstellung geht jedoch hervor, dass dies bisher durch manuelle Bearbeitung gemacht wurde.

**Idee und Vorgehen** Um die Genauigkeit der Klassifizierung zu erhöhen, wird der Top-3-Algorithmus eingeführt. Es werden keine Änderung an den Model vorgenommen. Der Top-3-Algorithmus gibt nicht nur die wahrscheinlichste ETIM Klasse zurück, sondern erweitert diese Rückgabe auf die wahrscheinlichsten drei ETIM Klassen. Dadurch entsteht eine Kombination aus maschineller Vorarbeit und der endgültigen Entscheidung durch einen Sachbearbeiter. Dadurch kann die theoretische Genauigkeit des Algorithmus verbessert werden.

**Daten** Die Datengrundlage bildet das in den jeweiligen Algorithmen beschriebene Schema. Es werden keine weiteren Trainingsdurchläufe unternommen, weshalb nur die Top-1-Accuracy und die Top-3-Accuracy präsentiert werden.

Als Validierungsset wurde auf das `prepared_data_validate` Set zurückgegriffen, wie bereits in der Validierung des DNNs und des Lookup Table.

Zusammenfassung	DNN_Top 1	DNN_Top 3	fastText_Top 1	fastText_Top 3
	0.921	0.951	0.963	0.985

**Test Accuracy** Das DNN erhöht seine Accuracy gegenüber dem Top-1 Resultat auf 0.951. Der fastText Algorithmus erhöht seine Accuracy gegenüber dem Top-1 Resultat auf 0.985.

**Zusätzliches** Zu beachten ist, dass hier die Accuracy den Fall beschreibt, dass unter den Top-3 `etim_id` die richtige dabei ist. Das bedeutet, dass ein Benutzer die falsche Entscheidung treffen kann und somit die wirkliche Accuracy verschlechtert. Jedoch bietet diese Idee die Möglichkeit unterschiedliche Entscheidungen vom User und vom Algorithmus miteinander zu vergleichen. Diese Erkenntnisse können genutzt werden, um die Algorithmen in Zukunft zu verbessern und sich der Auswahl des Benutzers anzupassen.

## 6.7.2 Lookup und DNN

**Idee und Vorgehen** Der Lookup Table hat eine sehr hohe Präzision auf den Inputkeys, welche jeweils nur eine ETIM Klasse zurückgeben. Die vorhergesagten `etim_id` stimmen in 0.988 der Fälle.

Der Lookup Table wird genutzt, um alle eindeutig zuweisbaren `etim_id` zu klassifizieren. Die nicht zuweisbaren Datensätze werden dem DNN übergeben.

**Daten** Die Datengrundlage bildet das in den jeweiligen Algorithmen beschriebene Schema. Es werden keine weiteren Trainingsdurchläufe unternommen, weshalb nur die Top-1-Accuracy und die Top-3-Accuracy präsentiert werden.

Als Validierungsset wurde auf das `prepared_data_validate` Set zurückgegriffen, wie bereits in der Validierung des DNNs und des Lookup Table.

Zusammenfassugn	Top 1	Top 3
	0.921	0.953

**Test Accuracy** Die Kombination des Lookup Table und des DNN erreicht eine Top-1-Accuracy von 0.921 und eine Top-3-Accuracy von 0.953.

## 7 Vergleich der Resultate

### 7.1 Analyse der Resultate

**Datenset** Die Analyse der Ergebnisse wurde auf dem `prepared_data_validate` durchgeführt. Dieses Set bildet das zu den Trainingsdaten komplementäre Validierungsset.

Metriken	Algorithmus	Accuracy
	<code>dnn_top_1</code>	0.911332385
	<code>lookup_dnn_top_1</code>	0.921289711
	<code>dnn_top_3</code>	0.951161688
	<code>lookup_dnn_top_3</code>	0.953532479
	<code>fastText_top_1</code>	0.963489806
	<code>fastText_top_3</code>	0.985775249

Tabelle 16: Algorithmen Vergleich Accuracy

Die obenstehende Tabelle weist den fastText Top-3-Algorithmus als genausten Algorithmus aus.

Es ist möglich einem Benutzer eine Auswahl von drei `etim_ids` zu präsentieren, in der die richtige Klassifizierung mit einer Wahrscheinlichkeit von 98.58% vorhanden ist.

**Ursachen** Dafür, dass fastText so viel besser abschneidet, als das DNN gibt es zwei Erklärungen:

1. Das DNN trainiert auf weniger Spalten als fastText es tut. Aufgrund der begrenzten Rechenkapazität, welche zur Verfügung stand, musste die Spaltenanzahl für das DNN auf ein Minimum reduziert werden.
2. Das Textembedding kann die Nähe zweier verschiedener Werte abbilden, während beim One-Hot-Encoding alle verschiedenen Werte äquidistant sind. Dies resultiert darin, dass das DNN Werte, welche beim Training nicht aufgetreten sind, auf `Null` abgebildet. Das Textembedding merkt hingegen, wenn zum Beispiel zwei `short_texts` bis auf ein Wort identisch sind.

**Metriken auf Klassen Stufe** Eine genauere Übersicht der Metriken sind im Excel Arbeit Sheet `Metriken.xlsx` enthalten. Dabei sind noch weitere Informationen bezüglich der Accuracy, des Recall und der Precision auf individuellen `etim_ids` vorhanden. Die Angaben geben Ausschluss darüber, welche Klassen im Einzelfall eine gute Zuteilung erfahren und bei welchen die Algorithmen Mühe haben.

etim_id	Lookup Accuracy	DNN Accuracy	fastText Accuracy
<code>EC010890</code>	0.0	0.0	0.0
<code>EC000406</code>	lookup undecided	0.0	0.0
<code>EC000474</code>	lookup undecided	0.0	0.0
<code>EC010042</code>	lookup undecided	0.0	0.0

<b>EC010081</b>	lookup undecided	0.0	0.0
<b>EC010170</b>	lookup undecided	0.0	0.0
<b>EC010394</b>	lookup undecided	0.0	0.0
<b>EC010786</b>	lookup undecided	0.0	0.0
<b>EC010827</b>	lookup undecided	0	0
<b>EC010391</b>	lookup undecided	0.5	0

Tabelle 17: Schlechteste etim\_id Accuracy

Die obenstehende Tabelle zeigt zum Beispiel die zehn `etim_ids`, die für die Algorithmen nicht zuordbar sind. Dabei handelt es sich im Normalfall um `etim_ids`, die ein sehr geringes Vorkommen aufweisen.

Nichtsdestotrotz wären diese Informationen bei einer definitiven Umsetzung von Seitens Geberit von Nutzen, da schwierig zu klassifizierende `etim_ids` als Spezialfall behandelt werden könnten.

## 7.2 Rücksprache mit Geberit

- Top-1 vs Top-3** Bei der Besprechung der erhobenen Metriken möchte Geberit sich nicht auf vollautomatisierte Klassifizierung verlassen, sondern immer auf noch einen Menschen in die Entscheidung einbeziehen. Aus diesem Grund macht es wenig Sinn, einen Algorithmus zu produzieren, welche nur eine Klasse ausgibt. Eine sinnvollere Lösung wäre ein Tool, welches die drei Wahrscheinlichsten Klassen ausgibt, von welchen dann ein Mensch die richtige auswählen kann.
- fastText** Von den Algorithmen, welche drei ETIM Klassen herausgeben ist fastText Top-3 mit 0.985 Accuracy klar der Beste.

---

## 8 Lösung

---

### 8.1 Klassifizierung

---

- Ziel** Das Ziel der Arbeit ist einen Algorithmus zu finden, der Geberit bei der Klassifizierung von `etim_ids` unterstützt und den grösstmöglichen Mehrwert bietet.
- Top-3-Algorithmus** Aus der Rücksprache mit Geberit ging hervor, dass vor allem ein Interesse am Top-3 Algorithmus besteht. `fastText` stellt für Geberit die eindeutige Wahl dar. `FastText` liefert in der Top-1 Version bereits die besten Resultate mit einer Genauigkeit von 96.35% und bietet in der Top-3 Version neben potenziell verbesserten Resultaten einen noch grösseren Mehrwert. Es erlaubt die Kombination von maschineller Vorarbeit und einem finalen menschlichen Entscheid mit einer obere Genauigkeitsgrenze von 98.57%.
- Empfehlung** Die Empfehlung der Autoren ist der `fastText` Top-3-Algorithmus.

### 8.2 Auswertung

---

- Lernen** Alle Hyperparameter wurden in den vorherigen Kapiteln trainiert und sind somit für die jeweiligen Sub Sets optimiert. Eine Überprüfung der momentanen Modelle der Algorithmen macht noch keinen Sinn, da es eine signifikante Anzahl von `etim_ids` im `final_test_set` gibt, auf denen durch die mehrfachen Splits niemals trainiert wurde. Deshalb werden alle Modelle mit denselben Parametern nochmal trainiert, aber auf allen Daten, welche nicht im `final_test_set` sind.
- Datenset** Die finale Auswertung wird auf dem unberührten `final_test_set` durchgeführt.
- Die neuen Modelle der jeweiligen Algorithmen werden mit den durchs Training optimierten Hyperparameter auf allem ausser `final_test_set` trainiert.
- Für den Lookup und den `fastText` Algorithmus wurde dafür das `final_prepared_and_formatted_data_train` Set zusammengestellt und für das DNN das `final_prepared_data_train` Set. Diese beiden Sets tragen dieselben Informationen und unterscheiden sich nur in der Daten Aufbereitung für die jeweiligen Algorithmen.
- EC000000** Da die Modelle bewusst 'EC000000' nicht ausgeben, sind diese Datensätze alle mit Sicherheit falsch. Die Anzahl Datensätze im `final_test_set` ist 250 und bilden somit einen Anteil von 9.679%
- Unbekannte etim\_ids** Abgesehen von 'EC000000' gibt es auch andere `etim_ids`, welche durch die zufällige Partition im `final_test_set` sind, aber beim Trainieren nie angetroffen werden. Dies sind 6 oder 0.232%, welche ebenfalls automatisch falsch klassifiziert werden.
- Metriken** Zur finalen Bewertung der Algorithmen wird jeweils die Top-1-Accuracy und die Top-3-Accuracy folgender Daten erhoben.

1. Das `final_test_set`: Dieses Set gibt pessimistische Metriken zurück, welche definitiv nicht zu Gunsten des Algorithmus manipuliert wurde.
2. Das `final_test_set` ohne `'EC000000'`: Diese Metriken und insbesondere die Top-3-Accuracy stellen der Meinung der Autoren nach das genaueste Analog dazu dar, wie gut der Algorithmus, in der Praxis, auf noch nicht klassifizierten Daten, abschneidet.
3. Das `final_test_set` ohne `'EC000000'` und unbekannte `etim_ids`: Dieses Set gibt einem die höchstmöglichen Metriken zurück für ungesehene Daten.

## Resultate

Der fastText Algorithmus gibt folgende Metriken:

Resultat	Top-1	Top-3
<code>final_test_set</code>	0.868	0.887
<code>final_test_set</code> ohne <code>'EC000000'</code>	0.961	0.982
<code>final_test_set</code> ohne <code>'EC000000'</code> und unbekannte <code>etim_ids</code>	0.963	0.985

Das DNN gibt folgende Metriken:

Resultat	Top-1	Top-3
<code>final_test_set</code>	0.835	0.872
<code>final_test_set</code> ohne <code>'EC000000'</code>	0.925	0.966
<code>final_test_set</code> ohne <code>'EC000000'</code> und unbekannte <code>etim_ids</code>	0.927	0.968

Der Kombinationsalgorithmus aus Lookup Table und DNN gibt folgende Metriken:

Resultat	Top-1	Top-3
<code>final_test_set</code>	0.839	0.872
<code>final_test_set</code> ohne <code>'EC000000'</code>	0.928	0.966
<code>final_test_set</code> ohne <code>'EC000000'</code> und unbekannte <code>etim_ids</code>	0.931	0.968

## 8.3 GUI

### 8.3.1 Funktionalität

- Allgemein** Das GUI wurde mit Hilfe des `tkinter` Python Paket erstellt. Es bietet die Möglichkeit die trainierten Modelle der Algorithmen in ihren Variationen auf geladenen Daten anzuwenden. Der angezeigte Datensatz widerspiegelt stets den resultierenden Datensatz nach einer ausgewählten Aktion.
- Laden** Das Tool bietet die Möglichkeit einen Datensatz als `csv` Datei direkt einzulesen und zu betrachten. Diese Funktionalität wird durch das `pandastable` Paket zur Verfügung gestellt.
- Überprüfung** Durch die Auswahl «check Classification» wird eine Überprüfung der geladenen Datensätze initiiert. `fastText` Top-1 übernimmt diese Aufgabe als Klassifizierungsalgorithmus. Der resultierende Datensatz besteht aus den Daten, bei denen die Klassifizierung eine andere `etim_id` erwartet hätte.

Die Idee hierbei ist, dass somit schnell und übersichtlich potenziell falsch klassifizierte Daten identifiziert werden können.

- Lookup** «Lookup» klassifiziert die Daten mit Hilfe des Lookups Algorithmus.
- DNN** «DNN» klassifiziert die Daten mit Hilfe des DNN Top-1-Algorithmus.
- fastText** «fastText» klassifiziert die Daten mit Hilfe des fastText Top-1-Algorithmus.
- Lookup\_DNN** «Lookup\_DNN» klassifiziert die Daten mit Hilfe des Lookup\_DNN Kombinations Algorithmus.
- DNN Top 3** «DNN Top 3» klassifiziert die Daten mit Hilfe des DNN Top-3-Algorithmus. Hierbei wird jeweils ein zusätzliches Fenster geöffnet, um die Benutzerauswahl entgegenzunehmen.
- fastText Top 3** «fastText Top 3» klassifiziert die Daten mit Hilfe fastText Top-3-Algorithmus. Hierbei wird jeweils ein zusätzliches Fenster geöffnet, um die Benutzerauswahl entgegenzunehmen
- Speichern** Nachdem die jeweilige Funktion durchgeführt wurde, gibt es die Möglichkeit den resultierenden Datensatz wieder abzuspeichern. «Save as...» speichert den sichtbaren Datensatz an beliebiger Stelle ab.
- Kontrolle** Der Prototyp erlaubt minimale Kontrolle über die Daten. Top-1-Algorithmen befüllen den Datensatz nach ihrem besten Wissen und Gewissen. Top-3-Algorithmus erlauben etwas mehr Interaktion. Dabei kann beispielsweise die jeweilige Reihe ausgewählt werden, auf der eine Klassifizierung gewünscht wird. Dies ist vor allem in Kombination mit «check Classification» nützlich, da man dabei gleich Vorschläge für potenziell falsch klassifizierte Daten erhalten und diese korrigieren kann.
- Log** Ein weiteres nützliches Feature ist die Log Funktionalität. Falls bei einem Top-3-Algorithmus der Benutzer nicht diejenige Auswahl trifft, die der Algorithmus als die Wahrscheinlichste ansieht, wird der gesamte Datensatz inklusive der durch den Benutzer ausgewählten und der wahrscheinlichsten `etim_id` in einem Log abgespeichert. Dies erlaubt eine Analyse, in welchen Daten sich der Algorithmus und ein Benutzer regelmäßig widersprechen und erlaubt damit mögliche Veränderungsansätze und kostbare Erkenntnisse.

## Showcase

ETIM Classification									
Open	check Classification	Lookup	Fattest	DNN Top 3	Welcome to the ETIM Classification Tool.				
Save as...	DNN	lookup_dnn	Fattest Top 3						
id	article_id	name	article_number	etim_description	etim_id	product_hierarchie	short_text	product_name	
1	ART_106295	21506 - Kreuzstück CSI 90G d15-22 abg	21506	Filling mit 4 Anschlüsse	EC000326	21.775.7723	Mapress Filings C-Stahl v	Kreuzstück CSI 90G d15-22 abg verz.	Gebert Mapress C-Stahl Kreuzstück
2	ART_106016	367.750.16.1 - GewStu.m/VerschKappe	367.750.16.1	Filling mit 2 Anschlüsse	EC000324	21.601.6058	PE-HD Formstücke d110	GewStu.m/VerschKappe PE-HD d110 H7	Gebert PE-Gewindestutzen mit Versch
3	ART_1524170	2691000 - Slim. Modo RF Set PL	2691000		EC000000	11.860.8701	Sets Installationselement	Slim. Modo RF Set PL	Kolo Set Slim Element für Wand/WC
4	ART_1847296	131.021.00.5 - MLith WC W-h Al-g Wau A	131.021.00.5	Aufbauspülkasten	EC010864	12.208.2082	Gebert Monolith Sanitär	MLith WC W-h Al-g Wau A	Gebert Monolith Sanitärmodul für War
5	ART_108283	371.745.00.1 - Losflansch Stahl d250 Ks	371.745.00.1	Feste Flansche	EC010338	21.802.8125	PE-HD Verbindungen Zub	Losflansch Stahl d250 Kbsch	Gebert PE Losflansch
6	ART_1825268	502.305.01.3 - US-F-1-WT Ge+Ch B90 2	502.305.01.3	Waschtischunterschrank	EC011382	31.820.8201	Unterschranke	US-F-1-WT Ge+Ch B90 2-Schbl v matt	Gebert (Con) Unterschrank für Wascht
7	ART_1796597	501.915.01.8 - US-1-MWT mit MWT Ge+Re	501.915.01.8	Waschtischunterschrank	EC011382	31.820.8202	Sets Waschtisch + Badez	US-1-MWT mit MWT Ge+RenP 960 Schbl	Gebert (Con) Unterschrank für Dopp
8	ART_101240	147.011.00.1 - Ausgleichspuffer erhöht B	147.011.00.1	Anschlagpuffer für WC-Sitz	EC012204	35.426.4296	Balena 4000 Ersatzteile	Ausgleichspuffer erhöht Balena 4000	Gebert Ausgleichspuffer erhöht für WC-Sitz
9	ART_352491	840237000 - S-Elm Ge+Ch B37H40 m	840237000	Badmöbelschrank	EC010013	31.820.8213	Seitenschränke	S-Elm Ge+Ch B37H40 m/SBoW v	Gebert (Con) Seitenelement mit Staub
10	ART_342763	20771320001 - SHOWERAMA 8-1 MON	20771320001		EC000000	32.860.8608	Duschabtrennungen Ersatz	SHOWERAMA 8-1 MONTERINGSSATS	Gebert Duschabtrennung Ersatz
11	ART_105384	361.559.16.1 - Reduktion PE-HD d50/40 z	361.559.16.1	Filling mit 2 Anschlüsse	EC000324	21.601.6053	PE-HD Formstücke d50m	Reduktion PE-HD d50/40 zentrisch	Gebert PE Reduktion zentrisch, kurz
12	ART_1424221	500.896.00.1 - SB BTW Ge+Fant 155 m	500.896.00.1	Bidelt	EC010126	31.920.9201	Ständindets wandbüding	SB BTW Ge+Fant 155 m/Ul.F. Bef.verd	Gebert Fantasia Ständindet wandbüding
13	ART_1825295	502.314.01.3 - US-1-1-DWT Ge+Ch B120	502.314.01.3	Waschtischunterschrank	EC011382	31.820.8201	Unterschranke	US-1-1-DWT Ge+Ch B120 2-Schbl v matt	Gebert (Con) Unterschrank für Dopp
14	ART_432977	500.627.01.3 - DWT PG-AcN B120 mULF	500.627.01.3	Waschtisch	EC011550	31.910.9103	Waschtische	DWT PG-AcN B120 mULF mHLo	Pozzi-Ginori Acanto Doppelwaschtisc
15	ART_111651	91045 - FlanschDicht Cent-HD3822 DN4	91045	Gummi Flanschdichtung (kein Norm)	EC010674	22.755.7555	Mapress Filings Zubehör	FlanschDicht Cent-HD3822 DN40 PN6	Gebert Flanschdichtung PN 6
16	ART_101280	147.248.EP.1 - Bedienfeld p-m B4000	147.248.EP.1	Zubehör für Dusch-WC	EC010204	35.426.4296	Balena 4000 Ersatzteile	Bedienfeld p-m B4000	Gebert Bedienfeld für Gebert Balena 4000
17	ART_112588	632.006.00.1 - Übergang MF IG Rg G12	632.006.00.1	Filling mit 2 Anschlüsse	EC000324	22.705.7155	Mapia Filings Zubehör	Übergang MF IG Rg G12/MF-1/2	Gebert Übergang mit Ausseingewinde
18	ART_103895	242.305.00.1 - Distanzscheiben weiss A	242.305.00.1	Zubehör für Dusch-WC	EC010204	35.426.4352	AquaClean 800/Opas Ersatz	Distanzscheiben weiss Acq/38/8000 - Set	Gebert Distanzscheiben für Gebert Aqua
19	ART_554222	3696101201 - SWC-T mAP-SPK a mWC-3696	3696101201	WC-Kombination	EC011318	31.900.9007	Sets Stand-WC + AP-SPK	SWC-T mAP-SPK a mWC-sz Id-Gw AB	IDO Glow Stand-WC Abgang multire
20	ART_345538	575100000 - JOOPI WC-Sitz mit Absenka	575100000	WC-Sitz	EC011196	31.955.9565	WC-Sitze aus Duroplast	JOOPI WC-Sitz mit Absenkautomatik	Gebert JOOPI WC-Sitz
21	ART_374848	243.661.00.1 - Befestigungsrahmen zu Pn	243.661.00.1	Zubehör Vorwand-Einbauelement Sanitär	EC011338	33.342.3423	Ersatzteile WC-Spualaufen	Befestigungsrahmen zu Pneumatikdrücker	Gebert Befestigungsrahmen für pneu
22	ART_1139325	595976000 - SchbIFr Ge+RenP sRf WT-U	595976000		EC000000	31.820.8210	Mobel Ersatzteile Keramik	SchbIFr Ge+RenP sRf WT-U B130 Ech	Gebert Set Schublendenfronten für Unt
23	ART_1774966	500.929.00.5 - Ft-1-Sw-Id-Design 90x200	500.929.00.5	Duschtür	EC010962	32.860.8601	Duschabtrennungen komple	Ft-1-Sw-Id-Design 90x200 s s-m G4r	IDO Design Pendeltür
24	ART_353164	597206000 - VerStl Ge+WbFlf WWC-U	597206000	Zubehör für Dusch-WC	EC010204	31.955.9564	Ersatzteile Keramik Integ	VerStl Ge+WbFlf WWC-U/WBL L 20.3c	Gebert Verankerungsstück für Wand
25	ART_436721	243.724.00.1 - Steuerung für AcQ Sela a	243.724.00.1	Zubehör für Dusch-WC	EC010204	35.426.4357	AquaClean Sela Ersatzteile	Steuerung für AcQ Sela ab 2017	Gebert Steuerung für Gebert AquaClean Sela
26	ART_108129	19455 - SysRoR MAP CSI d28x1 5 ia-ver	19455	Dünnwandiges Stahlrohr	EC010032	22.770.7704	Mapress Systemrotre C-Si	SysRoR MAP CSI d28x1 5 ia-ver	Gebert Mapress C-Stahl Systemrotre I
27	ART_108452	22965 - AnKreuzStk VLRL 2x CSI d15-12	22965	Passierstück	EC011717	22.775.7721	Mapress Filings C-Stahl v	AnKreuzStk VLRL 2x CSI d15-12 Kz	Gebert Mapress C-Stahl Anschlus-K
28	ART_106419	390.116.14.1 - Doppelsteckmuffe PP-MD	390.116.14.1	Filling mit 2 Anschlüsse	EC000324	21.622.6281	PP-MD Verbindungen d4	Doppelsteckmuffe PP-MD d40	Gebert Silent-PP Doppelsteckmuffe
29	ART_109886	51252 - Bogen mEEnd CSI 45G d88.9 F	51252	Filling mit 2 Anschlüsse	EC000324	22.775.7739	Mapress Filings C-Stahl v	Bogen mEEnd CSI 45G d88.9 FKM	Gebert Mapress C-Stahl Bogen mit E
30	ART_1424504	58203650000 - WWC-Top Ge+3Co W3S820	58203650000	WC	EC011289	21.905.9052	Wand-WCs	WWC-Top Ge+3Co W3S 8 T88.5	Gebert 300 Comfort Wand-WC Tiefsp
31	ART_1564903	501.598.00.1 - Sockel-EVWT Ge+SncCp	501.598.00.1	Zubehör für Badezimmerzubehör	EC012353	31.820.8206	Badzimmermöbel Zubehör	Sockel-EVWT Ge+SncCp B45 2-Lgr	Gebert Selnova Compact Ecksockel
32	ART_264141	116.053.00.1 - Reduktion VOL Ms d32-1	116.053.00.1	Filling mit 2 Anschlüsse	EC000324	22.708.7906	Volex Filings Ms 32 mm	Reduktion VOL Ms d32-16	Gebert Volex Reduktion
33	ART_351609	650580000 - ReBW Duo Ge+mD 1.180x	650580000	Badewanne	EC011609	32.840.8401	Badewannen aus Acryl	ReBW Duo Ge+mD 1.180x B80	Gebert DayRechtbadewanne E
34	ART_1316927	244.286.00.1 - Fixed Panel 150	244.286.00.1	Zubehör und Ersatzteil für Duschabtrenn	EC010068	32.860.8608	Duschabtrennungen Ersatz	Fixed Panel 150	Gebert Fixelement für Dusche GEO
35	ART_1417885	CG2161000 - WT-WFA Ge+RenB60 m	CG2161000		EC000000	31.910.9110	Sets Waschtisch + WT-A +	WT-WFA Ge+RenB60 mULF mHLo	Gebert Renovita Ceramid Renovita Set I
36	ART_104057	242.516.00.1 - Anschluss-Slitzen d50 P	242.516.00.1	Spülbogen	EC010411	11.155.1453	Ersatzteile Duplex	Anschluss-Slitzen d50 PVC Bogen	Gebert Anschluss-Slitzen d50
37	ART_338476	272411000 - HWB Ge+ReCp B40 uoULF	272411000	Waschtisch	EC011550	31.910.9103	Waschtische	HWB Ge+ReCp B40 uoULF oHLo vKa	Gebert Renovita Compact Handwasch
38	ART_1737555	501.828.00.1 - US-F-1-MWT Po+Ele B60 2	501.828.00.1	Waschtischunterschrank	EC011382	31.820.8201	Unterschranke	US-F-1-MWT Po+Ele B60 2.AZ vKa gr	Porsgrund Elegat Unterschrank für M
39	ART_1757758	550.929.00.1 - RedW Ge+Olo 170x90 m	550.929.00.1	Duschwanne	EC011443	32.875.8752	Duschwannen au Mineralred	RedW Ge+Olo 170x90 m Dicht	Gebert Redteckduschwanne Olona

Abbildung 10: GUI Vorschau

ETIM Classification									
Open	check Classification	Lookup	Fattest	DNN Top 3	Runs Fattest Top 3 Algorithm for classification with user selection				
Save as...	DNN	lookup_dnn	Fattest Top 3						
id	article_id	name	article_number	etim_description	etim_id	product_hierarchie	short_text	product_name	
1	ART_106295	21506 - Kreuzstück CSI 90G d15-22 abg	21506	Filling mit 4 Anschlüsse	EC000326	21.775.7723	Mapress Filings C-Stahl v	Kreuzstück CSI 90G d15-22 abg verz.	Gebert Mapress C-Stahl Kreuzstück
2	ART_106016	367.750.16.1 - GewStu.m/VerschKappe	367.750.16.1	Filling mit 2 Anschlüsse	EC000324	21.601.6058	PE-HD Formstücke d110	GewStu.m/VerschKappe PE-HD d110 H7	Gebert PE-Gewindestutzen mit Versch
3	ART_1524170	2691000 - Slim. Modo RF Set PL	2691000		EC000000	11.860.8701	Sets Installationselement	Slim. Modo RF Set PL	Kolo Set Slim Element für Wand/WC
4	ART_1847296	131.021.00.5 - MLith WC W-h Al-g Wau A	131.021.00.5	Aufbauspülkasten	EC010864	12.208.2082	Gebert Monolith Sanitär	MLith WC W-h Al-g Wau A	Gebert Monolith Sanitärmodul für War
5	ART_108283	371.745.00.1 - Losflansch Stahl d250 Ks	371.745.00.1	Feste Flansche	EC010338	21.802.8125	PE-HD Verbindungen Zub	Losflansch Stahl d250 Kbsch	Gebert PE Losflansch
6	ART_1825268	502.305.01.3 - US-F-1-WT Ge+Ch B90 2	502.305.01.3	Waschtischunterschrank	EC011382	31.820.8201	Unterschranke	US-F-1-WT Ge+Ch B90 2-Schbl v matt	Gebert (Con) Unterschrank für Wascht
7	ART_1796597	501.915.01.8 - US-1-MWT mit MWT Ge+Re	501.915.01.8	Waschtischunterschrank	EC011382	31.820.8202	Sets Waschtisch + Badez	US-1-MWT mit MWT Ge+RenP 960 Schbl	Gebert (Con) Unterschrank für Dopp
8	ART_101240	147.011.00.1 - Ausgleichspuffer erhöht B	147.011.00.1	Anschlagpuffer für WC-Sitz	EC012204	35.426.4296	Balena 4000 Ersatzteile	Ausgleichspuffer erhöht Balena 4000	Gebert Ausgleichspuffer erhöht für WC-Sitz
9	ART_352491	840237000 - S-Elm Ge+Ch B37H40 m	840237000	Badmöbelschrank	EC010013	31.820.8213	Seitenschränke	S-Elm Ge+Ch B37H40 m/SBoW v	Gebert (Con) Seitenelement mit Staub
10	ART_342763	20771320001 - SHOWERAMA 8-1 MON	20771320001		EC000000	32.860.8608	Duschabtrennungen Ersatz	SHOWERAMA 8-1 MONTERINGSSATS	Gebert Duschabtrennung Ersatz
11	ART_105384	361.559.16.1 - Reduktion PE-HD d50/40 z	361.559.16.1	Filling mit 2 Anschlüsse	EC000324	21.601.6053	PE-HD Formstücke d50m	Reduktion PE-HD d50/40 zentrisch	Gebert PE Reduktion zentrisch, kurz
12	ART_1424221	500.896.00.1 - SB BTW Ge+Fant 155 m	500.896.00.1	Bidelt	EC010126	31.920.9201	Ständindets wandbüding	SB BTW Ge+Fant 155 m/Ul.F. Bef.verd	Gebert Fantasia Ständindet wandbüding
13	ART_1825295	502.314.01.3 - US-1-1-DWT Ge+Ch B120	502.314.01.3	Waschtischunterschrank	EC011382	31.820.8201	Unterschranke	US-1-1-DWT Ge+Ch B120 2-Schbl v matt	Gebert (Con) Unterschrank für Dopp
14	ART_432977	500.627.01.3 - DWT PG-AcN B120 mULF	500.627.01.3	Waschtisch	EC011550	31.910.9103	Waschtische	DWT PG-AcN B120 mULF mHLo	Pozzi-Ginori Acanto Doppelwaschtisc
15	ART_111651	91045 - FlanschDicht Cent-HD3822 DN4	91045	Gummi Flanschdichtung (kein Norm)	EC010674	22.755.7555	Mapress Filings Zubehör	FlanschDicht Cent-HD3822 DN40 PN6	Gebert Flanschdichtung PN 6
16	ART_101280	147.248.EP.1 - Bedienfeld p-m B4000	147.248.EP.1	Zubehör für Dusch-WC	EC010204	35.426.4296	Balena 4000 Ersatzteile	Bedienfeld p-m B4000	Gebert Bedienfeld für Gebert Balena 4000
17	ART_112588	632.006.00.1 - Übergang MF IG Rg G12	632.006.00.1	Filling mit 2 Anschlüsse	EC000324	22.705.7155	Mapia Filings Zubehör	Übergang MF IG Rg G12/MF-1/2	Gebert Übergang mit Ausseingewinde
18	ART_103895	242.305.00.1 - Distanzscheiben weiss A	242.305.00.1	Zubehör für Dusch-WC	EC010204	35.426.4352	AquaClean 800/Opas Ersatz	Distanzscheiben weiss Acq/38/8000 - Set	Gebert Distanzscheiben für Gebert Aqua
19	ART_554222	3696101201 - SWC-T mAP-SPK a mWC-3696	3696101201	WC-Kombination	EC011318	31.900.9007	Sets Stand-WC + AP-SPK	SWC-T mAP-SPK a mWC-sz Id-Gw AB	IDO Glow Stand-WC Abgang multire
20	ART_345538	575100000 - JOOPI WC-Sitz mit Absenka	575100000	WC-Sitz	EC011196	31.955.9565	WC-Sitze aus Duroplast	JOOPI WC-Sitz mit Absenkautomatik	Gebert JOOPI WC-Sitz
21	ART_374848	243.661.00.1 - Befestigungsrahmen zu Pn	243.661.00.1	Zubehör Vorwand-Einbauelement Sanitär	EC011338	33.342.3423	Ersatzteile WC-Spualaufen	Befestigungsrahmen zu Pneumatikdrücker	Gebert Befestigungsrahmen für pneu
22	ART_1139325	595976000 - SchbIFr Ge+RenP sRf WT-U	595976000		EC000000	31.820.8210	Mobel Ersatzteile Keramik	SchbIFr Ge+RenP sRf WT-U B130 Ech	Gebert Set Schublendenfronten für Unt
23	ART_1774966	500.929.00.5 - Ft-1-Sw-Id-Design 90x200	500.929.00.5	Duschtür	EC010962	32.860.8601	Duschabtrennungen komple	Ft-1-Sw-Id-Design 90x200 s s-m G4r	IDO Design Pendeltür
24	ART_353164	597206000 - VerStl Ge+WbFlf WWC-U	597206000	Zubehör für Dusch-WC	EC010204	31.955.9564	Ersatzteile Keramik Integ	VerStl Ge+WbFlf WWC-U/WBL L 20.3c	Gebert Verankerungsstück für Wand
25	ART_436721	243.724.00.1 - Steuerung für AcQ Sela a	243.724.00.1	Zubehör für Dusch-WC	EC010204	35.426.4357	AquaClean Sela Ersatzteile	Steuerung für AcQ Sela ab 2017	Gebert Steuerung für Gebert AquaClean Sela
26	ART_108129	19455 - SysRoR MAP CSI d28x1 5 ia-ver	19455	Dünnwandiges Stahlrohr	EC010032	22.770.7704	Mapress Systemrotre C-Si	SysRoR MAP CSI d28x1 5 ia-ver	Gebert Mapress C-Stahl Systemrotre I
27	ART_108452	22965 - AnKreuzStk VLRL 2x CSI d15-12	22965	Passierstück	EC011717	22.775.7721	Mapress Filings C-Stahl v	AnKreuzStk VLRL 2x CSI d15-12 Kz	Gebert Mapress C-Stahl Anschlus-K
28	ART_106419	390.116.14.1 - Doppelsteckmuffe PP-MD	390.116.14.1	Filling mit 2 Anschlüsse	EC000324	21.622.6281	PP-MD Verbindungen d4	Doppelsteckmuffe PP-MD d40	Gebert Silent-PP Doppelsteckmuffe
29	ART_109886	51252 - Bogen mEEnd CSI 45G d88.9 F	51252	Filling mit 2 Anschlüsse	EC000324	22.775.7739	Mapress Filings C-Stahl v	Bogen mEEnd CSI 45G d88.9 FKM	Gebert Mapress C-Stahl Bogen mit E
30	ART_1424504	58203650000 - WWC-Top Ge+3Co W3S820	58203650000	WC	EC011289	21.905.9052	Wand-WCs	WWC-Top Ge+3Co W3S 8 T88.5	Gebert 300 Comfort Wand-WC Tiefsp
31	ART_1564903	501.598.00.1 - Sockel-EVWT Ge+SncCp	501.598.00.1	Zubehör für Badezimmerzubehör	EC012353	31.820.8206	Badzimmermöbel Zubehör	Sockel-EVWT Ge+SncCp B45 2-Lgr	Gebert Selnova Compact Ecksockel
32	ART_264141	116.053.00.1 - Reduktion VOL Ms d32-1	116.053.00.1	Filling mit 2 Anschlüsse	EC000324	22.708.7906	Volex Filings Ms 3		

## 9 Bewertung, Methodenreflexion, Empfehlungen

**Bewertung** Sowohl der Industriepartner als auch die Studenten sind mit dem resultierenden Algorithmus zufrieden. Neben einem anwendbaren Model wird ein Prototyp abgeliefert, der eine effiziente Klassifizierung ermöglicht. Durch die simple Abhandlung der Klassifizierung eines Produkts in wenigen Klicks wird der Aufwand eines Sachbearbeiters und die potenzielle menschliche Fehlerrate reduziert. Das Ziel der Verbesserung der ETIM Klassifizierung durch maschinelles Lernen ist erfüllt.

**Methodenreflexion**

Microsoft Teams	<p>Pandemie bedingt wurden Meetings jeweils über Microsoft Teams gehalten. Daraus ergab sich natürlich, dass Sharepoint für das Verwalten der Notebooks Modelle und Dokumentation benutzt wurde.</p> <p>Dies war insbesondere beim Schreiben des Berichts sehr angenehm, da man Änderungen der anderen Kollaboratoren in Real Time sieht.</p> <p>Ein Nachteil, dieses Systems ist Versionierung. Sicher zu stellen, dass alle Beteiligten Parteien, immer dieselben Versionen trainierter Modelle haben ist schwierig. Das Rekonstruieren älterer Versionen von Notebooks oder Modellen war mit Mehraufwand verbunden.</p>
Jupyter Notebooks	<p>Jupyter Notebooks ermöglichten es, im Python Code direkt mit Markdown zu dokumentieren, was die Zusammenarbeit stark erleichterte. Des Weiteren ermöglichen sie einzelne Codeblöcke zu ändern und laufen zu lassen, ohne dabei zeitaufwändige runtime Operationen wie zum Beispiel das Laden von Daten oder trainieren von Models zu wiederholen.</p>
Tensorflow/DNN	<p>Der Grund warum ein DNN mit One-Hot-Encoding als Ansatz gewählt wurde, war das Interesse der Studenten ein Neurales Netzwerk in der Arbeitswelt anzuwenden.</p> <p>Der Ansatz war rückblickend im Vergleich von fastText für die Problemstellung jedoch nicht effizient genug.</p>
fastText	<p>Nach der Datenanalyse war klar, dass ein Textembedding ein vielversprechender Ansatz war. Es gäbe diverse Alternativen zu fastText, welche voraussichtlich ähnlich gute Resultate produziert hätten. An fastText war ausschlaggebend, dass es 'light-weight', 'selfcontained' und 'opensource' ist.</p>

**Empfehlungen**

Gebrauch GUI	<p>Das fastText Model ist so gut trainiert, dass der Gerbrauch davon über das GUI empfohlen wird, um Artikel zu klassifizieren. Die anderen Algorithmen sind um einen so grossen Faktor schlechter, dass es sich nicht lohnt, etwas anderes als fastText einzusetzen.</p> <p>Zu beachten ist, dass das GUI die Resultate der Arbeit darstellen soll und daher alle Algorithmen und ihre Kombinationen enthält.</p>
EC000000	<p>Es wäre sinnvoll, mit dem GUI über alle Artikel der Datenbank zu gehen, welche momentan 'EC000000' zugeordnet sind.</p>

Verbesserung des  
fastText Models

Dadurch würden die Einträge eliminiert, welche nicht zugeordnet wurden und für den Algorithmus als Human Error angesehen werden.

Andererseits könnte das Model inklusive den 'EC000000' Einträgen neu trainiert werden. Somit würde der Algorithmus ebenfalls in Mapping auf 'EC000000' lernen und dieses vorschlagen, wenn die Wahrscheinlichkeit hoch ist, dass dem Produkt keine ETIM Klasse zugeordnet werden kann.

Die Logfunktion des GUIs ermöglicht es, Statistiken darüber zu erschliessen, ob das Model mit der Zeit schlechter wird. Zum Beispiel falls der endgültige Entscheid des Benutzers nicht der wahrscheinlichsten Vorhersage des Models entspricht. Diese Informationen könnten genutzt werden, um am Algorithmus Verbesserungen vorzunehmen.

Ein alternativer Ansatz wäre das periodische Trainieren der Modelle auf der vollumfänglichen Datenbank.

## 10 Literatur und Quellenverzeichnis

**Herkunft der Vorlage** Das Dokument wurde auf der Basis einer Vorlage für Technische Berichte erstellt. Die Vorlage ist ein Element des „Werkzeugkastens Technische Berichte“ der Hochschule für Technik Rapperswil. Sie orientiert sich an Prinzipien des Strukturierten Schreibens.

### 10.1 Quellenverzeichniss

[1] "ETIM Schweiz," [Online]. Available: <https://www.etim.ch/de/>. [Accessed 28 Mai 2021].

### 10.2 Abbildungen

Abbildung 1: ETIM Zusammensetzung, Quelle <a href="https://www.etim.ch/de/klassifizierung/modell-informationen">https://www.etim.ch/de/klassifizierung/modell-informationen</a> .....	6
Abbildung 2: ETIM Verteilung .....	10
Abbildung 3: ETIM Verteilung logarithmisch.....	11
Abbildung 4: ETIM Vorkommen .....	11
Abbildung 5: Prozentanteil der NaN Werte .....	14
Abbildung 6: Graphische Darstellung der Datenaufbereitung erstellt mit <a href="http://www.miro.com">http://www.miro.com</a> .....	22
Abbildung 7: Benchmark Train/Validate Accuracy.....	28
Abbildung 8: DNN Train/Validate Accuracy.....	34
Abbildung 9: DNN Optimierte Train/Validate Accuracy .....	36
Abbildung 10: GUI Vorschau .....	46
Abbildung 11: GUI Vorschau mit Auswahl .....	46
Abbildung 12 erstellt mit <a href="http://www.miro.com">http://www.miro.com</a> .....	52
Abbildung 13 erstellt mit <a href="http://www.miro.com">http://www.miro.com</a> .....	52
Abbildung 14 erstellt mit <a href="http://www.miro.com">http://www.miro.com</a> .....	53
Abbildung 15 erstellt mit <a href="http://www.miro.com">http://www.miro.com</a> .....	53
Abbildung 16 erstellt mit <a href="http://www.miro.com">http://www.miro.com</a> .....	54

### 10.3 Tabellen

Tabelle 1: Spaltenbeschreibung.....	9
Tabelle 2 Unique Values.....	12
Tabelle 3: Zuteilungen der Namen und Ids .....	14
Tabelle 4: Bayes Error mit einzelnen Komponenten .....	16
Tabelle 5: Bayes Error mit zwei Komponenten .....	16
Tabelle 6: Bayes Error mit drei Komponenten.....	17
Tabelle 7: Bayes Error mit vier Komponenten.....	17
Tabelle 8: Datenbearbeitung Aktionen .....	24
Tabelle 9: Datenbearbeitung numerische Details.....	25
Tabelle 10: Datenbearbeitung Zweck.....	25

Tabelle 11 Verwendete Metriken .....	26
Tabelle 12: Lookup Performance .....	30
Tabelle 13: Lookup Schlussfolgerung .....	31
Tabelle 14: Lookup Analyse.....	32
Tabelle 15: Bayes Error .....	32
Tabelle 16: Algorithmen Vergleich Accuracy.....	41
Tabelle 17: Schlechteste etim_id Accuracy.....	42

## Anhang

### I Auswertung der reservierten Daten für die Bewertung

- Allgemein** Hier ist die Auswertung mit den 1000 Datensätzen, welche am Anfang zur Bewertung der Arbeit entfernt wurden.
- Vorgehen** Das Vorgehen hier ist analog zum Testen mit dem `final_test_set` und es werden dieselben Metriken erhoben.

**Resultate** Der fastText Algorithmus gibt folgende Metriken:

Resultat	Top-1	Top-3
<code>final_test_set</code>	0.882	0.900
<code>final_test_set ohne 'EC000000'</code>	0.967	0.987
<code>final_test_set ohne 'EC000000' und unbekannte etim_ids</code>	0.969	0.989

Das DNN gibt folgende Metriken:

Resultat	Top-1	Top-3
<code>final_test_set</code>	0.843	0.876
<code>final_test_set ohne 'EC000000'</code>	0.924	0.961
<code>final_test_set ohne 'EC000000' und unbekannte etim_ids</code>	0.926	0.963

Der Kombinationsalgorithmus aus Lookup Table und DNN gibt folgende Metriken:

Resultat	Top-1	Top-3
<code>final_test_set</code>	0.846	0.875
<code>final_test_set ohne 'EC000000'</code>	0.928	0.959
<code>final_test_set ohne 'EC000000' und unbekannte etim_ids</code>	0.930	0.962

## II Bedienungsanleitung GUI

**Allgemein** Das Tool wurde nicht als eigenständige Applikation zur Verfügung gestellt. Es werden lediglich die gesamten Python Files und Skripts bereitgestellt.

**Startup** Um das graphische Tool zu starten:

1. Navigiere in den Ordner `prototyp_gui`
2. Installiere `requirements.txt`

```
pip install -r requirements.txt
```

3. Starte das Programm

```
python application.py
```

**Updates der Modelle** Das Tool lädt alle Modelle aus `prototyp_gui/algorithms/models`. Falls neue Modelle zur Verfügung stehen, genügt ein Austauschen der Modelle im oben genannten Ordner.

### III Abgabeordner

**Allgemein** Der Abgabeordner enthält jegliche Daten und Algorithmen, welche für die Arbeit benutzt wurden.

#### Top Level

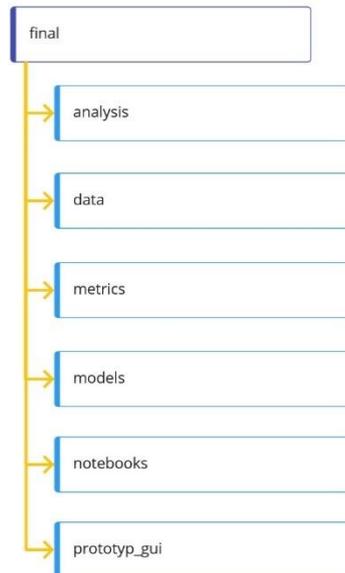


Abbildung 12 erstellt mit <http://www.miro.com>

**Analysis** Enthält Daten und Metriken aus der Datenanalyse.

#### Data

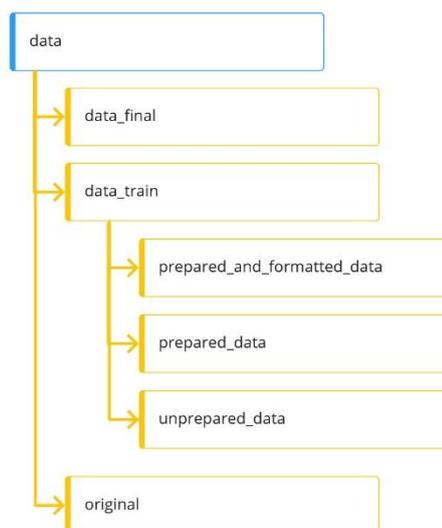


Abbildung 13 erstellt mit <http://www.miro.com>

Der Data Unterordner enthält alle benutzen Datensätze. Diese sind aufgeteilt nach `original`, `data_final` und `data_train`. Im `data_train` befinden sich die jeweiligen Unterteilungen der Daten nach der Datenaufbereitung.

## Metrics

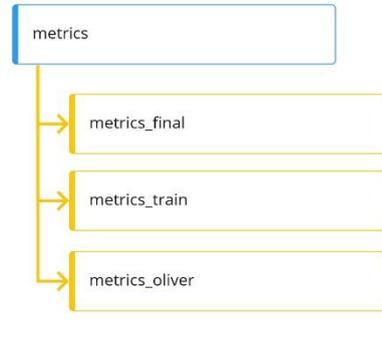


Abbildung 14 erstellt mit <http://www.miro.com>

Die Metriken enthalten alle Metriken der verschiedenen Algorithmen.

`metrics_train` enthält die Metriken der Trainingsalgorithmen auf dem `prepared_data_validate`.

`metrics_final` enthält die Metriken der finalen Algorithmen auf dem `final_test_set`.

`metrics_oliver` enthält die Metriken der finalen Algorithmen auf den Daten, die vor dem Beginn der Arbeit bereits abgetrennt wurden.

## Models

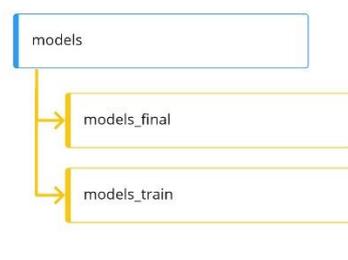


Abbildung 15 erstellt mit <http://www.miro.com>

Die Models enthalten die jeweiligen Versionen der Algorithmen.

`models_train` enthält die Modelle aller Algorithmen, die auf dem `prepared_data` oder einer Variation davon trainiert sind.

`models_final` enthält alle Modelle aller Algorithmen, die auf allen Daten ausser dem `final_test_set` trainiert sind.

**Notebooks**

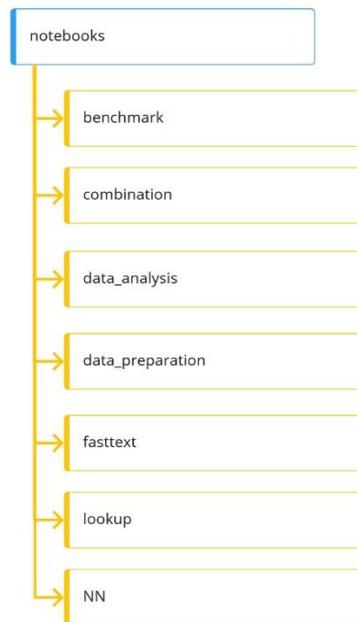


Abbildung 16 erstellt mit <http://www.miro.com>

Die Notebooks enthalten alle Jupyter Notebooks, die während der Arbeit erstellt wurden.

**Prototyp GUI**

Der Prototyp GUI Ordner enthält die Python Anwendung.

**IV Management Summary**

# ETIM Klassifikation mit ML-Approach

Vorhersage der ETIM-Klasse basierend auf Daten aus der Geberit Datenbank

Autoren:	Etienne Baumgartner, Nathanael Gall
Betreuer:	Prof. Oliver Augenstein
Experte:	Nico Schmid
Themengebiet:	Machine Learning
Industriepartner:	Geberit AG, Jona

**Ausgangslage** Das ETIM Klassifikationsmodell setzt sich auf dem internationalen Markt immer mehr durch. Die ETIM Klassifikation vereinfacht einerseits den Datenaustausch zwischen Händler Hersteller sowohl als auch die Klassifizierung eines Produkts im internationalen Elektrotechnik Fachbereich. Die Geberit AG nutzt diesen Standard für ihren Produktkatalog und die Klassifizierung wird manuell durch einen Sachbearbeiter vorgenommen.

Gesucht sind Möglichkeiten den Klassifizierungsprozess der Produkte mit Hilfe von maschinellem Lernen einerseits zu beschleunigen und andererseits potenzielle menschliche Fehler zu vermeiden.

**Vorgehen / Technologie** Die Suche nach passenden Lösungen verlangt die genaue Analyse der bereits vorhanden, manuell klassifizierten Produktdaten. Spezifisch wird eine Analyse auf den Daten durchgeführt, um allgemein gültige Charakteristiken zu extrahieren. Um die Eigenschaften und gewonnenen Informationen vollumfänglich auszunutzen, wird im zweiten Schritt eine Datenaufbereitung vorgenommen.

Auf der Basis der verarbeiteten Produktdaten werden zwei Modelle aufgebaut.

**One-Hot-Encoding:**

Die Daten werden für das Trainieren auf einem Deep Neural Network One-Hot-Encoded. Um die Anzahl Dimensionen des Netzwerks möglichst klein zu halten, wird anhand des Bayes Errors die optimale Spaltenkombination ermittelt. Auf dieser Auswahl wird ausserdem ein Lookup Table erstellt, um eindeutige Datensätze direkt zu klassifizieren

**Textembedding:**

Die hohe Anzahl und der hohe Informationsgehalt in den textbasierten Spalten ermöglichen den Aufbau einer Text Klassifizierung. Die Produkte werden durch den fastText Algorithmus und anhand von Textembeddings zugeordnet.

- 
- Ergebnisse** Nach der Optimierung dieser zwei Ansätze stellt sich heraus, dass der textbasierte fastText Algorithmus genauere Resultate liefert. Es wurde eine Pythonanwendung geschrieben, welche fastText benutzt, um dem Nutzer für jeden Datensatz eine Auswahl von drei potenziellen ETIM Klassen auszugeben. Die Wahrscheinlichkeit, dass sich die korrekte ETIM Klasse in dieser Auswahl befindet ist 98.2%.
- Ausblick** Wenn man diese Anwendung zur Klassifizierung benutzt, kann man das Model in Zukunft noch verbessern.  
Vom Algorithmus abweichende Entschlüsse des Sachbearbeiters können analysiert werden und bieten die Möglichkeit das Model basierend auf den Entscheidungen des Benutzers entsprechend anzupassen.