

Table of Contents

1. Intro
2. What is the tariff classification?
3. How machine learning works?
4. How texts can be classified?
5. How it could be implemented in our business processes?
6. How it works in our case?
7. How it performs on test data?
8. What is about results?
9. What can be improved?
10. Conclusion

1. Intro

Working as classification analyst at company_name, I meet on an everyday basis information dedicated to part numbers, part descriptions, HTS codes, failed notes and, of course, audit trails. I use SQL queries to search the existing classifications, similar parts and so on. In other words, I work with a huge amount of text data every day. company_name database has collected gigabytes of data over the years. It would be imprudent not to use such treasury. Some months ago, I started to dive deeper into this topic and now I want to represent a small research dedicated to approaches that can be implemented in our everyday work. The main goal of this research is to show the prospects of the machine algorithms and how they can enhance our job.

2. What is the tariff classification?

Everyday work of classification analysts is pretty routine. Generally speaking to classify part means to assign the appropriate HTS code depending on its properties. Firstly, there is a need to find a part to do. It could be found in the tracker if it is an urgent part, or in the xls-file that gathers work required part numbers from TradeVia database. When part is found the analyst does a research to determine a correct HS code. There are several methods to obtain required information like drawings from client's PLM systems, answers from engineers, US rulings, BTIs and other tools. Sometimes it is very useful to start analysis from searching for the part with similar part number or part description.

Next, when classification is done, part goes to technical advisors. They need to decide if the classification code and audit trail is correct or not. If the classification is incorrect, the analyst receives a failed review and part should be reworked. As analysts describe their logic of choosing a code in the audit trail, technical advisors also write a short failed review that explains why part hadn't passed the review. Technical advisors see HS code, part number, part description and audit trail. They use their experience and deep knowledge of the client's business to decide. It is significant that they work mainly with text data as well as classification analysts.

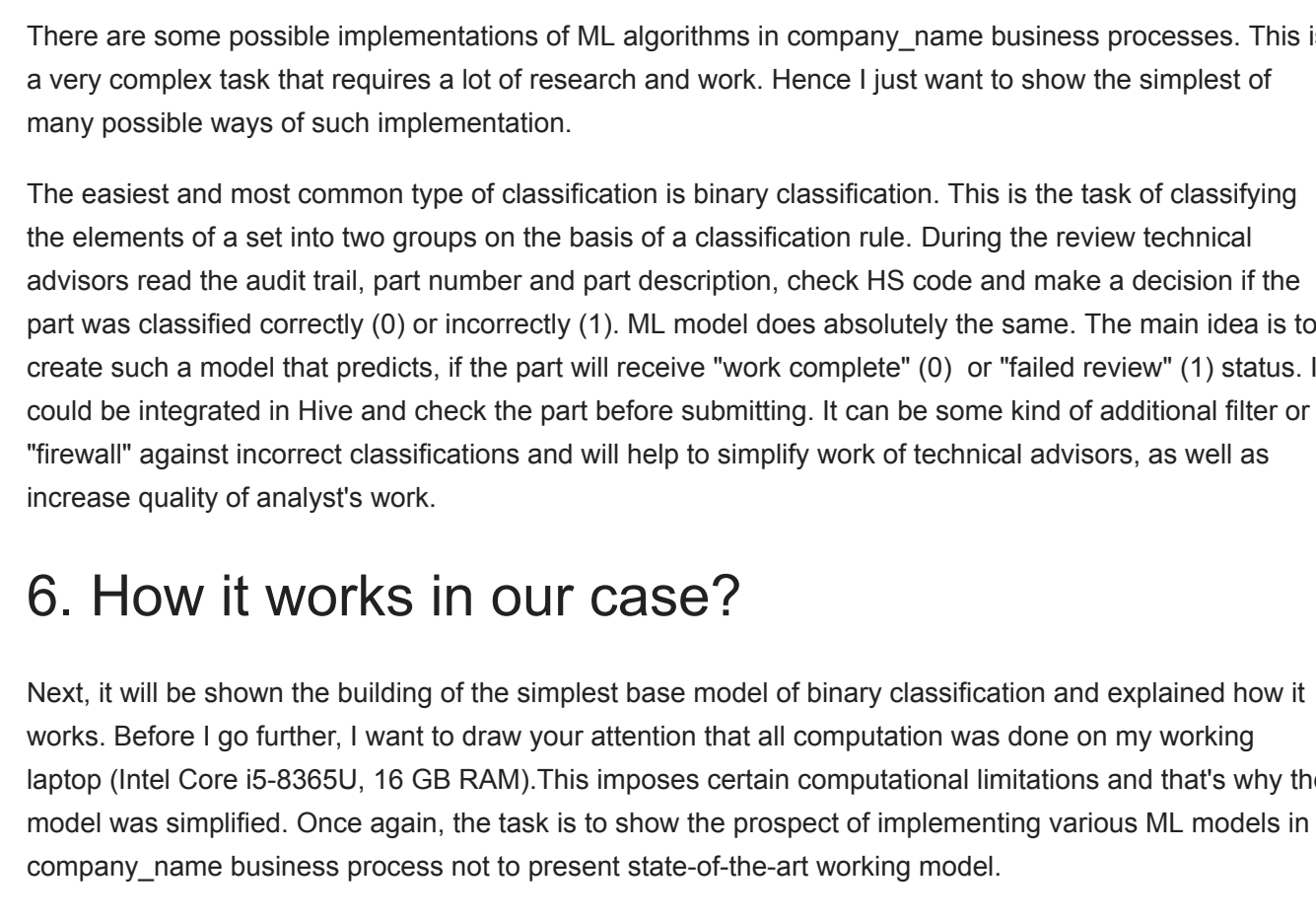
Mistakes always happen in our work and it is normal. Audit trail can be interpreted not clearly, part can have unusual part number, broker can add wrong part description and there are many other reasons. Of course, it should not be forgotten that analysts as well as technical advisors have a monthly target of part number to be completed. Analysts also have quality target. So there is not much time that can be spent on analysis of each part.

3. How machine learning works?

There is a great definition of machine learning (ML) was defined by Tom Mitchell, an industry pioneer:

Machine learning (ML) is the study of computer algorithms that allow computer programs to automatically improve through experience

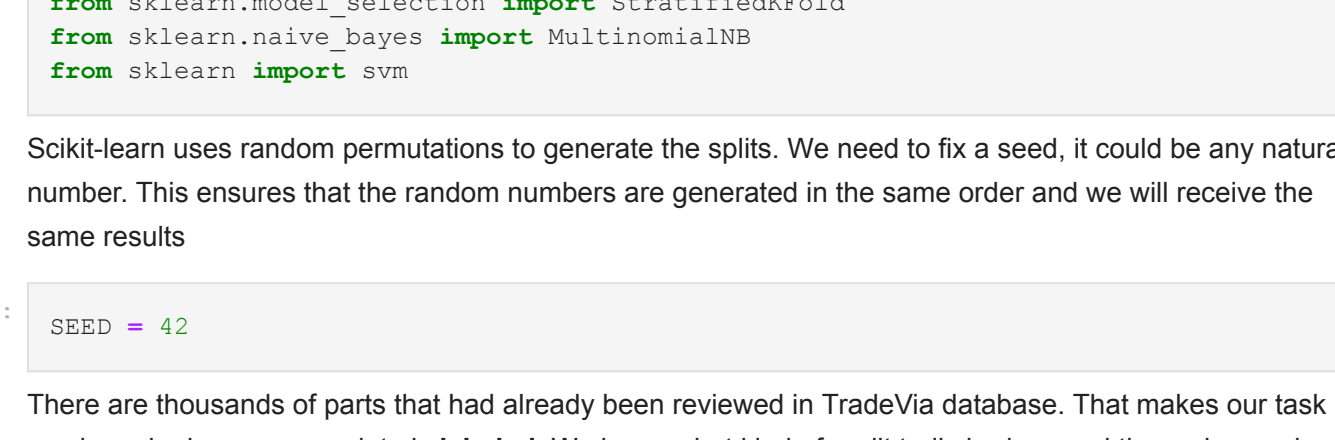
ML algorithms build a model based on sample data in order to make predictions or decisions without being explicitly programmed to do so. We meet these algorithms everyday, they have become an essential part of our life so we don't even notice how often we use them. Virtual Personal Assistants like Siri, personalized recommendations on Netflix, email spam and malware filtering in Gmail, all these technologies are based on ML algorithms. Basic types of algorithms are represented below



Machine learning algorithms use computational methods to "learn" information directly from data without relying on a predetermined equation as a model. They find natural patterns in data that generate insight and help you make better decisions and predictions.

4. How texts can be classified?

One of the big areas of machine learning is classification of texts. It is the process of assigning tags or categories to text according to its content. It is one of the fundamental tasks in natural language processing with broad applications such as sentiment analysis, topic labeling, spam detection, and intent detection. Usual baseline of text classification is shown on the figure.



Firstly, we extract our text data (corpus) from any kind of the source (databases, CRMs, ERPs, social networks and so on). Then text should be preprocessed and usually there is a main part of the work. Corpus need to be cleaned and transformed into a numerical representation in the form of a vector. Next important task is the creation of new features that will help to improve the quality of predictions. When data is ready then it could be applied various classification methods. There are a lot of different mathematical approaches that can solve the problem of classification and there is no universal solution. Various problems require various models and various approaches.

5. How it could be implemented in our business processes?

There are some possible implementations of ML algorithms in company_name business processes. This is a very complex task that requires a lot of research and work. Hence I just want to show the simplest of many possible ways of such implementation.

The easiest and most common type of classification is binary classification. This is the task of classifying the elements of a set into two groups on the basis of a classification rule. During the review technical advisors read the audit trail, part number and part description, check HS code and make a decision if the part was classified correctly (0) or incorrectly (1). ML model does absolutely the same. The main idea is to create such a model that predicts, if the part will receive "work complete" (0) or "failed review" (1) status. It could be integrated in Hive and check the part before submitting. It can be some kind of additional filter or "firewall" against incorrect classifications and will help to simplify work of technical advisors, as well as increase quality of analyst's work.

6. How it works in our case?

Next, it will be shown the building of the simplest base model of binary classification and explained how it works. Before I go further, I want to draw your attention that all computation was done on my working laptop (Intel Core i5-8365U, 16 GB RAM). This imposes certain computational limitations and that's why the model was simplified. Once again, the task is to show the prospect of implementing various ML models in company_name business process not to present state-of-the-art working model.

This model will try to predict if the part will pass the review or not based on its part description, part number, HTS code and audit trail.

Firstly, let's import necessary libraries. I use Python 3.8.3, **Pandas 1.1.4** and **NLTK 3.5** libraries for pre-processing of the data and **scikit-learn 0.23** for building model

```
In [1]: import numpy as np
import pandas as pd
import nltk
# Import warnings
warnings.filterwarnings('ignore')
nltk.download('stopwords')
nltk.download('wordnet')
from nltk import sent_tokenize, word_tokenize
from nltk.stem.snowball import SnowballStemmer
from nltk.stem.wordnet import WordNetLemmatizer
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
import string
from sklearn.feature_extraction.text import TfidfVectorizer
from scipy.sparse import coo_matrix, hstack
from sklearn.model_selection import train_test_split
from sklearn.linear_model import RidgeClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import f1_score, roc_auc_score, accuracy_score
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
from sklearn.naive_bayes import MultinomialNB
from sklearn import svm
```

Scikit-learn uses random permutations to generate the splits. We need to fix a seed, it could be any natural number. This ensures that the random numbers are generated in the same order and we will receive the same results

```
In [2]: SEED = 42
```

There are thousands of parts that had already been reviewed in TradeVia database. That makes our task much easier because our data is **labeled**. We know what kind of audit trails had passed the review and what had not. I used only part number from <client's name> business because I'm most familiar with that client. I used not all <client's name> businesses but only the following: <client's businesses>. In future it could be integrated for all clients and all business. Also due to the computational limitations I used limited quantity of parts that passed review. I took them since 2019-01-01 to 2019-06-30 for complete parts and whole 2019 for failed parts because there are fewer of them but our model must be trained on balanced dataset. Also I used parts reviewed by only several TAs because otherwise dataset would be too big. In additional I used parts that have strength ('8', '20', '55', '90', '91', '92', '95', '100'). To create a model that will find regularities in corpuses of audit trails and determine which of them leads to successful or unsuccessful classification, I will combine two parts of passed audit trails and failed audit trails. Firstly let's import data to Python from .csv file and throw out unnecessary columns.

First part of the model will have parts that had passed the review.

```
In [3]: passed = pd.read_csv('model_passed_review.csv', header = None)
passed = passed.drop([1, 5, 6, 8, 9], axis=1)
passed = passed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
passed = passed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)
passed_data = passed.fillna('No audit trail')
passed_data['Failed'] = 0
```

Then let's do the same with the parts that hadn't passed the review. Also, I concatenate two datasets into one and split the target variable as y

```
In [4]: failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)
```

```
In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number
```

```
Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64
```

```
In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)
```

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use **lemmatization** but in our case stemming works slightly better

```
In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)

In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number

Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64

In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use lemmatization but in our case stemming works slightly better

In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)

In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number

Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64

In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use lemmatization but in our case stemming works slightly better

In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)

In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number

Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64

In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use lemmatization but in our case stemming works slightly better

In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)

In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number

Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64

In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use lemmatization but in our case stemming works slightly better

In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)

In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number

Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64

In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use lemmatization but in our case stemming works slightly better

In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",
    4: "Audit_trail"
})
failed = failed.drop_duplicates(subset=['Audit_trail']).reset_index(drop=True)

In [5]: failed_data = failed.fillna('No audit trail').reset_index(drop=True)
failed_data['Failed'] = 1

In [6]: train_data = pd.concat([passed_data, failed_data])
train_data.groupby(['Failed']).count().Part_number

Out[6]: Failed    0    36674
          1    26106
          Name: Part_number, dtype: int64

In [7]: y = train_data['Failed']
train_data = train_data.drop(['Failed'], axis=1)

Our data is ready for further work. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length. That's why we need to preprocess it and prepare it for building the model. Here I will use the technique called "Bag-of-words". It is one of the most commonly used methods in text classifications and also it is easy to implement. In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity. The frequency of occurrence of each word is used as a feature for training a classifier. The following steps are:

1. reduce all words to lowercase
2. use regular expressions to represent all words and symbols as "tokens"
3. join tokens into one string separated by spaces
4. reduce words to their stem (e.g. "play" from "playing") using stemming algorithms. It also possible to use lemmatization but in our case stemming works slightly better

In [8]: TOKENIZE_RE = re.compile(r'[a-zA-Z]+[-?!\d+\'@,;~<\/p><\/div>
<div data-bbox="100 670 950 698" data-label="Text"><pre>In [9]: train_data = preparing_data(train_data)

In [10]: In a large text corpus, some words will be very present hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequency of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the <math>tf-idf</math> transform. We will throw away so-called "stopwords", like "and", "the", "per". Also, here we got two additional thresholds: <math>max\_df</math> and <math>min\_df</math>. <math>max\_df = 0.90</math> means "ignore terms that appear in more than 90% of the texts". <math>min\_df = 4</math> means "ignore terms that appear in less than 4 texts". It will vectorized our train data into a sparse two-dimensional matrix suitable for feeding into a classifier. <math>ngram\_range=(1, n)</math> means that it will be used a string of n words in a row. In our case <math>n = 3</math>

In [11]: vectorizer = TfidfVectorizer(analyzer='word',
                                stop_words=stop_words,
                                max_df=0.9,
                                min_df=4,
                                ngram_range=(1, 3))

In [12]: train_data_sparse = vectorizer.fit_transform(train_data)

In [13]: len(vectorizer.vocabulary_)

Out[13]: 186405</pre></div>
<div data-bbox="100 699 950 718" data-label="Text"><p>Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called <b>overfitting</b>. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a <b>holdout</b> part. Usually it is 20..30% of the whole train dataset.</p></div>
<div data-bbox="100 719 950 728" data-label="Text"><pre>In [14]: X_train, X_holdout, y_train, y_holdout = train_test_split(train_data_sparse, y,
                    failed = pd.read_csv('model_failed_review.csv', header = None)
failed = failed.drop([5, 6, 7, 8], axis=1)
failed = failed.rename(columns={
    0: "Business",
    1: "Part_number",
    2: "Part_description",
    3: "HS_code",

```