

Labeling Clinical Reports with Topic Modeling and Active Learning

*Uppmärkning av kliniska rapporter genom topic modeller och
active learning*

Simon Lindblad

Supervisor : Marco Kuhlmann
Examiner : Arne Jönsson

External supervisor : Mikael Nilsson

Upphovsrätt

Detta dokument hålls tillgängligt på Internet – eller dess framtida ersättare – under 25 år från publiceringsdatum under förutsättning att inga extraordinära omständigheter uppstår. Tillgång till dokumentet innebär tillstånd för var och en att läsa, ladda ner, skriva ut enstaka kopior för enskilt bruk och att använda det oförändrat för ickekommersiell forskning och för undervisning. Överföring av upphovsrätten vid en senare tidpunkt kan inte upphäva detta tillstånd. All annan användning av dokumentet kräver upphovsmannens medgivande. För att garantera äktheten, säkerheten och tillgängligheten finns lösningar av teknisk och administrativ art. Upphovsmannens ideella rätt innefattar rätt att bli nämnd som upphovsman i den omfattning som god sed kräver vid användning av dokumentet på ovan beskrivna sätt samt skydd mot att dokumentet ändras eller presenteras i sådan form eller i sådant sammanhang som är kränkande för upphovsmannens litterära eller konstnärliga anseende eller egenart. För ytterligare information om Linköping University Electronic Press se förlagets hemsida <http://www.ep.liu.se/>.

Copyright

The publishers will keep this document online on the Internet – or its possible replacement – for a period of 25 years starting from the date of publication barring exceptional circumstances. The online availability of the document implies permanent permission for anyone to read, to download, or to print out single copies for his/hers own use and to use it unchanged for non-commercial research and educational purpose. Subsequent transfers of copyright cannot revoke this permission. All other uses of the document are conditional upon the consent of the copyright owner. The publisher has taken technical and administrative measures to assure authenticity, security and accessibility. According to intellectual property law the author has the right to be mentioned when his/her work is accessed as described above and to be protected against infringement. For additional information about the Linköping University Electronic Press and its procedures for publication and for assurance of document integrity, please refer to its www home page: <http://www.ep.liu.se/>.

Abstract

Abstract.tex

Acknowledgments

Acknowledgments.tex

Contents

Abstract	iii
Acknowledgments	iv
Contents	v
List of Figures	vii
List of Tables	ix
1 Introduction	1
1.1 Motivation	1
1.2 Aim	2
1.3 Research questions	2
1.4 Delimitations	3
1.5 Structure of the Report	3
2 Background	4
3 Theory	6
3.1 Text Processing using Unsupervised Techniques	6
3.2 Text Classification	11
3.3 Active Learning	13
3.4 Evaluation Metrics	18
3.5 Related Work	20
4 Method	22
4.1 Frameworks, Tools and Implementation	22
4.2 Datasets	23
4.3 Pre-Processing and Text Representation	24
4.4 Exploratory Study	26
4.5 Experiments	28
4.6 Evaluating the Label Balance	31
5 Results	33
5.1 Exploratory Study	33
5.2 Filter Out Invalid Clinical Reports Using Topic Models and Clustering	34
5.3 Alternatives to Labeling at Random	37
5.4 Evaluating the Label Balance	42
6 Discussion	46
6.1 Results	46
6.2 Method	46

7 Conclusion	48
Bibliography	49

List of Figures

2.1	A screenshot of the interface used to label clinical reports at Sectra	5
3.1	Diagram of the LDA model.	8
3.2	(a) to (e) shows iterations of k-means until convergence. In (e) it can be seen that the new centroids capture the same documents as the previous iteration, and we have converged.	10
3.3	An overview of an active learning system	13
4.1	A sample report from the dataset provided by Sectra	24
4.2	The distribution over the labels in the initial set of labeled data provided by Sectra	25
4.3	The distribution over the labels in the Reuters data	25
4.4	A 2D plot of the text data, where each point is colored by topic with the highest probability	27
4.5	Wordclouds for a 75 topic LDA model	27
4.6	A way to visualize and analyze topics based on their relevance and frequency	28
5.1	A 2D plot of the full word2vec plot. Given the amount of terms used there is a lot to analyze.	34
5.2	A 2D plot of the zoomed in word2vec plot. Most of the values here are names. This represents the red box in from Figure 5.1	35
5.3	The perplexity scores for the different LDA models	35
5.4	The distribution of the topics with the highest probability for the invalid reports in the training set. Note that only topics that occurred at least once are shown in the histogram.	36
5.5	The distribution of the topics with the highest probability for the valid reports in the training set. Note that only topics that occurred at least once are shown in the histogram.	37
5.6	The distribution of the number of prominent topics assigned to the invalid reports in the training set.	38
5.7	The distribution of the number of prominent topics assigned to the valid reports in the training set.	39
5.8	The labeled data points plotted in 2D, and colored based on the first label of the report in alphabetical order.	39
5.9	The counts of the most likely topics for the four most common categories.	40
5.10	The categories of the different reports that are assigned a certain topic as the most likely one.	41
5.11	Accuracy of the models with initial sample size 25	41
5.12	Accuracy of the models with initial sample size 25	42
5.13	Micro recall of the models with initial sample size 25	42
5.14	Micro precision of the models with initial sample size 25	43
5.15	Time to query a new set of samples by the different strategies	43
5.16	The distribution of labels after random sampling	44

5.17	The distribution of labels after Binary Version Space Minimization	44
5.18	The distribution of labels after Binary Version Space Minimization with clustering	45
5.19	The distribution of labels after Adaptive Active Learning	45

List of Tables

3.1	Confusion matrix for explaining true positives, false positives, true negatives and false negatives	19
4.1	The different combination of topic model/k-means clusters that were evaluated. .	29
4.2	The different configurations of active learning strategies evaluated. BinMin stands for Binary Version Space Minimization, MMC stands for Maximum Loss Reduction with Maximum Confidence and Adaptive stands for Adaptive Active Learning.	32
5.1	The synonyms, misspellings and shorts found in the data that the author could with assert with confidence.	34
5.2	The results of the classification of the invalid reports. The manual column represents the use of manual interpretation of the LDA model.	37
5.3	The number of reports assigned a certain category, as well as the number of different topics assigned as the most likely one for reports with the given category. . . .	40



1 Introduction

The world's population is growing each year. Making healthcare more efficient and robust is of great importance in order to handle the challenges that arises with a growing population. One way of increasing the efficiency as well as the quality of healthcare is to create automated systems that can aid doctors in their process. As the population is growing it is of utmost importance to ensure that the quality of diagnoses remain high, and to minimize the risk of missing some critical piece of information. Taking advantage of the available medical information is key to create aforementioned systems.

The purpose of this thesis is to evaluate different techniques for labeling multi-label medical reports. It is done at Sectra Medical Imaging IT Solutions AB in conjunction with a research project where they investigate how they can use machine learning and text mining techniques on clinical reports to improve their products. By actively choosing the reports to be labeled, the goal is to make the set of labeled reports more useful in future systems.

1.1 Motivation

Information pertaining to a patient's diagnosis is often in the form of written clinical reports. This is one example of data that can be utilized in automated systems, by extracting information from old reports, the process of writing new ones can be enriched. A system could show cases with similar features as the one currently being written, the doctor could then compare the findings and check if they have obtained an abnormal result. Being able to perform such a comparison will result in extra quality assurance in the diagnostic flow. It could also provide doctors with more confidence in that their diagnosis is correct.

The problem systems like this would face is to categorize medical reports in order to make further suggestions. One approach that is commonly used for such problems is machine learning. Which is a field where a set of inputs is used to create a mapping to some output values [6]. This is done by using data to build a, usually statistical, model.

The task of predicting a category, or class, for a given text document is called text classification. Text classification is usually solved using supervised learning [1]. In supervised text classification, there exists a set of inputs, in this case text data, that already has a category assigned to it. This data is then used to fit the model so that it later make predictions for inputs that it has not yet been exposed to. One model that have been shown to be successful in text classification is Support Vector Machines (SVM) [19, 1, 40].

In order to fit a machine learning model to predict categories for clinical reports, we need a set of already categorized, or labeled, data. That is, we need to assign categories to an existing set of clinical reports. It is often the case that text data is widely available, but it is harder to come by data that is already labeled. Obtaining high quality data is important to use in machine learning systems, both in healthcare and other areas. Since the models require a sufficient amount of labeled reports, the task of labeling them can be cumbersome. Especially in the case of clinical data, since doctors and other clinicians time is both valuable and expensive. By improving the process and the quality of data to be labeled, and thereby reducing the number of reports that need to be categorized, they can spend more time doing their job.

A field within machine learning that is focused on the task of improving the data labeling process is called active learning. It is a form of semi-supervised learning. The algorithm queries an oracle (in this case a doctor) for labels for the data points that it think will help the model improve the most. This is used when there is plenty of readily available data, but assigning labels is expensive. Since the data points to be labeled are actively selected, the models can require fewer examples than if they were selected at random. The points can be selected by considering the certainty of the models, and request to label the documents that the model is less certain about. Another approach, which has not been given as much attention, is using the underlying structure of the data to select points. The goal with this approach is to capture the distribution of the categories.

If one of two classes is assigned to each document, you have a binary classification problem [6]. Problems where one of several classes is assigned to an instance is called a multi-class classification problem. Multi-labeled classification is when you assign one or more label to each document. This thesis is mainly concerned with multi-label classification. Assigning several classes to a document is more time consuming than in the cases where you only need to find one option. In those cases the process can be stopped when the appropriate label has been found. However, when a document can be assigned several classes the entirety of the text needs to be considered. For example, a news article be on several subjects, such as both economics and sport. Just because the category sport has been identified, the entire document has to be read in order to find any additional categories. This makes the use of active learning to enhance the labeling of documents even more useful in the multi-labeled case.

1.2 Aim

The purpose with this thesis project is to evaluate different solutions to increase the quality of labeled reports, and thereby reducing the amount of them needed for use within the project. Resulting from this will be a complete, standalone, system for labeling reports. The reports are interactively queried so a user can label the reports that are deemed most useful by the system. Labeling data to use in machine learning will probably be necessary for a long time ahead, but the aim here is to create a system that makes it easier. This will then be used together with an existing web interface that Sectra created for the purpose of labeling reports.

1.3 Research questions

The specific research questions that this thesis will treat are presented here. They will be the main focus of study.

1. *Is it possible filter out invalid clinical reports by using an unsupervised techniques such as topic modeling?*

In the dataset from Sectra, there are reports describing patients not showing up for or changing the time of their appointments, deceased patients or patients that have been ordered to another hospital. These reports does not contain any information of

value from a medical point of view and should not be considered in the report labeling process.

Unsupervised machine learning models such as topic modeling does by definition not require any labeled documents to train on. If it is possible to, without any such data, to capture the necessary variance and group these invalid reports together and they can be removed from the labeling process before a doctor is presented with them that would be an additional hurdle removed from the process.

2. *What active learning strategies are good alternatives to sample documents at random in a multi-label document labeling system? How well does these alternatives perform?*

How we are choosing the documents to be sampled is important. If the dataset that is being sampled is skewed, i.e. some categories are a more frequent than others, our labeled set will likely follow that distribution. This will result in models requiring a lot of labeled documents to gather a sufficient amount of reports that are of the less common categories. Without these, the model will only perform well on the frequently occurring categories.

If the decision boundaries of our models can be used to pick documents that would be more informative for the model, the number of labeled documents could be reduced and still obtain the same performance. Another approach to selecting the documents to sample is to take advantage of the underlying structure of the data.

The algorithms to be evaluated will be based on the models certainty, as well as taking advantage of the underlying structure of the data. When choosing the algorithm to use, there are several different factors that will affect the final results, and therefore need to be taken into consideration. How well the models perform on the data will be evaluated based on the accuracy, precision, recall and f1-score of the models. But they also need to be able to query documents in a reasonable time, if it is expensive to label reports it is likely to be expensive to wait for the reports to be queried. If the algorithm needs a big initial set of labeled reports is another factor that will be evaluated.

3. *How does the algorithms from question 2 effect the balance of labels in the labeled dataset? Another indication on the quality of the labeled reports is the balance between the classes. Based on the initial sampling, the underlying distribution of labels in the clinical data is far from uniform. There are certain categories, like the one describing that everything is okay with the patient, that is a lot more common than other more rare illnesses. Even though the original data may be imbalanced, selecting samples that contains a better balance between the different categories could improve the performance of the models. When a training set is imbalanced, the standard learning algorithms' performance can be significantly reduced [16] The goal here is to see which one of the different sampling techniques that will result in the best balance between the different categories in the resulting dataset.*

1.4 Delimitations

Even though the sampling strategies are evaluated objectively on the Reuters dataset, the applicability of the techniques on clinical data is only evaluated by one physician, on the one dataset provided by Sectra.

1.5 Structure of the Report

The next chapter covers the background theory that is relevant for the thesis. After the theory, the methodology used is described, which is followed by a chapter covering the results. The method and results are then discussed in Chapter 6. Finally, Chapter 7 presents the conclusions.




2 Background

Sectra AB is creating products both in medical IT and secure communications. It is a multinational corporation that currently has offices in 12 countries. Medical imaging is their biggest business, and radiology is the main area within that. This thesis is carried out at Sectra Medical Imaging IT Solutions AB, as a part of their research division that focuses on text analytics. They are currently pursuing a research project with Region Skåne. Region Skåne is responsible for the healthcare in Skåne, the southern most county of Sweden. The goal behind the project is to use machine learning and text mining techniques to improve the functionality of their products and aid the physicians in their work. Among other things, this can be suggesting categories to doctors while they are writing medical reports. Another case is to use the categories of documents to present doctors with medical reports that are handling similar cases from the past. With this information the doctors could get an extra quality assurance check in their diagnostic flow.

By using these techniques there are plenty of opportunities and ways to accomplish new interesting things with textual data. But in order to build these system, you need a substantial amount of clinical reports that are labeled. Therefore, the purpose of this thesis is to increase the quality and efficiency of the process of labeling these reports. This will be done by using unsupervised learning techniques such as topic models and clustering to first remove documents that are not supposed to be labeled. These invalid reports are documents that describe patients that did not get examined for some reason. Examples of this can be deceased patients and patients being moved to a different hospital.

For the process of labeling the reports, a complete system is needed. This system should use active learning in conjunction with the aforementioned unsupervised techniques to increase the quality of the labeled documents. In the work that they have done so far in the research project a doctor has done some initial labeling of reports. This was done by simply selecting the reports in the order they were on file. The doctor that primarily worked with the labeling of reports stated that the distribution over the labeled categories is very skewed. The vast majority of labeled documents were assigned to a small subset of the categories. In addition, a skewed dataset causes the number of clinical reports to be labeled to increase a lot. For a statistical model to be able to achieve good results with the less frequently occurring categories, a large number of reports needs to be labeled in order to obtain a good amount of reports with these categories.

Navigator

 victor

Antal gjorda: 1

Utlåtandets id: 24912948

Bekänt

Framåt

Ätergå

90 år

Progress av meningiom? Andra fokala förändringar?

Leukoencefalomalaci? Atrofi?

Anamnes

Kvinnu med typ 2 diabetes, B12-brist, astma, polyneuropati. Börjar bli glöms och har lite nedsatta kognitiva funktioner sedan några år tillbaka men klarar sig själv i huset utan någon hemhjälp utan klarar det mesta i hushålllet själv. Sedan tidigare känt falskt meningiom som vid kontroll 2013 respektive 2014 inte hade progredierat men här finns viss kompression av intilliggande hjärnvävnad, dock inget ödem eller medellinjeshöverskjutning. Lilly upplever tilltagande huvudvärksproblematik. Tackam snart CT för att se om här trots allt finns någon progress av denna förändring och för mer generell kartläggning vad gäller hennes begynnande kognitiva svårigheter.

Kreatininvärde tidigare i år 75 med estimerat GFR 50. Både baserat på kreatinin respektive cystatin C.

Med vänlig hälsning,

Leg. läk

Farmaka

Omnpaque Inj,lös 350 mg i/ml, 70 ml

Undersökningar

DT hjärna utan och med iv kontrast

Utlåtande

[NUM-SEQ]:56 Datortomografi av hjärna utan och med iv kontrast

Jämfört DT hjärna från [NUM-SEQ].

Känt meningiom bakre falx har oförändrat storlek och karaktär. Inget ödem i angränsande parenkym.

Ingen blödning. Ingen färsk eller gammal infarkt. Ingen intraparenkymal expansivitet. Inget ödem. Måttliga vitsubstansförändringar periventrikulärt. Lätt kortikal atrofi och måttlig atrofi av hippocampus bilateralt.

Normalt luftförande bihålor och mastoidceller.

Färsk Intrakraniell blödning

Gammal Intrakraniell blödning

Extrakraniell blödning

Färsk infarkt

Gammal infarkt

Intrakraniell infektion

Sinult

Tumör

Atrofi

Kärlsjukdom

Småkärlssjukdom

Likvorcirkulationsrubning

Postoperativa förändringar / resttillstånd

Fraktur

Annat

Normal

Nästa

Ej kategoriserbar

By using active learning, that is active selection of the samples to be labeled, the goal is to reduce the number of labeled samples needed to achieve sufficient results with the model. Sectra has a simple website that they use for labeling reports, which the active learning techniques will be incorporated into. The existing web interface can be seen in Figure 2.1.



3 Theory

In this section the theory behind the techniques used in this thesis work will be presented. The first part will go through the techniques used to process the data and to perform an exploratory analysis. After that text classifications, primarily with SVM, will be covered. Text classification, primarily using Support Vector Machines, will be covered in the section after that. Following this, there will be a section containing an overview of the field of active learning, as well as a comparison of some different active learning techniques for multi-labeled data. This chapter will then end by going through the different evaluation metrics that will be used to assess the methods used, as well as a section describing work that is related to this thesis.

3.1 Text Processing using Unsupervised Techniques

Techniques in machine learning that does not require you have a categorized or labeled set of data is called unsupervised learning. They use the structure of the data to obtain the information to use when processing it. These techniques may have different goals. Some are used to estimate a distribution, others are used to reduce the number of dimensions in the data by trying to find dimensions that capture as much as possible of the variance, and some are used for the purpose of identifying groups of similar points within the data [6]. When it comes to text data, there are a few common methods and techniques that are unsupervised, and can be used for different purposes. Examples of such techniques are *topic modeling* and *clustering*. Another interesting technique is word2vec, that is used to produce word embeddings.

When working with text, it needs to be represented in a way that allows the models to work with it effectively. *Bag-of-words* (BoW) is one of the more common representations when performing text analysis. Using BoW the text is represented as a multi-set. That is, a document is represented by the number of occurrences of the different words. Representing text in this way therefore becomes high-dimensional, there is one dimension for each word in the vocabulary. Like the name implies, the positions of the words are not taken into account, they are viewed as if they were taken from a bag. In written language a word can be used to express several different thoughts, and one thought can be expressed using several different words. This is something that cannot be captured with BoW. However, it is easy to work with, and is commonly used when performing topic modeling among other things.

One way to incorporate positional information into the representation is the use of n-grams. Instead of storing information pertaining to one term, information is stored with regards to n consecutive terms. Consider the text “Pattern Recognition and Machine Learning”. Using an n-gram with $n=2$, called a bigram, would result in the tokens: “Pattern Recognition”, “Recognition and”, “and Machine”, and “Machine Learning”.

Topic Modeling

A topic model is a statistical model for finding topics within text [12]. The topics build upon the probability that a certain word would occur in a text about a given topic, on the basis of terms occurring together. For example, if the topic represents United States politics, words such as “government”, “Trump”, “Reagan”, “Senate”, or “Medicaid” are more likely to appear than “sailboat” or “sweater”. Any given document can then contain a certain topic with some probability. This can be viewed as fuzzy clustering, and that the document has a degree of membership in a topic or cluster [12]. The most common topic model in use today is Latent Dirichlet Allocation (LDA) [12]. Another topic model that preceded LDA is Probabilistic latent semantic analysis (PLSA) [17]. However, PLSA has been shown to be more prone to overfitting than LDA [12].

In the rest of the report, the following notation will be used:

- D denotes a corpus of M documents: $D = \{w_1, w_2, \dots, w_M\}$.
- The number of topics is K . Each topic is indexed by i .
- N_d is the number of terms in document d .
- N_i is the number of terms in topic i .
- V denotes the number of words in the vocabulary.

Latent Dirichlet Allocation

Latent Dirichlet allocation (LDA) is a statistical model, where abstract topics in the model are defined as distributions over words [7]. LDA is based on a generative process, a model of which can be seen in Figure 3.1. The circles in this figure represent random variables. Dependencies between these random variables are shown with arrows, and if a variable is observed it is shaded in the figure. In this model, the only observed variable is the words in the document. Parts of the model are surrounded by a rectangle to show that the part is repeated several times.

The generation of a corpus is done with the following steps [12, 7]:

- **Draw a distribution over the words for each topic.** A sample ϕ_i is drawn from a symmetric Dirichlet distribution with parameter β . This sample represents the distribution of terms for the topic i .

$$\Phi_i \sim \text{Dir}(\beta) \quad (3.1)$$

$$p(\Phi_i | \beta) = \frac{\Gamma(V\beta)}{\Gamma(\beta)^V} \prod_{v=1}^V \phi_{iv}^{\beta-1} \quad (3.2)$$

Here, Γ is the gamma function.

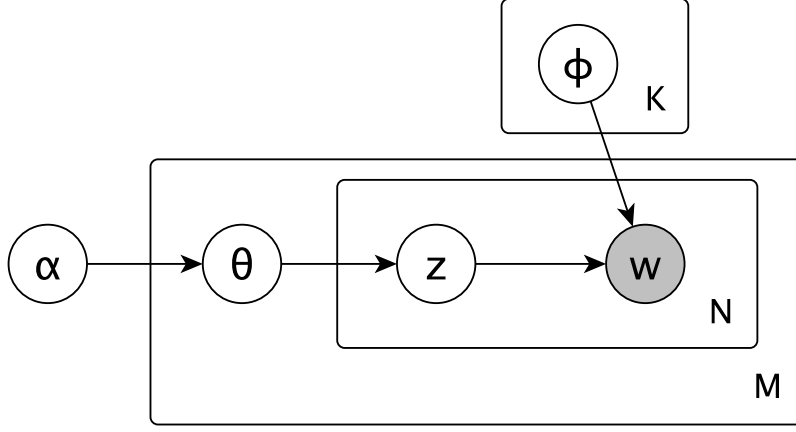


Figure 3.1: Diagram of the LDA model.

- **Draw a distribution over the topics for each document.** A sample θ_d is drawn from a Dirichlet distribution with parameters α . This sample represents the distribution of topics for document d .

$$\Theta_d \sim \text{Dir}(\alpha) \quad (3.3)$$

$$p(\Theta_d | \alpha) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{di}^{\alpha_i - 1} \quad (3.4)$$

- For each token with index n :
 - **Draw a topic assignment z_{dn} for the token index n .** z_{dn} is drawn from the distribution over topics for each document. That is, z_{dn} is drawn from a multinomial distribution using θ_d as a parameter.

$$z_{dn} \sim \text{Multinomial}(\Theta_d) \quad (3.5)$$

$$p(z_{dn} = i | \Theta_d) = \theta_{di} \quad (3.6)$$

- **Draw a token w_{dn} .** The token w_{dn} is drawn from the topic distribution assigned to the index n . That is, w_{dn} is drawn from a multinomial with parameter $\phi_{z_{dn}}$.

$$w_{dn} \sim \text{Multinomial}(\Phi_{z_{dn}}) \quad (3.7)$$

$$p(w_{dn} = v | z_{dn} = i, \Phi_i) = \phi_{iv} \quad (3.8)$$

The LDA model identifies topics from different terms that occur together. Consider the case where an LDA model has been used to learn a number of topics. Two terms that frequently occur together are then likely to be in the same topic. So, if the same word has been used to express different thoughts, and the word has the same probability in two topics, the words that it co-occurs with can be used to differentiate between the different thoughts.

The task of learning the LDA model is a Bayesian Inference problem. We have several variables that we cannot observe: the word distribution for the topics (ϕ_i), the topic assignments for the tokens (z), and the topic distribution for the documents (θ_d). The only observed variables are the words in the document. We have to approximate the posterior distribution using some sampling method, since it cannot be inferred automatically [7].

There exist a few algorithms that can be used to learn topics for the LDA model. Two of these that has shown to be able to extract useful topics from text are *collapsed Gibbs sampling* [15] and *variational Bayes* [7]. Variational Bayes works by using simple single-variable models to approximate the LDA. As a consequence, it disregards any dependencies between the variables.

Collapsed Gibbs Sampling

Gibbs sampling is a Markov chain Monte Carlo (MCMC) method that is often used to obtain a good estimate for the distribution of a probability model, when it is not feasible to sample the distribution directly [12]. With the help of a heuristic, or randomly, Gibbs sampling initializes the variables. During a large number of iterations the variables are then sampled. When variable is being sampled, it is conditioned on the others. In an MCMC fashion, a number of samples are rejected during an initial burn-in period. This is done in order to get to a state where the points are more representative of the distribution that is estimated.

Griffith et al. [15] came up with collapsed Gibbs sampling, where θ and ϕ are marginalized out. The only variable to then be repeatedly sampled is the topic assignment, z_{dn} , conditioned on the assignments of the other tokens.

Text Clustering

Cluster analysis is commonly defined as finding groups in a given dataset. The members of these groups are determined to be similar by a similarity measure [20, 2]. Text data is sparse, but yet have a very high dimensionality. With one dimension per term in the dictionary, it is not uncommon with dimensions in the order of 10^5 . For this reason, some of the more naive clustering algorithms does not work well for text data [2].

In distance-based clustering, a similarity function is used to measure the closeness between two text documents. For the purpose of measuring the similarity between text objects, the cosine similarity function is commonly used [2], as well as Euclidean distance. Two different approaches to distance-based clustering are distance-based partitioning, and agglomerative hierarchical clustering. For the distance-based approach, k-means and k-medoid are two frequently used algorithms.

K-means Clustering

When using the k-means clustering algorithm, the clusters are based upon an initial set of k representatives. A simple approach to k-means clustering can be seen as:

1. Select K seeds from the original dataset
2. Assign the rest of the documents to one of these seeds, based how how similar they are by the similarity function
3. Before each new iteration, select a new centroid for each cluster. This should be the point that is the best central point for the cluster.
4. Repeat step 2 and 3 until convergence.

A visualization of this can be seen in Figure 3.2. One advantage that k-means has over K-medoid is that it requires a small number of iterations, especially compared to K-medoid [2,

35]. However, k-means is rather sensitive to the selection of initial seeds. One approach is to just select them randomly, or selecting them based on the result of another lightweight clustering method. A frequently used method is k-means++, that has been shown to improve both the speed and accuracy of k-means clustering [4].

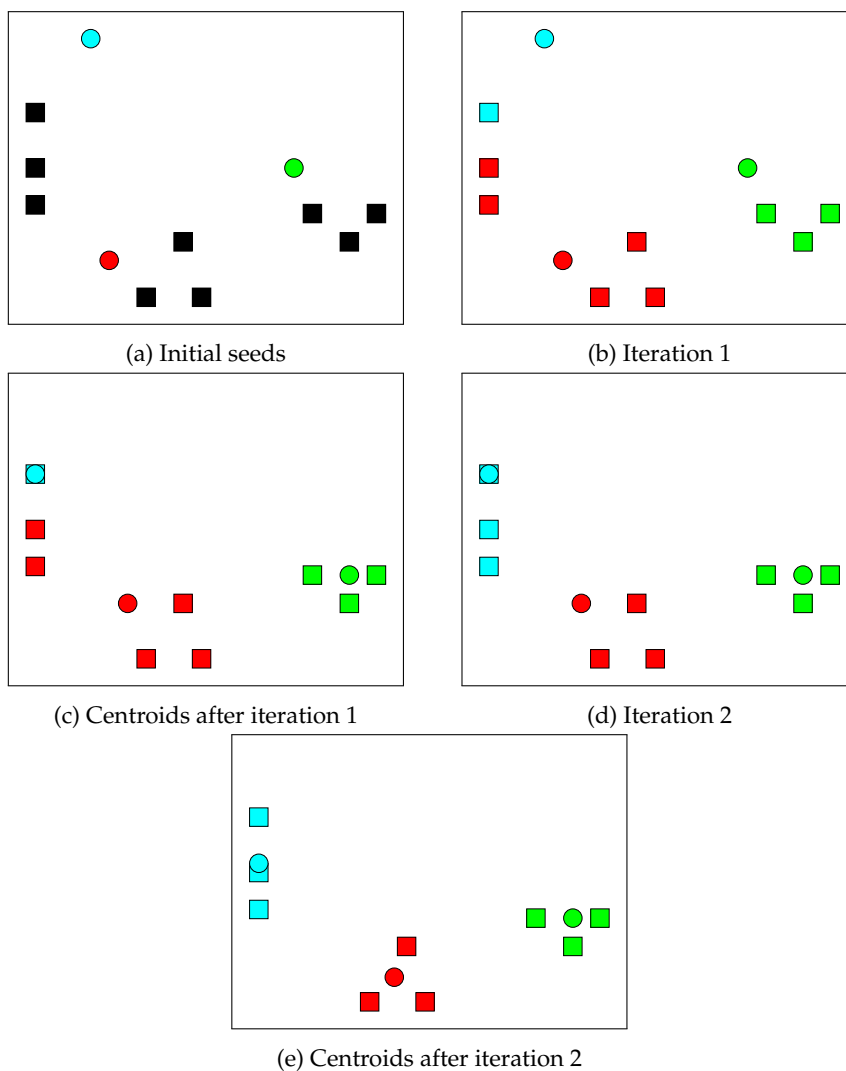


Figure 3.2: (a) to (e) shows iterations of k-means until convergence. In (e) it can be seen that the new centroids capture the same documents as the previous iteration, and we have converged.

k-means is commonly used with Euclidean distance, which is defined on two n dimensional points p_1 and p_2 as:

$$d(p_1, p_2) = \sqrt{(p_{1_1} - p_{2_1})^2 + (p_{1_2} - p_{2_2})^2 + \dots + (p_{1_n} - p_{2_n})^2} \quad (3.9)$$

Word Embeddings

Word embeddings can be done with several different techniques, and it is the process of representing a word as a real-valued vector instead of just an atomic unit. Viewing a word as a vector allows for doing interesting things with them, such as evaluating how similar to words are. Something that is hard to do when treating them as atomic units.

This can be accomplished by using a co-occurrence matrix to see how often certain words occur together, and then perform some dimensionality reduction on them [21, 22]. Another approach that shown to be very successful in producing high quality word embeddings is *word2vec*, which uses a neural network to accomplish this task [27]. In addition to being able to compute similarities between words, using simple algebraic some interesting relationships could be discovered. The example that Mikolav et al. [27] showed was that using the vector for “King”, subtracting the vector for “Man” and adding the vector for “Woman” resulted in a vector that was close to that for “Queen”.

One approach to this is using a continuous bag-of-words model [27]. A neural network is used to predict the middle word using both words occurring before and after it. The four words occurring before and after the middle word is used as inputs, but their internal order is not used, which is why the name bag-of-words is used.

The second approach that Mikolav et al. explored was a continuous skip-gram model [27]. Here a neural network with a continuous projection layer was used. With the current word as input, co-occurring words within a certain range are predicted. A bigger range is more computationally complex, but results in word vectors of higher quality.

Word2vec does not require labels to be provided with the data, but uses the data itself to generate targets. For this reason it is sometimes called a self-supervised technique.

3.2 Text Classification

Text classification is a widely studied field within Computer Science. It is an important problem in supervised machine learning, and it is the task of assigning one or more classes to a given text document [1]. The problem is mainly approached with supervised machine learning. That is, with a dataset that consists of a collection of text documents, where each document has one or more classes assigned to it. With the help of these labels, a classification model is fitted to the data. The goal of this is for the model to be able to correctly assign a class to a previously unseen text document. Some of these classification models can also produce a probability of a document being of a certain class. Other models are based on the concept of a margin that separates the classes, and the distance between a data point and a margin can be used to indicate how certain the model is of the assigned label [40]. Example of use cases for text classification is categorization of news articles, document retrieval and email filtering. There exists several different models for classifying text. Decision trees, neural networks and Support Vector Machines (SVM) are some have been previously applied to the text domain with successful results [2]. In this thesis, SVM are the main focus, since they have been studied extensively in the context of active learning. Logistic regression is used for the binary classification task of classifying invalid reports.

Support Vector Machines

SVMs work by implicitly map the training data to a feature space [6]. The goal is that the data should be linearly separable in the feature space, even if it is not in the input space. In the case of binary classification, a point is classified by the linear model:

$$y(x) = w^T \phi(x) + b \quad (3.10)$$

, where the sign of $y(x)$ determines which label will be assigned to x .

The goal of an SVM is to try to find the hyperplane that maximizes the margin. That is, the distance between any point and the decision boundary should be as large as possible. A subset of the data points will be used to determine where the decision boundary is, these

points are called the support vectors. The hyperplane that gives us the maximum margin can be found by:

$$\arg \min_{w,b} \frac{1}{2} \|w\|^2 \quad (3.11)$$

The parameter b is not apart of this equation. However it is implicitly set by the constraints determined by the data points relation to the decision boundary. In order to allow for better generalization, and for data that is not completely linearly separable, SVMs make use of slack variables, denoted ξ_n , to penalize points that are close to the decision boundary [6]. A parameter $C > 0$ controls how much effect the slack variables will have. The equation with the slack variables become:

$$\arg \min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_n \xi_n. \quad (3.12)$$

A smaller C -value allows more points to be misclassified, which is done in order to achieve better generalization.

Logistic Regression

Despite having regression in the name, logistic regression is a classification model. It is appropriate to use when there are categorical, binary, targets. Logistic regression works by determining the conditional probability of a class C_1 , given a feature vector ϕ . This probability can then be used with a certain threshold to determine whether or not a data point is a part of a certain class or not. Consider the case where we have two classes, C_1 and C_2 . Their probabilities are then calculated by [6]:

$$P(C_1|\phi) = \text{sigm}(w^T x) \quad (3.13)$$

, where *sigm* is the *logistic sigmoid* function defined as:

$$\text{sigm}(a) = \frac{1}{1 + e^{-a}} \quad (3.14)$$

The probability for C_2 is then obtained with:

$$P(C_2|\phi) = 1 - P(C_1|\phi) \quad (3.15)$$

From Equation 3.13 it is clear that the number of parameters in the model is the same as the number of dimensions in the feature vector. Maximum likelihood is used in order to determine the values of these parameters. For a dataset with the features x_n for $n = 1, \dots, N$ and targets $t_n \in \{0, 1\}$, the likelihood function can be written as:

$$P(t|w) = \prod_{n=1}^N P(C_1|\phi)^{t_n} (1 - y_n)^{1-t_n} \quad (3.16)$$

, where $t = (t_1, \dots, t_N)^T$.

This is traditionally used with the Iterative Reweighted Least Squares method [6].

Multi-Label Classification

Multi-label classification is the type of text classification where one instance can be associated with multiple labels. It is a generalized version of the multi-class classification problem, where you have more than 2 labels, but each document is only assigned one [41].

A common way of solving multi-label classification problems is the Binary Relevance method [30, 8, 26]. It is a way of transforming the multi-label classification problem into

several different binary ones. With Binary Relevance you fit one classification model per label in your data. Each of these classifiers are then predicting whether or not the document is associated with the corresponding label or not.

3.3 Active Learning

Conventional machine learning systems use a set of available data to find a hypothesis that can explain the patterns within the data. The purpose of active learning is to allow a system to *select* the data that it wants to be labeled, and therefore the data it wants the model to be trained on [36]. An active learning system samples a document to be labeled from a pool of unlabeled data, and then queries an oracle, which is often a human annotator, to get the label for that document. By being able to decide what data to label and use, the goal is that the system can achieve better results, and that the data will be of higher quality. To get a better understanding of how the process works, these are the steps commonly iterated over until enough samples have been labeled:

1. Evaluate the samples in the unlabeled pool based on a particular measure that is defined by the querying strategy.
2. The selected samples are presented to an oracle that is queried for the labels. This oracle is commonly in the form of a human annotating the instances.
3. The newly labeled samples are added to the labeled pool.
4. A machine learning model is trained on the labeled pool, this model is often used by the querying strategy in order to select samples to be label in step 2

A model of the active learning system can be seen in Figure 3.3.

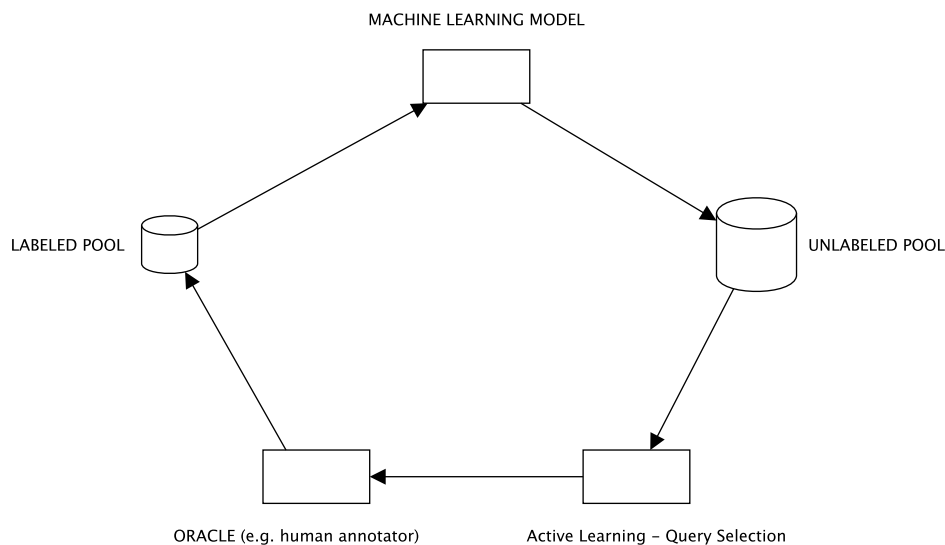


Figure 3.3: An overview of an active learning system

In several different domains, data is readily available and easy to obtain. But even if the data is abundant, labels for the data is often harder and more expensive to come by [36], especially when it comes to multi-label problems.

The next section will describe different ways to access the documents in active learning systems, followed by some theory of how the samples relate to the hypothesis space. After that some concrete methods for selecting the samples to queried are described and compared.

Pool-Based Sampling

The main focus in active learning is how to select the samples to be labeled. There are different sampling methods in use, and which one is more appropriate depends on how the data can be accessed. Pool-based sampling is motivated by the assumption that there exists a large ready pool of data, where only a small portion is labeled [23, 36]. The samples to be labeled are selected by evaluating the entire pool of unlabeled data, and choose the most appropriate ones based on a defined utility measure. If the entire pool is large, a subset could be used instead. For applied active learning, pool-based sampling seems to be the most popular choice [24], but there are some alternatives that have been used in theoretical settings such as stream-based selective sampling. The difference between stream-based selective sampling and pool-based sampling is the individuality of the decisions in the former. Where you draw one sample at a time from an input source and make the decision whether or not to query a label for it. For text applications, where a set of data is often readily available, pool-based sampling is often the more appropriate option since you can consider the entire dataset. Pool-based is therefore the sampling technique that is considered in this paper.

Searching the Hypothesis Space

In machine learning, a hypothesis is a specific configuration of a model, the purpose of which is to predict outputs on new instances of data by generalizing the training data. One hypothesis can, for example, be an SVM model with specific values for the parameters. The set of all possible hypothesis that we are working with is called the *hypothesis space* [32]. Following the SVM example, the hypothesis space would be the set of SVMs with the different values under consideration. The hypothesis space is defined as:

$$\mathcal{H} = \left\{ f \mid f(x) = \frac{w * \phi(x)}{|w|} \right\} \quad (3.17)$$

where \mathcal{W} is our parameter space and $w \in \mathcal{W}$. The subset of the hypothesis space that in the feature space separates the data is called the version space, which is defined as [36, 40]:

$$\mathcal{V} = \left\{ f \in \mathcal{H} \mid y_i f(x_i) > 0 \forall i \in \{1 \dots n\} \right\} \quad (3.18)$$

So the version space therefore represents the different hypotheses that make correct predictions on the training data. Under the assumption that one of the hypothesis can fully separate the data, the version space shrinks when more labeled data is acquired. So for new labeled instances the hypotheses in the version space will give better predictions for the training data. Based on this, an active learning algorithm should aim to reduce the size of the version space with each new sample, optimally make it half the size in each iteration.

There exists a useful relationship between the feature space \mathcal{F} and the hypothesis space \mathcal{H} called the *version space duality* [40, 43]. It states that hyperplanes in the hypothesis space corresponds to points in the feature space, and the other way around. So by selecting points to be labeled, constraints can be enforced on the hypothesis that form the version space.

One approach to this that has shown to be successful is *Uncertainty Sampling* [36]. This is commonly done with SVM, since the idea behind SVM is to find a hyperplane that separates two classes in a binary classification with the maximum margin. Out of the different hyperplanes in the hypothesis space, the version space contains those that can successfully separate the data. Uncertainty sampling aims to select the points in the feature space that will reduce the amount of the valid hypotheses the most. An SVM model tries to find the support vectors that maximizes the decision boundary in the feature space, separating the two classes. Considering this in \mathcal{H} , it will be analogous to the hypothesis in the center of the hypothesis space, which is encompassed by the constraints set by the labeled data. What Uncertainty Sampling

is doing is predicting the values for the unlabeled points, and then choose the one that it is most uncertain about, which will be the one closest to the decision boundary. Based on the version space duality, it is a good approximation for dividing the version space in two [36].

Binary Version Space Minimization

Binary version space minimization is a generalization of uncertainty sampling, that is designed to make it work with multi-label data. The approach taken is to decompose the multi-label problem to several binary tasks with the binary relevance method, as is discussed in Section 3.2. The unlabeled point that is chosen for labeling is then the one with the smallest SVM margin across all the binary classification tasks. By doing this, it does not incorporate the multiple labels into the decision process, but treats all classes individually and equally. Selecting a new data x_{new} to be labeled using binary version space minimization can formally be defined as:

$$x_{new} = \arg \min_{x \in D_U} \left(\min_{i \in \{1 \dots n\}} y_i(x) \right) \quad (3.19)$$

, where n is the number of labels, $y_i(x)$ is the function from Equation 3.10 for the binary SVM classifier that is used to predict label i , and D_U is the set of unlabeled data points.

Maximum Loss Reduction with Maximum Confidence

Maximum loss reduction with maximum confidence (MMC) was developed by Yang et al [44]. The goal of this technique is to find the samples that will reduce the expected model loss the most, and select this sample for labeling. These are the basic notations that will be used when explaining the MMC approach:

- The labeled dataset: D_L .
- The unlabeled dataset: D_U .
- Possible query set: D_S .
- Optimal query set: D_S^* .
- The classification function that is trained on dataset D_L : f_{D_L} .
- A data point: x , and its label: y .
- The loss of on data point x : $L(f_{D_L}(x))$.
- The expected loss of the model: $\widehat{\sigma_{D_L}}$.

The expected model loss that MMC is trying to reduce can be defined as follows [44]:

$$\widehat{\sigma_{D_L}} = \int_x \left(\sum_{y \in Y} L(f_{D_L}(y|x)) P(y|x) \right) P(x) dx \quad (3.20)$$

It is hard to estimate $P(x)$, so it is instead measured over all the samples in D_U . This results in the estimate:

$$\widehat{\sigma_{D_L}} = \frac{1}{|D_U|} \sum_{x \in D_U} \sum_{y \in Y} L(f_{D_L}(y|x)) P(y|x) \quad (3.21)$$

After a set of data points D_S has been labeled, the new dataset $D'_L = D_L + D_S$ is obtained. Under the assumption that any $x \in D_U - D_S$ has an equal effect on a model trained on the datasets D_L and D'_L , we get the following equation for the reduction of the expected loss [44]:

$$D_S^* = \arg \max_{D_S} (\widehat{\sigma_{D_L}} - \widehat{\sigma_{D'_L}}) = \arg \max_{D_S} \left(\sum_{x \in D_S} \sum_{y \in Y} (L(f_{D_L}(y|x)) - L(f_{D'_L}(y|x))) P(y|x) \right) \quad (3.22)$$

In their paper, Yang et al. considers the process of finding the greatest reduction in two steps: finding a good estimate for the conditional probability $p(y|x)$, and finding a way to assess the loss reduction of a multi-label classifier.

It is unfeasible for a query strategy to provide an estimation for all possible label combinations. If there are n different labels, there will be 2^n different label combinations. In order to estimate the conditional probability $p(y|x)$, MMC uses an approach that first estimates the number of labels for a given data point, and then uses that estimate to select the most probable labels. Consider the case where a data point has m labels. Since we can obtain the probability from our classification model for each label, we can sort them in descending order. The first m labels are then likely to have a high probability, while the rest a rather low probability.

Yang et al. [44] describes the process of estimating the number of labels as follows:

1. Use the classification model to obtain the probabilities for each label for all the data samples.
2. For each data sample, sort and normalize the probabilities for all the labels.
3. Using the labeled dataset, fit a logistic regression classifier with the sorted and normalized probabilities as features, and the number of labels as target.
4. With the fitted logistic regression model, predict the number of labels for the samples in the unlabeled pool.

After obtaining the predicted number of labels m for a sample x , the estimate for $p(y|x)$, denoted \hat{y} , is then obtained by selecting m the most probable labels based on the original classification models output.

The only task that's left is now to estimate the loss for the multi-label classifier. By using the model where there are k different binary classifiers for a problem with k labels, the model loss can be calculated by adding the loss for the different binary classifiers like:

$$L(f) = \sum_{i=1}^k L(f^i) \quad (3.23)$$

where the loss of a single binary classifier is denoted as $L(f^i)$. With this definition, it only remains to define the measure of loss on a single binary classifier. The measurement that is used by MMC is to estimate the model loss by the size of the version space of the SVM [40, 44]. The version space's size can be computed with Equation 3.18. However, computing this for each possible label set is expensive. By using the heuristic from [39], an approximation of the version space with the added label can be obtained from the current SVM classifiers' margin. The reduction rate after adding the data point (x, y^i) , where $y^i \in -1, 1$, can be expressed as follows [44, 39]:

$$\frac{L(f_{D_L+(x,y^i)}^i)}{L(f_{D_L}^i)} \approx \frac{V_{D_L+(x,y^i)}^i}{V_{D_L}^i} \approx \frac{1 + y^i f_{D_L}^i(x)}{2} \quad (3.24)$$

$L(f_{D_L}^i)$ does not involve the sample selected for labeling, so by writing the loss reduction as:

$$\begin{aligned} L(f_{D_L}) - L(f_{D_L'}) &= \sum_{i=1}^k (L(f_{D_L}^i) - L(f_{D_L'}^i)) \\ &= \sum_{i=1}^k (L(f_{D_L}^i) (1 - \frac{L(f_{D_L'}^i)}{L(f_{D_L}^i)})) \end{aligned} \quad (3.25)$$

it can be seen that focusing on the reduction rate is sufficient. By incorporating the result from Equation 3.24, the following approximation for the reduction rate is obtained:

$$\sum_{i=1}^k \left(\frac{1 - y^i f_{D_L}^i(x)}{2} \right) \quad (3.26)$$

The only thing that remains is now to combine the estimation of $p(x|y)$ with the estimate for loss reduction. The resulting equation, called maximum loss reduction with maximal confidence, is [44]:

$$D_S^* = \arg \max_{D_S} \left(\sum_{x \in D_S} \sum_{i=1}^k \left(\frac{1 - \hat{y}^i f_{D_L}^i(x)}{2} \right) \right) \quad (3.27)$$

The full algorithm for the approach can be seen in Algorithm 1.

Algorithm 1 Maximum Loss Reduction with Maximal Confidence procedure. Taken from Yang et al. [44], with some modifications to the notations used in order to make it coherent with the rest of the report.

Input: Labeled set D_L

Unlabeled set D_U

Number of classes k

Number of selected examples per iteration S

repeat

Train k binary SVM classifiers f^1, \dots, f^k based on training data D_L .

for each $x \in D_U$ **do**

Predict its label vector using the method described in Section 3.3.

Calculate the expected loss reduction with the most confident label vector \hat{y} , $score(x) =$

$$\sum_{i=1}^k \left(\frac{1 - \hat{y}^i f_{D_L}^i(x)}{2} \right)$$

end for

Sort $score(x)$ in decreasing order for all x in D_U .

Select a set of S examples D_S^* with the largest scores, and update the training set $D_L \leftarrow D_L + D_S^*$.

until enough instances are queried

Adaptive Active Learning

In their paper, Li et al. [24] presents two approaches to active learning:

- Max-Margin Uncertainty Sampling
- Label Cardinality Inconsistency

These two techniques are then combined in a weighted fashion into what they call Adaptive Active Learning.

Max-Margin Uncertainty Sampling

The idea behind *Max-Margin Uncertainty Sampling* comes from the observation that multi-label classification prediction is mainly about separating the positive labels from the negative labels [24]. That is, separating the labels are assigned to an instance and the ones that are not. In order to model the uncertainty of the prediction of the prediction on a datapoint, Li et al. suggest the usage of a global separation margin to separate the negative labels from the positive.

The positive labels for a data point x is defined as those where $\text{sign}(f'_{D_L}(x))$ is positive. The separation margin is defined as:

$$\begin{aligned} \text{sep_margin}(x) &= \min_{i \in \hat{y}^+} f'_{D_L}(x) - \min_{i \in \hat{y}^-} f'_{D_L}(x) \\ &= \min_{i \in \hat{y}^+} |f'_{D_L}(x)| + \min_{i \in \hat{y}^-} |f'_{D_L}(x)| \end{aligned} \quad (3.28)$$

where \hat{y}^+ denotes the set of predicted labels that are positive on the instance, and the \hat{y}^- denotes the negative ones.

The data point that the model is the most uncertain about is then the one with the smallest margin. Li et al. define their global measure, max-margin prediction uncertainty, as:

$$u(X) = \frac{1}{\text{sep_margin}(X)} \quad (3.29)$$

Label Cardinality Inconsistency

Label Cardinality Inconsistency is based on that the underlying distribution is the same for the labeled and unlabeled data. In a multi-label dataset, the *label cardinality* is defined as the average number of labels assigned to each class [41]. The selection strategy that Li et al. based on this measures the Euclidean distance between the number of assigned predicted labels on x , and the label cardinality of the labeled data:

$$c(x) = \left\| \sum_{i \in \hat{y}^+} 1 - \frac{1}{N_L} \sum_{y \in Y_L} \sum_{i \in y^+} 1 \right\|_2 \quad (3.30)$$

where N_L is the number of labeled samples, Y_L are the labels for those samples, and y^+ are the positive labels in y .

Integration - Adaptive Active Learning

Since *max-margin uncertainty sampling* and *label cardinality inconsistency* complement each other, an integration method is used:

$$q(x, \beta) = u(x)^\beta \cdot c(x)^{1-\beta} \quad (3.31)$$

where β is a parameter controlling the weight put on the two measures. This parameter is chosen by in each iteration evaluating a discrete set of values, for example $\{0, 0.1, 0.2, \dots, 1\}$. Then selecting β based on the most informative sample amongst them. Equation 3.31 shows the *approximate generalization error*, which is used to select the sample.

$$\epsilon(x) = \sum_{x \in D_U} \max_{i \in f(x)^+} (1 - f^i(x)) + \max_{i \in f(x)^-} (1 + f^i(x)) \quad (3.32)$$

where $f(x)^+$ and $f(x)^-$ are the predicted positive labels, and negative labels, respectively. So, the sample is then selected by:

$$x^* = \arg \min_{x \in D_U} \epsilon(x) \quad (3.33)$$

The complete algorithm for the integrated approach can be seen in Algorithm 2.

3.4 Evaluation Metrics

For classification or information retrieval systems, the typical evaluation metrics in use are *precision*, *recall* and *recall* [18]. We define the following metrics in terms of *true positives*, *false*

Algorithm 2 Adaptive Active Learning Procedure. Taken from Li et al. [24], , with some modifications to the notations used in order to make it coherent with the rest of the report.

Input: Labeled set D_L
 Unlabeled set D_U
 Parameter set B

```

repeat
  Train multi-label SVM classifier  $F^0$  on  $D_L$ .
  for each  $x_i \in D_U$  do do
    Compute  $u(x_i)$  and  $c(x_i)$ .
  end for
  for each  $\beta \in B$  do do
    Mark a candidate instance  $x = \arg \max_{x \in D_U} q(x, \beta)$ 
  end for
  Copy all marked candidate instances into a set  $S$ .
  for each  $x \in S$  do do
    Produce  $\hat{y}$  using classifier  $F^0$ .
    Retrain a new classifier  $F$  on  $(x, \hat{y}) \cup D_L$ .
    Compute  $e(x)$  using classifier  $F$  and Equation 3.32.
  end for
  Select instance  $x^*$  from  $S$  using Equation 3.33
  Remove  $x^*$  from  $D_U$ , query its label vector  $y^*$ .
  Add  $(x^*, y^*)$  to  $D_L$ .
until enough instances are queried
  
```

	Correct P	Correct N
Predicted P	True Positive	False Positive
Predicted N	False Negative	True Negative

Table 3.1: Confusion matrix for explaining true positives, false positives, true negatives and false negatives

positives, *true negatives* and *false negatives*. How they are defined can be seen in table 3.1. Data points that are correctly classified are then either *true positives* or *true negatives*.

Precision is the percentage of the results found by the system that are correct [33]. Recall is the percentage of correct results in the dataset that are found by the system. Precision and recall are defined as follows by Van Rijsbergen [31]:

$$Precision = \frac{tp}{tp + fp} \quad (3.34)$$

$$Recall = \frac{tp}{tp + fn} \quad (3.35)$$

F-score is the harmonic mean between recall and precision, and is defined as [33]:

$$F = \frac{2 * precision * recall}{precision + recall} \quad (3.36)$$

Another metric that is used for evaluation is *accuracy*, which is the percentage of prediction that matches the actual labels. Accuracy is defined as follows:

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn} \quad (3.37)$$

These metrics are designed to work when there are two labels. Given several labels, there are different average methods used to get a score for the system. Two of these are *micro* and *macro* averaging.

The number of true positives, false positives, true negatives and false negatives for an instance λ are here denoted as tp_λ , fp_λ , tn_λ and fn_λ . A binary evaluation measure on these is denoted as $B(tp_\lambda, fp_\lambda, tn_\lambda, fn_\lambda)$. Micro-average work by summing the individual true positives, false positives, true negatives and false negatives [42]. Then you use the sums to obtain the final score. Using the notation established above, the micro average is defined as [42]:

$$B_{micro} = B\left(\sum_{\lambda=1}^k tp_\lambda, \sum_{\lambda=1}^k fp_\lambda, \sum_{\lambda=1}^k tn_\lambda, \sum_{\lambda=1}^k fn_\lambda\right) \quad (3.38)$$

Macro-average on the other hand works by first calculating the binary measure, and then taking the average of all of them [42]. It is defined as:

$$B_{macro} = \frac{1}{k} \sum_{\lambda=1}^k B(tp_\lambda, fp_\lambda, tn_\lambda, fn_\lambda) \quad (3.39)$$

It is worth noting that for some measures, such as Accuracy, the result of the two averaging approaches is the same. However, it differs for *recall* and *precision*, and therefore also the *F1-score* [42].

Perplexity can be used in order to compare different probabilistic models. It is a measurement that determines how good a models predictions are, where a lower score means that the model is better at predicting. By evaluating the perplexity on a test set, it will give an indication of how well the model will generalize [7]. Perplexity for a set of M documents, on a dataset D is:

$$perplexity(D) = \exp \left\{ - \frac{\sum_{i=1}^M \log p(x_i)}{M} \right\} \quad (3.40)$$

3.5 Related Work

Using topic modeling for various clinical applications has been done before. Topic modeling have been a popular approach for this purpose since clinical data often is in the form of free text. The resulting topics can also be interpreted by humans, which allow doctors to get more insight into the system. Sarioglu et al. [34] to represent clinical reports with topic vectors in order to classify them. Chan et al. [10] used topic models to analyze patient records and clinical reports from cancer patients. In their paper, they found relationships between the content of the notes on the patients, with the data that was available on the patients' genetic mutations. The interpretability of topics generated from an LDA model was studied on clinical reports by Arnold et al [3]. They evaluated how interpretable topics were based on how many topics the model used.

Active learning has been researched in text classification with different approaches. They can be seen as two categories: searching through the hypothesis space by using the uncertainty of a model, or by exploiting the structure of the data through clustering [13].

One of the common baselines for active learning is uncertainty sampling [23]. That simply queries the label for the data point the model is most uncertain about. In [13] hierarchical clustering is used in an active learning system. The labels are queried from clusters where there is a lot of uncertainty when it comes to the majority label. By pruning the tree of clusters while querying for labels the goal is to obtain a pruning where each node mostly contains one label.

In [28] they also take advantage of a clustering to select the samples to be labeled in a two-class environment. They use that the data points closest to the centroids are the most important ones, and that most data points in one cluster have the same label. What this

approach has in common with a lot of the current research is that it is treating single-label or binary classification problems, which cannot be directly applied to a multi-label scenario.

Research in [9] is dealing with the multi-label problem. That is the paper that developed the *binary version space minimization* strategy that is described in section 3.3. It simply takes the instance with the smallest margin among the binary classifiers, using the binary relevance scheme. The MMC strategy [44] that is described in section 3.3, and the adaptive active learning strategy [24] in section 3.3 are also techniques for managing the multi-label problem. MMC tries to find the greatest reduction for the estimated loss. While adaptive active learning combines an uncertainty measure with a measure of how the label cardinality differs. Singh et al. [38] is another multi-label active learning approach that simply takes the minimum average of the margin among the classifiers for a data point. For image classification, there has been some methods develop, for example [25, 29]. In [25] the goal is to, after making predictions, selecting the sample with the biggest mean loss. However experiments have shown that this is not as suitable for text classification [44]. In [29], the approach is to use pairs of labels and samples to present to the annotator, and the aim is to minimize the Bayesian classification error. Due to the fact that labeling for text classification is more time consuming than image classification, since you have to read an entire text, this approach is not suitable for text classification [44].

Active learning has been used to deal with the problem of imbalanced datasets before. In a binary classification setting, Ertekin et al. [14] used uncertainty sampling with SVMs to get a more balanced dataset. In [5], Attenberg et al. uses density based active learning to improve the class balance. However, it does not attempt to apply this in a multi-label setting.



4 Method

The task of making a better system for labeling clinical reports was approached with different text mining techniques, support vector machines and three learning querying strategies. At first, the framework and tools used in the system are described, followed by a description of the provided dataset. Finally, the experiments used to answer the research questions are presented.

4.1 Frameworks, Tools and Implementation

The entire system was written in Python. The motivation behind this choice was mainly that, when it comes to machine learning and text mining, most of the existing infrastructure at Sectra is using Python. This, in combination with the fact that there exists several tools for these purposes in Python, such as *numpy*¹, *nltk*², *scikit-learn*³ and *gensim*⁴. Most of the plotting was done using the *seaborn*⁵ and *bokeh*⁶ libraries. *pyLDAvis*⁷ was used for some additional visualization purposes with regards to topic models.

However, when it comes to active learning, there does not seem to be a proven mainstream library that contains a set of readily available algorithms. In order to achieve better integration between the active learning system and the existing infrastructure at Sectra, as well as making adaptations such as the number of items queried at each iteration, an active learning framework was written from scratch. The ground for this framework were the algorithms presented in Section 3.3.

This framework consisted of three modules, called *model*, *dataset*, and *query strategy*. The model is a wrapper around different machine learning models. By providing an interface for a distance or certainty measure, any underlying model able to provide such an interface can be incorporated. For accessing the data pool a dataset wrapper was written, with an interface for accessing the labeled and unlabeled pools. Putting this in its own module opens up the

¹Numpy, <http://www.numpy.org/>

²Natural Language Toolkit, <https://www.nltk.org/>

³scikit-learn, <http://scikit-learn.org/stable/>

⁴Gensim, <https://radimrehurek.com/gensim/>

⁵Seaborn, <https://seaborn.pydata.org/>

⁶Bokeh, <https://bokeh.pydata.org/en/latest/>

⁷pyLDAvis, <https://github.com/bmabey/pyLDAvis>

possibility for using several different storage solutions, such as a database or files. The query strategy module contains the different active learning algorithms for selecting what sample to label next.

4.2 Datasets

Two different datasets were used in this thesis. They were the dataset of clinical reports provided by Sectra, as well as Reuters-21578⁸. The latter was used in order to be able to simulate a multi-label labeling process. Before being integrated into Sectra's system, the different strategies needed to be evaluated from an objective point of view, so that any tradeoffs were known beforehand. Since the vast majority of the dataset from Sectra was unlabeled, this could not effectively be done using only that.

The set of reports provided by Sectra contained 1068904 different entries, where 493 were initially labeled. The entries were spread out over several files and stored in the JSON format. However, those labels were subject to change, so they were mainly used to see if there was a correlation between the labels and clusters during the exploration phase. A sample report can be seen in Figure 4.1. The fields include:

1. **ExamId**: The ID of the exam.
2. **ReportText**: The text for the report written by the physician after the examination.
3. **Anamnesis**: The patient's medical history.
4. **PatientAlert**: Anything special about the patient.
5. **ExamComment**: Comments regarding the performed examination.
6. **Cancelled**: Whether or not the examination was Cancelled.
7. **ExamName**: Name of the exam.
8. **ExamCode**: Code for the exam.
9. **PatientSex**: The sex of the patient.
10. **PatientAge**: Age of the patient. This field is truncated if it is above 90 years.
11. **Urgent**: If the examination is urgent or not.
12. **Pharma**: List of administrated pharmaceuticals.

The work was mainly concerned with the ReportText field, since it contains the response to the result of the examination. But for the complete Active Learning system the Anamnesis was used as well. The labels that were initially assigned to these reports were: "Blödning", "Infektion", "Metabol", "Tumör", "Cysta", "Missbildning", "Syndrom", "Demens", "Hydrocefalus", "Infarkt", "Kärlsjukdom", "Trauma", "Systemsjukdom", "Inklämmning" and "Normal".

The distribution of labels among these initially labeled reports can be seen in Figure 4.2. Note that this is only a count of the individual labels, and the multi-label nature of the labeling is not taken into account in the plot.

The Reuters-21578 newswire dataset is widely used when it comes to text classification research, and provides a good multi-label benchmark that can be used to compare how well certain techniques perform to other papers. All experiments used the *ModApte* split of the dataset, which is commonly used and readily available. It splits the dataset into a predefined set of training and test documents, containing 7.769 and 3.019 entries respectively. This split

⁸Reuters-21578, <https://archive.ics.uci.edu/ml/datasets/reuters-21578+text+categorization+collection>


```

{
  'ExamId':      3302250,
  'ReportText':  '[NUM-SEQ]
                  Craniell datortomografi utan och med
                  intravenös kontrast:

                  Frontalt på höger sida finns ett c:a
                  4 x 3 cm stort lågattenuerande område, tolkas
                  representera rest efter genomgången
                  parenchymskada, troligen äldre kontusionsblödning.
                  Subcorticalt på ömse sidor om centralfåran på höger
                  sida finns ett några cm-stort lågattenuerande
                  område som kan vara ischemiskt. Lätt sänkt
                  attenuering av vit substans periventrikulärt
                  förenlig med leuko-araios av degenerativ natur. För
                  åldern normalstora ventriklar. Corticala sulci upp
                  mot konvexiteten är något smalare än förväntat
                  för åldern.

                  Någon tumörsuspekt förändring påvisas ej.',
  'Canceled':    False,
  'Question':    'Förändring vä temporalt?',
  'PatientAlert': 'Hepatit C-positiv.',
  'ExamComment': 'Alla kontrastfrågor: UA mnn',
  'ExamName':    'DT hjärna utan och med iv kontrast',
  'Anamnesis':   'Pat med skalltrauma på 60-talet. Kommer nu med nattliga
                  från-varoattacker. Skrikigt beteende som tolkats som
                  epilepsi. CT är aldrig gjort. HEPATIT C-positiv."',
  'ExamCode':    '81081',
  'PatientSex':  'MALE',
  'PatientAge':  59,
  'Urgent':      0,
  'Pharma':      [{"ExamId": 3200240, "Units": "100 ml",
                  "Pharma": "Omnipaque Inj.lösn 300 mg I/ml"}]
}

```

Figure 4.1: A sample report from the dataset provided by Sectra

contains a subset of the categories, specifically 90 different ones. Since the clinical dataset from Sectra only contained 15 different categories, this would not mirror that very well, so instead the 15 most common categories of those were taken out. The distribution of the top 15 Reuters-21578 categories can be seen in Figure 4.3. After filtering out the documents not labeled with any of the top 15 categories, there were 6880 documents left in the training set, and 2646 in the test set.

4.3 Pre-Processing and Text Representation

Before the data was used in the conducted experiments, several pre-processing steps were applied in order to clean the dataset and make it easier to work with. The steps were:

1. The first step was to extract the fields of interest. For the exploratory phase, and for the use of active learning techniques these were “ReportText” and “Anamnesis”. When it came to filter out invalid reports, the “ReportText” was the only field of concern. It

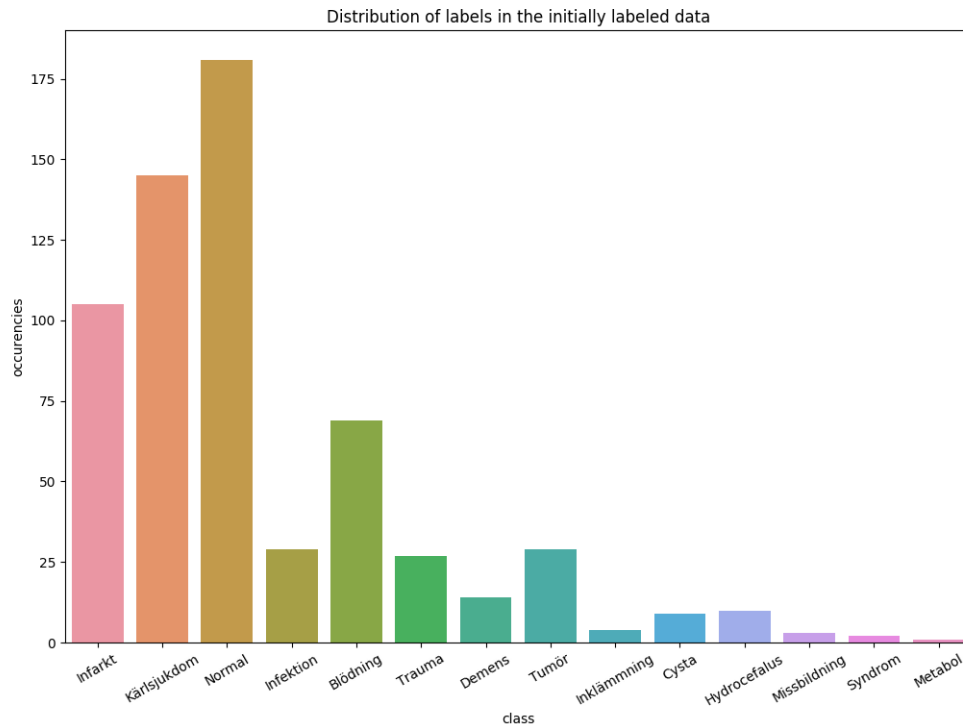


Figure 4.2: The distribution over the labels in the initial set of labeled data provided by Sectra

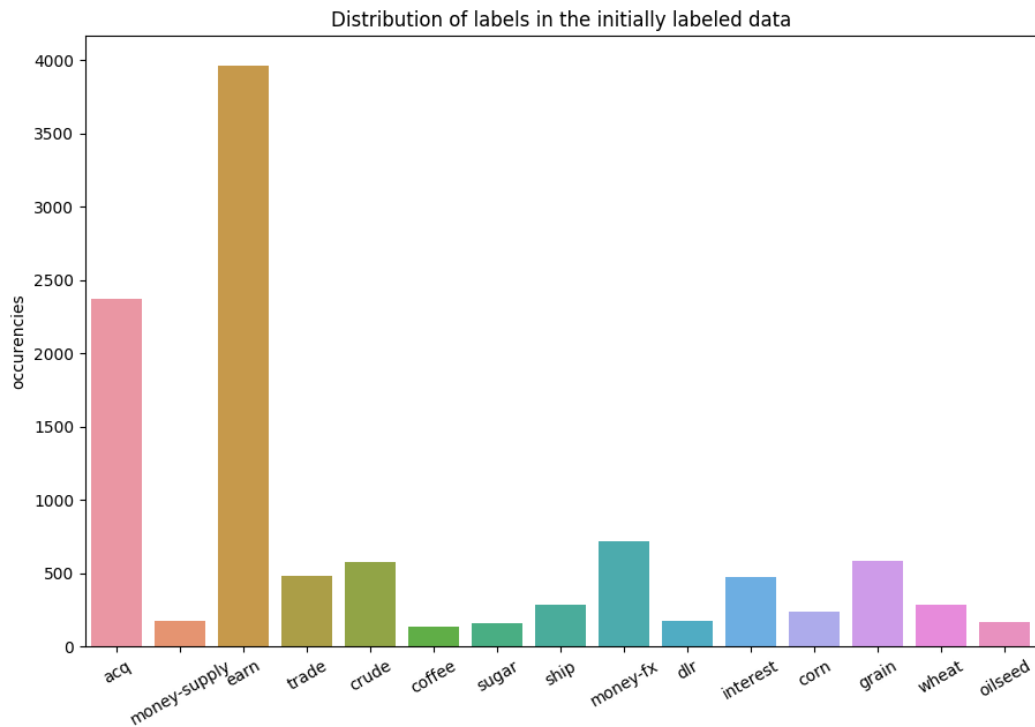


Figure 4.3: The distribution over the labels in the Reuters data

describes the results of the examination, and therefore if there was an examination at all.

2. White space and punctuation were stripped from the data.
3. All words were transformed into lowercase.
4. The most common words as well as very infrequent words were both filtered out. Specifically, words occurring in less than 1% of the documents, and words occurring in more than 90% were removed. The idea behind this is that these words would not contribute to differentiating different classes of documents. Removing of both frequent and infrequent words is commonly done when working with text and has been done in the context of classification, active learning or topic modeling before [40, 7, 9, 34].
5. A list of identified common stopwords were removed as well. This list of words was based on the Swedish nltk stopwords list. After iterating over the dataset words that occurred frequently but were not considered to be very informative for the models were identified. The list of stopwords was then extended to incorporate these words as dataset-specific information. For example, this included names of the doctors that had written the report. By removing names of doctors the idea is to make the system more applicable to new reports, written by other doctors.
6. Accents from the words were removed.
7. The text was tokenized and then stemmed using the swedish Porter2 stemmer ⁹.

Most of these steps have been performed in previous research dealing with text analysis in the form of classification or active learning [40, 7, 9, 34].

After transforming the text into a sequence of tokens, the final step before using it with the models was to create a representation that would be beneficial to work with. The representation chosen was bag of words, i.e. a matrix of tokens count. Each document is represented by the counts of each token, disregarding the order of the tokens. In order to get some positional information into the representation additional tokens are stored. The additional tokens are bigrams, which are pairs of tokens (i.e. processed words). By storing the frequency of how often such a pair occurs in the document, alongside the regular one word tokens, some positional information is retained.

4.4 Exploratory Study

For the exploratory study we used the representation described in Section 4.3, but without the bigrams. The main goal of this phase was to get to know and to better understand the dataset. A part of this goal was to go through the fields for the different reports to see how they worked and what values could be expected. In order to visualize the data in a 2D plot, t-distributed stochastic neighbor embedding (t-SNE) was used. It is a dimensionality reduction technique, that is able to transform high dimensional data into two dimensions, trying to retain as much variance as possible.

The first step was to fit an LDA model to the data. For the purpose of exploring the data, the number of topics were chosen to be a 100, with the hope of it not resulting in too granular topics that would be hard to manually analyze. A 100 topics is what is in the middle of the range Chang et al. [11] used for evaluating the interpretability of topics. Since the purpose of this is only to explore the data, it seemed like a sufficient starting point. The data points were plotted in a 2D plot after reducing the dimensions using t-SNE. Each data point was colored based on some trait that the specific data point had. For the exploratory study, the topic with the highest probability for a given data point was used to determine the color. A plot of this

⁹Swedish Porter Stemmer, <http://snowball.tartarus.org/algorithms/swedish/stemmer.html>

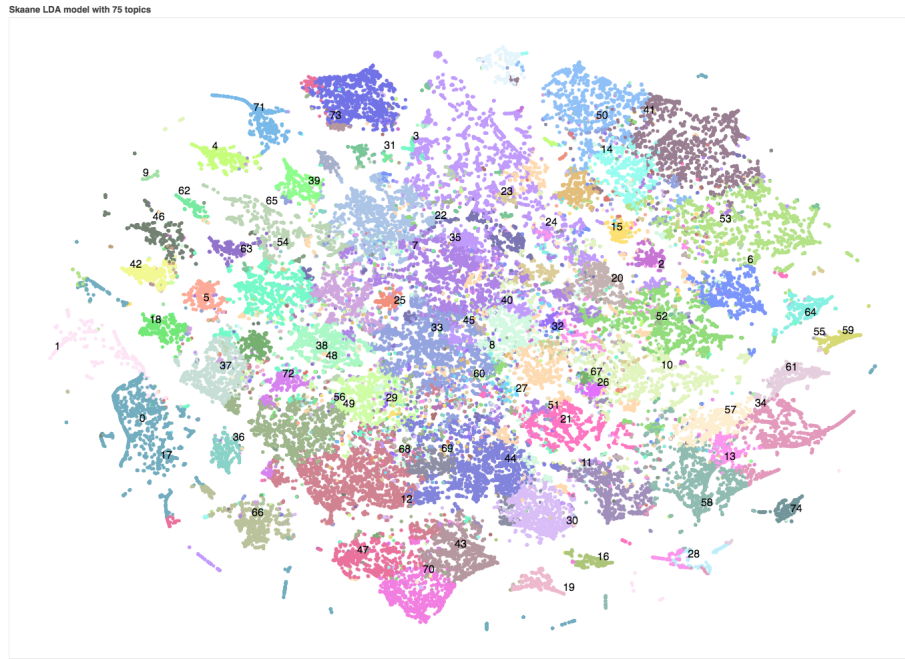
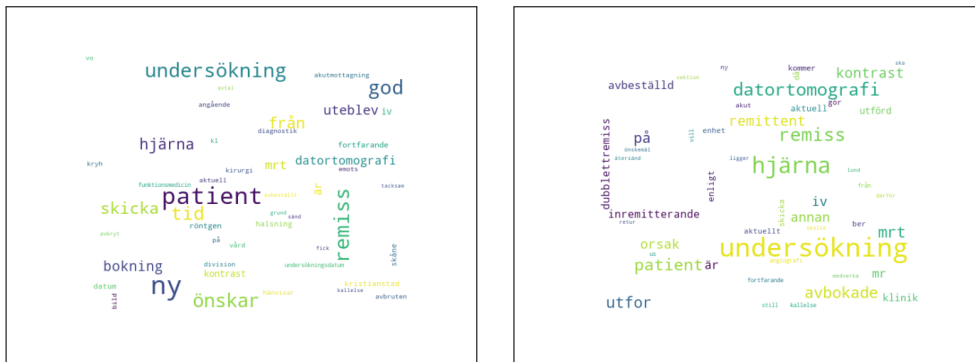


Figure 4.4: A 2D plot of the text data, where each point is colored by topic with the highest probability



(a) A wordcloud over the words occurring in topic 1 for an 75 topic LDA model.

(b) A wordcloud over the words occurring in topic 17 for an 75 topic LDA model

Figure 4.5: Wordclouds for a 75 topic LDA model

can be seen in Figure 4.4. Although it might be hard to interpret as a 2D plot in this report, the bokeh library allowed for the generation an interactive plot. Hovering over each data point would show the content of the report and the topics assigned to it, making it a convenient way to explore the data and the generated topics.

Samples of the generated topics can be seen in Figure 4.5. Another way to visualize the topics for inspection is using the techniques described by Sievert et al. [37]. They propose a *relevance* measure where the probability for a certain term within a topic is weighted against how common that topic is in the entire corpus. The interactive interface provided by pyLDAvis can be seen in Figure 4.6.

A word2vec model was used on the entire dataset to evaluate see the relationship between terms and find possible synonyms. In order to find synonyms, all words in the dataset that had a similarity over 95% were manually inspected. In addition to finding synonyms, this model was used to identify names and other identifiers in the reports. They would come

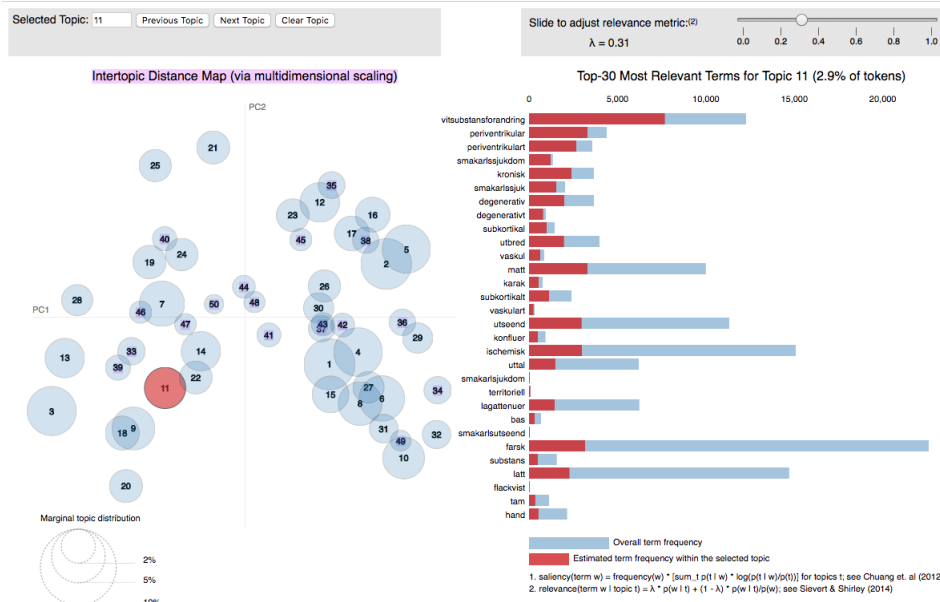


Figure 4.6: A way to visualize and analyze topics based on their relevance and frequency

up as similar entities from the model. Accomplishing this was done by exploring the data through an interactive plot, after using t-SNE to reduce the number of dimensions.

4.5 Experiments

In this section the experiments used to answer the research questions are described. There is one experiment designed for each question. The first experiment is trying to identify reports that are deemed to be invalid, the second one is a study to see what are the alternatives to labeling data points at random, as well as a comparison of how well they perform according to certain metrics. Finally, the third one aims to compare how the distribution of labels in the labeled dataset compared between the different methods.

Filter Out Invalid Clinical Reports Using Topic Models and Clustering

The task was to evaluate how well topic modeling and clustering could be used to filter out invalid reports. Invalid reports are considered to be reports that describe a situation where an examination never took place. This can be because of a deceased patient, a patient being moved to another hospital, a patient did not show up or for some reason did not want to go through with the examination.

Different topic model and cluster sizes were tested for this purpose. The topic model used was Latent Dirichlet Allocation, and the clustering algorithm used was k-means. The k-means clustering used the latent topic vectors produced by the topic model as representation when finding the clusters. As described in Section 3.1, in order to use the LDA model the number of topics has to be selected. The same applies to the number of clusters with k-means, which is described in Section 3.1. Selecting the best model was done by evaluating different number of topics k_{LDA} and the number of clusters k_c . The topic models were evaluated with k_{LDA} set to 25, 50, 75, 100, 150, 200. Evaluating the clusters was done by combining the different topic models with the different cluster sizes to find the best combination. The different combinations can be seen in Table 4.1.

ID	k_{LDA}	k_c	ID	k_{LDA}	k_c
1	25	25	19	100	25
2	25	50	20	100	50
3	25	75	21	100	75
4	25	100	22	100	100
5	25	150	23	100	150
6	25	200	24	100	200
7	50	25	25	150	25
8	50	50	26	150	50
9	50	75	27	150	75
10	50	100	28	150	100
11	50	150	29	150	150
12	50	200	30	150	200
13	75	25	31	200	25
14	75	50	32	200	50
15	75	75	33	200	75
16	75	100	34	200	100
17	75	150	35	200	150
18	75	200	36	200	200

Table 4.1: The different combination of topic model/k-means clusters that were evaluated.

In order to evaluate how well they performed on the clinical data, a set of reports had to be marked as valid/invalid. This was done by creating a script that presented a report to the user, and allowed it to be marked as valid or invalid. 5358 reports were labeled, in order to obtain 500 invalid reports. The labels were skewed, containing only 500 invalid reports, which makes up 9.3% of the labels. In order to get a more balanced dataset, 4358 of the valid reports were dismissed at random.

80 000 reports were used in the experiment, and they were selected at random. The reason for only using a subset is that the number of reports available would be too big to use in the final active-learning system due to time constraints. The models were fitted on 72 000, or 90%, of these reports, and the additional 10% were used as a held-out set to evaluate the models. To determine which topic model that should be used in the filtering of invalid reports, the perplexity was compared and the models with the lowest perplexity was chosen. Perplexity is used in the original LDA paper by Blei et al. [7] to compare different number of topics. Hofmann used it to evaluate the pLSI based topic models [17].

In order to make the models able to separate the invalid reports from the valid reports they had to be manually analyzed. They are both unsupervised methods and were therefore not fitted with a specific target. Approaching this in a way that would result in the models overfitted to the analyzed data could be hard since they are manually analyzed. To reduce the bias in the evaluation, the labeled reports were split into a training set and validation set, containing 80% and 20% of the reports respectively. The training set was used to be analyze and identify important topics, whereas the validation set was used to evaluate how well the model performed.

First, the LDA model was analyzed. This was done by inspecting the topics in the same way that was done in the exploratory study, Section 4.4. Based on the distribution of the most likely topics for the invalid reports in the training set, topics with a high indication of a report being invalid was selected for further analysis. A combination of these topics together with the length of the report text, as well as the number of topics assigned with a high probability to a report were used to determine whether or not a report was invalid or not. After some analysis, the number of topics that were assigned to a report with the probability of over 10%, referred to as prominent topics, was used as well.

Filtering out these reports with k-means results in a simpler method. K-means does not give any probability for its clusters, so filtering by the clusters must be done as a binary decision. Length restrictions were applied as before. The specific cluster that contained the most invalid reports was then used to determine if reports were invalid or not.

After some initial experiments with the topics generated it was clear that the topic vectors contained patterns ripe for exploitation. A topic vector is a vector with the topic probabilities for a document. The patterns are clear enough to motivate the manual identification of topics that are important to differentiate between invalid and valid reports. In order to compare this result with some more objective baseline, a logistic regression classifier was fitted and evaluated as well. A set of reports were already labeled with “invalid” or “valid” in order to evaluate the manual interpretation approach. This set was therefore used to fit the classifier. The topic vectors were used as features, and the targets were the labels indicating if a report is valid or not.

Both of these approaches were evaluated using four different metrics, recall, precision, F_1 -measure and accuracy. All of which are described in Section 3.4.

Alternatives to Label Reports at Random

This was done by first doing a literature study, and then exploring the relation between the initial set of labeled data with the structure of the data through clustering and topic analysis. At first, the labeled data was transformed using the LDA and k-means models. After that, they were plotted in the same 2D space as before. The color of the labeled data was set based on the first label, after an instance’s labels had been sorted.

Just as before this was an interactive plot, hovering over the data points revealed the report as well as the topics and the cluster assigned to the data point. The goal of this was to see if there existed a relationship between the topics/clusters and the labeled assigned to the data point. Even if the labels are not the same in the final system, knowledge of an existing relationship might still be exploited even if the specific labels change. Based on multi-label nature of the data and the results of this plot, active learning approaches were researched, with the goal of identifying methods that would be applicable in a multi-label setting. The research touched upon both methods that exploit the structure of the data, and methods that are purely uncertainty based.

After establishing the techniques that had some indication on providing a better labeling process than sampling documents at random, they were evaluated. In the end, the algorithms described in Section 3.3 were the ones selected for evaluation.

In order to provide a thorough evaluation of how well they techniques perform a set of already labeled documents was needed. For this reason, the Reuters-21578 dataset was used. The properties of the dataset, as well as a comparison between it and the clinical data provided by Sectra can be found in Section 4.2. The dataset is common in active learning research and has been used by Brinker et al. [9] and Yang et al. [44], among others. With this set of of labeled reports, a simulation could be used to compare the different strategies with different metrics. The metrics used were:

1. Accuracy
2. Micro recall
3. Macro recall
4. Micro precision
5. Macro precision
6. Micro F1-Score
7. Macro F1-Score

These are described in Section 3.4 and are frequently used to compare different active learning methods, for example by Yang et al. [44], Dasgupta et al. [13] and Li et al. [24]. In addition to these metrics, the time it took to query samples with each model was also compared.

With this dataset, the same pre-processing steps that were applied to the clinical dataset were applied to the Reuters data too. Some modifications of this includes the stopwords, instead of a curated list of words, the unmodified list of english stopwords provided by nltk was used. The main goal was to compare how the different techniques affected the labeled dataset, and how well an SVM model performed on it. Optimizing the process for the particular model and dataset was therefore not the focus of the study, but instead offering a more comprehensive comparison.

The strategies compared were:

- *Binary Version Space Minimization*: Described in Section 3.3
- *Maximum Loss Reduction with Maximum Confidence*: Described in Section 3.3
- *Adaptive Active Learning*: Described in Section 3.3

The strategies need a small initial set of labeled reports so that they can base the calculations on something. The techniques were evaluated both by selecting this initial set of points at random, as well as selecting them from the clusters generated by the k-means algorithm. Sampling from the clusters was done by iterating over the clusters and selecting an equal number of data points from each clusters. All samples selected from a given cluster were chosen randomly amongst the members of the clusters. In addition, the number of clusters selected was 25, in order to get an equal number of reports from each cluster in the different experimental settings. The topic model selected in Section 4.5 was used as input to the k-means algorithm.

Since the different models may depend on the initial samples in different ways, different initial sizes were evaluated. This is also done by Yang et al. [44]. In their paper they tried quite large initial sample sizes. Here, the sizes evaluated are: 25, 50, 100. The reason for this is that an large initial sample size would make it hard for the human annotator to see a difference in the class balance early on. The different active learning configurations that were tried is displayed in Table 4.2 In total there were 18 configurations, based upon the three different methods.

4.6 Evaluating the Label Balance

The goal with the last experiment is to evaluate how the labels in the produced labeled dataset is distributed. A set where the labels are more evenly, or uniform, distributed would be preferable. From a perspective of the person labeling, it could feel more productive not assigning the same labels most of the time. The more prosperous outcome from a balanced dataset would be that the models using the data in later stages could also benefit from this, and obtain better results.

The configurations used here are the same as in Table 4.2. Every iteration, the distribution of labels assigned are stored and analyzed. He et al. discusses the usage of ROC curves and measures such as g-means to compare multi-class imbalanced data. However, since the doctor involved at Sectra specifically requested that labels should be more uniformly distributed. The evaluation will therefore focus on measuring that instead of the models performance, which is done in 4.5. An evaluation of how the distribution progresses was done by comparing how the class imbalance is affected by the number of new samples obtained from the different methods.

Ertekin et al. [14] used a class imbalance ratio to compare how well an active learning strategy worked. This was only done for the binary case. In order to get a measure for the

ID	Active Learning Strategy	Initial Sampling	Initial Sample Size
1	BinMin	Random	25
2	BinMin	Random	50
3	BinMin	Random	100
4	BinMin	Sampled from clusters	25
5	BinMin	Sampled from clusters	50
6	BinMin	Sampled from clusters	100
7	MMC	Random	25
8	MMC	Random	50
9	MMC	Random	100
10	MMC	Sampled from clusters	25
11	MMC	Sampled from clusters	50
12	MMC	Sampled from clusters	100
13	Adaptive	Random	25
14	Adaptive	Random	50
15	Adaptive	Random	100
16	Adaptive	Sampled from clusters	25
17	Adaptive	Sampled from clusters	50
18	Adaptive	Sampled from clusters	100

Table 4.2: The different configurations of active learning strategies evaluated. BinMin stands for Binary Version Space Minimization, MMC stands for Maximum Loss Reduction with Maximum Confidence and Adaptive stands for Adaptive Active Learning.

multi-label problem, the evaluation of class imbalance was measured by the percentage of all total labels that were in the most common class, as well as in the top 3 most common classes.



5 Results

In this chapter the results are described. First, the outcome from the exploratory study is presented, followed by the different experiments. The first experiment, filtering out invalid reports, presents the evaluation of the topic model and k-means model used to filter out the reports, as well as the specific topics and clusters used in the process. In the second one, the methods considered and the decisions behind which ones that were appropriate are presented. Finally, the last section goes through the result of evaluating the different active learning techniques.

5.1 Exploratory Study

The goal with the exploratory study was to acquire a better understanding of the data, how it was structured and what kind of information might be extracted from it. Certain fields such as the cancelled field did not seem to be very reliable. Reports that clearly explained a situation where the patient had been transferred to another hospital, or for another reason not having performed an examination, still described a situation where the cancelled field was set to “false”. After further manual analysis it was clear that the vast amount of invalid reports were contained within a few topics, something that is used in the first research question. The evaluation of more concrete relationships were done within the context of that experiment, and is presented in Section 5.2.

The word2vec model produced results that allowed for synonyms to be detected. By doing this, 420 pairs were discovered. The vast majority of these were names, medical terms that (some of which the author was unable to evaluate, and therefore excluded) as well as words that are used in similar contexts, which includes opposites like “left” and “right”. Disregarding these, the synonyms and misspellings that were decided for use in the final system can be seen in Table 5.1. The original value was replaced with the new one in the final system.

In order to identify names from this word2vec model, it was plotted using the an interactive plot that allowed for exploring the data. Since names are commonly used in similar contexts, they would have similar attributes in the word embedding model. Figure 5.1 and Figure 5.2 show how this was done. Given that the names got similar coordinates in the plot, identifying the section with names allowed for identification of a lot of the names used in the reports.

Original	Replacement	Type
ordinärt	normalt	synonyms
ej	inte	synonyms
avbeställd	avbokad	synonyms
avebställd	avbokad	misspelling + synonym
belsutat	beslutat	misspelling
måttliga	lätta	synonyms
pat	patient	short
pt	patient	short
pateint	patient	misspelling
akuten	akutmottagningen	misspelling
us	undersökning	short

Table 5.1: The synonyms, misspellings and shorts found in the data that the author could with assert with confidence.

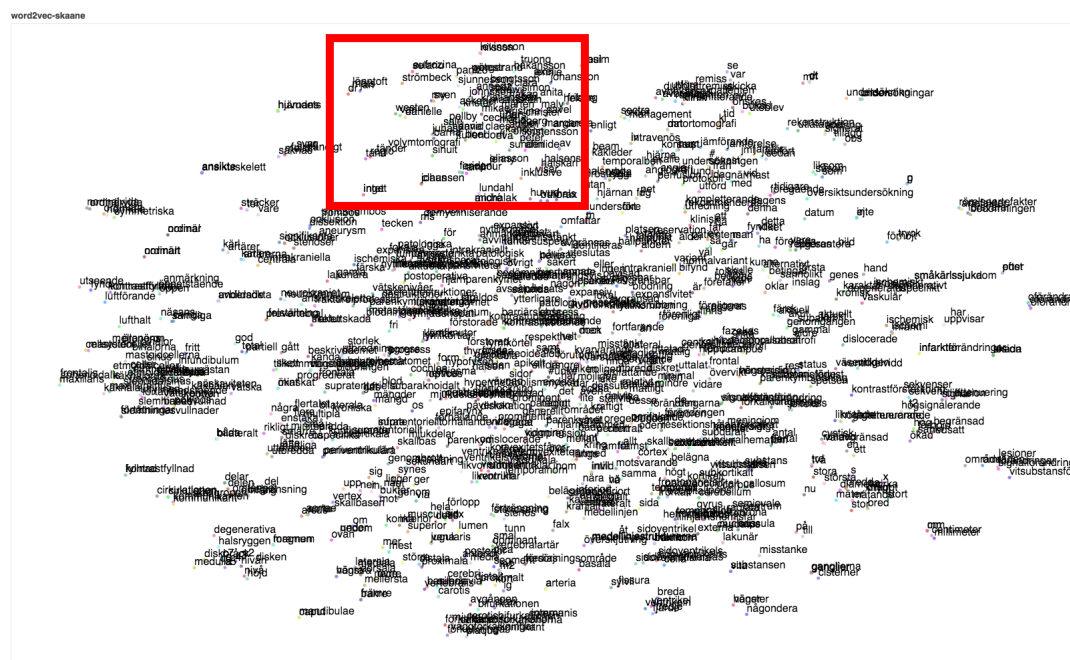


Figure 5.1: A 2D plot of the full word2vec plot. Given the amount of terms used there is a lot to analyze.

5.2 Filter Out Invalid Clinical Reports Using Topic Models and Clustering

The LDA models were evaluated by calculating the perplexity on the held-out set as described in Section 4.5. Perplexity for the evaluated models can be seen in Figure 5.3, with the lowest score being better. Based on this, the selected model was the LDA model with 75 topics.

The next step was identifying the topics that were assigned to the invalid reports. The distribution of the topics with the highest likelihood for the invalid reports in the training set can be seen in Figure 5.4. A couple of topics, 1 and 17 clearly stands out as the ones that most invalid reports gets assigned. Specifically, topic 1 had a count of 131 reports from the invalid reports in the training set, and topic 17 had 358. Some invalid reports have topics 0, 3, 9, 16, 47, 53, 58, 61, 62 as the most likely topic. 53 is the third most common, with a count of 3.

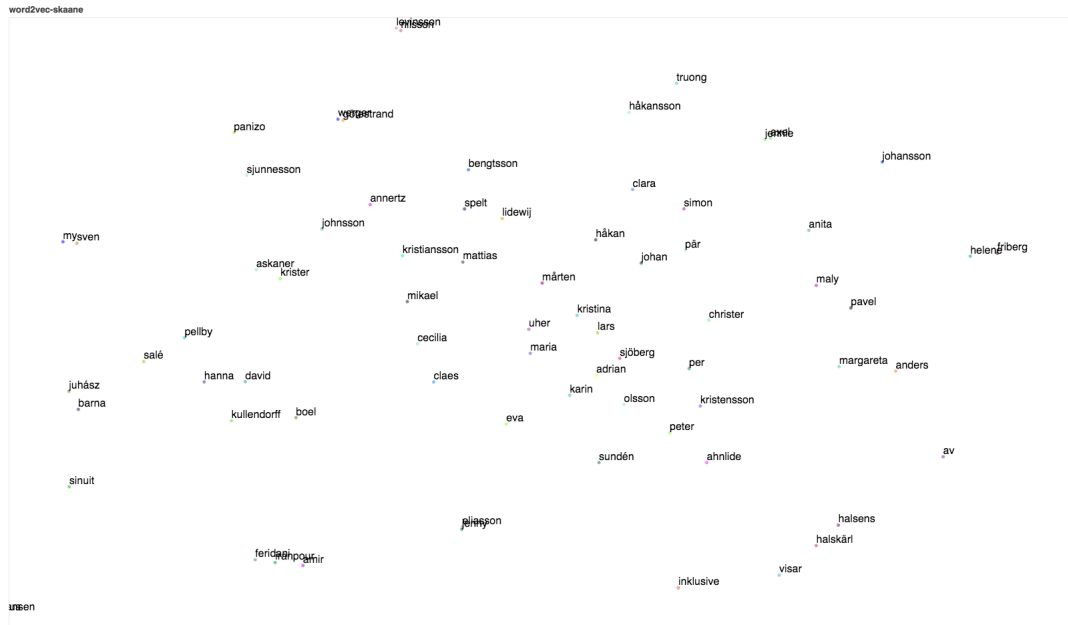


Figure 5.2: A 2D plot of the zoomed in word2vec plot. Most of the values here are names. This represents the red box in from Figure 5.1

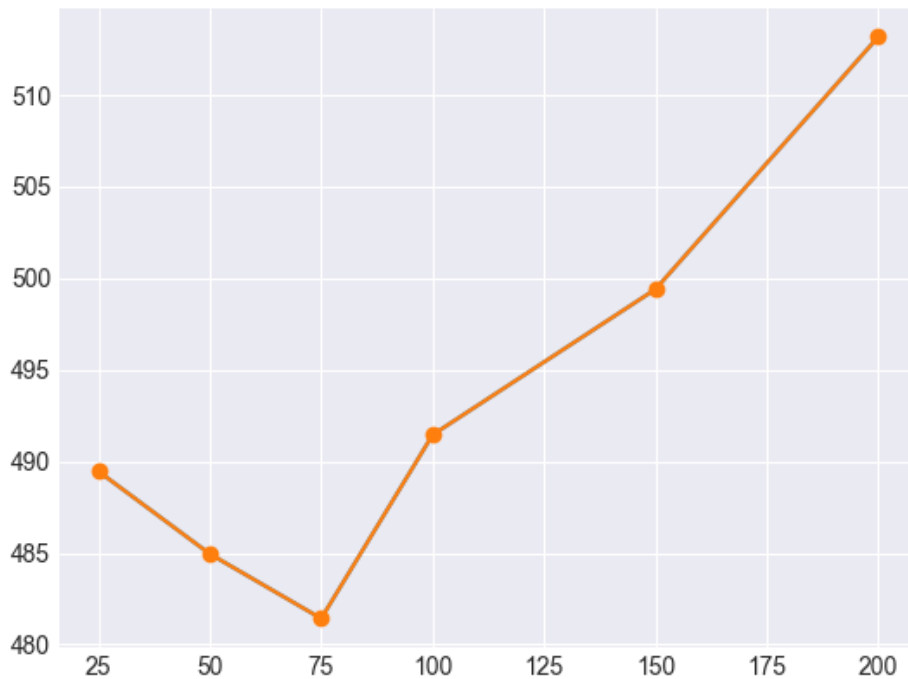


Figure 5.3: The perplexity scores for the different LDA models

The corresponding plot for the valid reports can be seen in Figure 5.5. There is a lot more variety among the most likely topics here. Topics 1 and 17 occurs very infrequently. Topic 1

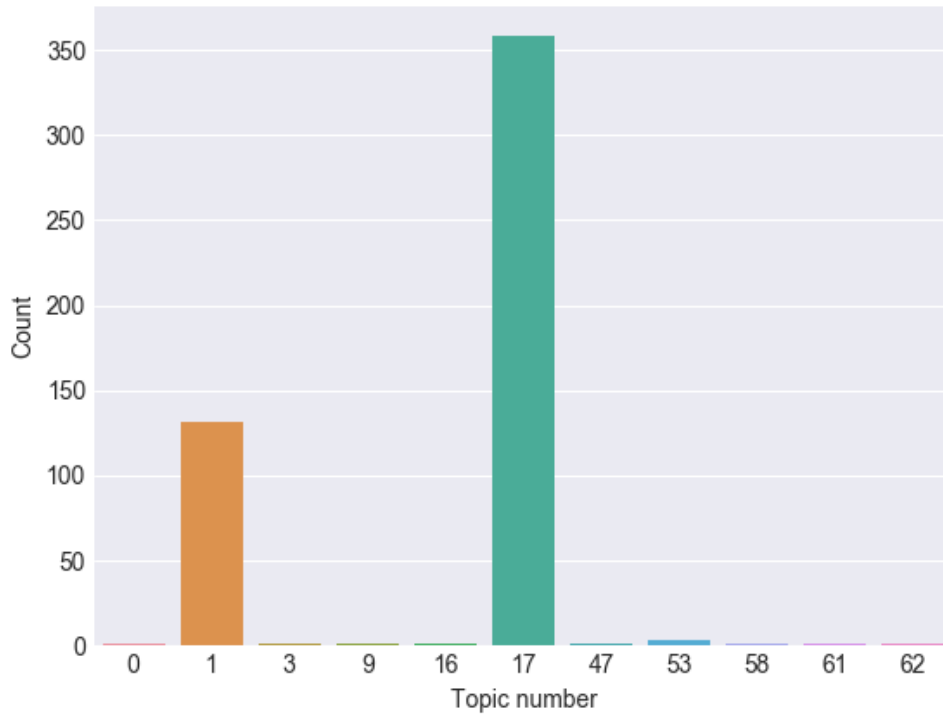


Figure 5.4: The distribution of the topics with the highest probability for the invalid reports in the training set. Note that only topics that occurred at least once are shown in the histogram.

occurs 0 times, while topic 17 occurs 2 times. The third most common topic from the invalid reports, 53, occurs 28 times. The ones having topic 17 had 4 and 8 prominent topics assigned to them.

Each topic that is assigned to a report with a probability above 10% was considered to be a prominent topic. The number of prominent topics for the invalid reports can be seen in Figure 5.6. 1, 2 and 3 number of prominent topics are the most common. And there are barely any reports having less than 6. The corresponding plot for the valid reports can be seen in Figure 5.7. In the case of valid reports, 6 is the most common number of prominent topics. The topics are a lot more spread out than in the case of the invalid reports.

Another idea was to evaluate whether or not 17 and 1 was among the most probable topics for the reports. A simple evaluation of this on the set of valid reports showed that it was not a good approach. For example, seeing if 17 or 1 had more than 10% probability returned 48 and 45 reports, respectively. Checking if both were above the threshold returned 11 reports. This threshold was tested with 1% increments from 5% to 20% without giving any good results.

Based on these findings, reports were determined to be invalid or not by the following criterion:

- Having either topic 1 or 17 as its most probable topic.
- Not having more than 6 prominent topics assigned to it.

The evaluation on the validation set can be seen in Figure 5.2. In the figure you can also see the results of the logistic regression classifier, which was fitted with the topic vectors as features and the invalid/valid labels as targets.

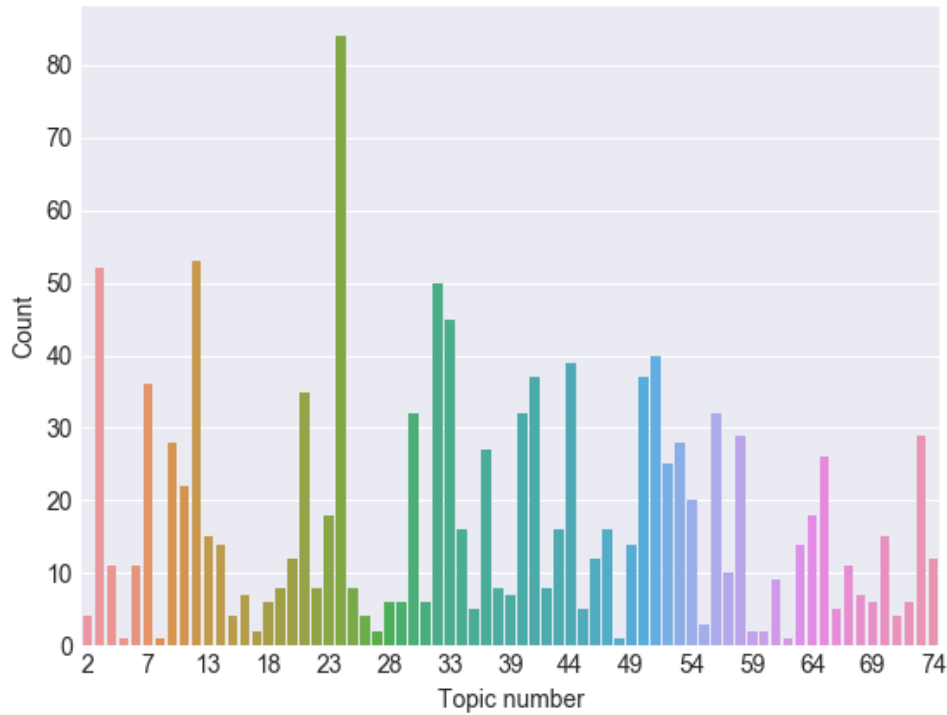


Figure 5.5: The distribution of the topics with the highest probability for the valid reports in the training set. Note that only topics that occurred at least once are shown in the histogram.

	Manual	Logistic Regression
Precision	97.2%	98.7%
Recall	100%	100%
F_1-measure	98.6%	99.4%
Accuracy	97.9%	99.1%

Table 5.2: The results of the classification of the invalid reports. The manual column represents the use of manual interpretation of the LDA model.

5.3 Alternatives to Labeling at Random

The first thing that was done with regards to this experiment was to explore and try to find a relationship between the initial set of labels and the inherent structure of the data. This was done by visualizing the LDA model again. In order to find any existing relationship between the topics and the labeled samples they were colored based on their assigned labels. Unlabeled samples were hidden from the plot. The resulting plot can be seen in Figure 5.8. From this plot it is clear that there is a grouping of labels. Certain labels are more likely to occur in documents assigned a specific topic. For example, the purple points in the lower part of the graph represents the “blödning” label, the gray labels in the middle represents “infarkt” and the light blue ones that are mostly concentrated in the bottom right corner represents “infektion”. It is clear that these labels are not evenly spread out over all the topics, but neither are they confined enough to make a mapping between topics and labels, like they were in Section 5.2.

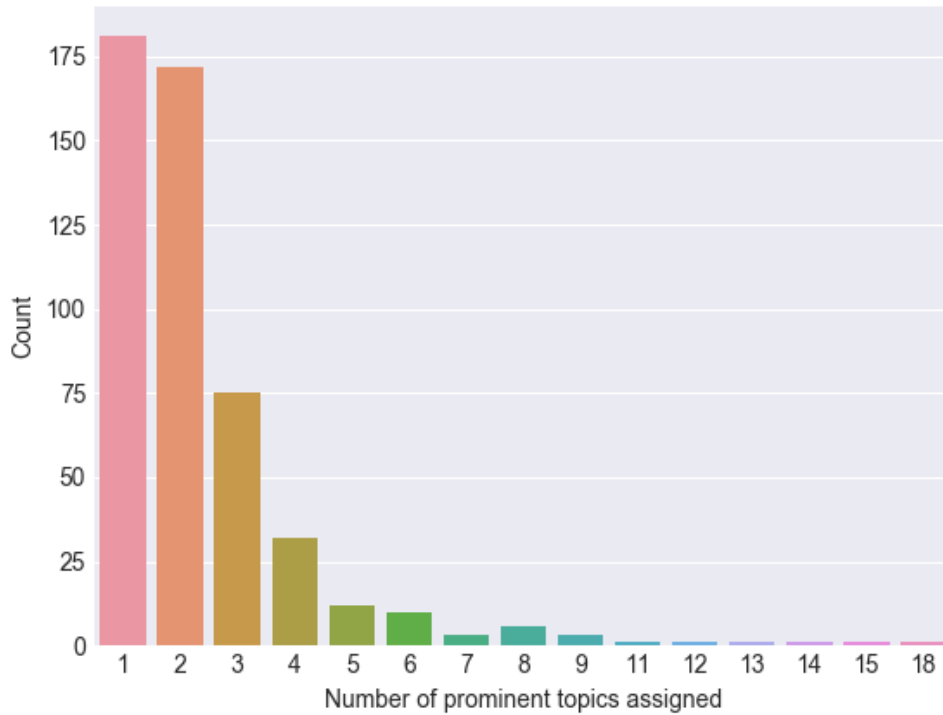


Figure 5.6: The distribution of the number of prominent topics assigned to the invalid reports in the training set.

In Figure 5.9 the counts of most likely topics for four most common categories are displayed. These categories are “infarkt”, “kärleksjukdom”, “normal” and “blödning”. From the histograms it is clear that documents of a certain category are more likely to be assigned certain topics, at least in these cases. Even though there exist a clear relationship, it is not exclusive enough to make any clear relation. The number of topics assigned and reports labeled for these 4 topics can be seen in Table 5.3.

In order to analyze these topics further, Figure 5.10 shows the different categories that has a certain topic assigned to it as the most likely one. Taking into account the information from Table 5.3, i.e. that some categories are a lot more common than others, there is not a clear enough pattern to distinguish between different categories based on the topics. This does not take the multi-label nature of the data into account. If a report has multiple labels assigned to it, both of the labels are counted separately.

Based on the knowledge that there exists a pattern, the initial goal was to find some methods that could exploit this. Some active learning approaches using different forms of clustering, such as Dasgupta et al’s approach using hierarchical clustering [13] would be good contenders. However, the method described by Dasgupta et al. is made for the single-label case with no obvious way of extending the technique into multi-label. The same applies to the density based technique suggested by Attenberg et al. [5].

Most of the active learning research seems to be focused on binary, or maybe multi-class classification. Therefore, the methods described in Section 3.3 where the ones decided on. Methods that are fully reliant on a models certainty, such as Binary Version Space Minimization [9] are tested. In addition to this, methods incorporating some information about the data in the form of label cardinality is included as well. These techniques are MMC and

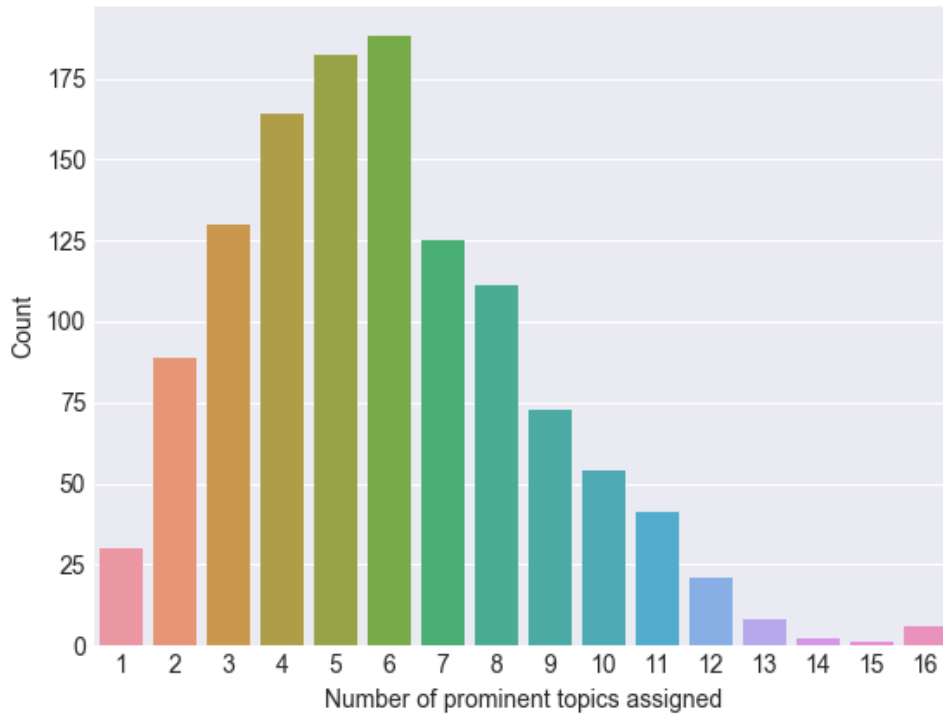


Figure 5.7: The distribution of the number of prominent topics assigned to the valid reports in the training set.

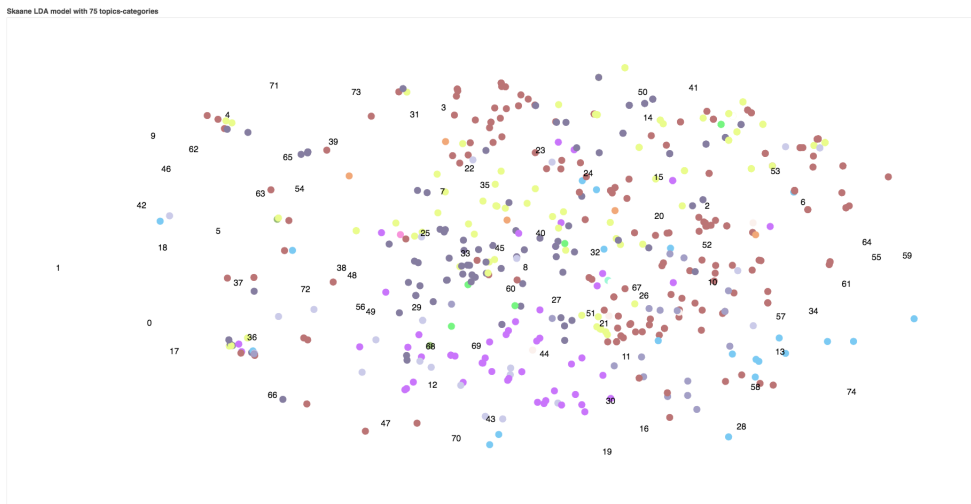


Figure 5.8: The labeled data points plotted in 2D, and colored based on the first label of the report in alphabetical order.

Adaptive Active Learning [44, 24]. An attempt to take advantage of the structure of the data is done by selecting the initial samples from different clusters, as described in Section 4.5.

The first metric for evaluation was Accuracy. A plot of this can be seen in Figure 5.11 to Figure ???. It contains the accuracy evaluations for the different number of initial sample sizes. It can easily be seen that the adaptive learner performs significantly better for the first 400

Label	No. reports	No. most likely topics
Normal	181	28
Tumör	29	15
Infarkt	105	27
Blödning	69	19
Kärlsjukdom	145	28
Hydrocefalus	10	5
Demens	14	9
Trauma	27	12
Cysta	9	7
Missbildning	3	3
Inklämmning	4	1
Infektion	29	15
Syndrom	2	2
Metabol	1	1

Table 5.3: The number of reports assigned a certain category, as well as the number of different topics assigned as the most likely one for reports with the given category.

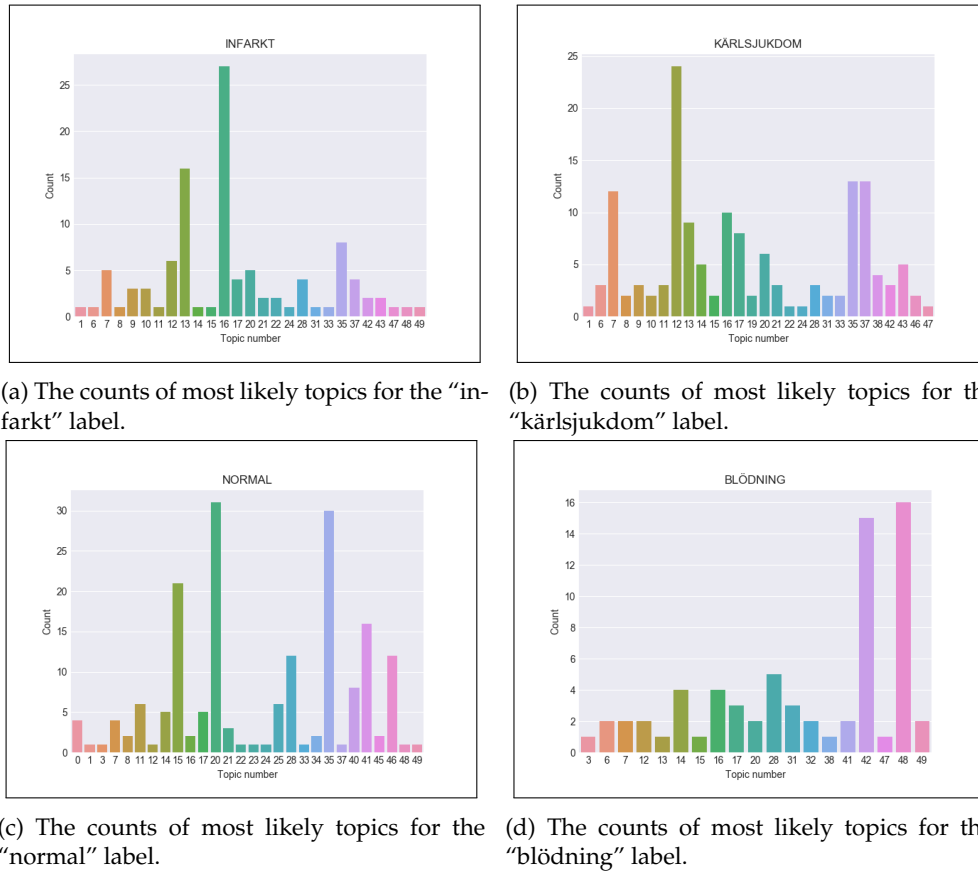


Figure 5.9: The counts of the most likely topics for the four most common categories.

added samples. However, after that MMC and Binary Version Space Minimization catches up to it, and surpasses it in terms of accuracy after around 800 labels. Both of them sustaining an significant improvement over labeling at random after round 200 labels. It takes sampling at random approximately 1600 labels to get to an accuracy of 78%. Adaptive Active Learning acquires this accuracy after around 300 labels, and BinaryMinimization and MMC after 450.

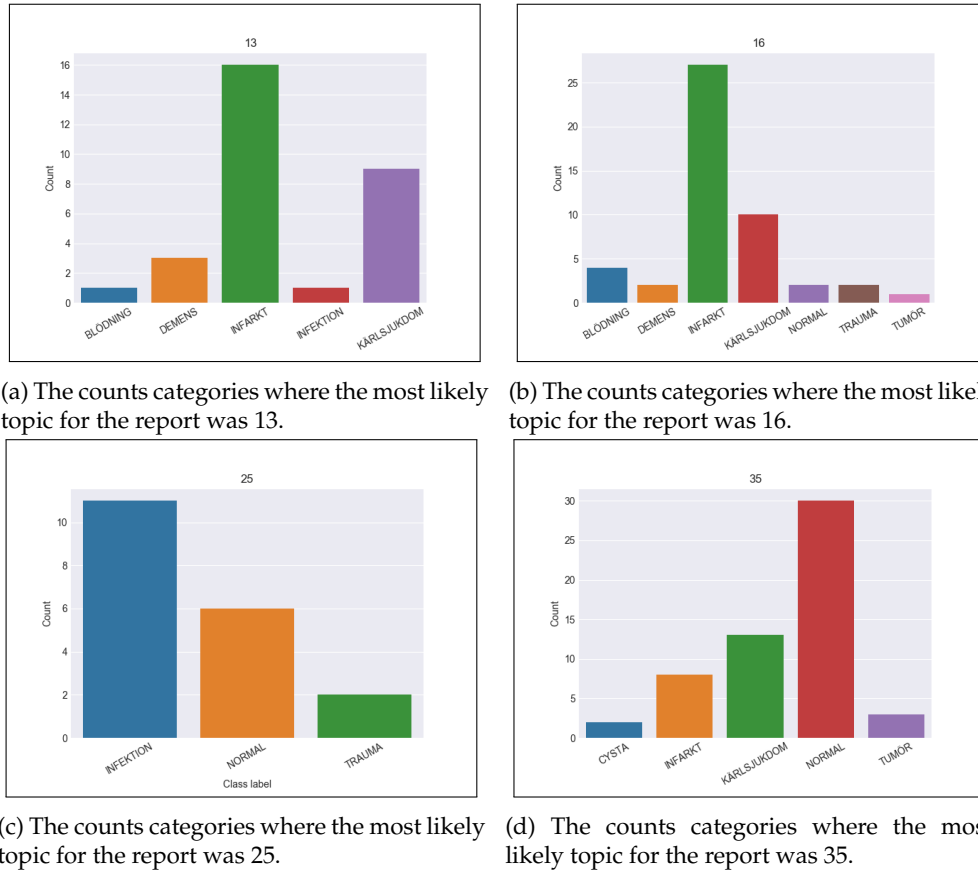


Figure 5.10: The categories of the different reports that are assigned a certain topic as the most likely one.

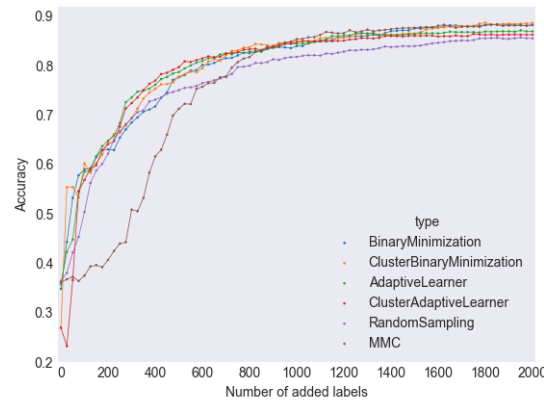


Figure 5.11: Accuracy of the models with initial sample size 25

When it comes to micro F1-Score, the results between the different methods are similar. Adaptive Active Learning performs better than the rest in the beginning, but the other active learning strategies obtain similar results after around 500 labels. One difference from the accuracy evaluation is that Binary Version Space Minimization does not quite surpass the adaptive learner in the same way. The gap to random sampling is smaller as well. The results can be seen in Figure 5.12 to Figure 5.12.

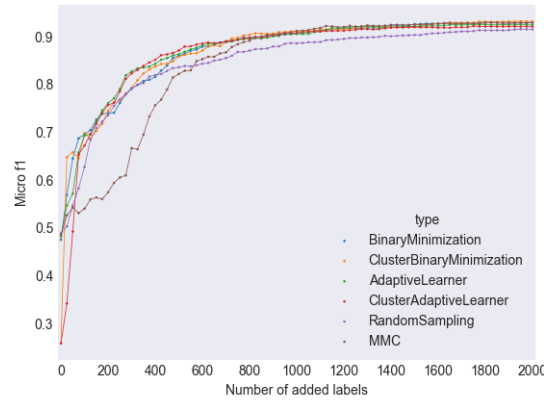


Figure 5.12: Accuracy of the models with initial sample size 25

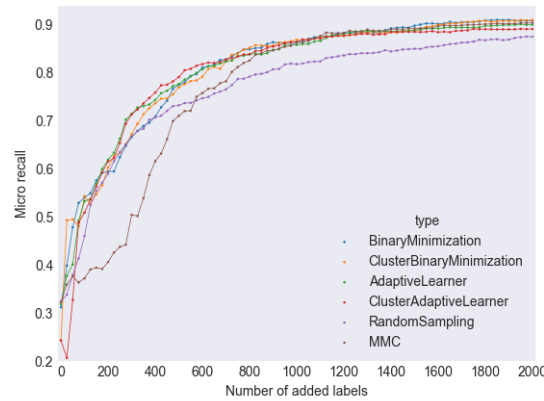


Figure 5.13: Micro recall of the models with initial sample size 25

The same applies to micro recall. The one difference is that the Adaptive Active Learner strategy is performing better than the rest for more iterations. It takes approximately 1000 added labels for the BinaryMinimization to obtain similar micro recall. This can be seen in Figure ?? . The random sampling performs considerably worse than the rest.

The micro precision of the models are very similar between all strategies. Algorithms using clustering obtains a higher precision between 200-400 added labels, but not by that much. From 200 added labels and onwards, Adaptive Active Learner is worse than Binary Version Space Minimization.

Another important aspect of the evaluation is the time it takes. This can be seen in Figure ?? . Time is constant for the random sampling, very close to 0. For the Binary Minimization, the time increases a bit every iteration, but stays under 200 seconds for the first samples. MMC gets similar results, but takes a bit longer. The Adaptive Active Learner takes far longer time than the other methods.

5.4 Evaluating the Label Balance

The last part to evaluate is how the different Active Learning techniques effected the balance of the labels. For the Reuters dataset, how the overall distribution of the labels is can be seen in Figure 4.3. After random sampling, the distribution can be seen in Figure 5.16. The

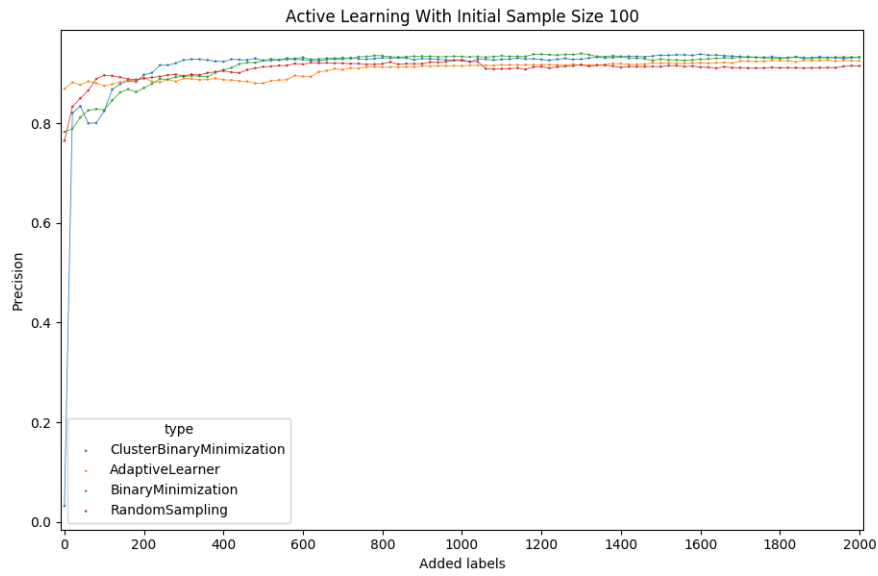


Figure 5.14: Micro precision of the models with initial sample size 25

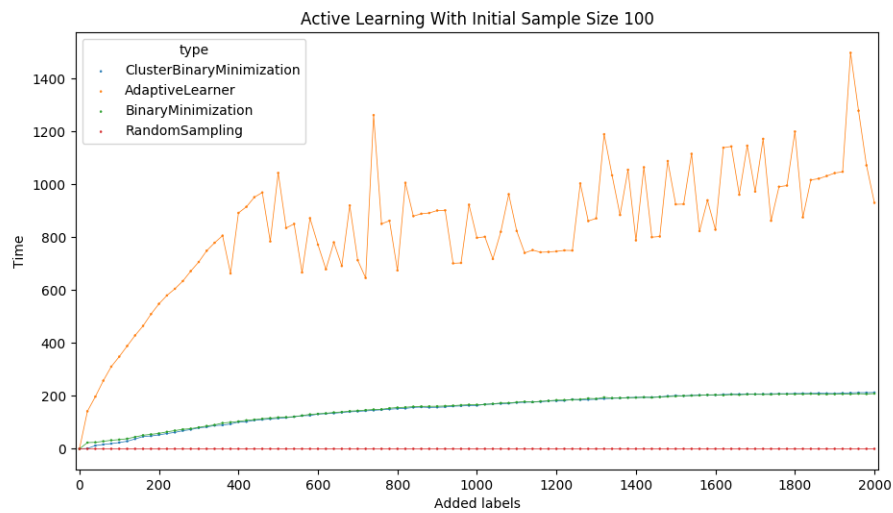
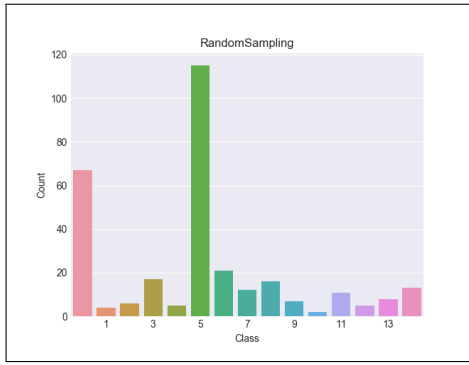
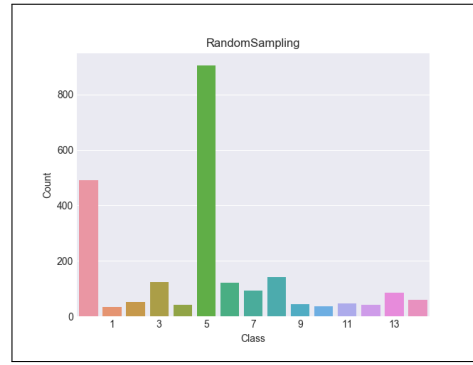


Figure 5.15: Time to query a new set of samples by the different strategies

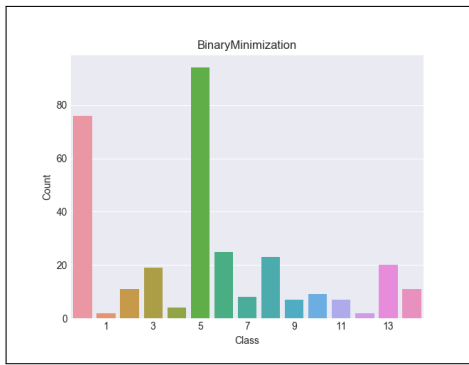


(a) The class distribution from random sampling after 200 labels

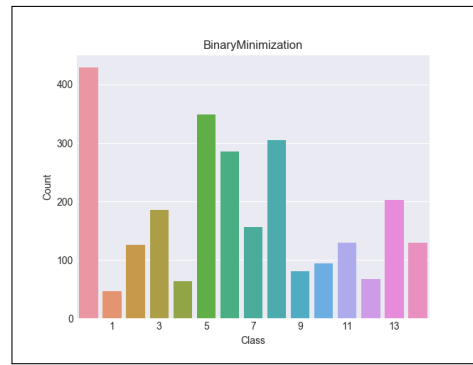


(b) The class distribution from random sampling after 2000 labels

Figure 5.16: The distribution of labels after random sampling



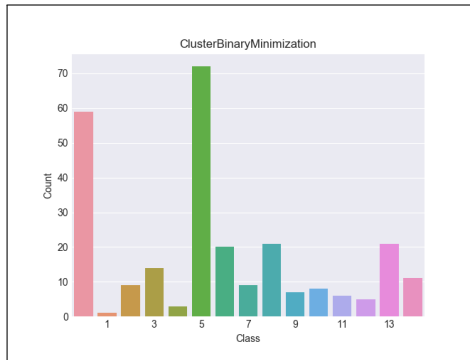
(a) The class distribution from Binary Version Space Minimization after 200 labels



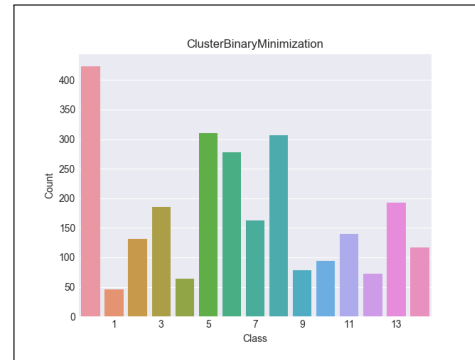
(b) The class distribution from Binary Version Space Minimization after 2000 labels

Figure 5.17: The distribution of labels after Binary Version Space Minimization

distribution after sampling with the original Binary Version Space Minimization can be seen in Figure 5.17, and with the initial samples taken from clusters in Figure 5.19. The ratio between the biggest and smallest class is here X_1 compared to X_2 for the random sampling. When the initial labeled set was selected from the clusters, the ratio was X . For MMC the distribution can be seen in Figure ??, and for Adaptive Active Learning in Figure ?. The imbalance ratio between the biggest and the smallest class here is X_1 and X_2 , respectively. Figure ?? and Figure ?? shows the distribution for the same methods, but with the initial samples taken from the clusters.

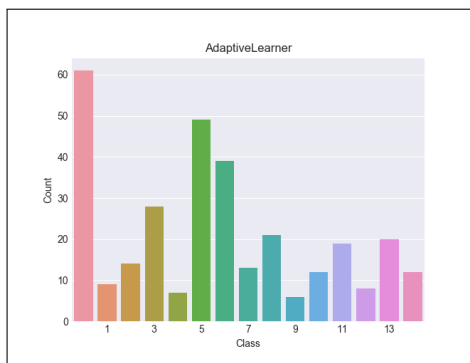


(a) The class distribution from Binary Version Space Minimization with clustering after 200 labels

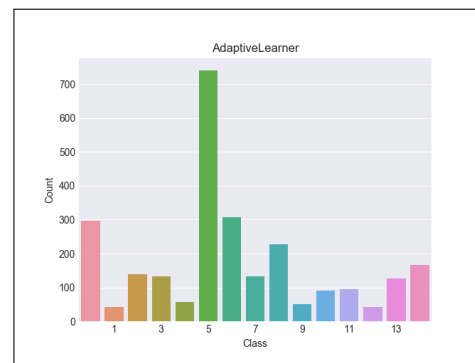


(b) The class distribution from Binary Version Space Minimization with clustering after 2000 labels

Figure 5.18: The distribution of labels after Binary Version Space Minimization with clustering



(a) The class distribution from Adaptive Active Learning after 200 labels



(b) The class distribution from Adaptive Active Learning after 2000 labels

Figure 5.19: The distribution of labels after Adaptive Active Learning



6 Discussion

6.1 Results

6.2 Method

Results that are based on a public dataset is naturally easier to reproduce. The method becomes more reliable in the sense that the same results can be expected by reproducing the concrete steps. However, the clinical dataset from Sectra is not publically available. Therefore any results that are derived from specific attributes of that dataset might not be reasonable to expect when applied in a new environment.

The comparisons of the active learning algorithms are using the Reuters dataset, which both public and a standard dataset for evaluation. Reproducing these results might therefore be an easier tasks. Any issues with the reliability in this case might come from the implementation of the presented in Chapter 3. The implementations of the Binary Version Space Minimization, Maximum Loss Reduction with Maximum Confidence and Adaptive Active Learning are based on the algorithms, but written independently for this thesis. They have not faced any public scrutiny and their correctness can not be guaranteed.

The usage and analysis of the LDA model was done in a rather way. During the exploration phase and the first experiment, the relationships and pattern found was in part a result of the authors intuition. For this reason, if another party would perform the same study they might come up with different results and find other patterns. However, the relationships found between the topics and the invalid reports are rather explicit, so the main difference between different researchers might be the specific thresholds. For example, the 10% threshold for prominent topics was chosen based on the author's experience with the dataset after seeing several iterations of topic models. By using logistic regression for comparison the subjectiveness of the thresholds can be compared with an objective baseline, which improves the reliability of the results. Any subsequent study is more likely to obtain similar results with this approach.

Another aspect that is questionable when it comes to the first experiment is the labeling of the invalid reports. This was done manually by the author, without any medical knowledge. However, the nature of the labeling is rather trivial. The medical knowledge required to understand the result of an examination is far greater than the one needed to see if an examination was performed. Which in most cases can be identified not despite of, but because

of the lack of medical information. The number of labeled reports, X , is probably sufficient to manually identify broad relationships to set some categorization rules based on the output of the LDA. However, the logistic regression based classification might have been able to achieve better results with more data to train on. This could have been done rather easily without the time constraints of this thesis work. In order to obtain a more balanced sample to train on, an active learning system could, as we have seen, be used. However, the system developed in this thesis was not finished when the labeling was done, and it is not focused on the binary classification task. Evaluating the categorization of invalid reports by accuracy, precision, recall and F1-score are fairly standard. The metrics have been used in a lot of text classification and information retrieval research [1, 6].



7 Conclusion

This chapter contains a summarization of the purpose and the research questions. To what extent has the aim been achieved, and what are the answers to the research questions?

The consequences for the target audience (and possibly for researchers and practitioners) must also be described. There should be a section on future work where ideas for continued work are described. If the conclusion chapter contains such a section, the ideas described therein must be concrete and well thought through.



Bibliography

- [1] Charu C Aggarwal and ChengXiang Zhai. “A survey of text classification algorithms”. In: *Mining text data*. Springer, 2012, pp. 163–222.
- [2] Charu C Aggarwal and ChengXiang Zhai. “A survey of text clustering algorithms”. In: *Mining text data*. Springer, 2012, pp. 77–128.
- [3] Corey W Arnold, Andrea Oh, Shawn Chen, and William Speier. “Evaluating topic model interpretability from a primary care physician perspective”. In: *Computer methods and programs in biomedicine* 124 (2016), pp. 67–75.
- [4] David Arthur and Sergei Vassilvitskii. “k-means++: The advantages of careful seeding”. In: *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics. 2007, pp. 1027–1035.
- [5] Josh Attenberg and Seyda Ertekin. “Class imbalance and active learning”. In: *imbalanced Learning: Foundations, Algorithms, and Applications* (2013), p. 101149.
- [6] Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [7] David M Blei, Andrew Y Ng, and Michael I Jordan. “Latent dirichlet allocation”. In: *Journal of machine Learning research* 3. Jan (2003), pp. 993–1022.
- [8] Matthew R Boutell, Jiebo Luo, Xipeng Shen, and Christopher M Brown. “Learning multi-label scene classification”. In: *Pattern recognition* 37.9 (2004), pp. 1757–1771.
- [9] Klaus Brinker. “On active learning in multi-label classification”. In: *From Data and Information Analysis to Knowledge Engineering*. Springer, 2006, pp. 206–213.
- [10] Katherine Redfield Chan, Xinghua Lou, Theofanis Karaletsos, Christopher Crosbie, Stuart Gardos, David Artz, and Gunnar Ratsch. “An empirical analysis of topic modeling for mining cancer clinical notes”. In: *Data Mining Workshops (ICDMW), 2013 IEEE 13th International Conference on*. IEEE. 2013, pp. 56–63.
- [11] Jonathan Chang, Sean Gerrish, Chong Wang, Jordan L Boyd-Graber, and David M Blei. “Reading tea leaves: How humans interpret topic models”. In: *Advances in neural information processing systems*. 2009, pp. 288–296.
- [12] Steven P Crain, Ke Zhou, Shuang-Hong Yang, and Hongyuan Zha. “Dimensionality reduction and topic modeling: From latent semantic indexing to latent dirichlet allocation and beyond”. In: *Mining text data*. Springer, 2012, pp. 129–161.