

NORWEGIAN UNIVERSITY OF SCIENCE AND  
TECHNOLOGY

PROJECT THESIS

---

# Improving Movie Recommendations with Unpersonal Social Media Data

---

*Author:*

Jonas MYRLUND

*Supervisors:*

Prof. Heri RAMAMPIARO

Prof. Helge LANGSETH

December 2013

NORWEGIAN UNIVERSITY OF SCIENCE AND TECHNOLOGY

## *Abstract*

Faculty of Information Technology, Mathematics and Electrical Engineering

Department of Computer and Information Science

Master of Technology

### **Improving Movie Recommendations with Unpersonal Social Media Data**

by Jonas MYRLUND

This is a project thesis about movie recommendations, specifically answering the question: *can movie recommendations be significantly improved by unpersonal social media data?* We look at how we can improve a set of recommendations – through filtering and annotation – by analyzing *the sentiment in social media*. By selecting a set of movies, and comparing the results to those found in open movie rating datasets, we consider whether the sentiment extracted from social media is of significant use in improving movie recommendations.

@TODO What do we find???

# *Acknowledgements*

@TODO: Acknowledge.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>Contents</b>	<b>iii</b>
<b>List of Figures</b>	<b>v</b>
<b>List of Tables</b>	<b>vi</b>
<b>Abbreviations</b>	<b>vii</b>
<b>Glossary</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.1.1 Data sparsity in collaborative filtering . . . . .	1
1.1.2 Twitter . . . . .	2
1.2 Research Questions . . . . .	2
1.3 Overview and Summary . . . . .	3
<b>2 Survey</b>	<b>4</b>
2.1 Relevant literature . . . . .	4
2.2 Similar applications . . . . .	5
2.3 Twitter data . . . . .	5
2.3.1 What makes a Tweet “popular”? . . . . .	6
2.3.2 Relevant parts of the Twitter API . . . . .	6
2.3.3 Novelty of available data . . . . .	7
2.3.4 Twitter as a Data Source . . . . .	8
2.4 The sentiment classifier . . . . .	8
<b>3 Design</b>	<b>9</b>
<b>4 Implementation</b>	<b>10</b>
4.1 Data retrieval . . . . .	11
4.2 Sentiment analysis . . . . .	11
4.3 Evaluating sentiment . . . . .	12

---

<b>5</b>	<b>Evaluation</b>	<b>13</b>
5.1	Evaluating against Netflix rating data . . . . .	13
<b>6</b>	<b>Conclusion</b>	<b>15</b>
6.1	Suggestions for further work . . . . .	15
<b>A</b>	<b>Evaluation Results</b>	<b>17</b>
	<b>Bibliography</b>	<b>18</b>

# List of Figures

3.1	The basic architecture of a system described in this thesis. We mainly describe the component labelled <i>A</i> . . . . .	9
3.2	The layout of the system within component <i>A</i> . . . . .	9

# List of Tables

2.1	Classification of microblogging behavior . . . . .	6
-----	--	---

# Abbreviations

**MSE**    **M**ean **S**quare **E**rror

**SaaS**    **S**oftware **a**s **a** **S**ervice

**SNS**    **S**ocial **N**etwork **S**ervice

**SVM**    **S**upport **V**ector **M**achine



# Glossary

<b>Tweet</b>	A Twitter message
<b>DatumBox</b>	A machine learning SaaS

# Chapter 1

## Introduction

### 1.1 Motivation

Many movie recommendation systems today use collaborative filtering techniques at their core. Although collaborative filtering has many advantages that have let it attain its position as one of the dominant algorithms in the field, there are still quite a few weaknesses left to remedy. Among others, two important challenges for collaborative filtering techniques left today are 1) their ability to handle data sparsity [1], and 2) their ability to explain predictions [2].

I will take a somewhat untraditional approach in an attempt to mitigate these issues: can data mined from one of today's largest sources of user-generated content, Twitter, contribute enough relevant information to solve both the problem of data sparsity and result explanation?

#### 1.1.1 Data sparsity in collaborative filtering

The data sparsity problem has several sides to it [1].

Here, we will take a closer look at the *cold start* problem – or more specifically: the *new item problem*. It occurs when a new item enters the system and there is no rating history to base similarity measures on, leaving a barebones collaborative filtering algorithm without anything to base predictions on – until, of course, some users rate it. If this is the case for a large number of items, we say that we have poor *coverage*.

Second,

### 1.1.2 Twitter

Twitter is one of the largest sources of user-generated content available today. It launched in 2006, was incorporated in 2007, and has seen active user growth ever since. At the time of writing, Twitter has more than 230 million registered users sending around 500 million Tweets per day.<sup>1</sup>

One of the most interesting things about Twitter is its simplicity. Each message is limited to 140 characters in length, for no other apparent reason than to force its author to formulate messages very concisely, as well as significantly lower the threshold for publishing content compared to traditional blogging services [3].

Furthermore, 76% of Twitter’s active users are on mobile – enabling use of the service from anywhere. Combined with Twitter’s well-established REST API, this provides us with a robust source of real-time data on almost any subject.

I’ll delve further into aspects of using Twitter as a data source in section 2.3.

## 1.2 Research Questions

*@REWORK*

In this thesis, I’ll mostly examine *improving a set of recommendations* taking social media data into account, and not so much try to generate new recommendations in their own right – although I might make a go of it if the data should prove agreeable.

I will rather look at ways of using copious<sup>2</sup>, unpersonal<sup>3</sup> social media data to *filter* sets of recommendations.

More specifically, the aim is to find answers to the following questions:

1. Can sentiment analysis of large quantities of unpersonal social media data be used to effectively filter or provide context to recommendations?

---

<sup>1</sup>Numbers from <https://about.twitter.com/company>.

<sup>2</sup>*Copious* in the sense that the size of the data source is arbitrarily large.

<sup>3</sup>*Unpersonal* in the sense that the data is not related to a single hypothetical user of the system.

2. Can unpersonal social media data in any way generate reliable ratings of its own?
3. How do we evaluate our efforts in order to answer the above questions?

### 1.3 Overview and Summary

I will look at improving an intermediate step in a hypothetical recommendation pipeline, filtering and/or annotating recommended content. For a more detailed look at the system design, and some reasoning around the parts of the system being touched upon, see chapter [3](#).

*@TODO Add more overview and summary information as it comes into existence.*

## Chapter 2

# Survey

### 2.1 Relevant literature

The task of improving recommendations through sentiment analysis of social media data requires digging into several fields of study.

Some people have attempted the same task as in this thesis, albeit with another type of social data. Singh et al. [4] investigated a “formulation, where [they] combined the content-based approach with a sentiment analysis task to improve the recommendation results.” Their approach is very similar to our approach, but differs in two important ways:

1. It uses user reviews from IMDB as content source<sup>1</sup>, and not a more general source of sentiment-carrying content – as in our case, with Twitter.
2. It is not designed to enhance presupplied recommendations, but rather to generate its own based on genre input.

As we’re looking at data from Twitter, the sentiment analysis task is a bit different than usual, as it needs to operate on texts that are all less than 140 characters long. More often than not, we will in fact need to work with texts that are merely one or two sentences long. Cho & Kang [5] “propose a method of classifying tendencies and opinions in texts of multiple sentence length extracted from social media and covering

---

<sup>1</sup>IMDB does not provide open API access at the time of writing.

both formal and informal vocabularies”. Among the things the more unusual things they consider when analysing content is position of each sentence and emotion icons, which is quite important in short, concise text like ones found on Twitter. We’ll be taking a closer look at this method for the actual sentiment analysis task.

*@TODO More.*

## 2.2 Similar applications

What spawned the idea of using an unpersonal social service like Twitter as a content source for filtering content is that Netflix recently rolled out personal social recommendations of their content. For this they use Facebook.

One huge limitation to the Facebook approach is that Facebook doesn’t expose how “close” you are to your various friends.

Content sharing patterns (Social influence and the diffusion of user-created content): [\[6\]](#).

## 2.3 Twitter data

*@TODO: Something about what terms we’re searching for, using only title etc. Ref. implementation with filtered hashtags etc.*

Micro-blogging services such as Twitter have enormous amounts of data on almost every topic imaginable. Content is limited in length, and users react to each others’ content by “re-tweeting”, “favoriting” or “replying to” it. This leaves us with a source of textual data that is:

**Instant** Users express reactions to events as they experience them.

**Weighted** Users weigh each others’ content by interacting with it.

**Concise** Due to limitations on content length, users must express themselves concisely.

Furthermore, the Twitter search API<sup>2</sup> supports returning both *popular* and *recent* content, or *a mix* of the two. This enables two interesting approaches to both filtering and annotating the recommended content, in that we can treat popular and recent comments separately.

### 2.3.1 What makes a Tweet “popular”?

As previously mentioned, “Tweets” can be *favorited*, *retweeted*, and *replied to*. Additionally, we can tell how big reach an author has by counting the number of *followers* he/she has, and use this as another indication of content popularity.

We want to be able to use the data as a source of implicit ratings. To be able to, we need to quantify the significance of these verbs. Oard and Kim [7, 8] and Kelly and Teevan [9] have developed a framework for classification of implicit feedback. They define three major categories for implicit feedback: examination, retention, and reference.

We adapt it to the domain of Twitter data, and wind up with table 2.1.

Original	Ours	Action
Examine	Consume	→ Follow
Annotate	Evaluate	→ Reply
Retain	Endorse	→ Favorite
Reference	Forward	→ Retweet

TABLE 2.1: Classification of microblogging behavior

Twitter does not, at the time of writing, allow querying of these qualities directly. However, they do have a concept of content popularity, and it’s based off this set of qualities.

### 2.3.2 Relevant parts of the Twitter API

The various endpoints in the Twitter REST API is divided into 16 categories:

1. Timelines
2. Tweets

---

<sup>2</sup><https://dev.twitter.com/docs/api/1.1/get/search/tweets>

3. Search
4. Streaming
5. Direct Messages
6. Friends Followers
7. Users
8. Suggested Users
9. Favorites
10. Lists
11. Saved Searches
12. Places Geo
13. Trends
14. Spam Reporting
15. OAuth
16. Help

Of these, only two are of relevance, namely *Search* and *OAuth*.

OAuth is only relevant because it enables us to search, so we won't elaborate further on how we use it.

### 2.3.3 Novelty of available data

As briefly mentioned, part of the hypothesis is that Twitter data is usable for predicting ratings for new content, where collaborative filtering systems perform worse than otherwise.

Novelty, however, is an inherently fundamental quality of Twitter data, even as a data source. The Twitter search API simply does not expose data more than about a week old. As it is put in their search API guidelines as of December 2013 [10], “The Search



API is not complete index of all Tweets, but instead an index of recent Tweets. At the moment that index includes between 6-9 days of Tweets.”

This quality makes Twitter data less than useful when it comes to predicting ratings for movies old enough to have calmed in public debate – a category which, not surprisingly, includes the vast majority of available movies.

#### **2.3.4 Twitter as a Data Source**

There are vast opportunities in managing to understand a data source with the qualities discussed above, but alas – the diversity and free nature of Twitter as a publishing platform comes at a price: *there is a lot of noise*. That is not to say that there is little relevant information, but when the service is designed to have such a low threshold for contributing content, a low signal-to-noise ratio seems inevitable. We will return to the issue of noise in chapter 5.1, where we try to evaluate correlations between movie titles in our test data and Tweets about the same titles.

This problem of finding content carrying relevant information turns out to perhaps be the biggest challenge of the entire study.

### **2.4 The sentiment classifier**

## Chapter 3

# Design

The basic architecture is outlined in figure 3.1. Our work will be aimed at the component labelled *A*.

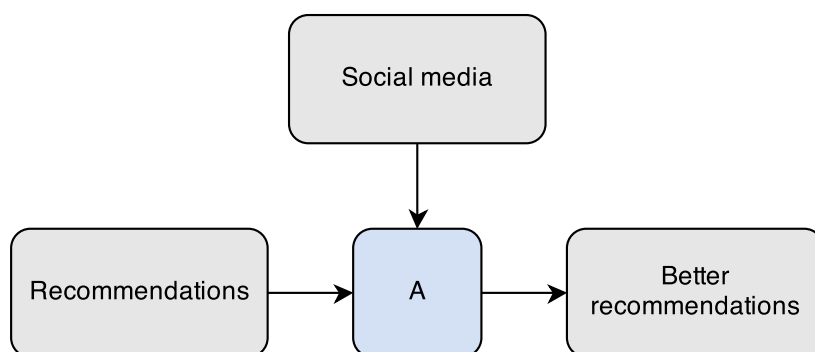


FIGURE 3.1: The basic architecture of a system described in this thesis. We mainly describe the component labelled *A*.

The *A* component is structured as illustrated in figure 3.2.

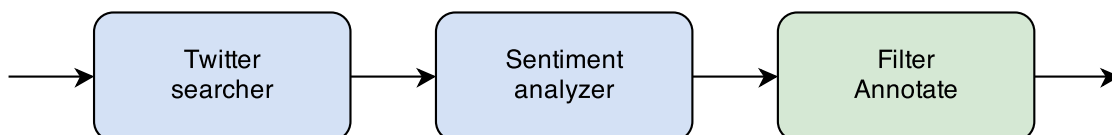


FIGURE 3.2: The layout of the system within component *A*.

## Chapter 4

# Implementation

The system does not directly employ any specific algorithms or data structures worth mentioning. The problem lies in determining the applicability of the data harvested from Twitter.

A central part of the system is the sentiment classification of the Twitter results. This sentiment analysis could well have been implemented locally, but for simplicity's sake it has been offloaded to an external service called DatumBox<sup>1</sup>, as it seems to employ a reasonable choice of algorithm, and performs well enough for our needs on the short Twitter messages. See section 4.2 for a more in-depth look at the techniques the DatumBox service utilizes.

The system is implemented as a pipeline of sorts, consisting of two main steps, and one evaluation step:

1. Data retrieval – retrieves, cleans and packages the Twitter data into simple classes for easier subsequent use.
2. Sentiment analysis – analyzes each tweet, classifying it as either positive, neutral, or negative.
3. Evaluation – maps sentiment to a final score, and compares these results to those found in external datasets.

---

<sup>1</sup>[datumbox.com](http://datumbox.com)

## 4.1 Data retrieval

As discussed in section 2.3.4, Twitter has a disturbingly low signal-to-noise ratio, so we had to examine every opportunity to refine the data retrieval step. Luckily, the Twitter API, at the time of writing, has a quite extensive search interface, with many ways of tweaking the results in a desired direction.

After much trial and failure, the following settings seem to yield the best results:

- Ensure that the title is searched for in its entirety, not the individual words.
- Exclude tweets containing the following terms<sup>2</sup>: “download”, “stream”, “nw”, “nowwatching”, and “RT”.

Then remained the choice between the two modes of search: “popular”, “recent”, or the optional combination of the two.

## 4.2 Sentiment analysis

As mentioned above, the task of sentiment analysis is offloaded to a SaaS called [DatumBox](#). They outline the techniques applied in an article on their service blog [11].

With the approach taken, texts are classified as either positive, neutral, or negative. A training set of 1.2 million tweets were tokenized “by extracting their bigrams and by taking into account the URLs, the hash tags, the usernames and the emoticons”, and were subsequently fed into a Mutual Information algorithm for feature selection. The classifier in use is the Binarized Naïve Bayes, after having outperformed SVM, Max Entropy and others on the test set.

With a 10-fold cross-validation, the best performing classifier allegedly achieves an accuracy of 83.26%.

When calling the DatumBox API, the best results were achieved when removing the movie title itself from the query. Quite a lot of titles have sentiment-carrying words in their titles<sup>3</sup>, and this obviously confused the classifier quite a bit.

---

<sup>2</sup>Any time a set of irrelevant results shared a common term, it would be added to the list. There are probably many ways of fine-tuning this further.

<sup>3</sup>“Breaking Bad” consequently scoring way below “Cheers” was a rather clear cut case.

## 4.3 Evaluating sentiment

## Chapter 5

# Evaluation

- When rating popular and well-known movies<sup>1</sup> the predicted ratings achieve a correlational coefficient of 0.75 with regard to average Netflix ratings for the same movie.
- B-movies or older less-known movies rarely collect enough Twitter search results to warrant any further analysis.
- Movies with titles that are fairly common words or expressions in their own right achieve very low precision, and often return only noise. This is hard to detect without manual interference. Need to perform some sort of Named Entity Disambiguation, maybe something like the techniques outlined in Cucerzan [12] or Sarmiento [13] (is elaboration needed?).

### 5.1 Evaluating against Netflix rating data

To find out how the Twitter-based predictions fare, we will use the average of the available Netflix ratings of the same titles as a benchmark.

With predictions  $p$ , and benchmark ratings  $r$ , we will compute the MSE (Mean Square Error) of  $N$  sample movies in the following way:

---

<sup>1</sup>The sample in question consisted of “Pulp Fiction”, “The Shining”, “Mission: Impossible”, “The Matrix”, “The Godfather”, “Forrest Gump”, and “A Clockwork Orange”.

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (p_i - r_i)^2 \quad (5.1)$$

As a baseline metric, we'll compare against the MSE of the average of all the benchmark ratings:

$$\text{MSE}_{\text{baseline}} = \frac{1}{N} \sum_{i=1}^N (\bar{r} - r_i)^2 \quad (5.2)$$

*@TODO: Reasonable at all? Twitter loses heavily every time.*

## Chapter 6

# Conclusion

I thought that the data would contain too much noise to be usable in any other way than to annotate and – in the best of cases – adjust ratings of content where the social sentiment disagreed strongly with the proposed rating.

However, for certain kinds of content, the predictions generated by Twitter were more or less spot on the same as the ones from Netflix, as shown in chapter 5.

This leads me to believe that Twitter as a data source for recommendations can have a greater role than that of augmenting recommendations coming from elsewhere. (@TODO More specifically....)

Moving towards a conclusion:

1. Twitter’s strongest suite lies in its abundance of novel content.
2. Some of the traditional CF systems’ weakest points relate to recommending novel content.
3. Augmenting new content, with extremely sparse user ratings, might well be a good application.

### 6.1 Suggestions for further work

Suggestions to further work:



- Applying NED to Twitter entities.
- Further improving sentiment analysis of informal texts.
- CRF og sentimentanalyse? – klassifisere bort tekster? filtrere i et første steg.
- Nyere filmer – tråle IMDB for nye filmer.
- Gi kontekst til Netflix-ratings: (1. legge til adjektiver?, 2. kontrastifisere positive/negative adjektiver. ordsky?)
- *Hva vil jeg jobbe med til våren?*

## Appendix A

# Evaluation Results

# Bibliography

- [1] Xiaoyuan Su and Taghi M. Khoshgoftaar. A survey of collaborative filtering techniques. *Adv. in Artif. Intell.*, 2009:4:2–4:2, January 2009. ISSN 1687-7470. doi: 10.1155/2009/421425. URL <http://dx.doi.org/10.1155/2009/421425>.
- [2] Jonathan L. Herlocker, Joseph A. Konstan, and John Riedl. Explaining collaborative filtering recommendations. In *Proceedings of the 2000 ACM Conference on Computer Supported Cooperative Work, CSCW '00*, pages 241–250, New York, NY, USA, 2000. ACM. ISBN 1-58113-222-0. doi: 10.1145/358916.358995. URL <http://doi.acm.org/10.1145/358916.358995>.
- [3] Akshay Java, Xiaodan Song, Tim Finin, and Belle Tseng. Why we twitter: Understanding microblogging usage and communities. In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 Workshop on Web Mining and Social Network Analysis*, WebKDD/SNA-KDD '07, pages 56–65, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-848-0. doi: 10.1145/1348549.1348556. URL <http://doi.acm.org/10.1145/1348549.1348556>.
- [4] Vivek Kumar Singh, Mousumi Mukherjee, and Ghanshyam Kumar Mehta. Combining a content filtering heuristic and sentiment analysis for movie recommendations. In K.R. Venugopal and L.M. Patnaik, editors, *Computer Networks and Intelligent Computing*, volume 157 of *Communications in Computer and Information Science*, pages 659–664. Springer Berlin Heidelberg, 2011. ISBN 978-3-642-22785-1. doi: 10.1007/978-3-642-22786-8\_83. URL [http://dx.doi.org/10.1007/978-3-642-22786-8\\_83](http://dx.doi.org/10.1007/978-3-642-22786-8_83).
- [5] Sang-Hyun Cho and Hang-Bong Kang. Statistical text analysis and sentiment classification in social media. In *Systems, Man, and Cybernetics (SMC), 2012*

- IEEE International Conference on*, pages 1112–1117, 2012. doi: 10.1109/ICSMC.2012.6377880.
- [6] Eytan Bakshy, Brian Karrer, and Lada A. Adamic. Social influence and the diffusion of user-created content. In *Proceedings of the 10th ACM conference on Electronic commerce*, EC '09, pages 325–334, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-458-4. doi: 10.1145/1566374.1566421. URL <http://doi.acm.org/10.1145/1566374.1566421>.
- [7] Douglas Oard and Jinmook Kim. Implicit feedback for recommender systems. In *in Proceedings of the AAAI Workshop on Recommender Systems*, pages 81–83, 1998.
- [8] Douglas W. Oard and Jinmook Kim. Modeling information content using observable behavior, 2001.
- [9] Diane Kelly and Jaime Teevan. Implicit feedback for inferring user preference: a bibliography. *SIGIR Forum*, 37(2):18–28, September 2003. ISSN 0163-5840. doi: 10.1145/959258.959260. URL <http://doi.acm.org/10.1145/959258.959260>.
- [10] Twitter Inc. Using the twitter search api, October 2013. URL <https://dev.twitter.com/docs/using-search>.
- [11] Vasilis Vryniotis. How to build your own twitter sentiment analysis tool, September 2013. URL <http://blog.datumbox.com/how-to-build-your-own-twitter-sentiment-analysis-tool/>.
- [12] Silviu Cucerzan. Large-scale named entity disambiguation based on Wikipedia data. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 708–716, Prague, Czech Republic, June 2007. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/D/D07/D07-1074>.
- [13] Luís Sarmiento, Alexander Kehlenbeck, Eugénio Oliveira, and Lyle Ungar. An approach to web-scale named-entity disambiguation. In *Proceedings of the 6th International Conference on Machine Learning and Data Mining in Pattern Recognition*, MLDM '09, pages 689–703, Berlin, Heidelberg, 2009. Springer-Verlag. ISBN 978-3-642-03069-7. doi: 10.1007/978-3-642-03070-3\_52. URL [http://dx.doi.org/10.1007/978-3-642-03070-3\\_52](http://dx.doi.org/10.1007/978-3-642-03070-3_52).