

# Hiërarchische zelflerende modellen voor muziekbeveiling

Erik Vandeputte

Promotor: prof. dr. ir. Benjamin Schrauwen  
Begeleiders: Sander Dieleman, Philémon Brakel

Masterproef ingediend tot het behalen van de academische graad van  
Master in de ingenieurswetenschappen: computerwetenschappen

Vakgroep Elektronica en Informatiesystemen  
Voorzitter: prof. dr. ir. Jan Van Campenhout  
Faculteit Ingenieurswetenschappen en Architectuur  
Academiejaar 2012-2013





# Hiërarchische zelflerende modellen voor muziekbeveiling

Erik Vandeputte

Promotor: prof. dr. ir. Benjamin Schrauwen  
Begeleiders: Sander Dieleman, Philémon Brakel

Masterproef ingediend tot het behalen van de academische graad van  
Master in de ingenieurwetenschappen: computerwetenschappen

Vakgroep Elektronica en Informatiesystemen  
Voorzitter: prof. dr. ir. Jan Van Campenhout  
Faculteit Ingenieurwetenschappen en Architectuur  
Academiejaar 2012-2013



# Voorwoord

Muziek en computers zijn altijd 2 grote passies van mij geweest. Vandaar dat ik tijdens het beschikbaar komen van de lijst met thesisonderwerpen begonnen ben met zoeken op het trefwoord 'muziek'. Ongeveer een jaar later kan ik zeggen dat de combinatie van computerwetenschappen en muziek voor mij een leuke en verrassende uitdaging is gebleken.

Ik wil eerst en vooral mijn promotor prof. dr. ir. Benjamin Schrauwen bedanken. Daarnaast wil ik in het bijzonder ook mijn begeleiders Philemon Brakel, Sander Dieleman, Pieter-Jan Kindermans en Aäron van den Oord bedanken voor hun advies en feedback.

Ten slotte wil ik ook mijn ouders bedanken voor hun steun en vertrouwen.

# Toelating tot bruikleen

De auteur geeft de toelating deze masterproef voor consultatie beschikbaar te stellen en delen van de masterproef te kopiëren voor persoonlijk gebruik. Elk ander gebruik valt onder de beperkingen van het auteursrecht, in het bijzonder met betrekking tot de verplichting de bron uitdrukkelijk te vermelden bij het aanhalen van resultaten uit deze masterproef.

Erik Vandeputte, mei 2013

# Music recommendation based on listening history and acoustic models

Erik Vandeputte

Supervisor(s): prof. dr. ir. Benjamin Schrauwen, ir. Sander Dieleman, ir. Philémon Brakel

**Abstract**—In this article, we try to come up with a successful solution to the cold start problem by building a hybrid recommender system that combines the benefits of social and content-based recommendation algorithms.

**Keywords**—Music Recommendation, Collaborative Filtering, Cold start problem

## I. INTRODUCTION

SHIFTING from physical sales to digital sales during recent years, the online music industry is booming. Users are virtually overwhelmed with the available music content. Both from the perspective of the users, who want to discover relevant new music, as from the perspective of the content providers, who want to match their content with the right consumer group, the interest in recommendation engines has risen. Social-based recommendation algorithms, also known as collaborative filtering techniques (CF) are based on analyzing the large amount of information on user behavior such as rating information or listening/viewing history. Based on this information they are able to learn new relations between users and items so that they are able to predict what users will like based on their similarity to other users or items they've liked. CF is known to suffer from the cold start problem, where a system cannot recommend items for which it doesn't have any ratings for. Content-based techniques attempt to overcome this problem by extracting features out of the actual content and matching them with known user preferences for those features. We attempt to build a hybrid music recommender system that performs well on songs for which we have a lot of information for, but is also able to recommend new relevant songs to users.

### A. Dataset

The dataset consists of a subset of the Million Song Dataset [1] containing listening information for 20.000 users and 10.000 songs. This means that the input for the recommender system is a form of implicit feedback. We can represent this data in an user-item matrix  $R$ . One dimension corresponds to a user  $u$  while the other dimension corresponds to the song  $i$ . An rating  $r_{ui}$  represents the number of times user  $u$  has listened to song  $i$ . Additionally, there is a 30s audio file available that will be used for extracting features out of the audio.

## II. COLLABORATIVE FILTERING

Our CF approach is based on a latent factor model that takes implicit feedback into account. Latent factor models try to discover latent features that explain the observed ratings. Our model tries to decompose an user-item matrix  $R$  into a product of an user-factor matrix  $X$  and a song-factor matrix  $Y$ . If the

number of latent features is  $k$  then the observed score  $r_{ui}$  in the user-item matrix is approximated by an innerproduct between a user-factor vector  $x_u \in \mathbb{R}^k$  and a song-factor vector  $y_i \in \mathbb{R}^k$  i.e.,  $r_{ui} \cong x_u^T y_i$ . The main difference between recommender systems based on explicit feedback datasets and those based on implicit feedback datasets is that a score of 0 in an implicit feedback dataset doesn't mean that the user doesn't like the item. In our case it just means that he hasn't listened to that song before. A suitable loss function for a recommender system for implicit feedback datasets can look like this:

$$\min_{x^*, y^*} = \sum_{u,i} c_{ui} (p_{ui} - x_u^T y_i)^2 + \lambda \left( \sum_u \|x_u\|^2 + \sum_i \|y_i\|^2 \right) \quad (1)$$

$c_{ui}$  is a confidence parameter for a certain rating that is 1 when user  $u$  hasn't listened to item  $i$  but that increases lineary or logarithmically with the number of play counts.  $p_{ui}$  is a binarized form of the user-item matrix where  $p_{ui} = 1$  means that user  $u$  has listened to item  $i$ .  $\lambda$  is used for regularizing the model. Because the user-item matrix cannot be considered sparse, Alternating Least Squares (ALS) optimization is used for equation 1. We alternate between re-computing the user-factors and the song-factors by considering the former as fixed and solve for the latter and vice versa. In each step we are guaranteed to lower the value for the cost function. [2] describes an efficient implementation of this process, scaling lineary in the size of the dataset.

### A. Results

The results were compared to two trivial recommendation algorithms and one neighborhood-based model. A first baseline algorithm, called *popularity*, recommends each song in descending order of its popularity, regardless of the user's preferences. The *same artist* algorithm recommends songs in descending order of popularity as well, taking into account only the artists that the user has played before. The neighborhood-based model is based on [3].

The evaluation metrics used are the mean average precision (mAP) and the recall when 50 recommendations are made (see table I).

As expected, the neighborhood-based model and the latent factor model outperform the trivial recommender systems. An important advantage of the latent factor model is that we obtain a representation of each song in  $k$  factors.

| Recommender system              | mAP@50  | Rec@50  |
|---------------------------------|---------|---------|
| <i>Popularity</i>               | 0.01295 | 0.02739 |
| <i>Same artist</i>              | 0.03063 | 0.04723 |
| <i>Neighborhood-based model</i> | 0.12405 | 0.25469 |
| <i>Latent factor model</i>      | 0.12830 | 0.27362 |

TABLE I

THE EVALUATION FOR DIFFERENT RECOMMENDER SYSTEMS EVALUATED ON THE SUBSET OF THE MSD

### III. CONTENT-BASED FILTERING

Motivated by the idea that some of the song factors correlate with various acoustical attributes of the song, an alternative approach is for songs for which we don't have any listening information to predict the song factors. Predicting the song factors directly out of the raw audio files will not work. We present an alternative approach based on the clustering of Mel Frequency Cepstral Coefficients.

#### A. Feature Learning

MFCCs [4] originated in the context of speech recognition, but have already successfully been used in the field of music data mining, f.e. in genre recognition. MFCCs present a compact representation of how humans perceive sound. Each raw audio file was converted to a sequence of 2905 MFCC-vectors. As input for a feature learning algorithm, a sequence of  $l$  MFCC-vectors was grouped into one frame. Motivated by [5], k-means clustering was chosen as feature learning algorithm, operating directly on all the frames. K-means is extremely fast, has no hyper-parameters to tune other than the number of clusters and is very easy to implement. After learning a number of cluster means, we consider two choices for the feature mapping function  $f$ . The standard formula for cluster assignment is *hard clustering*:

$$f_k(x) = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|c^{(j)} - x\|_2^2 \\ 0 & \text{else} \end{cases} \quad (2)$$

An alternative approach is to assign each frame to multiple cluster means by formula:

$$f_k(x) = \max\{0, \mu(z) - z_k\} \quad (3)$$

The latter we will refer to as *soft clustering*.

#### B. Predicting Song Factors

After performing the feature learning algorithm, we can now represent a song by its clustered MFCC-representation. To be able to predict the song-factors out of this representation, we make use of regression analysis. First we splitted the song-factormatrix into a 80/20 training- and testset. By using ridge regression we estimated the song factors, based on the training data. The hyperparameter  $\lambda$  is obtained through 10-fold cross-validation. The evaluation on the testset for different frame sizes  $l$ , different number of features and for *hard* and *soft clustering* is shown in figure 1.

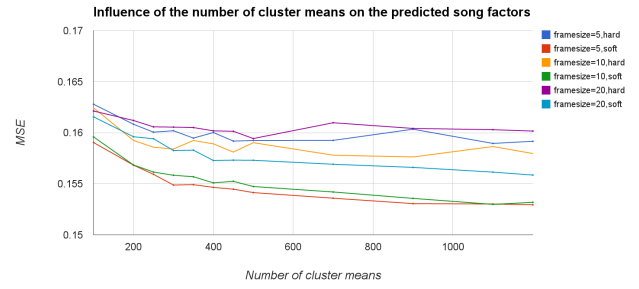


Fig. 1. Evolution of the number of cluster means on the MSE between the predicted and the true song factors for different types of clustering and frame sizes.

The goal is to estimate the predicted song factors for new songs or songs that still have very few listeners. These song-factors can then be plugged into the song-factor matrix. To test whether this content-based approach is feasible, different experiments were conducted.

### IV. EXPERIMENTS

Starting from the original user-item matrix, the dataset is divided into two parts as described in figure 2. A latent factor model is trained on  $X_{train}$ , resulting in a song-factormatrix  $V_{train}$  and a user-factormatrix  $U$ . The song-factormatrix  $V_{train}$  is now used to train a linear model that estimates song factors based on the clustered MFCC-representation described in III-A. This model is now used to predict the song factors for the songs in  $X_{test}$ . We will refer to those song factors as  $V_{test_{audio}}$ .

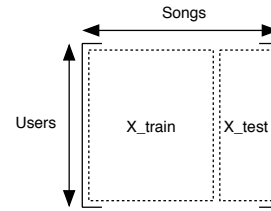


Fig. 2. Initial splitting of the dataset.

To simulate an actual cold start problem, we now remove listening information for the songs in  $X_{test}$ . Concretely, we split the matrix  $X_{test}$  randomly into 2 disjunct collections called  $X_{test_A}$ , that will serve as our new testset, and  $X_{test_B}$ , that will be used to calculate a new song-factormatrix  $V_{test_{mf}}$  using the user-factormatrix  $U$  and one iteration out of the latent factor model described in III where we take the user factors  $U$  to be constant and solve for the song factors. The parameter  $\alpha$  is used to determine the percentage listening information that is removed for each song. The matrix  $X_{test_A}$  contains the removed listening information. The matrix  $X_{test}$  consisted out of 264.767 nonzero entries. When  $\alpha = 0.5$ , both matrices  $X_{test_A}$  and  $X_{test_B}$  contain about 132.000 nonzero entries.

For different values of  $\alpha$ , we can generate our prediction scores for each song for the two types of recommender systems as follows:

$$pred_{audio} = U * (V_{test_{audio}})^T \quad (4)$$

$$pred_{mf} = U * (V_{test_{mf}})^T \quad (5)$$

Furthermore, the content-based technique (4) and the social-based technique (5) could easily be combined into a hybrid recommender system:

$$pred_{hybrid} = pred_{audio} + pred_{mf} \quad (6)$$

The results of this experiment for multiple values of  $\alpha$  are presented in tables II and III.

|                 | $\alpha$      |               |               |               |               |               |               |
|-----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                 | 0.5           | 0.8           | 0.9           | 0.95          | 0.99          | 0.995         | 0.999         |
| $pred_{mf}$     | <b>0.8580</b> | <b>0.8354</b> | <b>0.8129</b> | <b>0.7844</b> | <b>0.7022</b> | 0.6669        | 0.6451        |
| $pred_{audio}$  | 0.6630        | 0.6629        | 0.6628        | 0.6625        | 0.6626        | 0.6625        | 0.6580        |
| $pred_{hybrid}$ | 0.8520        | 0.8259        | 0.7987        | 0.7668        | 0.7019        | <b>0.6877</b> | <b>0.6808</b> |

TABLE II

AUC FOR DIFFERENT RECOMMENDER SYSTEMS AND DIFFERENT VALUES OF  $\alpha$

|                 | $\alpha$      |               |               |               |               |               |               |
|-----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                 | 0.5           | 0.8           | 0.9           | 0.95          | 0.99          | 0.995         | 0.999         |
| $pred_{mf}$     | <b>0.1367</b> | <b>0.1446</b> | <b>0.1353</b> | <b>0.1174</b> | <b>0.0572</b> | <b>0.0368</b> | 0.0228        |
| $pred_{audio}$  | 0.0174        | 0.0203        | 0.0214        | 0.0219        | 0.0224        | 0.0224        | 0.0226        |
| $pred_{hybrid}$ | 0.1278        | 0.1292        | 0.1091        | 0.0792        | 0.0365        | 0.0302        | <b>0.0277</b> |

TABLE III

MAP@500 FOR DIFFERENT RECOMMENDER SYSTEMS AND DIFFERENT VALUES OF  $\alpha$

We can tune the importance of each technique in the hybrid approach by using an extra parameter  $\beta$ . Equation 6 then becomes:

$$pred_{hybrid} = (1 - \beta) * pred_{audio} + (\beta) * pred_{mf} \quad \beta \in [0, 1] \quad (7)$$

$\beta = 1$  means that we only take the social-based recommender system into account. We now can search for the optimal value of  $\beta$  for different values of  $\alpha$ :

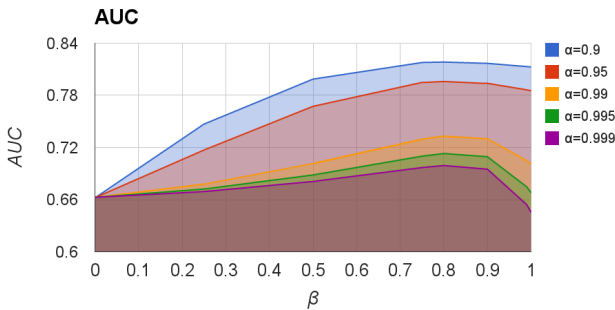


Fig. 3. AUC for different values of  $\beta$

We can clearly see that for different values of  $\alpha$ , the optimal value for  $\beta$  lies between 0.8-0.9, but then decreases when we

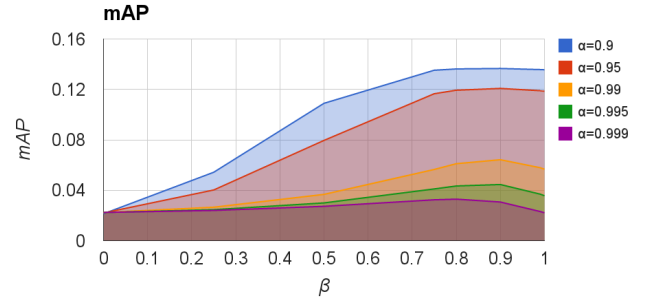


Fig. 4. mAP for different values of  $\beta$

only take the social-based system into account. We can conclude that in the tuned hybrid recommender system the social-based technique has a significant bigger contribution than our own content-based technique, but that the correct combination of both can lead to an improved recommender system in the case of a cold start.

## V. CONCLUSION & FUTURE WORK

In this paper, a solution to the cold start problem in music recommendation is proposed by using a latent factor model that allows to predict song factors based on MFCC-vectors when none or few listening history is available. If no listening information is available for certain songs, this approach works better than randomly recommend those songs, but the latent factor model tends to be superior even when few listening information is available. The main reason why the prediction of the song factors is poor is that not all song factors can be predicted out of MFCCs because they can correlate to other non-acoustical features such as home country of the artist. A possible solution to this problem is to incorporate more, non-acoustical attributes while predicting song factors.

Furthermore, a hybrid recommendation system that correctly combines the social-based and content-based technique outperforms the latent factor model when few listening information is available.

## REFERENCES

- [1] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere, "The million song dataset," in *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- [2] Yifan Hu, Yehuda Koren, and Chris Volinsky, "Collaborative filtering for implicit feedback datasets," in *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, Washington, DC, USA, 2008, ICDM '08, pp. 263–272, IEEE Computer Society.
- [3] Fabio Aielli, "A preliminary study on a recommender system for the million songs dataset challenge," 2012.
- [4] S. Davis and P. Mermelstein, "Experiments in syllable-based recognition of continuous speech," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 28, pp. 357 – 366, Aug. 1980.
- [5] Adam Coates, Andrew Y. Ng, and Honglak Lee, "An analysis of single-layer networks in unsupervised feature learning," *Journal of Machine Learning Research - Proceedings Track*, vol. 15, pp. 215–223, 2011.



# Inhoudsopgave

|  |             |
|--|-------------|
| <b>Voorwoord</b>   | <b>i</b>    |
| <b>Toelating tot bruikleen</b>                             | <b>ii</b>   |
| <b>Lijst met afkortingen</b>                               | <b>viii</b> |
| <b>1 Inleiding</b>   | <b>1</b>    |
| 1.1 Wat is een aanbevelingssysteem? . . . . .              | 1           |
| 1.1.1 Verschillende soorten aanbevelingssystemen . . . . . | 2           |
| 1.1.1.1 Inhouds-gebaseerde aanbevelingssystemen . . . . .  | 3           |
| 1.1.1.2 Sociaal-gebaseerde aanbevelingssystemen . . . . .  | 3           |
| 1.1.1.3 Hybride aanbevelingssystemen . . . . .             | 4           |
| 1.2 De lange staart . . . . .                              | 4           |
| 1.3 Evaluatie van een aanbevelingssysteem . . . . .        | 5           |
| 1.3.1 mean Average Precision . . . . .                     | 6           |
| Precisie (P) . . . . .                                     | 6           |
| Gemiddelde Precisie (AP) . . . . .                         | 7           |
| Afgeknotte Gemiddelde Precisie (mAP) . . . . .             | 7           |
| 1.3.2 Recall . . . . .                                     | 7           |
| 1.3.3 Area Under Curve . . . . .                           | 8           |
| 1.3.4 Novelty . . . . .                                    | 10          |
| 1.4 Datasets . . . . .                                     | 11          |
| 1.4.1 Million Song dataset (MSD) . . . . .                 | 11          |
| 1.4.2 MusiXmatch dataset . . . . .                         | 12          |
| 1.4.3 Last.fm dataset . . . . .                            | 12          |
| <b>2 Collaboratieve Filtertechnieken</b>                   | <b>13</b>   |
| 2.1 Buurgebaseerde modellen . . . . .                      | 13          |
| 2.1.1 Algoritme . . . . .                                  | 14          |
| 2.2 Verborgene Factormodellen . . . . .                    | 14          |
| 2.2.1 Algoritme . . . . .                                  | 16          |
| 2.3 Resultaten . . . . .                                   | 17          |
| <b>3 Het voorspellen van emotie in muziek</b>              | <b>21</b>   |
| 3.1 Data voorverwerking en feature extractie . . . . .     | 21          |
| 3.1.1 Echonest features . . . . .                          | 23          |
| 3.1.2 Songtekst features . . . . .                         | 26          |
| 3.2 Classificatie en Experimenten . . . . .                | 27          |

---

|          |   |           |
|----------|---|-----------|
| 3.2.1    | Classificatie . . . . .                               | 27        |
| 3.2.2    | Experimenten . . . . .                                | 28        |
| <b>4</b> | <b>Het voorspellen van muziekfactoren</b>             | <b>31</b> |
| 4.1      | Mel Frequency Cepstral Coefficients (MFCCs) . . . . . | 31        |
| 4.2      | Feature Learning . . . . .                            | 34        |
| 4.2.1    | Feature encoding . . . . .                            | 34        |
| 4.3      | Voorspellen van muziekfactoren . . . . .              | 37        |
| 4.3.1    | Ridge-regression . . . . .                            | 37        |
| 4.3.2    | Resultaten . . . . .                                  | 39        |
| 4.4      | Het genereren van aanbevelingen . . . . .             | 40        |
| 4.4.1    | Resultaten . . . . .                                  | 41        |
| 4.4.1.1  | Experiment 1 . . . . .                                | 41        |
| 4.4.1.2  | Experiment 2 . . . . .                                | 44        |
| 4.4.1.3  | Experiment 3 . . . . .                                | 45        |
| 4.4.1.4  | Experiment 4 . . . . .                                | 48        |
| <b>5</b> | <b>Besluit</b>  | <b>50</b> |
|          | <b>List of Figures</b>                                | <b>54</b> |
|          | <b>List of Tables</b>                                 | <b>57</b> |

# Lijst met afkortingen

|             |   |
|-------------|---|
| <b>ALS</b>  | Alternating Least Squares, alternerende kleinste kwadratenmethode |
| <b>AUC</b>  | Area Under the Curve, de oppervlakte onder een curve              |
| <b>API</b>  | Application Programming Interface                                 |
| <b>CF</b>   | Collaboratieve Filtertechnieken                                   |
| <b>DFT</b>  | Discrete Fouriertransformatie                                     |
| <b>DCT</b>  | Discrete Cosinustransformatie                                     |
| <b>FPR</b>  | Fals Positive Rate, vals positief ratio                           |
| <b>MSD</b>  | Million Song Dataset  |
| <b>mAP</b>  | mean Average Precision, gemiddelde precisie                       |
| <b>MFCC</b> | Mel Frequency Ceptral Coefficients                                |
| <b>MSE</b>  | Mean Squared Error, gemiddelde kwadratische fout                  |
| <b>PCA</b>  | Principal Component Analysis                                      |
| <b>RSS</b>  | Residual Sum of Squares   |
| <b>SVD</b>  | Singular Value Decomposition, singuliere waarden decompositie     |
| <b>TPR</b>  | True Positive Rate, echt positief ratio                           |

# Hoofdstuk 1

## Inleiding

In dit hoofdstuk wordt het begrip aanbevelingssysteem uitgelegd waarna de verschillende klassen van aanbevelingssystemen kort toegelicht worden. Vervolgens worden de problemen van hedendaagse aanbevelingssystemen aangekaart en de aanpak die deze masterproef voorstelt om deze te verhelpen. Daarna worden enkele populaire metrieken aangehaald die doorheen de masterproef gebruikt werden om verschillende aanbevelingssystemen te gaan evalueren. Tenslotte worden de verschillende datasets besproken die voorkomen in deze masterproef.

### 1.1 Wat is een aanbevelingssysteem?

Mede door de explosie aan beschikbare informatie hebben aanbevelingssystemen de laatste jaren enorm aan populariteit gewonnen op het internet. Deze systemen gaan mensen helpen door gepersonaliseerde informatie te tonen waar de gebruiker in eerste instantie nog niet noodzakelijk aan heeft gedacht. Contentleveranciers proberen vaak een zo groot mogelijk publiek te bereiken zodat de eindgebruiker vaak overspoeld wordt door een enorme hoeveelheid aan informatie. In de plaats van de gebruiker zelf te laten zoeken wat hij precies nodig heeft uit die informatie, filtert het systeem de informatie en presenteert het enkel de informatie waarvan vermoed wordt dat die van belang kan zijn voor de gebruiker. Aanbevelingssystemen vormen dus eigenlijk een brug tussen contentleveranciers en eindgebruikers en zijn interessant voor beide partijen. Enerzijds moet de eindgebruiker niet meer op zoek gaan naar de juiste informatie en de aanbieder kan automatisch de juiste doelgroep bereiken. De lijst van aanbevolen video's op Youtube, een webshop zoals Amazon die steeds relevante producten probeert te promoten of de online

streaming videodienst Netflix die relevante films en series aanraadt zijn allemaal voorbeelden van aanbevelingssystemen.

Terwijl vroeger nog naar de platenzaak werd gegaan om een CD is het tegenwoordig mogelijk om met enkele klikken een enorme hoeveelheid aan muziek terug te vinden en te beluisteren via het internet. Ook hier is er nog steeds met een kloof tussen de eindgebruiker met zijn specifieke voorkeuren en de muziekindustrie die natuurlijk zoveel mogelijk van zijn muziek wil verkopen. Wat een muziekaanbevelingssysteem doet is het voorzien van de juiste muziek voor de juiste mensen op het juiste moment. Één van de bekendste online muziekaanbevelingssystemen werd ontwikkeld door de communitywebsite Last.fm. Gebruikers kunnen telkens wanneer ze een muzieknummer op hun computer of mobiel toestel beluisterd hebben deze informatie 'scrobblen' naar een server. Dit betekent dat het nummer wordt toegevoegd aan het muziekprofiel van de gebruiker. Uit al deze informatie kan Last.fm o.a. gepersonaliseerde aanbevelingen maken of de gebruiker in contact brengen met mensen die een een gemeenschappelijke muzieksmaak hebben.

Aanbevelingssystemen kunnen dus pas gebouwd worden wanneer er reeds een hoeveelheid invoerdata voorhanden is. Deze informatie wordt typisch voorgesteld in een gebruiker-item matrix waar één dimensie correspondeert met de gebruikers en de andere dimensie met de verschillende items. Een item kan bijvoorbeeld een film of een muzieknummer zijn. Wanneer een gebruiker een score heeft toegekend aan een bepaald item, spreken we van een systeem met expliciete feedback. Enkele voorbeelden van expliciete feedback zijn de scores die een gebruiker toekent aan een bepaalde film op Netflix of de "ik vind dit leuk" en "ik vind dit niet leuk" knoppen van Youtube. Een gebruiker heeft typisch slechts een klein percentage van alle items beoordeeld, wat resulteert in een vrij ijle gebruiker-item matrix.

Wanneer er geen expliciete scores voorhanden zijn kan de invoerdata bestaan uit de zoekgeschiedenis, luistergeschiedenis, zoekopdrachten,... . Dit zijn allemaal vormen van impliciete feedback. Wanneer de invoerdata bestaat uit de luistergeschiedenis van verschillende gebruikers, kan elke positie in de gebruiker-item matrix corresponderen met het aantal keer dat een gebruiker een bepaald nummer heeft afgespeeld.

### **1.1.1 Verschillende soorten aanbevelingssystemen**

In het algemeen kunnen we 3 grote klassen van aanbevelingssystemen onderscheiden, namelijk inhouds-gebaseerde, sociaal-gebaseerde en hybride systemen [1].

#### 1.1.1.1 Inhouds-gebaseerde aanbevelingssystemen

Inhouds-gebaseerde aanbevelingssystemen vormen de meest intuïtieve klasse. Hierbij wordt geprobeerd een bepaald profiel van de eindgebruiker en het item samen te stellen aan de hand van externe informatie. Een profiel van een gebruiker kan bestaan uit verschillende attributen zoals leeftijd, woonplaats, voorkeur voor een bepaald muziekgenre,... Wanneer men muzieknummers wil aanbevelen kan men attributen verzamelen zoals artiest, tempo, emotie, muziekgenre,... Het aanbevelingsproces bestaat er dan in om beide profielen op elkaar af te stemmen. De internetradio Pandora is een commerciële toepassing van dit type aanbevelingssystemen die voortgevloeid is uit het Music Genome Project. In dit project heeft men geprobeerd om elk muzieknummer voor te stellen als een gen. Een gen kan op zijn beurt worden voorgesteld als een vector met ongeveer 400 componenten die corresponderen met o.a. het geslacht van de artiest, toonbereik van het refrein,... De luisteraar kan specifieke voorkeuren opgeven waarna beide profielen op elkaar afgestemd worden. Deze informatie hoeft niet noodzakelijk manueel ingevoerd te worden. Bij een inhouds-gebaseerd aanbevelingssysteem voor muziek kan bijvoorbeeld door akoestische analyse van het muzieknummer automatisch een goeie schatting gemaakt worden van attributen zoals tempo, genre of emotie. Een voorbeeld van deze automatische detectie wordt besproken in hoofdstuk 2.

#### 1.1.1.2 Sociaal-gebaseerde aanbevelingssystemen

In de praktijk moeten bedrijven zoals Amazon, IMDB en Last.fm het echter vaak stellen met een bepaald gedrag dat de gebruikers vertonen op hun website zoals bijvoorbeeld zoekgeschiedenis of luistergeschiedenis. Dit brengt ons bij de sociaal-gebaseerde aanbevelingssystemen. Het basisidee van sociaal-gebaseerde aanbevelingssystemen of collaboratieve filtertechnieken (CF) is vrij eenvoudig. Wanneer een gebruiker de items A,B,C en D goed vond en er veel gebruikers zijn die items A,B,C,D en E goed vonden dan kan met vrij hoge zekerheid gesteld worden dat die gebruiker item E ook goed zal vinden. De relaties tussen gebruikers en items worden geanalyseerd om nieuwe relaties te leren tussen gebruikers en items. De ontwikkeling van sociaal-gebaseerde aanbevelingssystemen kende een sterke groei onder invloed van 'The Netflix prize' [2]. Dit was een wedstrijd die het bedrijf Netflix heeft uitgeschreven in 2006 met als doel om de score te voorspellen die een bepaalde gebruiker zal toekennen aan een film. Er werd een prijs van 1 miljoen dollar uitgereikt aan het team dat als eerste een techniek ontwikkelde die een verbetering van

10% opleverde ten opzichte van het eigen aanbevelingssysteem van Netflix. Een overzicht van de belangrijkste CF wordt gegeven in [3].

Een belangrijk voordeel dat collaboratieve filtertechnieken hebben ten opzichte van inhouds-gebaseerde technieken is dat ze domeinvrij zijn. Dit betekent dat dezelfde technieken kunnen worden gebruikt om verschillende soorten items aan te bevelen. Daar tegenover staat dan wel dat deze technieken onderhevig zijn aan problemen zoals *het koude startprobleem*, waarbij een aanbevelingssysteem er niet in slaagt om items aan te bevelen waarvoor nog geen gebruikersinformatie beschikbaar is. Indien er in het bovenstaande voorbeeld geen enkele gebruiker item E reeds geconsumeerd heeft, dan kan het aanbevelingssysteem onmogelijk weten of dat item van nut zal zijn voor een bepaalde gebruiker. Een ander probleem waarmee CF te maken hebben is dat ze de voorkeur geven om globaal populaire items aan te bevelen ten koste van niche items (zie 1.2).

#### 1.1.1.3 Hybride aanbevelingssystemen

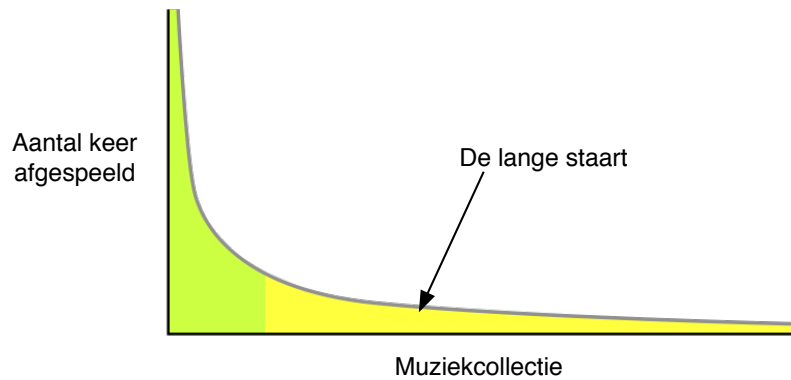
In een poging om de beste eigenschappen van de eerste twee klassen te gaan benutten kan ook worden geprobeerd om inhouds-gebaseerde en sociaal-gebaseerde systemen te gaan combineren. In dit geval spreken we van een hybride aanbevelingssysteem. Een voorbeeld van een hybride aanbevelingssysteem dat films aanraadt kan deels gebaseerd zijn op het zoekgedrag van gelijkaardige gebruikers (sociaal-gebaseerd), maar ook films aanraden die gemeenschappelijke karakteristieken (genre, regisseur,...) bevatten met films waarvoor een gebruiker reeds een hoge score heeft toegekend (inhouds-gebaseerd).

## 1.2 De lange staart

Door de populariteit van het internet is ook de manier waarop mensen muziek consumeren sterk veranderd. Technologische vooruitgang in netwerken, opslagcapaciteit en afspeelmedia zorgt ervoor dat mensen nu vaak een grotere muziekcollectie dan vroeger bezitten. Dit leidt echter niet noodzakelijk tot het feit dat mensen meer verschillende muziknummers gaan beluisteren. Vaak is het zo dat veel muziknummers zelden of nooit worden afgespeeld. Dit verschijnsel heet *de lange staart* [4]. Huidige aanbevelingssystemen focussen zich voornamelijk op een hoge graad van nauwkeurigheid. Dit leidt vaak tot het probleem dat deze systemen de voorkeur geven aan populaire items, terwijl items die zich in de lange staart bevinden misschien meer

nut hebben voor de gebruiker. Een goed aanbevelingssysteem moet hiermee rekening houden en moet gepersonaliseerde aanbevelingen genereren die de gebruikers helpen om nieuwe, relevante nummers te ontdekken. Een bekend probleem bij sociaal-gebaseerde aanbevelingssystemen is dat deze moeite hebben om items aan te bevelen die zich in de lange staart bevinden [5].

In deze masterproef zullen we proberen om een oplossing voor dit probleem te bedenken. Concreet gaan we kijken of een hybride aanpak op basis van een sociaal-gebaseerde en een inhouds-gebaseerde methode het koude start- en lange staartprobleem kan verhelpen.



FIGUUR 1.1: Een voorbeelddistributie van de populariteit van nummers in een muziekcollectie

### 1.3 Evaluatie van een aanbevelingssysteem

Vooraleer men aan muziekaanbeveling kan doen is het belangrijk om een geschikte manier te vinden om de aanbevelingen te evalueren. Voor een aanbevelingssysteem dat expliciete feedback bevat, is het mogelijk om de gemiddelde kwadratische fout (MSE) te bepalen tussen de score die het systeem voorspelt en de effectieve score die een gebruiker gaf:

$$MSE = \frac{1}{n} \sum_{i=1}^n (p_i - a_i)^2 \quad (1.1)$$

Waarbij  $n$  het aantal voorspelde scores bedraagt,  $p_i$  de voorspelde score voor de gebruiker op item  $i$ , en  $a_i$  de effectieve score die de gebruiker toekende aan item  $i$ .



In deze masterproef werd echter met impliciete feedback gewerkt, namelijk het aantal keer dat een gebruiker een muzieknummer heeft afgespeeld (zie 1.4.1). Het aantal afspreekbeurten correspondeert niet noodzakelijk met de score die een gebruiker zou geven aan het muzieknummer. Dit zorgt ervoor dat andere evaluatiemetriecken dienen te worden gebruikt.

Een dataset wordt in een eerste fase altijd opgesplitst in een trainingset en een testset. De trainingset bestaat uit de informatie die gebruikt wordt om de CF-modellen te bouwen waarna deze in staat moeten zijn om aanbevelingen voor elke gebruiker te genereren. Een goed aanbevelingssysteem zal erin slagen om aanbevelingen te genereren die effectief in de testset voorkomen. Het systeem zal als het ware de testset moeten voorspellen hoewel het enkel de informatie uit de trainingset te verwerken kreeg.

Allereerst kan worden opgemerkt dat een goede evaluatiemaat voor een aanbevelingssysteem met impliciete feedback de volgende principes in acht moet nemen:

- Een realistisch aanbevelingssysteem zal maximum  $x$  items aanbevelen.
- Bij het genereren van aanbevelingen zullen ook telkens  $x$  items worden aanbevolen. Indien foute aanbevelingen typisch niet worden afgestraft, ze beperken enkel het aantal 'goede' aanbevelingen die kunnen voorkomen in de lijst van  $x$  items.
- De volgorde van aanbevelingen is van belang, dus het systeem zal eerst de items aanbevelen waar het het meest zeker van is.

### 1.3.1 mean Average Precision

Een eerste evaluatiemaat waarvan gebruikt gemaakt werd in deze masterproef is de gemiddelde precisie (mAP). De mAP is een metriek die de nadruk legt op de eerste items die aanbevolen worden. Er kan bovendien een drempelwaarde opgegeven worden voor het maximum aantal aanbevelingen dat voor elke gebruiker gegenereerd mag worden zodat we kunnen spreken van een afgeknotte mAP. We veronderstellen dat de drempelwaarde  $\tau = 50$ .

**Precisie (P)** Om tot de mAP te komen wordt vertrokken vanuit formule 1.2. De matrix  $M \in \{0, 1\}^{m \times n}$  is een binaire gebruiker-item matrix, waarbij  $M_{u,i} = 1$  betekent dat de gebruiker  $u$  naar item  $i$  geluisterd heeft. De lijst van aanbevelingen wordt voorgesteld door  $y$  waarbij  $y(j) = i$  betekent dat item  $i$  voorkomt op positie  $j$  in de lijst. Verder veronderstellen we dat er

in een lijst  $y_u$  voor een gebruiker  $u$  geen nummers voorkomen die hij reeds beluisterd heeft. We willen immers nieuwe nummers aanbevelen voor elke gebruiker.

$$P_k(u, y) = \frac{1}{k} \sum_{j=1}^k M_{u,y(j)} \quad (1.2)$$

In bovenstaande formule wordt de precisie op positie  $k$  voor een gebruiker  $u$  en een lijst van aanbevelingen  $y$  berekend. Dit kan gezien worden als de proportie van juist voorspelde items binnen de top- $k$  van de aanbevelingen.

**Gemiddelde Precisie (AP)** Voor elke gebruiker wordt nu de gemiddelde precisie berekend:

$$AP(u, y) = \frac{1}{n_u} \sum_{k=1}^{\tau} P_k(u, y) \cdot M_{u,y(j)} \quad (1.3)$$

Hierbij stelt  $n_u$  het minimum voor van het aantal aangerade muzieknnummers en het aantal nummers voor die gebruiker die in de testset voorkomen.  $M_{u,y(j)}$  is opnieuw een binaire matrix die waarde 1 bevat als gebruiker  $u$  geluisterd heeft naar het item dat zich bevindt in de aanbevelingslijst op positie  $k$ . Een voorbeeld van de berekening van de gemiddelde precisie wordt weergegeven in tabel 1.1.

**Afgeknotte Gemiddelde Precisie (mAP)** In een laatste stap wordt dan uitgemiddeld over alle  $m$  gebruikers om de mAP te bekomen:

$$mAP = \frac{1}{m} \sum_u AP(u, y_u) \quad (1.4)$$

### 1.3.2 Recall

Een tweede metriek die van pas komt bij aanbevelingssystemen is de recall. Recall is de verhouding tussen het aantal relevante aangeraden items en het totaal aantal relevante items voor die gebruiker. Wanneer we voor iedere gebruiker telkens  $\tau$  aanbevelingen genereren kunnen we de recall als volgt formuleren:

| Gebruiker u                                  |           |             |
|--|-----------|-------------|
| j  | relevant? | $P_k(u, y)$ |
| 1  | ja        | 1.00        |
| 2  | nee       | 0.5         |
| 3  | nee       | 0.33        |
| 4  | nee       | 0.25        |
| 5  | ja        | 0.40        |
| 6  | nee       | 0.33        |
| 7  | nee       | 0.29        |
| 8  | nee       | 0.25        |
| 9  | ja        | 0.33        |
| 10   | nee       | 0.30        |
| $AP(u, y) = \frac{1.00+0.40+0.33}{3} = 0.58$ |           |             |

TABEL 1.1: Berekening van de gemiddelde precisie voor een bepaalde gebruiker wanneer  $\tau = 10$  nummers worden aangeraden en 3 nummers in de testset zitten.

$$Rec = \frac{1}{m} \sum_u \frac{1}{r} \sum_{j=1}^{\tau} M_{u,y(j)} \quad (1.5)$$

Hierin stelt  $r$  het aantal nummers voor gebruiker  $u$  die in de testset voorkomen.

### 1.3.3 Area Under Curve

Een alternatieve manier om aanbevelingen te evalueren is om niet te kijken naar de  $\tau = 50$  beste aanbevelingen alleen, maar de volgorde van aanbevelingen te gaan bestuderen. In deze masterproef werd gebruik gemaakt van de *Receiver Operating Characteristic curve* (ROC-curve). De ROC-curve vergelijkt de true positive rate (TPR) met de false positive rate (FPR) bij verschillende drempelwaarden. Om een ROC-curve op te bouwen wordt vertrokken van een *confusiematrix*. Een confusiematrix is een tabel met 2 rijen en 2 kolommen die het aantal echt positieve, vals positieve, echt negatieve en vals negatieve gevallen meet tijdens testfase.

|              |                    | echte waarde       |  |
|--------------|--------------------|--------------------|--|
| voorspelling | Echt positief (EP) | Vals positief (VP) |  |
|              | Vals negatief (VN) | Echt negatief (EN) |  |

TABEL 1.2: Een confusiematrix

Ook de evaluatie van een lijst met aanbevelingen kan worden opgedeeld in 4 disjuncte gevallen:

| voorspelling | echte waarde                 |                                   |
|--------------|------------------------------|-----------------------------------|
|              | aanbevelen + in testset      | aanbevelen + niet in testset      |
|              | niet aanbevelen + in testset | niet aanbevelen + niet in testset |

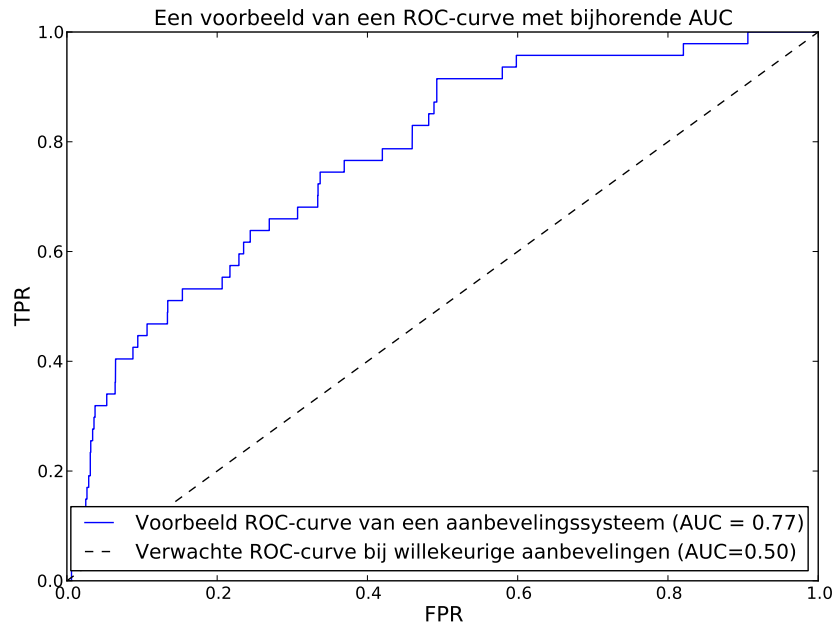
TABEL 1.3: Een confusiematrix toegepast op een aanbevelingssysteem

De true positive rate en false positive rate wordt bekomen aan de hand van volgende formules:

$$TPR = \frac{EP}{EP + VN} \quad (1.6)$$

$$FPR = \frac{VP}{VP + EN} \quad (1.7)$$

Bij evaluatie van het aanbevelingssysteem worden nu alle aanbevelingen in beschouwing genomen. Er wordt gestart met de hoogste score als drempelwaarde te zetten waarna men stap voor stap de drempelwaarde verlaagt tot uiteindelijk alle scores aangeraden worden. In elke stap worden de TPR en FPR berekend en tegenover elkaar uitgezet. Wanneer men deze punten nu met elkaar verbindt, bekomt men een ROC-curve. Een ROC-curve begint altijd in het punt (0,0) en eindigt in het punt (1,1). Een goed aanbevelingssysteem zal snel een hoge TPR bereiken in verhouding tot de FPR. Vandaar dat de oppervlakte onder de curve (AUC) een goede manier is om een aanbevelingssysteem te evalueren. De AUC kan dus bepaald worden door de ROC-curve te integreren en kan een maximale waarde van 1.0 aannemen. Wanneer alle nummers in een willekeurige volgorde worden aanbevolen is de verwachte  $AUC = 0.5$ .



FIGUUR 1.2: Een voorbeeld van een ROC-curve

### 1.3.4 Novelty

Men kan zich de vraag stellen of metrieken enkel gebaseerd op relevantie de beste metrieken zijn. Als het nieuwste album van een zeer bekende popgroep op Amazon wordt opgezocht dan bevatten de aanbevelingen op de resultatenpagina ook vaak de andere albums van diezelfde popgroep. Voor een webwinkel zoals Amazon zal dit wellicht zijn nut hebben maar de kans is echter vrij groot dat de gebruiker reeds weet heeft van deze albums. Bij het luisteren van muziek is men echter ook vaak op zoek naar nieuwe muziek. Een zelf ontworpen metriek geeft de *nieuwheid* van de aanbevelingen weer:

$$Nov = \frac{1}{m} \sum_u \sum_{i \in y(i)} \frac{\tau}{\log(total\_plays_i)} \quad (1.8)$$

$total\_plays_i$  geeft het totaal aantal afspelerbeurten van een nummer  $i$  terug, berekend over de testset.

Een metriek die gebaseerd is op gebruikerstevredenheid is uiteraard de beste optie. Het nadeel van deze soort metrieken is echter dat deze vrij omslachtig zijn aangezien een eindgebruiker telkens de aanbevelingen moet evalueren.

## 1.4 Datasets

### 1.4.1 Million Song dataset (MSD)

In deze masterproef werd vertrokken van een subset van een dataset die toegankelijk werd gesteld in het kader van de Million Song Dataset Challenge [6]. Deze subset bestaat uit 10000 muzieknummers en 20000 gebruikers die al dan niet geluisterd hebben naar een muzieknummer. Een record uit de dataset bestaat uit een gebruikersid, een muzieknummerid en een getal dat correspondeert met het aantal afspeelbeurten. Elk muzieknummerid kan gelinkt worden aan offline metadata die beschikbaar werd gesteld door The Echo Nest [7]. Deze metadata bestond o.a. uit titel, nationaliteit van de artiest, akoestische informatie zoals toonaard en pitch informatie en een verwijzing naar een kort audiofragment van het nummer. Deze audiofragmenten worden gebruikt als basis om de MFCC-vectoren te berekenen. Verder kan elk muzieknummerid gekoppeld worden aan een trackid. Dit trackid wordt gebruikt in de MusiXmatch dataset (zie 1.4.2) en de Last.fm dataset (zie 1.4.3).

|  |                     |    |
|--|---------------------|----|
| 10cbcd627472477dfbec90fb75017f8df6ce84ec | S0GPLBE12A58A80442  | 1  |
| 10cbcd627472477dfbec90fb75017f8df6ce84ec | S0BWSGV12AB018B5E0  | 1  |
| 10cbcd627472477dfbec90fb75017f8df6ce84ec | S0WINIH12AB018CC51  | 1  |
| 10cbcd627472477dfbec90fb75017f8df6ce84ec | S0XHVRT12A81C2320D  | 1  |
| e9dc6b4c2b22aa6dc8260e1963021567728055b2 | S0UIGCD12AB0186713  | 1  |
| e9dc6b4c2b22aa6dc8260e1963021567728055b2 | S0YOMRA12A6D4F9975  | 11 |
| e9dc6b4c2b22aa6dc8260e1963021567728055b2 | S0LMIUU12A58A79C99  | 3  |
| e9dc6b4c2b22aa6dc8260e1963021567728055b2 | S0M0BBZ12A8C144831  | 8  |
| e9dc6b4c2b22aa6dc8260e1963021567728055b2 | S0PXXYD12A6D4FA876  | 6  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0TQBKC12A8C13E961  | 1  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0AIBYI12AB0185C5B  | 1  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0S0M0NK12A8C139059 | 1  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0MYDVM12A6D4F7935  | 1  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0HXBKV12A6D4F820D  | 2  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0DZ0ER12A8C1360FB  | 1  |
| 30e4a688e6fc9c8bfe55998af3996a909ae34449 | S0QGYDJ12A6310F359  | 1  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0KWSEA12A8C141C9D  | 1  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0PGCXT12A8C138AD1  | 5  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0DDDW12AB0188D17   | 1  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0TVJCB12A8C136E46  | 1  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0UFTBI12AB0183F65  | 6  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0SXLTC12AF727F54   | 9  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0SNJIT12A8159E8DB  | 2  |
| 18ce1da0e1017e31baaa5f80afa64ee3c7fab379 | S0KENKR12AB01828F7  | 7  |
| 0b254c684efb08fd04933add2d1e4191d2a87bac | S00PUTL12A8C143381  | 1  |
| 0b254c684efb08fd04933add2d1e4191d2a87bac | S0POANU12A8AE48C9B  | 2  |
| 0b254c684efb08fd04933add2d1e4191d2a87bac | S0UDEUC12A6D4F95A7  | 1  |
| 0b254c684efb08fd04933add2d1e4191d2a87bac | S0EB0WM12AB017F279  | 1  |

FIGUUR 1.3: Een voorbeeldfragment uit de MSD. Links bevindt zich het gebruikersid, in het midden het muzieknummerid en rechts het aantal afspeelbeurten.

### 1.4.2 MusiXmatch dataset

De musixmatch dataset is een verzameling van songteksten in bag-of-words formaat waarbij enkel gekeken werd naar de 5000 populairste woorden, die we topwoorden noemen. De eerste regel bevat deze topwoorden in dalende volgorde van populariteit. Elke daaropvolgende regel uit de dataset bestaat uit een trackid, een id van musixmatch en een lijst koppels van de vorm  $\langle idx : cnt \rangle$ . Dit correspondeert met het aantal voorkomens  $cnt$  van het woord met index  $idx$  in de lijst van topwoorden. Deze dataset is gelinkt aan de MSD via het trackid.

```
TRAAEDY12903CBBFA7,4193161,1:5,23:11,24:4,42:6,72:1,90:7,156:1,190:1,229:1,243:4,251:1,474:7,482:12,530:1,544:2,553:1,603:1,632:15,650:15,695:1,710:8,826:2,978:12,1022:1,1246:2,1276:1,1328:2,1349:6,1380:1,1452:30,1641:10,1673:1,2105:1,2414:19,2533:1,2700:2,2747:12,2806:2,2825:1,2843:4,3055:1,3178:1,3271:1,3359:2,3665:1,3989:5,4189:2,4199:1,4213:1,4502:2,4633:2,4684:5,4815:2
TRAAEEQ128F42180B2,3561951,47:1,408:1,665:1,1057:1
TRAAEHW128F9344FD3,5967283,1:5,2:7,5:3,6:2,8:3,10:2,11:5,12:8,13:1,18:1,19:2,22:1,23:2,24:6,28:1,30:2,32:1,39:2,40:1,43:6,51:1,55:3,67:1,75:6,82:1,87:1,95:6,125:1,133:1,146:2,163:1,180:1,195:1,283:1,291:1,337:1,338:1,403:1,455:1,491:1,497:1,503:2,535:1,548:1,569:6,584:1,596:2,676:1,753:1,784:1,877:1,1035:1,1097:1,1114:1,1223:1,1525:1,1888:1,2069:1,2614:1,2774:1,3175:6,4089:1,4525:1
TRAAEJH128E0785506,1018402,1:9,2:12,3:2,4:3,5:1,7:1,9:2,11:8,16:1,17:4,20:4,21:1,25:2,26:1,29:1,30:1,31:2,32:10,34:1,35:1,37:1,44:1,51:1,53:1,55:2,57:1,59:1,61:1,62:1,67:1,72:2,73:3,74:2,89:1,108:1,113:1,120:1,123:1,153:2,154:3,155:1,164:1,166:10,173:1,199:1,229:1,240:1,245:1,254:1,275:1,303:1,310:1,324:1,337:1,347:2,352:11,444:2,497:1,536:1,602:3,635:1,668:1,914:1,3314:8
TRAAEJV128F423CF04,1861819,1:11,3:19,6:3,7:1,10:1,12:7,13:1,17:1,22:1,23:1,25:2,26:1,30:3,33:9,36:3,43:3,56:7,61:1,68:1,82:3,84:2,85:1,88:1,89:1,103:1,133:1,186:1,211:1,247:1,264:3,360:1,419:1,499:8,617:3,680:1,2766:1
```

FIGUUR 1.4: Een voorbeeldfragment uit de musixmatch dataset.

### 1.4.3 Last.fm dataset

De Last.fm dataset is een dataset die tag-informatie van muzieknummers en similariteit tussen verschillende muzieknummers bevat. Deze informatie kan ook online opgevraagd worden met behulp van de Last.fm API<sup>1</sup> wanneer de methoden 'Track.getTopTags' en 'Track.getSimilar' gebruikt worden. De dataset bestaat uit een verzameling van json-gecodeerde tekstbestanden. Elk muzieknummer correspondeert met 1 bestand. Ook deze dataset is gelinkt aan de MSD via het attribuut 'track.id'.

```
{"artist": "The Bar-Kays",
 "timestamp": "2011-08-03 15:26:43.246805",
 "similar": [],
 "tags": [{"soul", "100"}, {"funk", "40"}, {"60s", "20"}, {"Rhytm and Blues", "20"}, {"Saturday Night", "20"}, {"REAL soul music", "20"}, {"rock top", "20"}, {"rock top funky soul", "20"}, {"title is a full sentence", "20"}, {"groove", "0"}, {"you cant sit down lp version", "0"}],
 "track_id": "TRTOKHS128F147ED44", "title": "You Can't Sit Down (LP Version)"}
```

FIGUUR 1.5: Een voorbeeldbestand uit de Last.fm dataset.

Alle aanbevelingssystemen baseren zich op de data die zich in de MSD bevindt. In hoofdstuk 3 wordt gebruik gemaakt van de Last.fm dataset en de MusiXmatch dataset.

<sup>1</sup>[www.last.fm/api](http://www.last.fm/api)

## Hoofdstuk 2

# Collaboratieve Filtertechnieken

In dit hoofdstuk wordt de klasse van collaboratieve filtertechnieken uitvoeriger besproken. We bespreken de buurbaseerde modellen en de verborgen factormodellen. Vervolgens wordt een voorbeeld van beide modellen samen met twee eenvoudige aanbevelingsalgoritmes vergeleken tijdens het genereren van aanbevelingen wanneer een deel van de Million Song dataset (zie 1.4.1) gebruikt wordt als invoerdata. Het overige deel vormt dan de testset.

### 2.1 Buurbaseerde modellen

Een eerste soort van collaboratieve filtertechnieken zijn de buurbaseerde modellen. De voorspelde score voor een bepaalde gebruiker en een bepaald item wordt bepaald op basis van de score van de *buren*. Deze burens kunnen zowel items als gebruikers zijn. In het eerste geval spreekt men van een gebruikers-georiënteerd model. In deze masterproef werd echter een item-georiënteerd model geïmplementeerd. Dit houdt in dat voor een bepaalde gebruiker items worden aanbevolen die als goede burens beschouwd kunnen worden van eerder positief bevonden items. Dit vereist natuurlijk dat men eerst a priori voor elk item zijn belangrijkste burens kan definiëren. Dit gebeurt door eerst tussen alle items de similariteit te meten op basis van de beschikbare invoerdata met behulp van een similariteitsmaat. Er zijn verschillende gekende similariteitsmaten zoals de Pearson correlatie of de cosinus similariteit. Het zoeken naar een geschikte similariteitsmaat is een belangrijk probleem op zich. Wanneer nu voor een bepaalde gebruiker voorspellingen dienen te worden gemaakt wordt gekeken naar de similariteit tussen nieuwe items en de items waarvoor de gebruiker reeds een rating heeft toegekend. Dit noemt wordt omschreven als de



scorefunctie  $h_{ui}$ . Ook hier bestaan opnieuw vele varianten. Een voorbeeld van een scorefunctie voor een gebruiker  $u$  en een item  $i$  kan bestaan uit de gemiddelde score van de 5 beste burens. Wanneer de voorspelling gebaseerd is op alle reeds beoordeelde items, kan een hoger gewicht worden toegekend aan de belangrijkste burens.

### 2.1.1 Algoritme

De implementatie die in deze masterproef werd gebruikt steunt grotendeels op de winnende oplossing van de Million Song Dataset Challenge [8]. Stel als gegeven een set  $U$  van  $m$  gebruikers en een set  $I$  van  $n$  items. De waarden van de gebruikers-item matrix  $R = \{r_{ui}\} \in \mathbb{R}^{n \times m}$  worden voor de eenvoud voorgesteld door een binaire waarde waarvoor geldt dat  $r_{ui} = 1$  als gebruiker  $u$  reeds geluisterd heeft naar nummer  $i$ . De scorefunctie  $h_{ui}$  ziet er als volgt uit:

$$h_{ui} = \sum_{j \in I(u)} f(w_{ij}) \quad (2.1)$$

De scorefunctie is dus proportioneel met de similariteiten  $w_{ij}$  van het item  $i$  en de items die eerder al beluisterd werden door de gebruiker  $u$  ( $j \in I(u)$ ).

Om verschillende items met elkaar te vergelijken werd gekozen voor de cosinus similariteitsfunctie. De cosinus similariteit is symmetrisch en door het feit dat enkel binaire waarden in de gebruikers-item matrix voorkomen kan deze dus bovendien snel berekend worden als volgt:

$$w_{ij} = \frac{|U(i) \cap U(j)|}{|U(i)|^{\frac{1}{2}} |U(j)|^{\frac{1}{2}}} \quad (2.2)$$

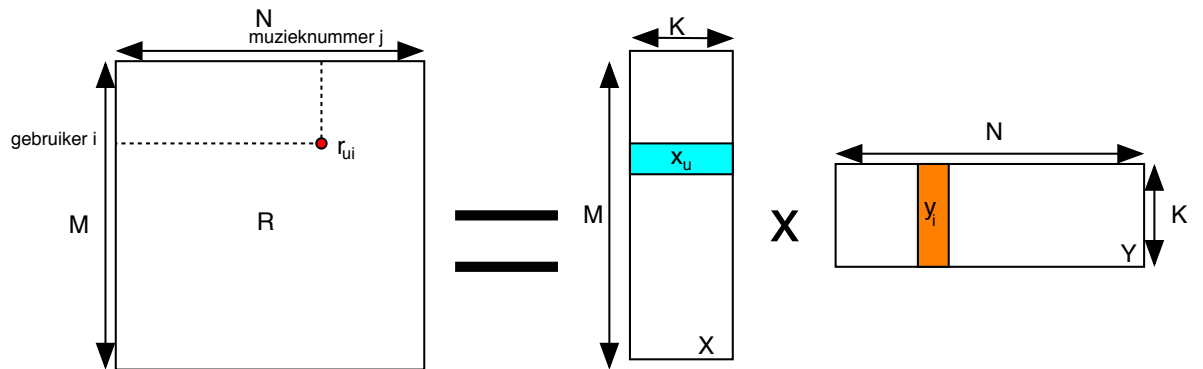
Voor elke gebruiker wordt nu voor de nummers die hij nog niet beluisterd heeft een score bepaald met behulp van vergelijking 2.1. Deze scores worden in dalende volgorde gesorteerd waarna de nummers met de hoogste scores eerst worden aangeraden.

## 2.2 Verborgene Factormodellen

Een alternatieve aanpak ten opzichte van buurtegebaseerde modellen is gebruik te maken van factormodellen. Bij deze wordt geprobeerd de scores tussen gebruikers en items te modelleren door verborgen variabelen bloot te leggen die deze scores verklaren. Binnen deze klasse bestaan

verschillende technieken, maar in deze masterproef werd de focus gelegd op modellen die een matrixfactorisatie uitvoeren op de gebruiker-item matrix. Een eerste eenvoudige aanpak kan bestaan uit het genereren van een singuliere waarden decompositie (SVD) uit de gebruikers-item matrix. Standaard SVD-technieken werken echter in op een matrix waarvan alle waarden gekend zijn. In een realistische dataset zijn er echter veel ontbrekende waarden aangezien de gebruikers niet elk item beoordeeld hebben. Dit zorgt ervoor dat er enige aanpassingen nodig zijn. Bovendien zou een SVD-oplossing erg te leiden hebben onder overfitting. Overfitting is een gekend probleem binnen het domein van machinaal leren en statistiek. Bij overfitting gaat het model relaties leren die specifiek zijn voor de trainingsdata maar die zich niet veralgemenen over alle data. Een typisch verschijnsel dat optreedt bij overfitting is dat het model zeer goed presteert op de trainingsdata. Wanneer hetzelfde model echter geëvalueerd wordt op nieuwe data zal de prestatie plots veel slechter zijn.

Een alternatieve methode vertrekt ook vanuit de gebruiker-item matrix  $R$  die wordt nu wordt benaderd als het product van 2 andere matrices  $X$  en  $Y$  die respectievelijk de gebruikersfactormatrix en de item- of muziekfactormatrix worden genoemd. Indien de gebruikers-item matrix dimensie  $M \times N$  heeft, dan komen de dimensies van  $X$  en  $Y$  overeen met  $M \times K$  en  $K \times N$ , waarbij  $K$  dus vrij te bepalen is. De waarde van  $K$  correspondeert met het aantal verborgen variabelen of factoren dat gebruikt zal worden om de observeerbare scores te verklaren.



$$r_{ui} = \sum_k x_{uk} y_{ik}$$

### 2.2.1 Algoritme

Een score is telkens gebaseerd op een inwendig product van gebruikersfactoren en muziekfactoren. Wanneer de gebruikersfactoren van een gebruiker  $u$  als  $x_u$  worden voorgesteld en de muziekfactoren van een muzieknummer  $y$  als  $y_i$  worden voorgesteld, dan bedraagt de voorspelde score  $x_u^T * y_i$ .

Indien met expliciete scores wordt gewerkt, bijvoorbeeld een systeem waar de gebruiker zelf een score toekent aan een muzieknummer, dan kan de kostfunctie als volgt worden neergeschreven:

$$\min_{x^*, y^*} = \sum_{r_{ui} \text{ is known}} (r_{ui} - x_u^T * y_i)^2 + \lambda(||x_u||^2 + ||y_i||^2) \quad (2.3)$$

Deze kostfunctie minimaliseert het verschil tussen de voorspelde scores en de geobserveerde scores. De parameter  $\lambda$  is de regularisatieparameter en voorkomt dat het model overfit op de trainingdata. Deze kostfunctie neemt dus enkel maar de observeerbare scores in beschouwing en kan geoptimaliseerd worden met behulp van een stapgebaseerde methode zoals bijvoorbeeld stochastische gradient descent.

Vertrekkend vanuit deze kostfunctie kan nu een kostfunctie worden afgeleid die overweg kan met een dataset waar enkel impliciete feedback aanwezig is.

Allereerst observeren we dat een waarde  $r_{ui} = 0$  nu niet overeenkomt met een zeer slechte score, maar gewoon met het feit dat de gebruiker het betreffende nummer nog niet beluisterd heeft. Indien een gebruiker een nummer nog niet heeft afgespeeld, kunnen we maar moeilijk bepalen of de gebruiker het nummer goed zal vinden. Hier zijn verschillende redenen voor. De gebruiker kan nog nooit gehoord hebben van het muzieknummer of hij kan het nummer bijvoorbeeld onmogelijk afspelen. Bovendien is het slechts enkele malen afspelen van een muzieknummer geen sterke indicatie dat de gebruiker het een goed nummer vindt. We weten wel dat naarmate de gebruiker het nummer meerdere keren heeft afgespeeld, we met grotere zekerheid kunnen stellen dat hij het een goed nummer vindt.

De kostfunctie wordt nu uitgebreid met een nieuwe variabele  $c_{ij}$  die het vertrouwen in een bepaalde score modelleert. Deze zal laag zijn wanneer het nummer nog niet werd afgespeeld door de gebruiker maar zal lineair of logaritmisch stijgen volgens het aantal afspeelbeurten.  $r_{ui}$  zelf wordt dan vervangen door de variabele  $p_{ui}$  die 0 of 1 zal aannemen alnaargelang de gebruiker het item al dan niet heeft afgespeeld. Dit brengt ons uiteindelijk tot de volgende kostfunctie:

$$\min_{x^*, y^*} = \sum_{u,i} c_{ui} (p_{ui} - x_u^T * y_i)^2 + \lambda \left( \sum_u \|x_u\|^2 + \sum_i \|y_i\|^2 \right) \quad (2.4)$$

Een belangrijk verschil van deze kostfunctie ten opzichte van vergelijking 2.3 is dat onze nieuwe kostfunctie nu  $M \times N$  termen heeft wat in een typische dataset problemen kan geven m.b.t. de schaalbaarheid van ons model.

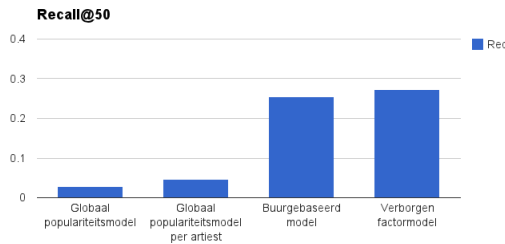
Vergelijking 2.4 kan men opnieuw proberen te optimaliseren met stochastic gradient descent, maar dit zal veel te traag convergeren naar een optimum. Een alternatieve oplossingsmethode is om gebruik te maken van Alternating Least Squares(ALS). Wanneer we de gebruikersfactoren of muziekfactoren berekend hebben, dan wordt de kostfunctie kwadratisch zodat we het globale optimum snel kunnen berekenen. Wat we in één iteratie van het ALS-algoritme doen is 1 van beide types factoren als constant beschouwen en telkens het andere type optimaliseren in functie van die constante en vice versa. In elke stap zullen we immers een lagere totale kostfunctie verkrijgen.

In [9] wordt er een efficiënte manier voorgesteld om dit probleem op te lossen waarbij het optimalisatieproces lineair schaalt in functie van de grootte van de dataset.

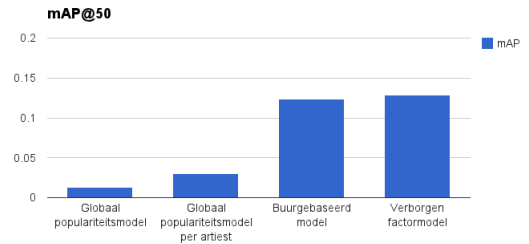
Het product van de muziekfactormatrix en de gebruikersfactormatrix vormt een benadering voor de originele gebruikers-item matrix  $R$ . We kunnen deze benadering  $R'$  nu gebruiken om relevante aanbevelingen te genereren. Concreet zullen we voor elke gebruiker  $i$  de nummers aanbevelen met de hoogste score in  $R'$  voor die gebruiker, op voorwaarde dat hij dat nummer nog niet eerder heeft afgespeeld.

## 2.3 Resultaten

In een eerste fase werden enkele eenvoudige aanbevelingssystemen gebouwd. Het eerste aanbevelingssysteem was gebaseerd op de globale populariteit van alle nummers. Het gaat voor elke gebruiker de nummers aanraden die het meest afgespeeld werden over de volledige dataset op voorwaarde uiteraard dat de gebruiker nog niet naar dat nummer geluisterd heeft. Een iets geavanceerder aanbevelingssysteem gaat eerst gaan kijken welke artiesten een gebruiker reeds beluisterd heeft, en zal enkel de meest globaal populaire nummers aanraden waarvoor de gebruiker minstens al 1 nummer van die artiest eerder heeft beluisterd.



FIGUUR 2.1: De recall op de testset voor de verschillende aanbevelingssystemen wanneer 50 aanbevelingen werden gegenereerd



FIGUUR 2.2: De mAP op de testset voor de verschillende aanbevelingssystemen wanneer 50 aanbevelingen werden gegenereerd

Zoals weergegeven wordt in figuren 2.1 en 2.2, presteerden de geavanceerde aanbevelingssystemen duidelijk beter.

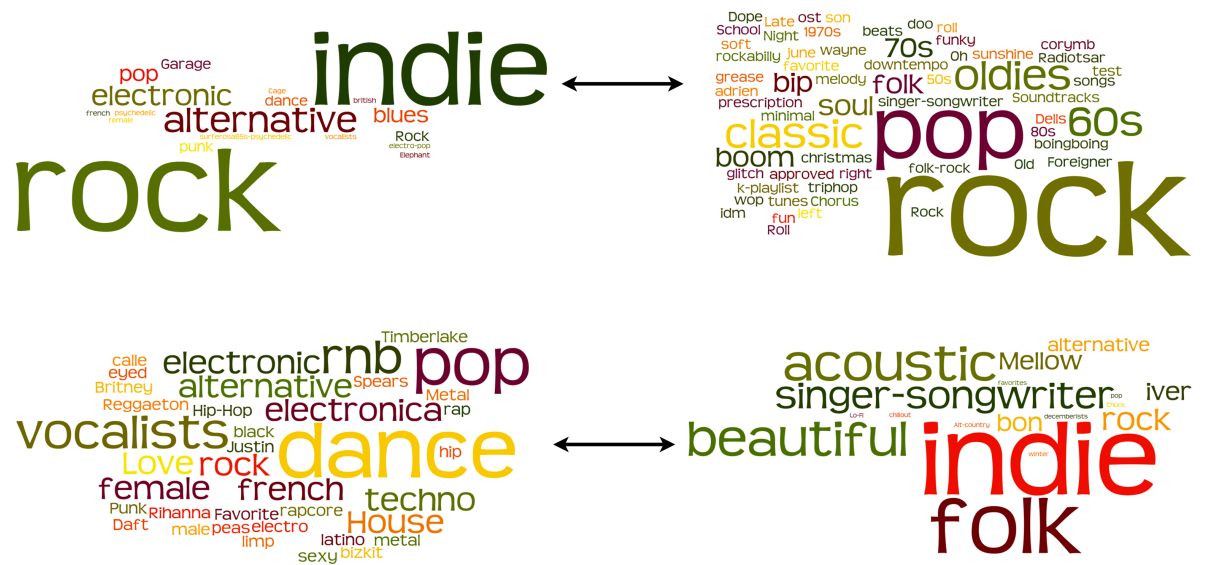
| Aanbevelingssysteem                    | mAP@50  | Rec@50  | Nov     |
|--|---------|---------|---------|
| Globaal populariteitsmodel             | 0.01295 | 0.02739 | 3.90522 |
| Globaal populariteitsmodel per artiest | 0.03063 | 0.04723 | 4.28761 |
| Buurgebaseerd model                    | 0.12405 | 0.25469 | 4.36076 |
| Verborgen Factormodel                  | 0.12830 | 0.27362 | 4.37244 |

TABEL 2.1: De evaluatie van de verschillende aanbevelingssystemen

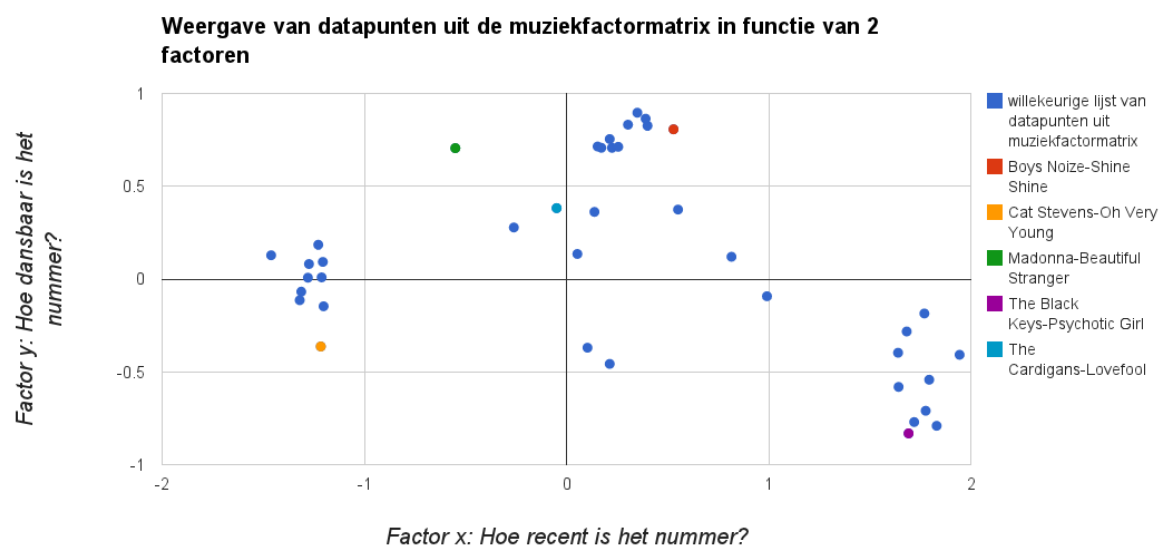
Het grote voordeel van het verborgen factormodel is dat we, vertrekkend vanuit de gebruiker-item matrix, nu een abstracte voorstelling hebben van elk muzieknummer en elke gebruiker in enkele factoren. Deze muzikfactoren kunnen we zien als een soort automatisch gegenereerd profiel gelijkaardig aan wat wordt bekomen bij een inhoudsgebaseerde aanpak. Een bepaalde dimensie of factor kan bijvoorbeeld corresponderen met de geografische locatie van de artiest, de luidheid van een bepaald nummer, dansbaarheid van het nummer,... . De gebruikersfactoren geven aan wat de precieze voorkeuren voor die dimensies zijn. Om dit te verifiëren kan met behulp van regressie-analyse worden nagaan of de muzikfactoren te voorspellen zijn uit externe informatie, bijvoorbeeld uit de akoestische informatie die meegeleverd werd met de initiële dataset (zie 1.4.1). Een verdere besprekening hiervan volgt in hoofdstuk 4. Eenvoudige lineaire regressie waarbij de te voorspellen muzikfactoren uit die beschikbare akoestische kenmerken werden geschat, leerde ons reeds heel wat. Het viel op dat de kenmerken of features die de muziek het meest verklaren, zoals timbre-features en pitch-informatie, de grootste coëfficiënten bevatten in absolute waarde na toepassing van lineaire regressie.

Om nog een betere indicatie van die muzikfactoren te krijgen werden de nummers met de meest extreme waarden voor een factor bepaald waarna bijhorende tag informatie werd opgehaald uit

een Last.fm dataset (zie 1.4.3). Enkele factoren kunnen duidelijk gelinkt worden aan een bepaald kenmerk van het muzieknummer zoals genre of tijdsgeest. Het resultaat hiervan is gevisualiseerd in figuur 2.3. In figuur 2.4 werden de waarden voor deze 2 factoren ten opzichte van elkaar uitgezet. Een belangrijke opmerking is dat de muziektactoren niet absoluut maar relatief moeten geïnterpreteerd worden. Wanneer een muzieknummer A bijvoorbeeld een grotere waarde bezit voor factor  $x$  dan muzieknummer B, dan kan verwacht worden dat nummer A recenter werd uitgebracht dan nummer B. Om een exact tijdsverschil tussen beide nummers te bepalen is de factor echter te onnauwkeurig.



FIGUUR 2.3: Taginformatie van muziknummers die een sterk negatieve waarde versus een sterk positieve waarde hadden voor 2 factoren x en y. Factor x geeft een indicatie van hoe recent het muziknummer is terwijl factor y kan geïnterpreteerd worden als een maat voor dansbaarheid van het nummer.



FIGUUR 2.4: Weergave van enkele datapunten uit de muziekfactormatrix in functie van 2 factoren.

## Hoofdstuk 3

# Het voorspellen van emotie in muziek

Muziek en emotie zitten dicht bij elkaar. Mensen luisteren vaak naar verschillende muziekstijlen op basis van hun stemming. Automatische detectie van emotie van een bepaald nummer kan dus een belangrijke component in muziekaanbeveling zijn. Naast de traditionele collaboratieve filtertechnieken kan de herkenning van emotie in muziek een oplossing bieden voor het koude startprobleem. Emotieherkenning van muziek kan gebaseerd zijn op akoestische features, songteksten of metadata zoals taginformatie van een bepaald nummer. In dit hoofdstuk wordt onderzocht of het mogelijk is om aan de hand van akoestische features en songtekstinformatie emotie te kunnen herkennen in muziek. Er wordt vertrokken van een gelabelde dataset. Dit is een dataset met al deze informatie over de muzieknnummers en telkens de aanduiding van de stemming. Uit deze dataset gaat het algoritme emoties leren herkennen. Het resultaat is dat het algoritme in staat is om emotie te herkennen in nieuwe muzieknnummers.

### 3.1 Data voorverwerking en feature extractie

Een eerste overweging die werd gemaakt is hoeveel emoties er in rekening gebracht moeten worden. Er bestaan vele soorten emoties, de ene al gemakkelijker te herkennen dan de andere. Hoe meer emoties we in rekening brengen, hoe geavanceerder het algoritme zal moeten zijn om het onderscheid te kunnen maken. Daarnaast is het aantal emoties ook een belangrijk aspect tijdens het opstellen van de gelabelde dataset. Het is eenvoudig om het verschil te herkennen

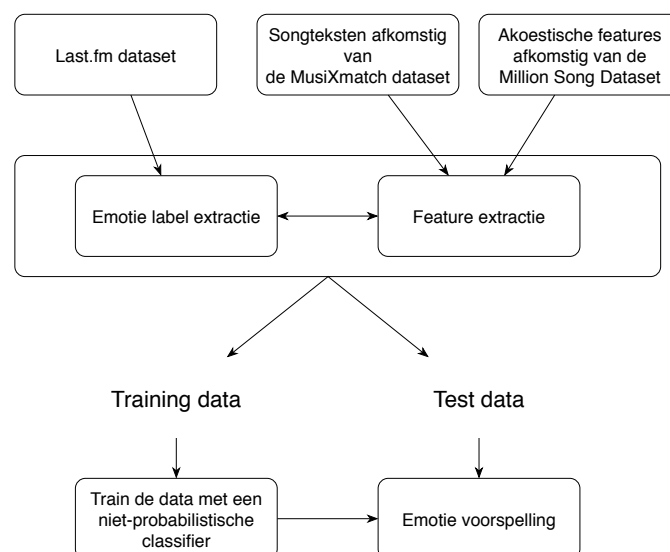


tussen blijdschap en verdriet, maar het verschil tussen verdriet en woede is vaak veel subjectiever en dus moeilijker om te classificeren. In een eerste fase werd enkel onderscheid gemaakt tussen een positieve emotie en een negatieve emotie. In een later stadium werd een negatieve emotie nog opgesplitst tussen verdriet en woede.

Aangezien het manueel annoteren van emotie bij elk nummer tijdsrovend en subjectief kan zijn en in diverse papers [10–12] aangetoond wordt dat sociale taginformatie een goede semantische descriptor is van het muziknummer werden de datasets gelabeld aan de hand van taginformatie. Concreet werd taginformatie van de Last.fm dataset gebruikt (zie 1.4.3). Enkel de muziknummers waarvan de taginformatie woorden bevatte die corresponderen met een positieve of negatieve emotie werden geselecteerd.

De 1e dataset werd bekomen door de 10.000 muziknummers uit de subset van de MSD (zie 1.4.1) te linken aan de Last.fm dataset. De 2e dataset werd gevormd door een ander deel ( $\sim 50.0000$  muziknummers) uit de volledige MSD te koppelen aan de Last.fm dataset. Veel van deze muziknummers kenden echter geen emotiegevoelige taginformatie. Daardoor telden de 2 verschillende datasets uiteindelijk respectievelijk 783 nummers en 5142 nummers. Er waren ongeveer evenveel positief als negatief gelabelde nummers in elke dataset aanwezig. Songtekst-informatie werd verkregen uit de MusiXmatch dataset (zie 1.4.2).

Figuur 3.1 toont een overzicht van dit proces.



FIGUUR 3.1: Opbouw van beide datasets en verdere verwerking vooraleer men aan emotieherkenning kan doen. De akoestische features waren bij de eerste dataset afkomstig van de eigen subset van de MSD. In de tweede dataset werd een groter deel van de MSD opgenomen.

### 3.1.1 Echonest features

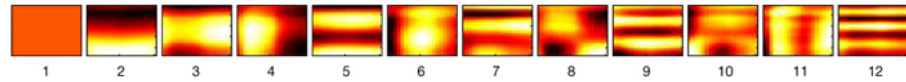
De akoestische features zijn afkomstig van de MSD (zie 1.4.1). In het kader van de MSD-challenge<sup>1</sup> was heel wat informatie reeds offline beschikbaar. Tijdens het bestuderen van deze informatie werd vastgesteld dat bepaalde features zoals *danceability* en *energy* ontbrekende waarden bevatten. Aangezien deze features wellicht ook enige correlatie vertonen met de emotie van een muziknummer, werd besloten om deze features nog eens afzonderlijk op te halen m.b.v. de Echonest API. Dit is een webservice die toegang geeft tot een enorme hoeveelheid aan akoestische informatie over muziknummers. Helaas waren er nog steeds verschillende ontbrekende waarden voor deze features. Een eenvoudige oplossing kan erin bestaan om deze nummers gewoon te verwijderen uit de dataset. De ontbrekende waarden vervangen door de gemiddelde waarde voor de feature bleek echter een beter alternatief.

Echonest deelt een muziknummer op in verschillende segmenten. Een segment is een verzameling van opeenvolgende geluiden die uniform zijn in timbre en harmonie. De duur van een segment ligt typisch onder 1 seconde. Aan elk segment zijn vervolgens verschillende features gekoppeld. De *segment pitches* en *segments timbre* features worden door Echonest beschreven als toongebaseerde en MFCC-gebaseerde features (zie 4.1). Een segment pitch feature bestaat uit een 12-dimensionele vector waarvan de verschillende componenten waarden bevatten tussen 0 en 1. Deze beschrijven de intensiteit van een bepaalde grondtoon in een bepaald segment van het nummer, er wordt dus geen onderscheid gemaakt tussen de verschillende harmonischen van deze grondtoon. De exacte implementatie van de segment timbre feature is onbekend, de enige informatie die gevonden kon worden is dat deze feature correspondeert met een projectie van het spectro-temporele oppervlak op een lagere deelruimte van 12 dimensies. De verschillende basisfuncties van deze projectie worden weergegeven in figuur 3.2. De horizontale as correspondeert met het tijdsverloop terwijl de verticale as overeenstemt met het frequentieverloop. De kleuren variëren van zwart tot wit en corresponderen met de amplitude. Het is moeilijk om de exacte relatie tussen de basisfunctie en het waargenomen geluid te bevatten, maar we kunnen bijvoorbeeld de eerste component van de timbre vector omschrijven als de gemiddelde luidheid over het volledige segment. De vierde component van de vector zal een grote waarde bevatten wanneer de segmenten een snellere aanzet hebben. De aanzet of *attack* van een toon is de wijze waarop een hoorbaar geluid begint en zijn maximale luidheid bereikt. Geluiden met een snelle aanzet zijn bijvoorbeeld geweerschoten of het toeslaan van een deur terwijl het traag openen van

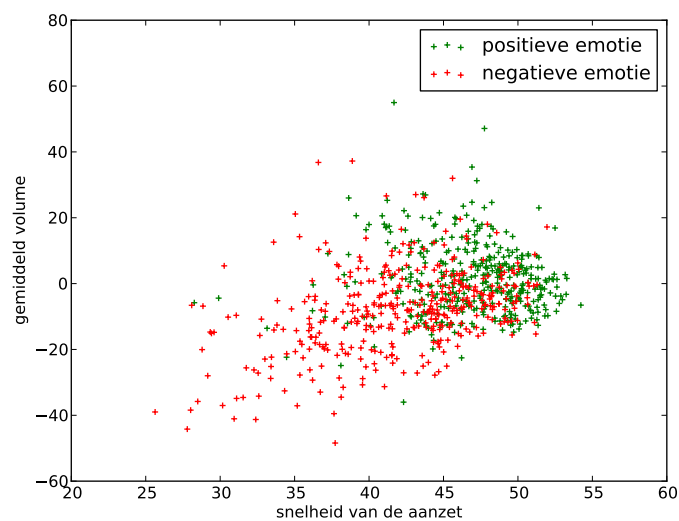
---

<sup>1</sup><http://www.kaggle.com/c/msdchallenge>

een piepende deur dan weer overeenkomt met een trage aanzet. Bij wijze van voorbeeld is in figuur 3.3 de 1e dataset gevisualiseerd in functie van de gemiddelde eerste en vierde component van elk nummer. We kunnen duidelijk zien dat muzieknummers met een positieve emotie eerder een luider volume en een snellere aanzet zullen hebben dan muzieknummers met een negatieve emotie.



FIGUUR 3.2: De 12 basis functies van de timbre vector



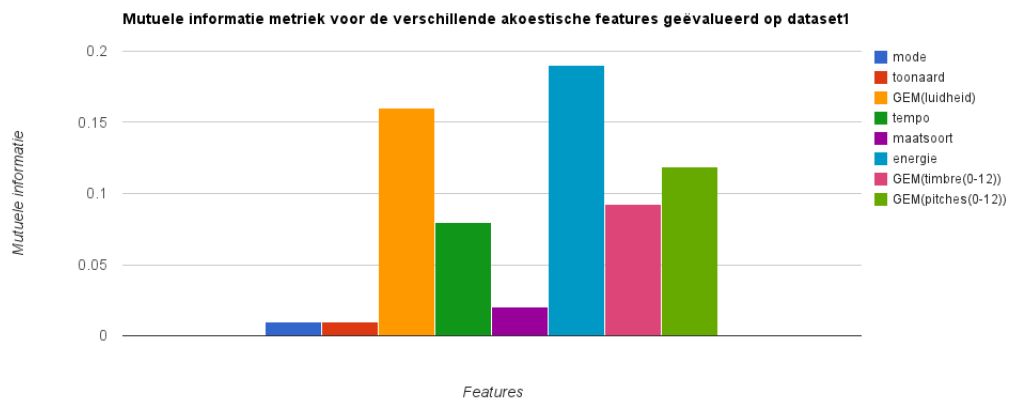
FIGUUR 3.3: De relatie tussen de gemiddelde luidheid en de snelheid van aanzet voor nummers met een positieve en negatieve emotie.

Verschillende muzieknummers hebben typisch een verschillende duur. Het aantal segmenten is dus verschillend per nummer. Om de verschillende features om te zetten naar een tijds-invariante representatie werd beslist om verschillende statistische metrieken van de segments timbre en segment pitches features op te nemen. In de praktijk bleken het gemiddelde, de variantie, de mediaan en de minimale en maximale waarden het beste resultaat op te leveren.

Van de toonaard en het tempo van een muzieknummer zou men ook kunnen verwachten dat ze in sterke mate de emotie van het nummer beïnvloeden. Het correct bepalen van de toonaard van een muzieknummer is echter geen eenvoudige taak, en er werd al snel vastgesteld dat vele waarden niet klopten. Er werden ook verschillende features die te maken hadden met de luidheidsverdeling doorheen een muzieknummer in rekening gebracht.

Vooraleer er verschillende combinaties van features werden uitgetest leek het interessant om in een eerste fase een idee te krijgen welke features het belangrijkste zijn bij emotieherkenning van muziek. Daarvoor werd gebruik gemaakt van de mutuele informatie. De mutuele informatie is de reductie in onzekerheid over een bepaalde variabele  $X$ , of de verwachte reductie in het aantal ja-nee vragen die gesteld moeten worden om de variabele  $X$  te bepalen na het observeren van  $Y$ . In dit geval kan de variabele  $X$  gelinkt worden aan de emotie van het nummer en de variabele  $Y$  komt dan overeen met de feature.

Features die een hoge mutuele informatie hebben m.b.t. de emotie zullen dus leiden tot een betere classificatie.



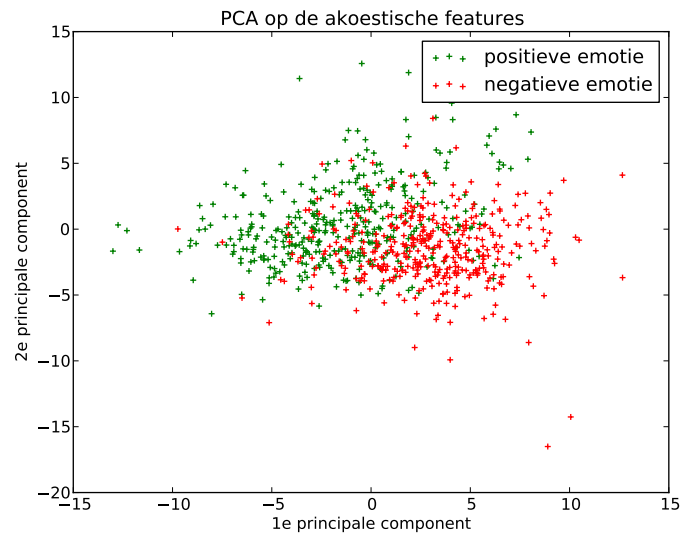
FIGUUR 3.4: Overzicht van de mutuele informatie tussen verschillende features en de emotielabels.

Zoals in bovenstaande figuur zichtbaar is, bestaat de kans dat het incorporeren van de toonaard en de toonsoort (mode) wellicht niet zal leiden tot een betere classificatie. Dit bleek na controle ook te kloppen. Tabel 3.1 bevat de kleinst mogelijke combinatie van features die resulteerden in de beste nauwkeurigheid. Het uitbreiden van deze combinatie met extra features ([7]) ging gepaard met extra rekentijd en resulteerde bovendien niet in een hogere nauwkeurigheid.

Tenslotte werd gebruik gemaakt van principale-componentenanalyse (PCA) om de hoogdimensionale feature vector te kunnen visualiseren. PCA is een dimensionaliteitsreductietechniek die toelaat om een hoeveelheid gegevens te beschrijven met een kleiner aantal relevante features. Dit zorgt ervoor dat een muziknummer met 125 features geschreven kan worden in functie van 2 features, die we nu principale componenten noemen, zodat het visueel geïnterpreteerd kan worden. Het resultaat is te zien in onderstaande figuur, waarbij reeds kan worden afgeleid dat het bouwen van een classifier die emoties kan onderscheiden een complex probleem is.

| No.     | feature                              |
|---------|--------------------------------------|
| 0-11,   | Gemiddelde van elke timbre component |
| 12-23   | Variantie van elke timbre component  |
| 24-35   | Mediaan van elke timbre component    |
| 36-47   | Minimum van elke timbre component    |
| 48-59   | Maximum van elke timbre component    |
| 60-71,  | Gemiddelde van elke pitch component  |
| 72-83   | Variantie van elke pitch component   |
| 84-95   | Mediaan van elke pitch component     |
| 96-107  | Minimum van elke pitch component     |
| 108-119 | Maximum van elke pitch component     |
| 120     | Tempo                                |
| 121     | gemiddelde luidheid                  |
| 122     | variantie luidheid                   |
| 123     | hoogste-laagste luidheid             |
| 124     | Energy                               |

TABEL 3.1: De verzameling van akoestische features die gebruikt werd bij classificatie.



FIGUUR 3.5: De visualisatie van de 1e dataset met behulp van PCA.

### 3.1.2 Songtekst features

Het specifieke formaat van de songtekst dataset staat beschreven in 1.4.2. De songteksten werden voorgesteld in een bag-of-words formaat: elk nummer werd beschreven aan de hand van het aantal keer dat een topwoord voorkwam. Een woord is een topwoord indien het tot de 5000 populairste woorden behoort, berekend over alle songteksten. Helaas waren de meeste topwoorden niet significant ('the, and, we, I') wanneer we aan emotieherkenning willen doen.

Daarom werd een subset gevormd van deze topwoorden. Deze subset ontstond uit de doorsnede van de topwoorden met verschillende gevoelswoord datasets zoals General Inquirer [13] en AFINN [14]. Dit zorgde ervoor dat er slechts 2049 topwoorden overbleven. Daarna werd geëxperimenteerd met verschillende feature representaties zoals binair (komt een gevoelswoord voor in de songtekst?) of een term frequency-inverse document frequency (tf-idf) representatie [15]. Deze waarde bepaalt hoe belangrijk het voorkomen van elk woord is in de dataset. De tf-idf waarde stijgt proportioneel volgens het aantal keer dat het topwoord voorkomt in 1 songtekst, maar is begrensd door de frequentie van het topwoord over de volledige dataset.

De beste resultaten werden echter verkregen wanneer er gewoon met de frequentie van elk topwoord in het muzieknnummer werd gewerkt.

## 3.2 Classificatie en Experimenten

### 3.2.1 Classificatie

Voor muziekclassificatietaken worden vaak k-nearest neighbour (K-NN) en support vector machines (SVM) gebruikt [16]. Mede door hun recent toenemende populariteit werd gekeken of emotieherkenning aan de hand van Random Forests ook goeie resultaten opleverde. De implementatie van de verschillende classifiers is afkomstig uit scikit-learn 0.12.1 [17], een bibliotheek die verschillende machine-learning algoritmes bevat voor de Python programmeertaal. Er werd telkens gebruik gemaakt van 10-fold kruisvalidatie in combinatie met grid search om de hyperparameters voor de verschillende classifiers te bepalen. K-fold kruisvalidatie houdt in dat de data wordt opgesplitst in K verschillende gelijke delen waarna er vervolgens telkens 1 deel afgesplitst wordt om de classifier op te testen, terwijl de overige K-1 delen gebruikt worden tijdens het trainen. Dit proces wordt K keer herhaald waarna de uiteindelijke nauwkeurigheid wordt bepaald door het gemiddelde te nemen van de K nauwkeurigheden op de testsets. Grid search is een manier om de hyperparameters van een classifier te bepalen. Bij grid search wordt gewoon exhaustief gezocht in een subset van de hyperparameterterruimte. Tabel 3.2 bevat een overzicht van de verschillende hyperparameters en hun waarden voor de classifiers die enkel gebruik maken van akoestische features en de classifiers die enkel gebruik maken van songtekstinformatie.

| hyperparameters bij classificatie op basis van akoestische informatie |  |                            |                            |                |                |
|---|--|----------------------------|----------------------------|----------------|----------------|
| classifier  | hyperparameters                                | dataset 1                  |                            | dataset 2      |                |
| K-Nearest Neighbors (KNN)   | Aantal buren                                   | k = 6                      | k = 8                      |                |                |
| Support Vector Machines (poly)  | Straf van de foutterm<br>graad van de polynoom | C = 100<br>graad = 3       | C = 10<br>graad = 1        |                |                |
| Support Vector Machines (RBF)   | Straf van de foutterm<br>kernel coëfficiënt    | C = 1<br>$\gamma$ = 0.0001 | C = 10<br>$\gamma$ = 0.001 |                |                |
| Random Forest   | Aantal bomen                                   | n = 400                    | n = 500                    |                |                |
| hyperparameters bij classificatie op basis van songtekstinformatie    |  |                            |                            |                |                |
|   | hyperparameters                                | alle woorden               | gevoelswoorden             | alle woorden   | gevoelswoorden |
| Lineair SVM   | Straf van de foutterm                          | C = 0.00001                | C=1.5                      | C = 0.001      | C= 0.5         |
| Naive Bayes   | Smoothing parameter                            | $\alpha$ = 0.01            | $\alpha$ = 5               | $\alpha$ = 1.5 | $\alpha$ = 5   |

TABEL 3.2: De hyperparameters van de verschillende classifiers

### 3.2.2 Experimenten

In een eerste experiment werd de nauwkeurigheid van de verschillende classifiers geëvalueerd. De nauwkeurigheid is gedefiniëerd als het percentage juist voorspelde emoties, gegeven de features van elk nummer. In eerdere studies zoals [18] werd aangetoond dat SVMs de beste resultaten opleveren m.b.t. emotieherkenning in muziek. Na het uitvoeren van het experiment bleek dat Random Forests erin slagen om een gelijkaardige nauwkeurigheid te realiseren. In onderstaande tabel is een overzicht terug te vinden van de resultaten van de verschillende classifiers die bekomen werden na 10-fold kruisvalidatie op beide datasets.

| classifier                     | nauwkeurigheid op dataset 1 | nauwkeurigheid op dataset 2 |
|--------------------------------|-----------------------------|-----------------------------|
| K-Nearest Neighbors (KNN)      | 0.7842                      | 0.8213                      |
| Support Vector Machines (poly) | 0.7906                      | 0.8394                      |
| Support Vector Machines (RBF)  | <b>0.797</b>                | 0.8431                      |
| Random Forest                  | 0.7906                      | <b>0.8454</b>               |

TABEL 3.3: De nauwkeurigheid van verschillende classifiers bij akoestische classificatie.

In een volgende experiment werd onderzocht hoe sterk de nauwkeurigheid zou afnemen indien er een extra emotie werd toegevoegd. Er werd vertrokken vanuit de eerste dataset waarbij een negatieve emotie verder werd opgesplitst in woede en verdriet. Zoals verwacht is er een sterke reductie in nauwkeurigheid.

| classification | nauwkeurigheid |
|----------------|----------------|
| 2 klassen      | 0.7906         |
| 3 klassen      | 0.6926         |

TABEL 3.4: Binomiale en multiklasse classificatie met behulp van Random Forests.

Om een verklaring te vinden voor deze sterke afname werd er onderzocht wat de voorspelling was voor elk muziknummer. Bij dit experiment werd gebruik gemaakt van de Random Forest

classifier. De tabellen geven het percentage voorspelde emoties terug in vergelijking met de echte emoties die gekoppeld zijn aan de muzieknnummers.

| Echte emotie | Voorspelde emotie |          |
|--------------|-------------------|----------|
|              | postief           | negatief |
|              | postief           | 0.8111   |
|              | negatief          | 0.1684   |

TABEL 3.5: De voorspelde en echte emoties bij binomiale classificatie

| Echte emotie | Voorspelde emotie |         |        |
|--------------|-------------------|---------|--------|
|              | postief           | droevig | woede  |
|              | postief           | 0.6045  | 0.2773 |
|              | droevig           | 0.1676  | 0.8113 |
|              | woede             | 0.3633  | 0.1490 |

TABEL 3.6: De voorspelde en echte emoties bij multiklasse classificatie

Op basis van deze resultaten blijkt dat het moeilijker is om muziektracks gekoppeld aan een positieve emotie te onderscheiden dan muziektracks gekoppeld aan een woede emotie dan om het onderscheid te maken tussen een droevige en een woede emotie. Dit lijkt op het eerste gezicht een verrassend resultaat, maar als rekening gehouden wordt met het feit dat de voorspelling enkel gebeurde op basis van akoestische informatie is dit verklaarbaar. Droevige muziek is vaak sereen terwijl nummers die een positief gevoel of woede opwekken vaak veel meer actie bevatten. Verder onderzoek van de gelabelde dataset leidde tot de vaststelling dat muzieknnummers die het woedelabel opgeplakt kregen vaak gerelateerd waren aan elektronische muziek of het metal genre.

Als laatste experiment werd onderzocht of emotieherkenning op basis van songtekstinformatie mogelijk is en of de combinatie van akoestische en songtekstinformatie tot een betere classifier kan leiden. Helaas werd er niet voor elk nummer een bijhorende songtekst gevonden wat ervoor zorgde dat de dataset 1 gereduceerd werd van 783 naar 537 nummers, en dataset 2 van 5142 naar 2886 nummers.

| methode        | aantal woorden | classifier  | nauwkeurigheid op dataset 1 | nauwkeurigheid op dataset 2 |
|----------------|----------------|-------------|-----------------------------|-----------------------------|
| alle woorden   | 5000           | Naive Bayes | 0.6115                      | <b>0.7278</b>               |
|                |                | Linear SVM  | 0.589                       | 0.7174                      |
| gevoelswoorden | 2039           | Naive Bayes | <b>0.6723</b>               | 0.7211                      |
|                |                | Linear SVM  | 0.6462                      | 0.6393                      |

TABEL 3.7: De nauwkeurigheid van verschillende classifiers bij classificatie op basis van songtekstinformatie.

Deze resultaten tonen duidelijk aan dat muziekclassificatie op basis van akoestische features veel beter werkt dan op basis van songteksten. Enkel gevoelswoorden in beschouwing nemen lijkt alleen voordelig te zijn bij een kleinere dataset. Het classificeren van een nummer enkel op basis van zijn songtekst blijkt in de praktijk ook vaak moeilijk te zijn. Allereerst kan een



nummer bestaan uit meerdere emoties, denk maar aan muzieknnummers die droevig beginnen maar eindigen met een positieve noot. Bovendien kunnen songteksten vaak vrij poëtisch zijn waarbij een gevoel of emotie beschreven wordt zonder expliciet gebruik van gevoelswoorden. Een laatste reden die aangewend kan worden is het feit dat in sommige nummers de melodie en songtekst tegenstrijdig met elkaar zijn. Een nummer over liefdesverdriet kan gekoppeld zijn aan een oppeppende melodie.

Een classifier die zowel gebruik maakte van de songteksten en akoestische features haalde een nauwkeurigheid van 82.29% op dataset 2. Dit is dus lager dan een classifier enkel gebaseerd op akoestische features. In [16] wordt echter aangetoond dat deze hybride aanpak wel tot een betere classificatie kan leiden.

## Hoofdstuk 4

# Het voorspellen van muziekfactoren

In dit hoofdstuk wordt de methode beschreven die de muziekfactoren van een verborgen factor-model (zie hoofdstuk 2) nu rechtstreeks uit het muziknummer zelf voorspelt. Dit gebeurt door elk muziknummer eerst naar een compacte voorstelling om te zetten waarna de muziekfactoren voorspeld worden uit die voorstelling. De voorspelde muziekfactoren worden in verschillende experimenten vergeleken met de muziekfactoren die berekend werden met het verborgen factor-model en er wordt nagegaan of deze aanpak het koude startprobleem en het lange staartprobleem kan verhelpen.

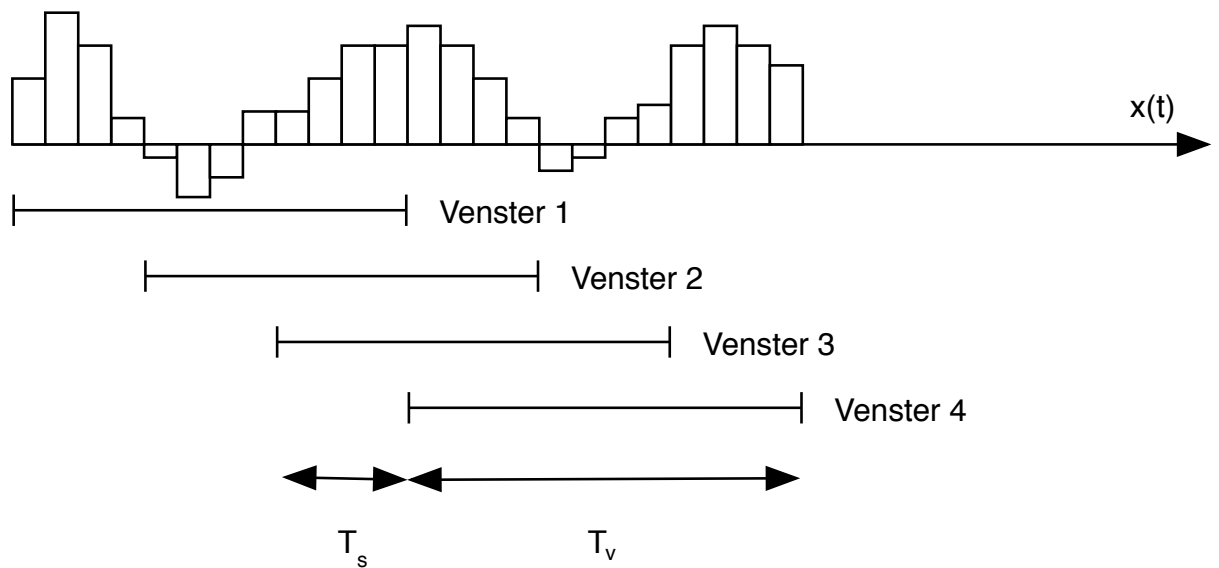
### 4.1 Mel Frequency Cepstral Coefficients (MFCCs)

Geluid is een akoestische golf die correspondeert met drukveranderingen die zich in de lucht voortplanten. Wanneer geluid wordt opgenomen met behulp van een microfoon, dan bekomt men een elektrisch signaal dat kan worden beschreven in functie van de tijd. Men noemt deze functie ook vaak het akoestisch signaal. Dit signaal vertoont op bepaalde ogenblikken een quasi-periodiek en op andere ogenblikken een grillig verloop. Wanneer het signaal periodiek is komt de frequentie ruwweg overeen met de waargenomen toonhoogte. Hoe hoger de frequentie, hoe hoger de toon. Een computer kan echter geen continue functie voorstellen en zal dus dit akoestisch signaal gaan benaderen als een serie van gequantiseerde waarden. Een digitale CD heeft bijvoorbeeld een sampling rate van 44,100 Hz en een bereik van 16 bits. Dit betekent dat elke seconde geluid voorgesteld kan worden door 44,100 getallen of samples die elk voorgesteld kunnen worden door 16 bits. Rechtstreeks muziekfactoren gaan voorspellen op basis van deze getallen is

uiteraard onbegonnen werk. Vandaar dat eerst geprobeerd werd om deze getallen om te zetten naar een compactere representatie.

MFCCs [19] kennen hun oorsprong in spraakverwerking, maar zijn in het verleden reeds succesvol toegepast bij het modelleren van muziekfragmenten [20]. Een MFCC-decompositie probeert een voorstelling te maken van het geluid m.b.t. hoe het menselijk gehoor dit geluid ervaart. Om over te gaan van een digitaal audiosignaal naar een MFCC-voorstelling zijn verschillende stappen nodig.

Eerst wordt het akoestisch signaal  $x(t)$  opgesplitst in verschillende overlappende vensters die elk afzonderlijk zullen worden geanalyseerd. In deze masterproef werd de lengte van een venster  $T_V$  vastgelegd op 25 ms en de verschuivingstijd  $T_S$  bedroeg 10ms.



FIGUUR 4.1: Opsplitsing van het akoestisch signaal in vensters

Daarna kan het signaal door een eenvoudige hoogdoorlaatfilter gestuurd worden. Deze filter modelleert het feit dat het gehoor lage tonen minder goed waarneemt dan hogere tonen. Elk venster wordt vervolgens omgezet vanuit het tijdsdomein naar het frequentiedomein via een discrete fouriertransformatie (DFT). We bezitten nu voor het signaal de intensiteitwaarden bij verschillende frequenties, maar dit komt niet overeen met de intensiteit die mensen waarnemen. Mensen kunnen o.a. een klein verschil in frequentie of toonhoogte beter waarnemen bij lage frequenties dan bij hoge frequenties. Om dit probleem op te lossen wordt de output van de DFT geconverteerd naar de mel schaal. De mel schaal relateert de waargenomen frequentie

van een zuivere toon aan de echte frequentie van die toon. Een populaire methode om van het frequentiedomein over te gaan naar de mel schaal is:

$$mel_{freq} = 2595 * \log_{10} \left[ 1 + \frac{f}{700} \right] \quad (4.1)$$

Een andere belangrijke eigenschap van het gehoor is dat het eerder gevoelig is voor relatieve dan absolute verschillen in intensiteit. Om het waargenomen volume van een geluid te verdubbelen moeten we de geluidsbron van 8 keer zoveel meer energie voorzien. Om dit te modelleren in de MFCC-decompositie is het aangewezen om de logaritme van de mel schaal te nemen. Op die manier kan men ook bij het vergelijken van 2 MFCC-vectoren gebruik maken van Euclidische afstanden.

Om uiteindelijk een discrete voorstelling te verkrijgen in een beperkt aantal coëfficiënten wordt in de laatste stap een discrete cosinustransformatie (DCT) uitgevoerd. Door de specifieke eigenschappen van de DCT zullen de eerste coëfficiënten de belangrijkste frequentie-eigenschappen van het signaal bevatten en de hogere coëfficiënten eerder de details van het spectrum. In deze masterproef werden enkel de eerste 12 coëfficiënten overgehouden om het log-mel-spectrum voor te stellen en vormen ze een MFCC-vector  $C = \langle c_0, c_1, \dots, c_{11} \rangle$ .

Net zoals de opeenvolgende samples waarvan vertrokken werd vaak met elkaar gecorreleerd zijn, zijn ook opeenvolgende MFCC-vectoren met elkaar gecorreleerd. Omdat we in een later stadium enkel MFCC's afzonderlijk gaan behandelen dreigt deze correlatie verloren te gaan. Een mogelijke oplossing bestaat erin om een MFCC van een venster  $n$  uit te breiden met informatie over de MFCC's van zijn omliggende vensters. Dit wordt gemodelleerd door de eerste-orde differenties  $\Delta C_{nk}, k = 0, \dots, 11$ :

$$\Delta C_{nk} = \frac{\sum_{m=1}^{m=2} m(X_{n+m,k} - X_{n-m,k})}{\sum_{m=1}^{m=2} m^2} \quad (4.2)$$

Elke MFCC-vector bestaat nu uit 24 componenten. In deze masterproef werd gebruik gemaakt van de bibliotheek TuneR [21] om de MFCC-vectoren voor verschillende audiosignalen te verkrijgen.

## 4.2 Feature Learning

Rechtstreekse voorspelling van de muziekfactoren op basis van de MFCC's zal niet goed werken omdat deze voorstelling nog steeds te ruw is om verbanden uit te leren. Daarom werd gekeken om MFCC's op een hoger niveau te gaan beschrijven. Dit is een typische voorbeeldtoepassing van *feature learning*, een verzameling van technieken binnen machinaal leren waarbij de features eerst omgezet worden naar een nieuwe ruimte die beter geschikt is voor een gesuperviseerde taak zoals bijvoorbeeld lineaire regressie.

Na het lezen van [22] werd besloten om het K-means algoritme te gebruiken als feature learning algoritme. [22] toont namelijk aan dat andere factoren zoals *whitening* van de data en het aantal nieuwe features een grotere impact hebben op de prestatie van de gesuperviseerde taak dan de complexiteit van het feature learning algoritme zelf. Een bijkomend voordeel van K-means ten opzichte van andere feature learning algoritmes is dat het zeer snel is, enkel het aantal clusters heeft als hyperparameter en dat het eenvoudig te implementeren is.

### 4.2.1 Feature encoding

Vertrekkende vanuit een verzameling MFCC-vectoren wordt nu geprobeerd om een meer compacte voorstelling van elk muziknummer te verkrijgen. Elk muziknummer wordt voorgesteld door een sequentie van 2905 MFCC-vectoren met elk dimensie 24. Een allereerste stap in een machinaal leren algoritme is vaak dat er normalisatie van de invoerdata plaats vindt. Dit betekent dat elke component in de MFCC-vector vervangen wordt door de gemiddelde waarde van die component, berekend over de volledige invoerdata, af te trekken van de originele waarde en vervolgens te delen door de standaardafwijking op deze component. Op die manier heeft elke feature een gemiddelde waarde 0 en bedraagt de standaardafwijking van elke feature 1. Vergelijking 4.3 toont een voorbeeld van input normalisatie:

$$\begin{bmatrix} 1 & -1 & 2 \\ 2 & 0 & 0 \\ 0 & 1 & -1 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 & -1,22... & 1,33... \\ 1,22... & 0 & -0,26... \\ -1,22... & 1,22... & -1,06... \end{bmatrix} \quad (4.3)$$

Verder wordt in [22] ook aangetoond dat whitening van de data resulteert in betere resultaten. Een whitening transformatie transformeert een set van variabelen met correlatiematrix  $M$  in een

set van nieuwe variabelen waarvan de correlatiematrix de eenheidsmatrix is, zodat alle variabelen ongecorreleerd zijn en allemaal variantie 1 hebben.

Een populaire methode om de geobserveerde data te whitenen is om gebruik te maken van de eigenwaardedecompositie van de correlatiematrix van  $X$  [23].

$$\Sigma = XX^T = VDV^T \quad (4.4)$$

$D$  is een diagonaalmatrix die de eigenwaarden van de correlatie bevat en  $V$  is een orthogonale matrix met de eigenvectoren. Whitening gebeurt vervolgens aan de hand van de volgende formule:

$$X_w = VD^{-\frac{1}{2}}V^TX \quad (4.5)$$

Dat de correlatiematrix van  $X_w$  gelijk is aan 1 kan als volgt aangetoond worden [23]:

$$\begin{aligned} \Sigma &= X_w X_w^T \\ &= (VD^{-\frac{1}{2}}V^TX)(VD^{-\frac{1}{2}}V^TX)^T \quad (\text{zie 4.5}) \\ &= (VD^{-\frac{1}{2}}V^T)(XX^T)(VD^{-\frac{1}{2}}V^T)^T \\ &= (VD^{-\frac{1}{2}}V^T)(VDV^T)(VD^{-\frac{1}{2}}V^T)^T \quad (\text{zie 4.4}) \\ &= I \end{aligned} \quad (4.6)$$

In de laatste stap werd gebruik gemaakt van het feit dat  $V$  een orthogonale matrix is zodat  $VV^T = V^TV = I$  en  $D^{-\frac{1}{2}}DD^{-\frac{1}{2}} = I$ .

Nadat deze voorverwerkingsstappen uitgevoerd werden, was het tijd om nieuwe, geschiktere features te leren die de muziekfactoren kunnen voorspellen. Het K-means clusteralgoritme deelt de invoerdata op in verschillende clusters. Elke cluster bezit 1 centrum. Dit centrum correspondeert met de gemiddelde waarde berekend over alle elementen die in de cluster aanwezig zijn. Het algoritme wordt gebruikt om centra te leren uit de invoerdata waarna elke MFCC-vector toegewezen wordt aan een of meerdere centra. Het centrum van de cluster waaraan de MFCC-vector werd toegewezen kan nu worden gebruikt als nieuwe voorstelling. De enige parameter van het algoritme is het aantal clusters  $K$  dat gebruikt wordt om de dataset te clusteren.

De trainingsfase van het K-means algoritme werkt als volgt:

1. Initialiseer  $K$  verschillende clustercentra  $\mu_k$  door bijvoorbeeld  $K$  willekeurige inputvectoren te kiezen.
2. Ken elke inputvector  $x_n$  toe aan het dichtstbijzijnde clustercentrum volgens deze formule:

$$r_{nk} = \begin{cases} 1 & \text{als } k = \operatorname{argmin}_j \|c^{(j)} - x_n\|_2^2 \\ 0 & \text{anders} \end{cases} \quad (4.7)$$

$r_{nk}$  is 1 enkel indien inputvector  $x_n$  werd toegekend aan cluster  $k$  met gemiddelde  $\mu_k$

3. Update de clustercentra:

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}} \quad (4.8)$$

4. Indien er geen wijzigingen meer waren in stap 3 of aan een andere stopvoorwaarde voldaan is, stop. In het andere geval wordt teruggekeerd naar stap 2.

In de praktijk bleek het algoritme snel te convergeren in de eerste iteraties waarna convergentie veel trager verliep zoals werd aangetoond in [24]. 10 iteraties bleek voldoende te zijn om reeds een goede clustering te verkrijgen. Er was dus aan de stopvoorwaarde uit stap 4 voldaan indien het algoritme 10 iteraties had uitgevoerd.

Wanneer de verschillende centra bepaald zijn, kunnen we de MFCC-vectoren nu toekennen aan de centra. Dit kunnen we op meerdere manieren doen. Enerzijds is er de eenvoudige *harde* toekenning, waarbij we elke vector toekennen aan juist 1 clustercentrum:

$$f_k(x) = \begin{cases} 1 & \text{als } k = \operatorname{argmin}_j \|c^{(j)} - x\|_2^2 \\ 0 & \text{anders} \end{cases} \quad (4.9)$$

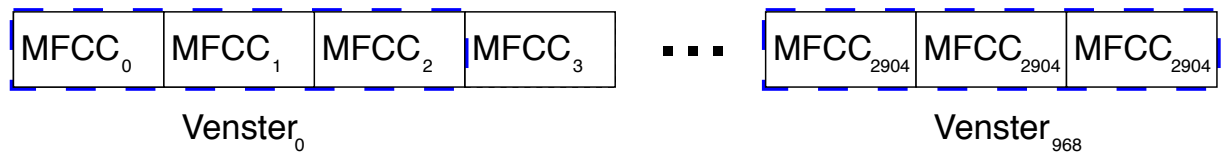
De functie  $f_{hard}(x) : \mathbb{R}^{24} \rightarrow \mathbb{R}^k$  bevat dus  $k - 1$  nulwaarden.

Een andere mapping die in [22] voorgesteld werd, is een *zachte* clustering:

$$f_k(x) = \max \{0, \mu(z) - z_k\} \quad (4.10)$$

waarbij  $z_k = \|x - c^{(k)}\|_2$  en  $\mu(z)$  het gemiddelde is van de elementen van  $z$ . Zachte clustering laat ons toe om iets meer informatie op te slaan dan enkel de dichtstbijzijnde cluster.

Bij zachte clustering is het dus zo dat, afhankelijk van het aantal clustercentra dat gekozen wordt, de dimensie van de geclusterde MFCC-vector kan toenemen terwijl bij harde clustering de geclusterde MFCC-vector telkens kan worden voorgesteld door 1 getal, namelijk de dichtstbijzijnde cluster. Om niet te veel data te genereren bij zachte clustering werd vaak een sequentie van  $l$  opeenvolgende MFCC-vectoren eerst samengebracht in een venster alvorens deze aan te bieden aan het clusteralgoritme. Wanneer er bijvoorbeeld telkens 3 MFCC-vectoren samen genomen worden, dan daalt het aantal vensters per muzieknummer van 2905 naar 969 (zie figuur 4.2).



FIGUUR 4.2: Het samennemen van MFCC-vectoren in vensters.

## 4.3 Voorspellen van muziekfactoren

### 4.3.1 Ridge-regression

Vertrekkende vanuit onze geclusterde MFCC-vectoren werd onderzocht of het mogelijk is om de muziekfactoren te voorspellen. Dit werd gedaan aan de hand van regressie-analyse. Regressie-analyse bestudeert of er een verband bestaat tussen een onafhankelijke variabele en een afhankelijke variabele, waarna de waarde van de onafhankelijke variabele gebruikt kan worden om de afhankelijke variabele te voorspellen:

$$Y = f(X) + U \quad (4.11)$$

$X$  vormt de onafhankelijke variabele en  $Y$  de afhankelijke variabele.  $U$  stelt de storingsterm voor, die onafhankelijk van  $X$  is. De functie  $f$  is onbekend, maar na toepassing van regressie-analyse op een verzameling trainingsvoorbeelden  $\{(x_i, y_i); i = 1, \dots, n\}$  is het mogelijk om  $f$  te karakteriseren aan de hand van enkele parameters.  $X$  zelf kan ook uit meerdere variabelen bestaan. In het bijzonder geval van lineaire regressie kan de functie  $f$  als volgt voorgesteld worden:



$$f(X) = \beta_0 + \sum_{j=1}^d \beta_j X_j \quad (4.12)$$

De onafhankelijke variabelen zijn nu  $X_1, \dots, X_d$ . De onbekende parameter vector  $\beta = (\beta_0, \dots, \beta_d)$  moet geschat worden uit de trainingsvoorbeelden  $\{(x_i, y_i); i = 1, \dots, n\}$ . Een populaire methode maakt gebruik van de kleinste kwadratenmethode, die ervoor zorgt dat de waarden van  $\beta$  bepaald worden door een foutterm te minimaliseren. In het geval van de kleinste kwadratenmethode wordt deze foutterm de residual sum of squares (RSS) genoemd:

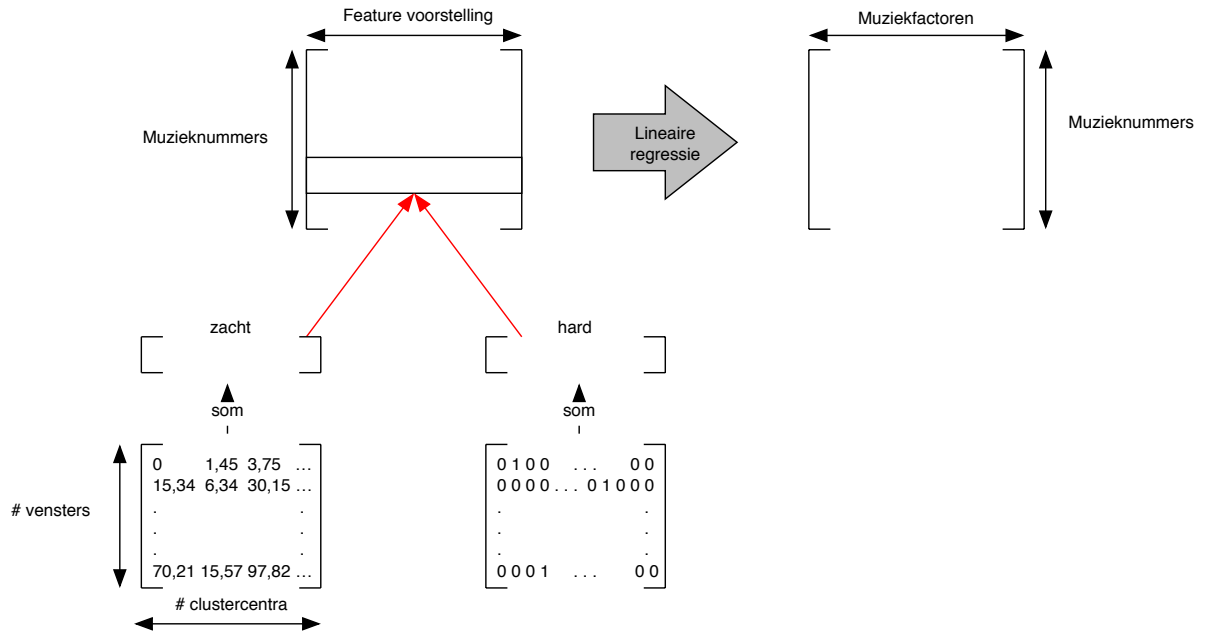
$$RSS(\beta) = \sum_{i=1}^n \|y_i - f(x_i)\|_2^2 \quad (4.13)$$

Deze regressie coëfficiënten zijn optimaal voor de trainingsdata, maar het uiteindelijke doel is om aan de hand van deze regressiecoëfficiënten, nieuwe voorspellingen te genereren. Om het probleem van overfitting op de trainingsdata te vermijden kunnen we gebruik maken van verschillende regularisatiemethoden. Ridge regression is een van de bekendste regularisatiemethoden waarbij de regressiecoëfficiënten worden geregulariseerd door een strafterm toe te kennen die gerelateerd is aan de grootte van de coëfficiënten. In het geval van ridge regression wordt een straf opgelegd aan de kwadratensom van de regressiecoëfficiënten ( $L_2$ -straf). De objectieffunctie kan er nu als volgt uitzien:

$$\beta_{ridge} = \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^n \|y_i - f(x_i)\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2 \quad (4.14)$$

De regularisatieparameter  $\lambda$  wordt bepaald aan de hand van kruisvalidatie.

In figuur 4.3 wordt nog eens een overzicht gegeven van het volledige proces dat gepaard gaat met het trainen van een model dat muziekfactoren kan voorspellen uit MFCC-vectoren. De geclusterde MFCC-vectoren van 1 muzieknummer worden via een somfunctie ingevoegd in een designmatrix  $X$ .



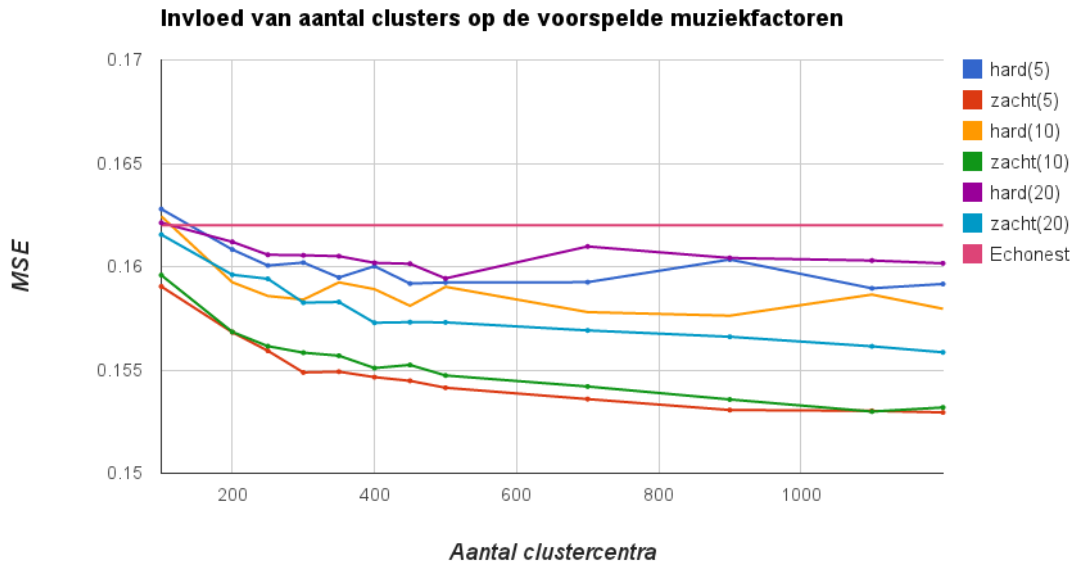
FIGUUR 4.3: overzicht van het trainingsproces van de muziekfactoren uit een geclusterde MFCC-voorstelling

### 4.3.2 Resultaten

De resultaten op de testset werden geëvalueerd aan de hand van de gemiddelde kwadratische fout (MSE). Als  $\hat{y}_i$  de voorspelde muziekfactoren van nummer  $i$  zijn en  $y_i$  de corresponderende echte muziekfactoren, dan wordt de MSE gegeven door:

$$MSE = \frac{1}{n} \sum_{i=1}^n \|\hat{y}_i - y_i\|_2^2 \quad (4.15)$$

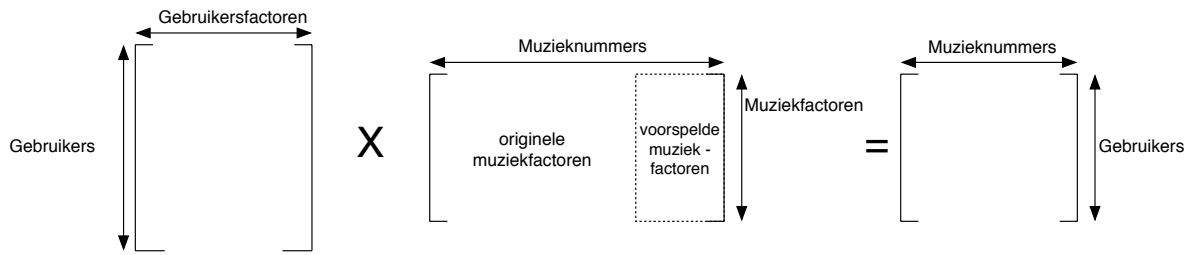
De echte muziekfactoren werden bekomen door het verborgen factormodel te trainen op de initiële dataset waarbij het aantal gebruikers en muziekfactoren  $K = 50$  bedroeg. In figuur 4.4 wordt de invloed van het aantal clustercentra weergegeven in functie van de bekomen MSE voor harde en zachte clustering bij verschillende venstergroottes. De MSE die bekomen werd door de verschillende echonestfeatures uit tabel 3.1 te gebruiken werd als referentiewaarde genomen. We kunnen duidelijk zien dat de nauwkeurigheid van onze voorspellingen stijgt indien we kleinere vensters gebruiken en meer clustercentra. Er werd ook gekeken of het combineren van verschillende venstergroottes in een model tot een betere benadering zou leiden. De kleine reductie in MSE woog echter niet op tegen de extra rekenkracht die gepaard ging met het trainen van dit model.



FIGUUR 4.4: De invloed van het aantal clustercentra op de MSE voor verschillende venstergroottes

#### 4.4 Het genereren van aanbevelingen

Het uiteindelijke doel van het voorspellen van de muziekfactoren is dat ze nauwkeurig genoeg moeten zijn om ook relevante aanbevelingen te maken. Concreet is het de bedoeling om voor nummers waarvoor geen luistergeschiedenis voorhanden is de muziekfactoren te laten voorspellen uit de geclusterde MFCC-vectoren. Om dit te testen werden de voorspelde muziekfactoren achteraf opnieuw ingevoegd in de originele muziekfactormatrix en werden opnieuw aanbevelingen gegenereerd. Om aanbevelingen te genereren nemen we nog steeds het product van de gebruikersfactormatrix en de muziekfactormatrix. De matrix  $R'$  ontstaan uit dit product correspondeert met een benadering voor de originele gebruikers-item matrix  $R$ . Voor elke gebruiker  $i$  worden vervolgens de nummers aangeraden volgens dalende score in  $R'$ , op voorwaarde dat de gebruiker nog niet eerder geluisterd heeft naar dat nummer. Bij het genereren van deze aanbevelingen werd steeds vertrokken vanuit zachte clustering van de MFCC-vectoren met 700 clustercentra en een verborgen factormodel met 50 factoren. De resultaten met deels voorspelde muziekfactoren werden vergeleken met de originele matrixfactorisatiemethode en een geval met willekeurige muziekfactoren. Om goede referentiewaarden te hebben bij willekeurige muziekfactoren worden deze gevormd door de voorspelde muziekfactoren te gaan permuteren. Op die manier hebben de willekeurige muziekfactoren vrij realistische waarden.



FIGUUR 4.5: Het aanbevelingsproces met voorspelde muziekfactoren.

#### 4.4.1 Resultaten

##### 4.4.1.1 Experiment 1

In het eerste experiment werd de prestatie van een aanbevelingssysteem gebaseerd op een volledig verborgen factormodel vergeleken met hetzelfde systeem waar een deel (25%) van de factoren voorspeld werden uit de audio en met een semi-willekeurig systeem waar de voorspelde factoren eerst nog een willekeurige permutatie ondergingen alvorens ze in de muziekfactormatrix werden ingevoegd.

Eerst werden de resultaten geëvalueerd volgens de mAP (zie 1.3). Er werden telkens 50 aanbevelingen gegenereerd waarna de mAP berekend werd over alle gebruikers. Er werd vastgesteld dat de mAP gevoelig daalde in het geval de muziekfactoren volledig of zelfs deels voorspeld werden. Verder bleken deze resultaten niet significant beter te zijn dan in het geval van willekeurige aanbevelingen. Na verder onderzoek werd het duidelijk dat de nummers waarvoor de muziekfactoren voorspeld werden veel minder voorkwamen in de top 50 aanbevelingen in vergelijking met het originele aanbevelingsproces.

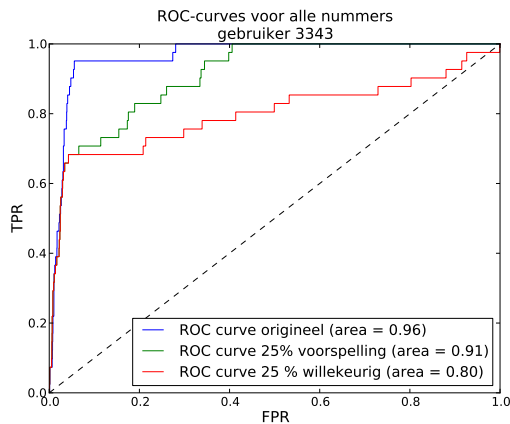
|                        | mAP@50  | aantal voorkomens in de top 50 |
|------------------------|---------|--------------------------------|
| verborgen factormodel  | 0.11706 | 246305                         |
| voorspeld uit de audio | 0.08016 | 3532                           |
| willekeurig voorspeld  | 0.08002 | 0                              |

TABEL 4.1: Een eerste vergelijking van de prestatie van het verborgen factormodel met hetzelfde model waar 25% van de factoren voorspeld werden uit de audio en een semi-willekeurig geval.

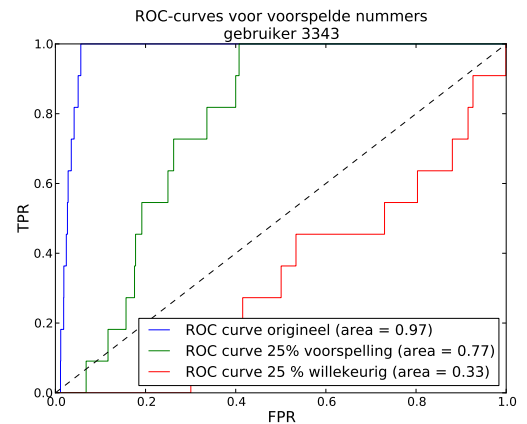
Het vermoeden waarom dit zo is ligt bij het feit dat een regressietechniek elke waarde even goed gaat proberen te benaderen. Muzieknummers waarvan de muziekfactoren veel extreme waarden bevatten hebben een grotere kans om snel aanbevolen te worden. Indien een bepaalde component van de muziekfactor dus een extreme waarde bevat, is het veel belangrijker om die

correct te kunnen voorspellen dan het geval waar die waarde rond de gemiddelde waarde voor die component schommelt.

Een alternatieve metriek die werd voorgesteld in hoofdstuk 1 is de ROC-curve en de bijhorende AUC. In figuren 4.6 en 4.7 worden de ROC-curves weergegeven voor de 3 gevallen wanneer alle nummers in de testset in beschouwing worden genomen en wanneer enkel gekeken wordt naar de nummers waarvan de muziekfactoren voorspeld werden en zich in de testset bevonden.



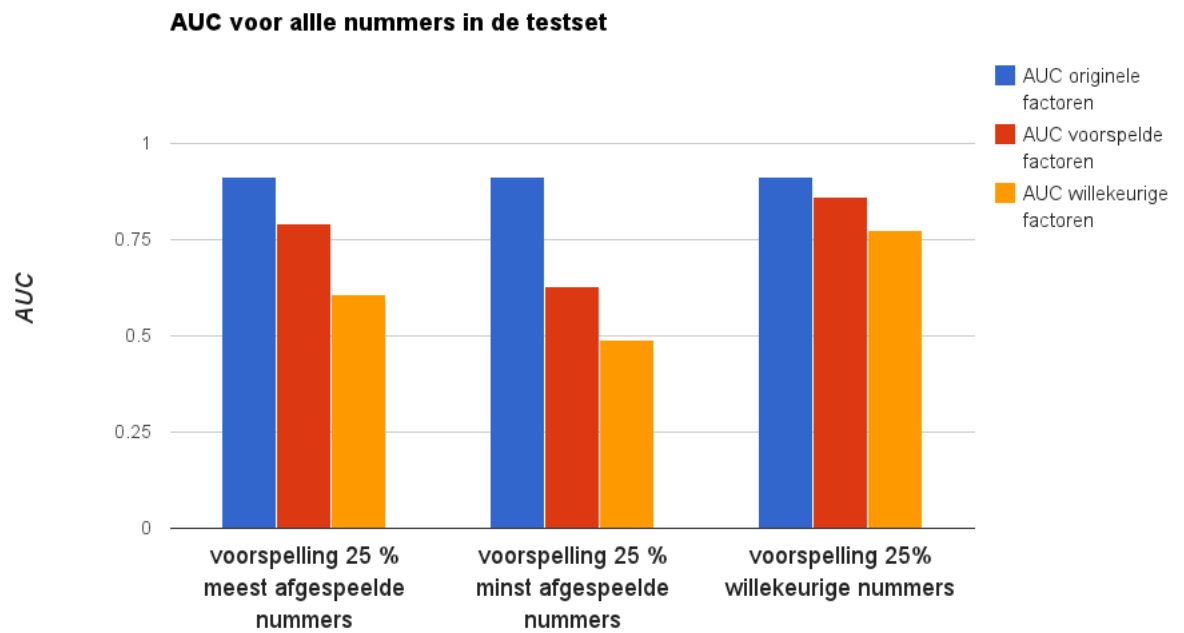
FIGUUR 4.6: ROC-curves voor alle nummers in de testset voor een bepaalde gebruiker.



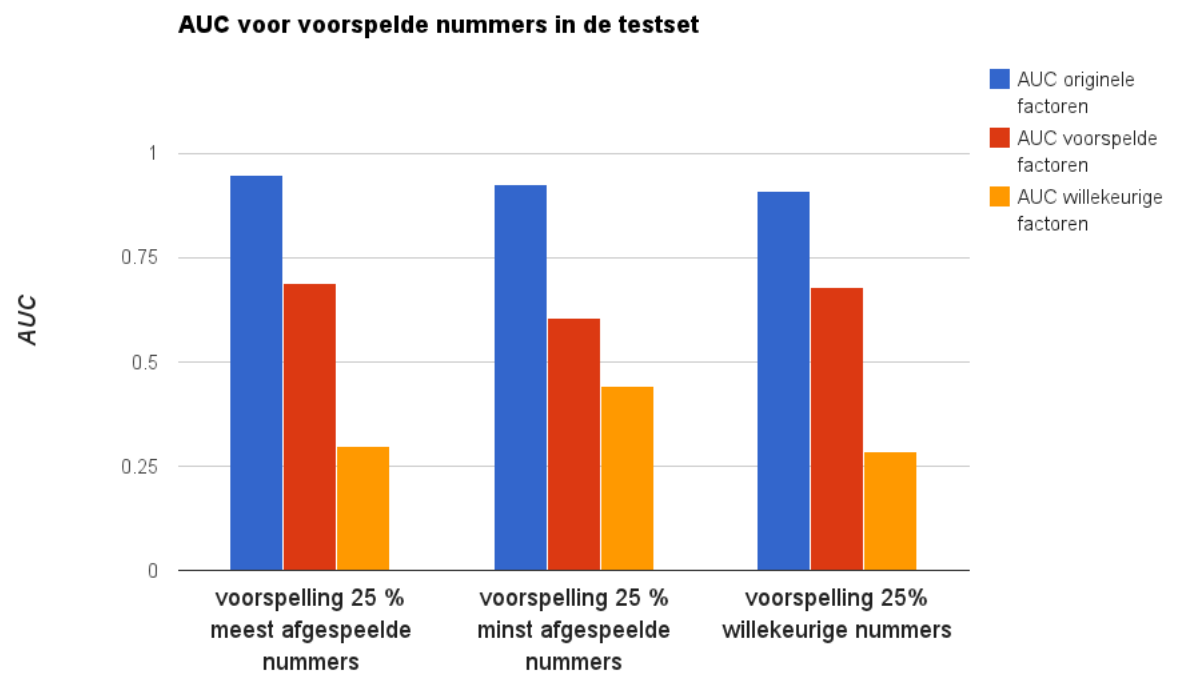
FIGUUR 4.7: ROC-curves voor enkel de voorspelde nummers in de testset voor een bepaalde gebruiker.

Wanneer we de AUC berekenen voor elke gebruiker en vervolgens uitmiddelen over alle gebruikers bekomen we de gemiddelde AUC voor het totale aanbevelingsproces. We kunnen ook bestuderen wat er gebeurt met de gemiddelde AUC indien we de 25% meest populaire nummers of de 25% minst populaire nummers voorspellen uit de audio. Zoals in figuren 4.8 en 4.9 te zien is, daalt de gemiddelde AUC sterk in alle gevallen wanneer we een deel van de muziekfactoren gaan voorspellen. De voorspelling is echter in alle mogelijke gevallen nog steeds beter dan indien de factoren achterna willekeurig gepermuteerd worden (het semi-willekeurig geval) en is bovendien nog steeds hoger dan 0.5, de verwachte waarde wanneer de nummers in totaal willekeurige volgorde aangeraden zouden worden.

Uit figuren 4.8 en 4.9 kan dus worden afgeleid dat onze voorspellingsmethode voor de muziekfactoren wel degelijk zijn nut kan hebben indien er absoluut geen luistergeschiedenis voorhanden is.



FIGUUR 4.8: Gemiddelde AUC waarden voor alle nummers.

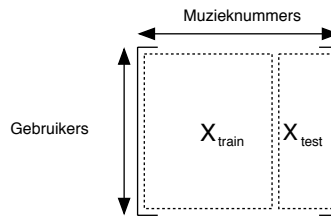


FIGUUR 4.9: Gemiddelde AUC waarden voor voorspelde nummers.

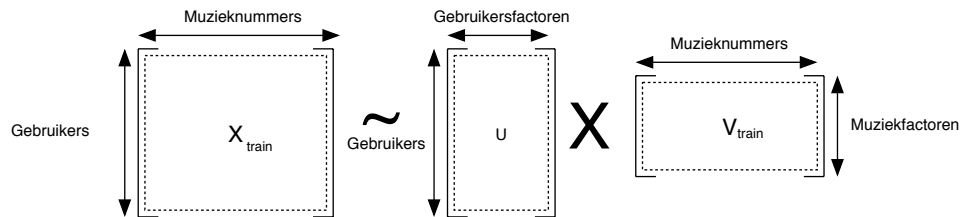
#### 4.4.1.2 Experiment 2

In dit experiment werd onderzocht of er een punt is voor de hoeveelheid luistergeschiedenis waar het voorspellen van de muziekfactoren uit de audio tot betere aanbevelingen zal leiden dan wanneer aanbevelingen worden gegenereerd enkel op basis van een collaboratieve filtertechniek.

Er werd opnieuw vertrokken van de initiële dataset (zie figuur 4.10) die wordt opgesplitst in een training en testset. De trainingset bevat informatie over de luistergeschiedenis van 8000 nummers terwijl de luistergeschiedenis van de overige 2000 nummers zich in de testset bevinden. Daarna werd een verborgen factormodel getraind op de trainingset  $X_{train}$  zoals beschreven staat in hoofdstuk 2. Dit resulteert in een muziekfactormatrix  $V_{train}$  en een gebruikerfactormatrix  $U$  (zie figuur 4.11). Deze gebruikerfactormatrix bleef steeds dezelfde gedurende het vervolg van het experiment.



FIGUUR 4.10: Opdeling van de data in training en testset.



FIGUUR 4.11: Opdeling van de data in training en testset.

De muziekfactormatrix  $V_{train}$  werd dan gebruikt om aan de hand van de geclusterde voorstelling van de MFCC-vectoren een lineair model te trainen dat in staat is om de muziekfactoren te voorspellen voor de muzieknummers die zich in de testset  $X_{test}$  bevinden. De uit de audio voorspelde muziekfactoren noemen we  $V_{test_{audio}}$ .

Nu vertrekken we vanuit  $X_{test}$  en gaan we telkens informatie over de luistergeschiedenis verwijderen. Voor elk nummer wordt een percentage van de gebruikers die naar het nummer geluisterd hebben verwijderd. Dit percentage wordt bepaald door de parameter  $\alpha$ . Concreet betekent dit

dat in de gebruiker-item matrix  $X_{test}$  het aantal luisterbeurten op 0 gezet wordt. Dit zorgt ervoor dat de gebruiker-item matrix opnieuw in een nieuwe training- en testset opgesplitst wordt. Het deel van de informatie dat wordt verwijderd noemen we  $X_{test_A}$  en is de nieuwe testset. Deze matrix wordt omgezet naar een binaire gebruikers-item matrix omdat enkel zal gekeken worden of de voorspellingen effectief ooit al eens beluisterd werden door de gebruiker. Het aantal luisterbeurten doet er niet echt toe. Het andere deel  $X_{test_B}$  zal opnieuw gebruikt worden om een nieuwe muziekfactormatrix  $V_{test_{mf}}$  te bouwen met behulp van de gebruikerfactormatrix  $U$  en 1 iteratie uit het verborgen factormodel. In deze iteratie worden de gebruikersfactoren constant beschouwd en zullen de muziekfactoren geoptimaliseerd worden in functie van de gebruikersfactoren en de scores (zie vgl 2.4). Doordat deze matrixfactorisatie voor grote waarden van  $\alpha$  zeer weinig informatie heeft, wordt verwacht dat de resulterende muziekfactoren vrij slecht zullen zijn. Wanneer dan voorspellingen gegenereerd worden en deze geëvalueerd worden op  $X_{test_A}$  moeten de resultaten vrij slecht zijn en wordt er hopelijk een waarde voor  $\alpha$  gevonden waarbij het beter is om de muziekfactoren te laten voorspellen uit de audio, dan op basis van de beperkte luistergeschiedenis. De oorspronkelijke matrix  $X_{test}$  bevatte 264.767 niet-nulwaarden. Wanneer  $\alpha = 0.5$  dan bevat  $X_{test_B}$  dus ongeveer 132.000 niet-nulwaarden. De nieuwe scores waaruit dan aanbevelingen worden gegenereerd worden opnieuw als volgt bekomen:

$$pred_{audio} = UV_{test_{audio}}^T \quad (4.16)$$

$$pred_{mf} = UV_{test_{mf}}^T \quad (4.17)$$

De resultaten kan men terugvinden in tabel 4.2 en 4.3. Deze werden bekomen door dit experiment 3 maal uit te voeren met telkens een verschillende  $X_{test_A}$  en  $X_{test_B}$  en vervolgens de resultaten uit te middelen.

#### 4.4.1.3 Experiment 3

Een andere vraag die gesteld kan worden stellen is of een combinatie van een sociaal-gebaseerde techniek en een inhouds-gebaseerde techniek in bepaalde situaties kan leiden tot een beter aanbevelingssysteem. Een eenvoudige manier om dit na te gaan bestaat erin om de scores uit 4.16 en 4.17 op te tellen en vervolgens de resulterende scores te gebruiken om aanbevelingen te genereren. Deze techniek zal in het bijzonder goed werken indien de 2 aanbevelingstechnieken



$$\begin{bmatrix} 0 & 0 & 6 & 0 & 5 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 4 & 0 & \mathbf{40} & 0 & \mathbf{13} & 0 & \mathbf{1} & \mathbf{11} & 3 & \mathbf{6} \\ 0 & 0 & 0 & 0 & \mathbf{4} & 1 & 0 & 0 & 0 & \mathbf{1} \\ 0 & 1 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & \mathbf{10} & 0 \\ \mathbf{1} & \mathbf{2} & 0 & 3 & 5 & 20 & 0 & 0 & 0 & 0 \\ 2 & 1 & 2 & 0 & 0 & 1 & 0 & 6 & 0 & 5 \\ 0 & \mathbf{3} & \mathbf{1} & 0 & 1 & 0 & 1 & 3 & 6 & 0 \\ 0 & 0 & 0 & 6 & \mathbf{2} & \mathbf{1} & 0 & 0 & \mathbf{3} & 0 \\ \mathbf{1} & 0 & 6 & \mathbf{7} & 0 & 2 & \mathbf{9} & 0 & 0 & 7 \\ 10 & \mathbf{5} & 0 & 0 & 0 & 0 & 0 & \mathbf{4} & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 6 & 0 & 5 & 0 & 0 & 0 & 0 & 0 \\ 4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 5 & 20 & 0 & 0 & 0 & 0 \\ 2 & 1 & 2 & 0 & 0 & 1 & 0 & 6 & 0 & 5 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 3 & 6 & 0 \\ 0 & 0 & 0 & 6 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 6 & 0 & 0 & 2 & 0 & 0 & 0 & 7 \\ 10 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & \mathbf{1} & 0 & \mathbf{1} & 0 & \mathbf{1} & \mathbf{1} & 0 & \mathbf{1} \\ 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & \mathbf{1} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & \mathbf{1} & 0 \\ \mathbf{1} & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \mathbf{1} & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \mathbf{1} & \mathbf{1} & 0 & 0 & \mathbf{1} & 0 \\ \mathbf{1} & 0 & 0 & \mathbf{1} & 0 & 0 & \mathbf{1} & 0 & 0 & 0 \\ 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 \end{bmatrix}$$

FIGUUR 4.12: Een voorbeeld van hoe de dataset  $X_{test}$  wordt opgesplitst in  $X_{test_A}$  en  $X_{test_B}$  wanneer  $\alpha = 0.5$ . Waarden in het **vet** worden verwijderd uit  $X_{test}$ .

complementair zijn. Dit betekent dat de inhouds-gebaseerde techniek erin slaagt om voorspellingen te maken die de sociaalgebaseerde techniek niet doet en vice versa:

$$pred_{hybrid} = pred_{audio} + pred_{mf} \quad (4.18)$$

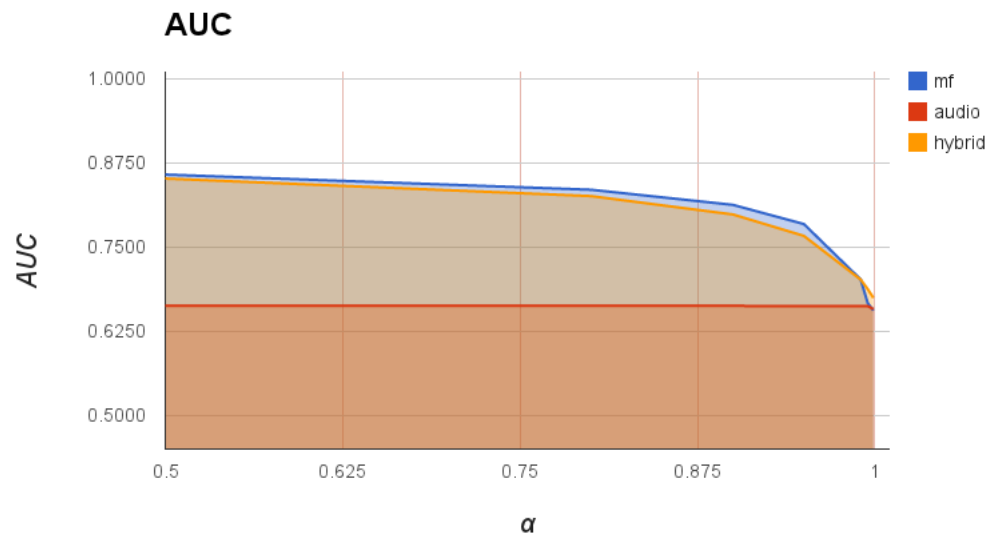
|                 | $\alpha$      |               |               |               |               |               |               |
|-----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                 | 0.5           | 0.8           | 0.9           | 0.95          | 0.99          | 0.995         | 0.999         |
| $pred_{audio}$  | 0.6630        | 0.6629        | 0.6628        | 0.6625        | 0.6626        | 0.6625        | 0.6580        |
| $pred_{mf}$     | <b>0.8580</b> | <b>0.8354</b> | <b>0.8129</b> | <b>0.7844</b> | <b>0.7022</b> | 0.6669        | 0.6451        |
| $pred_{hybrid}$ | 0.8520        | 0.8259        | 0.7987        | 0.7668        | 0.7019        | <b>0.6877</b> | <b>0.6808</b> |

TABEL 4.2: AUC voor de verschillende technieken ontwikkeld in experiment 2 en 3.

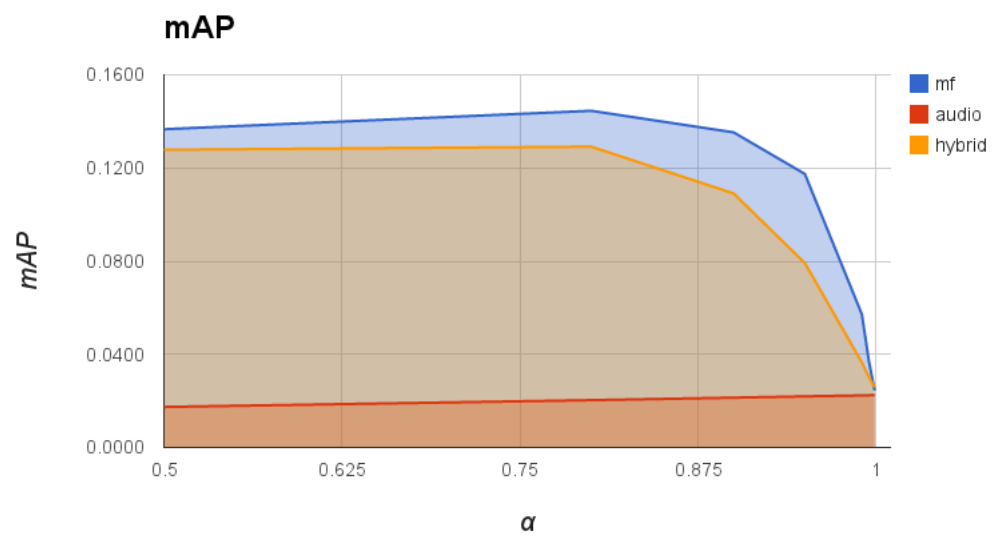
|                 | $\alpha$      |               |               |               |               |               |               |
|-----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                 | 0.5           | 0.8           | 0.9           | 0.95          | 0.99          | 0.995         | 0.999         |
| $pred_{audio}$  | 0.0174        | 0.0203        | 0.0214        | 0.0219        | 0.0224        | 0.0224        | 0.0226        |
| $pred_{mf}$     | <b>0.1367</b> | <b>0.1446</b> | <b>0.1353</b> | <b>0.1174</b> | <b>0.0572</b> | <b>0.0368</b> | 0.0228        |
| $pred_{hybrid}$ | 0.1278        | 0.1292        | 0.1091        | 0.0792        | 0.0365        | 0.0302        | <b>0.0277</b> |

TABEL 4.3: mAP voor de verschillende technieken ontwikkeld in experiment 2 en 3.

Matrixfactorisatie blijkt superieur te zijn ten opzichte van voorspelling van de muziekfactoren op basis van de MFCC-vectoren. Enkel wanneer zeer weinig luisterinformatie beschikbaar is, presteert het inhouds-gebaseerde aanbevelingssysteem op basis van MFCC-vectoren beter. Er werd een evenwichtspunt gevonden wanneer  $\alpha = 0.999$  was. In deze situatie bevatte de trainingset nog slechts een kleine 2000 niet-nulwaarden. Het gebruik van een hybride aanbevelingssysteem blijkt echter sneller de prestatie van de sociaalgebaseerde techniek te evenaren en zelfs te verbeteren wanneer weinig luistergeschiedenis voorhanden is.



FIGUUR 4.13: AUC van de aanbevelingen voor verschillende waarden van  $\alpha$  bij experiment 2 en 3



FIGUUR 4.14: mAP van de aanbevelingen voor verschillende waarden van  $\alpha$  bij experiment 2 en 3

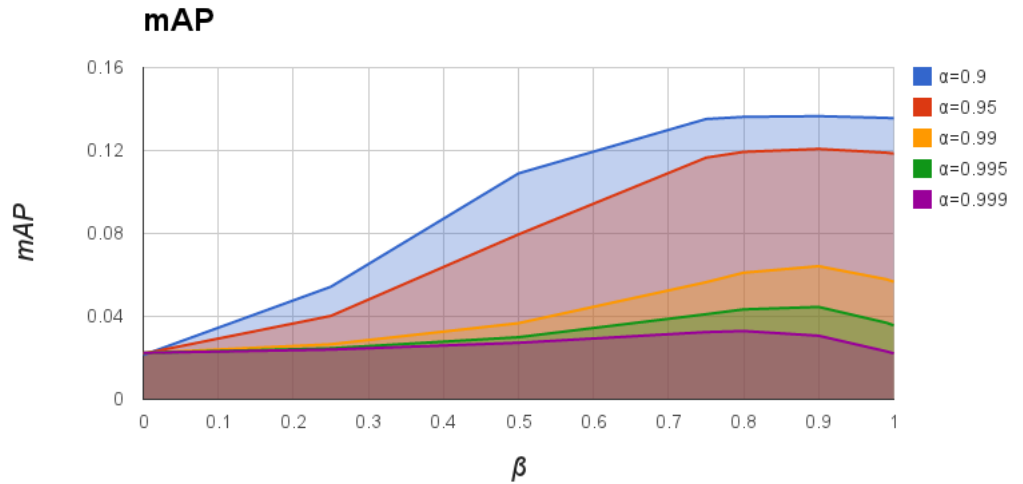
#### 4.4.1.4 Experiment 4

Uit experiment 3 konden we afleiden dat de sociaalgebaseerde techniek op basis van een verborgen factor model beter is. Een vraag die we kunnen stellen is of we door beide technieken correct te combineren in een hybride aanbevelingssysteem, toch nog een betere prestatie kunnen realiseren. We kunnen het belang van elke techniek in het hybride aanbevelingssysteem wijzigen door een extra parameter  $\beta$  in te voeren. Vergelijking 4.18 wordt dan uitgebreid tot:

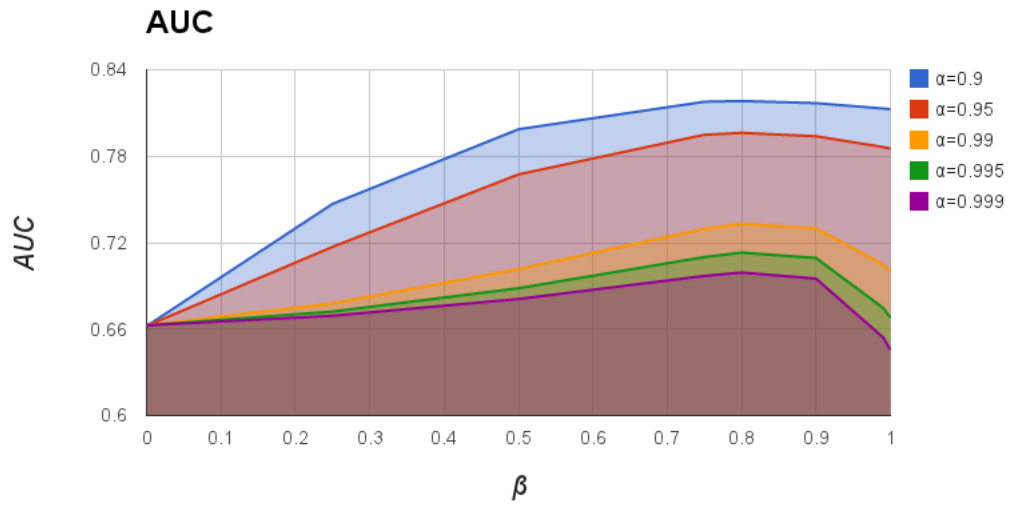
$$pred_{hybrid} = (1 - \beta) * pred_{audio} + (\beta) * pred_{mf} \quad \beta \in [0, 1] \quad (4.19)$$

$\beta = 0$  betekent dat we enkel de inhouds-gebaseerde techniek in beschouwing nemen terwijl  $\beta = 1$  enkel rekening houdt met de sociaalgebaseerde techniek. We kunnen nu opnieuw voor elke waarde van  $\alpha$  kijken wat de optimale waarde is voor  $\beta$ .

Opnieuw werd dit experiment meerdere malen uitgevoerd op een verschillende training- en test-set. De resultaten zijn gevisualiseerd in figuren 4.15 en 4.16.



FIGUUR 4.15: mAP voor verschillende waarden van  $\beta$ .

FIGUUR 4.16: AUC voor verschillende waarden van  $\beta$ .

We kunnen vaststellen dat voor verschillende waarden van  $\alpha$  het optimum  $\beta_{opt}$  rond 0.8-0.9 ligt. Dit zorgt ervoor dat bij het aanbevelingssysteem de sociaalgebaseerde techniek op basis van matrixfactorisatie zwaar doorweegt, maar dat men door in mindere mate de inhouds-gebaseerde scores mee te nemen in de uiteindelijke predictie tot een beter aanbevelingssysteem kan komen. In tabellen 4.4 en 4.5 worden de resultaten van het optimale hybride aanbevelingssysteem weergegeven samen met de sociaal-gebaseerde en inhoudsgebaseerde techniek. Deze resultaten werden opnieuw bekomen door dit experiment 3 maal uit te voeren met telkens een verschillende  $X_{test_A}$  en  $X_{test_B}$  en vervolgens de resultaten uit te middelen.

|                              | $\alpha$      |              |               |               |               |
|------------------------------|---------------|--------------|---------------|---------------|---------------|
|                              | 0.9           | 0.95         | 0.99          | 0.995         | 0.999         |
| $pred_{audio}(\beta = 0)$    | 0.6620        | 0.6622       | 0.6626        | 0.6626        | 0.6626        |
| $pred_{mf}(\beta = 1)$       | 0.8126        | 0.7853       | 0.7004        | 0.6680        | 0.6456        |
| $pred_{hybrid}(\beta_{opt})$ | <b>0.8183</b> | <b>0.796</b> | <b>0.7331</b> | <b>0.7131</b> | <b>0.6993</b> |

TABEL 4.4: AUC voor de sociaal-gebaseerde techniek, inhouds-gebaseerde techniek en optimale combinatie van beide.

|                              | $\alpha$      |               |               |               |               |
|------------------------------|---------------|---------------|---------------|---------------|---------------|
|                              | 0.9           | 0.95          | 0.99          | 0.995         | 0.999         |
| $pred_{audio}(\beta = 0)$    | 0.0214        | 0.0220        | 0.0402        | 0.0224        | 0.0246        |
| $pred_{mf}(\beta = 1)$       | 0.1356        | 0.1186        | 0.0568        | 0.0357        | 0.0222        |
| $pred_{hybrid}(\beta_{opt})$ | <b>0.1366</b> | <b>0.1208</b> | <b>0.0642</b> | <b>0.0446</b> | <b>0.0330</b> |

TABEL 4.5: mAP voor de sociaal-gebaseerde techniek, inhouds-gebaseerde techniek en optimale combinatie van beide.

## Hoofdstuk 5

### Besluit

In deze masterproef werd onderzocht of het mogelijk is om het *koude startprobleem* deels te verhelpen door een hybride aanbevelingssysteem te ontwikkelen. Dit systeem is gebaseerd op een collaboratieve filtertechniek die toelaat om data deels te voorspellen op basis van akoestische informatie in de plaats van luistergeschiedenis. Met behulp van verborgen factormodellen kan een score voor een gebruiker en een muziknummer geschreven worden als het product van een muziekfactor en een gebruikerfactor. Voor nieuwe nummers zijn de muziekfactoren nog onbekend en kunnen deze aan de hand van geclusterde MFCC-vectoren met behulp van regressie-analyse geschat worden. De resultaten toonden aan dat deze aanpak beter werkt dan willekeurige muziekfactoren, maar deze benaderingen scoorden beduidend slechter dan wanneer de correcte muziekfactoren gebruikt worden. Wanneer deze inhoudsgebaseerde techniek echter correct gecombineerd wordt met de sociaalgebaseerde techniek kan het aanbevelingssysteem in vele gevallen reeds betere aanbevelingen genereren.

Een mogelijke oplossing zou zijn om naast MFCC-vectoren die enkel akoestische informatie opnemen, heel wat externe informatie over het nummer ook in beschouwing te nemen. Meta-informatie zoals de tijdsgeest van het nummer of de taal waarin gezongen wordt draagt wellicht ook in belangrijke mate bij tot vorming van de verschillende muziekfactoren. Bij nieuwe muziknummers is deze informatie echter zelden beschikbaar. Ook een verdere akoestische analyse waarbij zaken zoals ritme en tempo ook in beschouwing worden genomen kan de voorspelling nauwkeuriger maken.

Verder kan ook worden nagegaan of het mogelijk is om de muziekfactoren beter te kunnen voorspellen met andere technieken dan lineaire regressie. De kostfunctie zou ook aangepast

kunnen worden zodat de nadruk meer ligt op het correct voorspellen van extreme waarden voor muziekfactoren.

We kunnen besluiten dat het voorspellen van muziekfactoren een complex probleem is en dat deze voorspelling nuttig kan zijn wanneer absoluut geen luisterinformatie voorhanden is, maar dat collaboratieve filtertechnieken op basis van matrixfactorisatie al snel superieur zijn indien er een beperkte luistergeschiedenis beschikbaar is.

# Bibliografie

- [1] Gediminas Adomavicius and Alexander Tuzhilin. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. on Knowl. and Data Eng.*, 17(6):734–749, June 2005. ISSN 1041-4347. doi: 10.1109/TKDE.2005.99. URL <http://dx.doi.org/10.1109/TKDE.2005.99>.
- [2] James Bennett, Stan Lanning, and Netflix Netflix. The netflix prize. In *In KDD Cup and Workshop in conjunction with KDD*, 2007.
- [3] 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>.
- [4] Chris Anderson. *The Long Tail: Why the Future of Business Is Selling Less of More*. Hyperion, 2006. ISBN 1401302378.
- [5] Ò. Celma and P. Cano. From hits to niches? or how popular artists can bias music recommendation and discovery. In *2nd Workshop on Large-Scale Recommender Systems and the Netflix Prize Competition (ACM KDD)*, Las Vegas, USA, 24/08/2008 2008. URL <http://mtg.upf.edu/files/publications/Celma-ACM-Netflix-KDD2008.pdf>.
- [6] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- [7] Million Song Dataset. *Example track description*, 2013 (accessed May 18, 2013). URL <http://labrosa.ee.columbia.edu/millionsong/pages/example-track-description>.
- [8] Fabio Aioli. A preliminary study on a recommender system for the million songs dataset challenge. 2012.

- [9] Yifan Hu, Yehuda Koren, and Chris Volinsky. Collaborative filtering for implicit feedback datasets. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining, ICDM '08*, pages 263–272, Washington, DC, USA, 2008. IEEE Computer Society. ISBN 978-0-7695-3502-9. doi: 10.1109/ICDM.2008.22. URL <http://dx.doi.org/10.1109/ICDM.2008.22>.
- [10] Gijs Geleijnse. The quest for ground truth in musical artist tagging in the social web era.
- [11] Mark Levy and Mark Sandler. A semantic space for music derived from social tags. In *ISMIR*, pages 411–416, 2007.
- [12] Ling Chen, Phillip Wright, and Wolfgang Nejdl. Improving music genre classification using collaborative tagging data. In *Proceedings of the Second ACM International Conference on Web Search and Data Mining, WSDM '09*, pages 84–93, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-390-7. doi: 10.1145/1498759.1498812. URL <http://doi.acm.org/10.1145/1498759.1498812>.
- [13] Philip J. Stone, Dexter C. Dunphy, Marshall S. Smith, and Daniel M. Ogilvie. *The General Inquirer: A Computer Approach to Content Analysis*. MIT Press, 1966. URL <http://www.webuse.umd.edu:9090/>.
- [14] F. Å. Nielsen. Afinn, mar 2011. URL <http://www2.imm.dtu.dk/pubdb/p.php?6010>.
- [15] Gerard Salton and Michael J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill, Inc., New York, NY, USA, 1986. ISBN 0070544840.
- [16] C. Laurier. *Automatic Classification of Musical Mood by Content-Based Analysis*. PhD thesis, Universitat Pompeu Fabra, 10/2011 2011. URL [files/publications/PhD\\_Cyril\\_Laurier\\_2011\\_Music\\_Mood\\_Classification.pdf](files/publications/PhD_Cyril_Laurier_2011_Music_Mood_Classification.pdf).
- [17] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [18] Kerstin Bischoff, Claudiu S. Firan, Raluca Paiu, Wolfgang Nejdl, Cyril Laurier, and Mohamed Sordo. Music mood and theme classification - a hybrid approach. In *ISMIR*, pages 657–662, 2009.



- 
- [19] S. Davis and P. Mermelstein. Experiments in syllable-based recognition of continuous speech. *IEEE Trans. Acoust., Speech, Signal Processing*, 28:357 – 366, Aug. 1980.
- [20] Oscar Celma. *Music Recommendation and Discovery: The Long Tail, Long Fail, and Long Play in the Digital Music Space*. Springer Publishing Company, Incorporated, 1st edition, 2010. ISBN 3642132863, 9783642132865.
- [21] Uwe Ligges. *tuneR: Analysis of music*, 2011. URL <http://r-forge.r-project.org/projects/tuner/>.
- [22] Adam Coates, Andrew Y. Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. *Journal of Machine Learning Research - Proceedings Track*, 15:215–223, 2011. URL <http://dblp.uni-trier.de/db/journals/jmlr/jmlrp15.html#CoatesNL11>.
- [23] A. Hyvärinen and E. Oja. Independent component analysis: algorithms and applications. *Neural Netw.*, 13(4-5):411–430, May 2000. ISSN 0893-6080. doi: 10.1016/S0893-6080(00)00026-5. URL [http://dx.doi.org/10.1016/S0893-6080\(00\)00026-5](http://dx.doi.org/10.1016/S0893-6080(00)00026-5).
- [24] Adam Coates and Andrew Y. Ng. Learning feature representations with k-means.

# Lijst van figuren

|     |  |    |
|-----|--|----|
| 1.1 | Een voorbeelddistributie van de populariteit van nummers in een muziekcollectie  | 5  |
| 1.2 | Een voorbeeld van een ROC-curve . . . . .  | 10 |
| 1.3 | Een voorbeeldfragment uit de MSD. Links bevindt zich het gebruikersid, in het midden het muzieknummerid en rechts het aantal afspeelbeurten. . . . .   | 11 |
| 1.4 | Een voorbeeldfragment uit de musixmatch dataset. . . . .   | 12 |
| 1.5 | Een voorbeeldbestand uit de Last.fm dataset. . . . .   | 12 |
| 2.1 | De recall op de testset voor de verschillende aanbevelingssystemen wanneer 50 aanbevelingen werden gegenereerd . . . . .   | 18 |
| 2.2 | De mAP op de testset voor de verschillende aanbevelingssystemen wanneer 50 aanbevelingen werden gegenereerd . . . . .  | 18 |
| 2.3 | Taginformatie van muzieknummers die een sterk negatieve waarde versus een sterk positieve waarde hadden voor 2 factoren x en y. Factor x geeft een indicatie van hoe recent het muzieknummer is terwijl factor y kan geïnterpreteerd worden als een maat voor dansbaarheid van het nummer. . . . . | 19 |
| 2.4 | Weergave van enkele datapunten uit de muziekfactormatrix in functie van 2 factoren. 20   |    |
| 3.1 | Opbouw van beide datasets en verdere verwerking vooraleer men aan emotieherkenning kan doen. De akoestische features waren bij de eerste dataset afkomstig van de eigen subset van de MSD. In de tweede dataset werd een groter deel van de MSD opgenomen. . . . .                                 | 22 |
| 3.2 | De 12 basis functies van de timbre vector . . . . .  | 24 |
| 3.3 | De relatie tussen de gemiddelde luidheid en de snelheid van aanzet voor nummers met een positieve en negatieve emotie. . . . .   | 24 |
| 3.4 | Overzicht van de mutuele informatie tussen verschillende features en de emotielabels. 25   |    |
| 3.5 | De visualisatie van de 1e dataset met behulp van PCA. . . . .  | 26 |
| 4.1 | Opsplitsing van het akoestisch signaal in vensters . . . . .   | 32 |
| 4.2 | Het samennemen van MFCC-vectoren in vensters. . . . .  | 37 |
| 4.3 | overzicht van het trainingsproces van de muziekfactoren uit een geclusterde MFCC-voorstelling . . . . .  | 39 |
| 4.4 | De invloed van het aantal clustercentra op de MSE voor verschillende venstergroottes . . . . .   | 40 |
| 4.5 | Het aanbevelingsproces met voorspelde muziekfactoren. . . . .  | 41 |
| 4.6 | ROC-curves voor alle nummers in de testset voor een bepaalde gebruiker. . . . .  | 42 |
| 4.7 | ROC-curves voor enkel de voorspelde nummers in de testset voor een bepaalde gebruiker. . . . .   | 42 |
| 4.8 | Gemiddelde AUC waarden voor alle nummers. . . . .  | 43 |
| 4.9 | Gemiddelde AUC waarden voor voorspelde nummers. . . . .  | 43 |

---

|      |  |    |
|------|--|----|
| 4.10 | Opdeling van de data in training en testset. . . . .   | 44 |
| 4.11 | Opdeling van de data in training en testset. . . . .   | 44 |
| 4.12 | Een voorbeeld van hoe de dataset $X_{test}$ wordt opgesplitst in $X_{test_A}$ en $X_{test_B}$<br>wanneer $\alpha = 0.5$ . Waarden in het <b>vet</b> worden verwijderd uit $X_{test}$ . . . . . | 46 |
| 4.13 | AUC van de aanbevelingen voor verschillende waarden van $\alpha$ bij experiment 2 en 3   | 47 |
| 4.14 | mAP van de aanbevelingen voor verschillende waarden van $\alpha$ bij experiment 2 en 3   | 47 |
| 4.15 | mAP voor verschillende waarden van $\beta$ . . . . .   | 48 |
| 4.16 | AUC voor verschillende waarden van $\beta$ . . . . .   | 49 |

# Lijst van tabellen

|     |  |    |
|-----|--|----|
| 1.1 | Berekening van de gemiddelde precisie voor een bepaalde gebruiker wanneer $\tau = 10$ nummers worden aangeraden en 3 nummers in de testset zitten. . . . .                               | 8  |
| 1.2 | Een confusiematrix . . . . .   | 8  |
| 1.3 | Een confusiematrix toegepast op een aanbevelingssysteem . . . . .  | 9  |
| 2.1 | De evaluatie van de verschillende aanbevelingssystemen . . . . .   | 18 |
| 3.1 | De verzameling van akoestische features die gebruikt werd bij classificatie. . . . .   | 26 |
| 3.2 | De hyperparameters van de verschillende classifiers . . . . .  | 28 |
| 3.3 | De nauwkeurigheid van verschillende classifiers bij akoestische classificatie. . . . .   | 28 |
| 3.4 | Binomiale en multiklasse classificatie met behulp van Random Forests. . . . .  | 28 |
| 3.5 | De voorspelde en echte emoties bij binomiale classificatie . . . . .   | 29 |
| 3.6 | De voorspelde en echte emoties bij multiklasse classificatie . . . . .   | 29 |
| 3.7 | De nauwkeurigheid van verschillende classifiers bij classificatie op basis van songtekstinformatie. . . . .  | 29 |
| 4.1 | Een eerste vergelijking van de prestatie van het verborgen factormodel met hetzelfde model waar 25% van de factoren voorspeld werden uit de audio en een semi-willekeurig geval. . . . . | 41 |
| 4.2 | AUC voor de verschillende technieken ontwikkeld in experiment 2 en 3. . . . .  | 46 |
| 4.3 | mAP voor de verschillende technieken ontwikkeld in experiment 2 en 3. . . . .  | 46 |
| 4.4 | AUC voor de sociaal-gebaseerde techniek, inhouds-gebaseerde techniek en optimale combinatie van beide. . . . .   | 49 |
| 4.5 | mAP voor de sociaal-gebaseerde techniek, inhouds-gebaseerde techniek en optimale combinatie van beide. . . . .   | 49 |

