

Inhaltsbasierte Musikempfehlung mit Convolutional Neuronalen Netzwerken

WEIDHAS PHILIPP

Matr.nr: 123456

philipp.weidhas@st.oth-regensburg.de

WILDGRUBER MARKUS

Matr.nr: 123456

markus.wildgruber@stud.oth-regensburg.de

Zusammenfassung

Hier kommt die Zusammenfassung...

1. EINLEITUNG

Im ersten Halbjahr des Jahres 2017 wurden 62% der Einnahmen der amerikanischen Musikindustrie durch Streaming Plattformen¹ erzielt. Im Vergleich zu Vorjahr erhöhten sich dadurch die Einnahmen um 48% auf 2.5\$ Milliarde [1]. Dieser Erfolg basiert nicht nur auf einer guten Verfügbarkeit der Lieder und einem günstigen Preis sondern auch auf automatischen Musikempfehlungsdiensten, welche dem Nutzer ein angenehmeres Konsumverhalten ermöglichen. Obwohl Empfehlungsdienste in den letzten Jahren viel erforscht wurden, ist das Problem der Musikempfehlung sehr komplex. Neben einer große Anzahl an verschiedenen Stile und Genres, beeinflussen sowohl soziales- und geographisches Umfeld, sowie der aktuelle Gemütszustand die Vorliebe eines Hörers. [2]

TO DO

In der Musik Information Retrieval (MIR) gibt es vier Kategorien [3] die einen Einfluss auf die Wahrnehmung von ähnlicher Musik haben.

Musikmerkmale sind Eigenschaften, welche aus dem Audiosignal eines Liedes extrahiert werden. Dazu zählen Aspekte wie der Rhythmus,

¹wie Spotify, Apple Music, Pandora etc.

die Melodie, die Harmonie oder die Stimmung eines Stücks.

Als *Musikkontext* versteht man alle Aspekte, die nicht aus dem Audiosignal abgeleitet werden, sondern Informationen die über ein Musikstück bekannt sind. Beispielsweise Metadaten wie der Titel eines Lieds, das Genre, Name des Künstlers oder das Erscheinungsjahr.

Die *Benutzereigenschaften* beziehen sie auf Persönlichkeitsmerkmale, wie Geschmack, musikalisches Wissen und Erfahrung oder den demographischen Hintergrund.

Im Unterschied dazu steht der *Benutzerkontext*, der sich auf die aktuelle Situation des Hörers bezieht. Dabei wird er durch seine Umgebung, seiner Stimmung oder der aktuellen Aktivität beeinflusst. [4]

Bislang werden Informationen über den Hörer durch ein Benutzerprofil repräsentiert. Das Profil enthält nur wenig Hintergrund Informationen des Hörers und beschränkt sich auf Lieder, die ein Benutzer angehört und bewertet hat [2]. Das Nutzen dieser Daten, um Musikvorschläge abzugeben wird als kollaboratives Filtern (KF) bezeichnet. In der Studie von Vigliensoni und Fujinaga [5] zeigt sich ein deutlicher Unterschied zwischen herkömmlichen Benutzerprofilen und das Einfügen von Zusatzinformationen. Durch das Hinzufügen der Features demographischen Hintergrund und Entdeckergeist des Hörers konnte im Vergleich zu einem herkömmlichen Profil eine 12% besser Genauigkeit erreicht werden.

Der weitere Verlauf der wissenschaftlichen Arbeit ist wie folgt organisiert. Im 2. Abschnitt werden verschiedenen Ansätze in den jeweiligen Methodenbereichen vorgestellt. Im 3. Kapitel werden die erfolgreichsten Ansätze miteinander verglichen. Teil 4 zeigt ein eigenes Experiment zu dem Thema. Abschnitt 5 schließt diese Arbeit ab und diskutiert zukünftige Forschungsrichtungen. //TO DO

2. METHODEN ZUR MUSIKEMPFEHLUNG

Es gibt verschiedene Methoden, die in Musikempfehlungssystemen verwendet werden: kollaboratives -, merkmalsbasiertes -, kontextbasiertes Filtern und die hybride Methode. Diese werden genutzt, um Informationen aus der in der Einleitung genannten Eigenschaften zu gewinnen und diese für Empfehlungen an den Nutzer zu verarbeiten. [6]

2.1 Kollaborativer Filter

Kollaboratives Filtern prognostiziert Vorlieben eines Hörers, indem es aus unterschiedlichen Benutzer-Lied Verhältnissen lernt. Es basiert auf der Annahme, dass Verhalten und Bewertungen andere Nutzer auf eine vernünftige Vorhersage für den aktiven Benutzer schließen lassen [7]. Durch explizite² und implizite³ Rückmeldung eines Hörers an das Empfehlungssystem empfiehlt dieses neue Lieder, indem es Gemeinsamkeiten auf Basis der Bewertungen vergleicht [8].

Im diesen Verfahren wird der Ansatz verfolgt, dass Lieder einem Nutzer auf Grundlage von Nutzungsverhalten anderer Anwender der gleichen Plattform vorgeschlagen werden. In der praktischen Umsetzung bedeutet dies: hört ein Anwender ein bestimmtes Musikstück, werden ihm von der Empfehlungsplattform, Lieder vorgeschlagen welche Nutzer in Zeitraum zuvor nach diesem Stück hörten. Dieses Verfahren

² Bewertungen eines Nutzers

geht davon aus, dass durch die Verbindung der Lieder durch vorhergehende Aufrufe eine gute Aussage darüber getroffen werden kann wie gut diese Stücke zusammen passen. Werden Lieder häufig nacheinander gehört, wird diese Verbindung höher bewertet und die Empfehlung häufiger ausgesprochen. Auch wird das Verhalten und der Musikgeschmack des Kunden selbst durch ein System analysiert, um so über Ähnlichkeiten der Kundenpräferenzen mit derer anderer, diesen wiederum bessere Empfehlungen aussprechen zu können. So werden Lieder einem Musikstil zugeordnet und so zielgerichtet dem Nutzer nahegelegt.

Verschiedene Studien ([8][9]) zeigen, dass KF alternative Methoden in der Genauigkeit übertrifft, weshalb es nicht nur im Bereich der Musikempfehlung als die erfolgreichste gilt.

Trotz der Popularität des KF gibt es Probleme, die bei der Verwendung dieser Methode beachtet werden müssen. Das Cold-Start Problem besteht darin, dass noch keine Bewertungen für ein Lied vorliegen, wodurch es auch nicht vorgeschlagen werden kann. Dasselbe Problem gibt es bei einem neuen Benutzer: diesem kann kein guter Vorschlag gemacht werden, da es an Information mangelt welche Art von Musik ihm gefällt. Neben dem Cold-Start Problem gibt es noch weitere Probleme. [7]

2.2 Merkmalsbasierter Filter

2.3 Kontextbasierter Filter

2.4 Hybride Methode

Bei hybriden Methoden werden kollaborative, merkmalsbasierte und kontextbasierter Filter miteinander verknüpft, wodurch ein besseres Empfehlungsergebnis mit weniger Nachteilen der einzelnen Methode zu erzielen. Meistens wird ein kollaborativer Filter mit einem der beiden anderem kombiniert.

Als *gewichtet* wird eine hybride Methode bezeichnet, bei der Empfehlungswerte der ein-

³Beobachten des Konsumverhalten

zelnen Methoden durch eine Linearkombination zusammengerechnet wird. Das Ergebnis der Linearkombination stellt den Empfehlungswertes eines Liedes dar. Durch unterschiedliche Gewichtung der Methoden kann das Empfehlungsergebnis optimiert werden. Der *wechselnde* Ansatz benutzt ein bestimmtes Kriterium anhand dessen es die Methode zur Vorschlagsbestimmung wechselt. Dies kann beispielsweise dann der Fall sein, wenn der erste Filter kein zuverlässiges Ergebnis⁴ liefert. Dann wechselt das System den Filter und kann ein besseres Empfehlungsergebnis bekommen. Bei *gemischten* hybriden Empfehlungen werden unterschiedliche Techniken direkt miteinander vermischt. Dadurch kann für ein System mit inhaltsbasierten Filter das Cold-Start Problem vermieden werden.

// TODO Hybride Methoden können einige Nachteile von kollaborativen Filtern entfernen. Allerdings stehen auch sie vor dem Neuen Benutzerproblem. Dennoch sind hybride Methoden sehr beliebt, da Information über einen neuen Benutzer schnell herausgefunden werden oder durch Profilangaben bereits nach der Registrierung vorhanden sind. [10]

3. CNN FÜR AUDIOSIGNALE

CNN sind durch das biologische Sehen inspiriert und konnten den ersten großen Erfolg im Bereich der Bildklassifizierung [11] verzeichnen. Trotzdem werden CNN auch in verschiedenen Audiobereich, wie der Spracherkennung [12] und der MIR mehr genutzt und erforscht.

In der MIR nutzen die ersten Forschungen CNNs, um die Aufgabe der Musikgenre-Klassifizierung [13] zu untersuchen. Die Ergebnisse⁵ zeigen, dass eine automatisierte Klassifizierung die herkömmliche Methode MFCC deutlich übertrifft. Das erste CNN für inhaltsbasierte Musikempfehlung [2] benutzt zunächst eine Matrix-Faktorisierung um Eigen Vektoren

⁴semantische Unterschiede

für alle Lieder zu erhalten. Anschließend wird das Neuronale Netz für die Zuordnung der Audio-Inhalte zu den Eigen Vektoren genutzt. [6]

Im nachfolgenden Absatz wird der Aufbau, das Training und die Optimierung eines CNN beschrieben.

3.1 Convolutional Neuronale Netze

Im Unterschied zu regulären DNN verwendet das CNN Neuronen, die drei Dimensionale angeordnet sind. Durch diese Anordnung ist es möglich größere Inputdaten in derselben Geschwindigkeit zu verarbeiten wie zuvor [14]. Um eine CNN Architektur zu erstellen werden drei Haupttypen von Schichten verwendet: Convolutional Layer (CL), Pooling Layer (PL) und ein Fully-Connected Layer (FCL).

3.1.1 Schichten eines Convolutional Neuronale Netzwerks

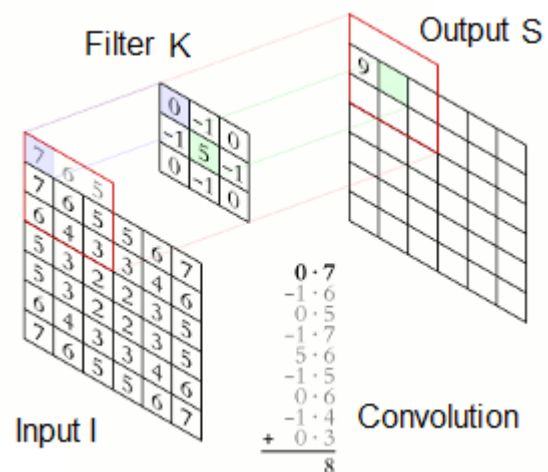


Abbildung 1: Faltung einer 6x6 Matrix mit einem 3x3 Filter [16]

⁵richtigen Klassifizierung

Convolutional Layer

In einem CL findet eine Faltung der Eingangsdaten, in Form einer Matrix, und einem oder mehreren Filtern statt. Ein Filter dient beispielsweise zur Glättung oder zur Verkleinerung der Daten. Eine Verkleinerung der Eingangsmatrix findet statt, wenn ein Filter ohne Zero-padding⁶ verwendet wird. Die Parameter eines Filters werden zufällig initialisiert, können aber mit Hilfe eines Optimierungsverfahrens (3.1.3) angepasst werden. Werden mehrere Filter auf die Eingangsdaten angewendet, ändert sich die Tiefe der gesamten Ausgangsmatrix entsprechend der Anzahl der Filter. [14]

In Abbildung 1 ist die Eingabematrix I eine 6x6 Matrix und K ein 3x3 Filter. Die Ausgangsmatrix S wird an den Stellen (i,j) durch die Gleichung (1) berechnet. Eine genauere Herleitung der Gleichung findet der Leser u. a. bei [15](328f).

$$S(i, j) = (I \star K)(i, j) \quad (1)$$

$$(I \star K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n) \quad (2)$$

Pooling Layer

Ein PL wird zwischen zwei CL eingefügt. Ihre Funktion besteht darin, die Größe der Daten zu reduzieren und damit die Anzahl der Parameter für das nächste CL. Durch die Reduzierung wird die Berechnung des gesamten Netzwerkes beschleunigt. [14]

Ein PL wandelt die Ausgabe eines CL, durch eine statistische Zusammenfassung von nebeneinander liegenden Ausgängen, um. Verschiedene Methoden für ein PL sind: Max Pooling [17], eine Übergabe der größten Zahl in einem rechteckigen Umfeld; die Durchschnittsberechnung des Umfeldes oder ein gewichteter Durchschnitt basierend auf der Entfernung eines zentralen Punktes [15](355).

⁶Eine Matrix wird am Rand um Nullen erweitert.
Bsp. aus einer 7x7 Matrix wird eine 9x9 Matrix

Abbildung 2 zeigt einen 2x2 Max-Filter, der auf eine 4x4 Datenmatrix angewandt wird. Die Verschiebung oder Stride des Filters ist 2 d.h. der Filter wird zunächst auf der y-Achse verschoben. Erreicht er dort das Ende wird er um eine Stride auf der x-Achse verschoben und beginnt wieder mit der y-Verschiebung.

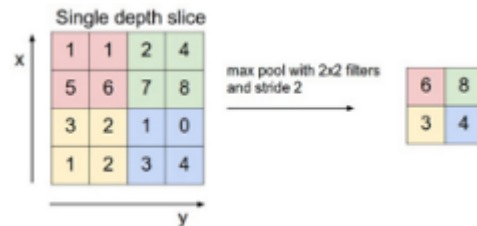


Abbildung 2: Maxpooling mit einem 2x2 Filter[14]

Fully-Connected Layer

Neuronen in einer FCL haben Verbindungen zu allen Knoten der vorherigen Schicht. Ihre Aktivierung wird durch eine Matrixmultiplikation und einem Bias-Offset berechnet [14]. Die FCL wird als Ausgangsschicht verwendet um aus der Eingangsmatrix einen Vektor zu erzeugen.

3.1.2 Training

Cross entropie

3.1.3 Optimierung

Dropout mini batch verfahren

3.2

3.3 Hybride Musikempfehlung mit einem Neuronalen Netzwerk

Im Unterschied zu der zuvor dargestellten Forschung (3.2) wird in der jetzigen ein Deep Belief Netzwerk(DBN) verwendet, um ein hybride

des inhaltsbasiertes Musikempfehlungssystem zu entwickeln. Bisherige inhaltsbasiertes Systeme verfolgen typischerweise einem zweistufigen Ansatz: zunächst extrahieren sie aus Audioinhalte den MFCC Koeffizienten; anschließend prognostizieren sie Musikpräferenzen eines Nutzers. Das nachfolgende Modell führt dieses beiden Schritte simultan und automatisch aus. [6]

Das hybride Modell basiert auf einem hierarchischen linearen Modell mit einem Deep Belief Netzwerk(HLDBN), dass zunächst erläutert wird, um anschließend die Funktionsweise des hybriden Systems darzustellen.

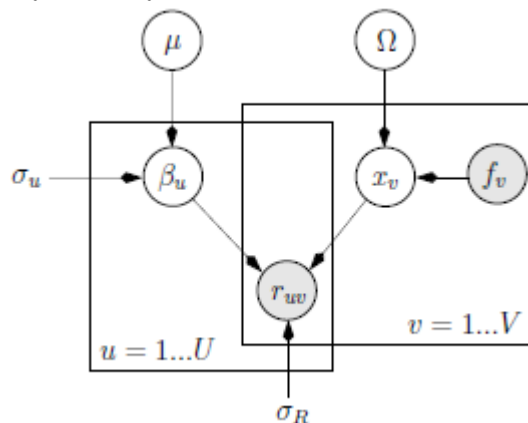


Abbildung 3: Hierarchisches lineares Modell eines Deep Belief Netzwerks [6]

3.3.1 Hierarchisches lineares Modell mit einem Deep Belief Netzwerk

Das in Abbildung 3 gezeigte Modell ist wie folgt definiert: f_v sind Musikmerkmale eines Liedes v , die durch den Eigenvektor x_v automatisch errechnet werden. Die bevorzugte Musik eines Benutzer u wird als Vektor β_u bezeichnet. Ω bezeichnet die Parameter, die das DBNs lernt. Die Bewertung, die u einem Lied v gibt, ist ein Skalarprodukt von r_{xv} und β_u . Durch σ_R wird die Varianz aller Bewertungen des Nutzers betrachtet. μ repräsentiert den allgemeinen Musikgeschmack aller Benutzer,

⁷ $\mathcal{N}(a,b)$ ist die Normalverteilung mit Mittelwert a und Varianz b . $x \sim p$ zeigt, dass x die Verteilung p erfüllt

wobei σ_u die Varianz des einzelnen Nutzers definiert. Alle Benutzer und Lieder Paare werden als I bezeichnet. Für eine Regularisierung der Werte wird die Gaußsche Normalverteilung \mathcal{N} verwendet.⁷ [6]

Das Modell wird wie folgt formuliert:

$$\begin{aligned} r_{xv} &\sim \mathcal{N}(\beta'_u x_v, \sigma^2_R) \\ \beta &\sim \mathcal{N}(\mu, \sigma^2_u I) \\ x_v &= \text{DBN}(f_v; \Omega) \end{aligned}$$

Für das Training des Systems wird die Maximum Likelihood-Funktion oder auch Cross-Entropy verwendet. Als Optimierungsmethode wird das stochastische Mini-Batch Verfahren genutzt, um ein Overfitting der Parameter zu vermeiden. Nach der Lernphase kann r_{xv} geschätzt werden, wodurch auch neue Lieder empfohlen werden können. [6]

3.3.2 Hybrides Modell mit einem Deep Belief Netzwerk

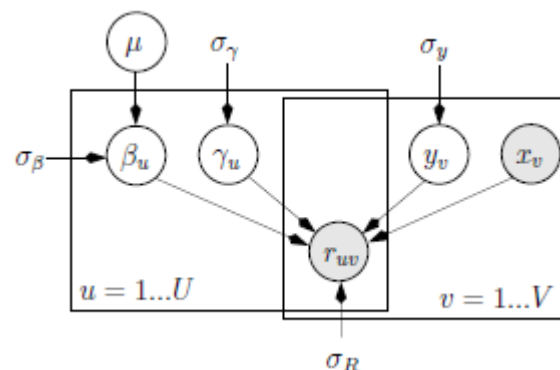


Abbildung 4: Hybrides Empfehlungs Modell [6]

4. VERGLEICH MIT STAND DER FORSCHUNG

5. DISKUSSION DER ZUKÜNFTIGEN FORSCHUNGSTRENDS

LITERATUR

- [1] Joshua P. Friedlander. News and notes on 2017 mid-year riaa revenue statistics. RIAA, 2017.
- [2] Aäron van den Oord, Sander Dieleman, and Benjamin Schrauwen. Deep content-based music recommendation. *Advances in Neural Information Processing Systems* 26, 2013.
- [3] Markus Schedl, Arthur Flexer, and Julián Urbano. *The neglected user in music information retrieval research*, volume 36. Springer, 2013.
- [4] Peter Knees and Markus Schedl. *Music Similarity and Retrieval*, volume 41. Springer, 2016.
- [5] Gabriel Vigliensoni and Ichiro Fujinaga. Automatic music recommendation systems: Do demographic, profiling, and contextual features improve their performance? *Proceedings of the 17th International Society for Music Information Retrieval Conference*, pages 94–100, 2016.
- [6] Xinxi Wang and Ye Wang. Improving content-based and hybrid music recommendation using deep learning. *Proceedings of the ACM International Conference on Multimedia*, pages 627–636, 2014.
- [7] Òscar Celma. *Music Recommendation and Discovery - The Long Tail, Long Fail, and Long Play in the Digital Music Space*. Springer, 2010.
- [8] B. McFee, T. Bertin-Mahieux, D. P. W. Ellis, and G. R. G. Lanckriet. The million song dataset challenge. *21st International Conference Companion on World Wide Web*, pages 909–916, 2012.
- [9] Luke Barrington, Reid Oda, and Gert Lanckriet. Smarter than genius? human evaluation of music recommender systems. *Proceedings of the 10th International Society for Music Information Retrieval Conference*, pages 357–362, 2009.
- [10] Robin Burke. Hybrid recommender systems: Survey and experiments. *User Model. User-Adapt. Interact.*, pages 331–370, 2002.
- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [12] Alex Graves, Abdel-Tahman Mohamed, and Geoffrey E. Hinton. Speech recognition with deep recurrent neural networks. *Acoustics, Speech and Signal Processing, IEEE International Conference on*, pages 6645 – 6649, 2013.
- [13] Honglak Lee, Yan Largman, Peter Pham, and Andrew Y. Ng. Unsupervised feature learning for audio classification using convolutional deep belief networks. *Advances in Neural Information Processing Systems*, 2009.
- [14] Andrej Karpathy. *Convolutional Neural Networks for Visual Recognition*. Stanford University, 2017. <https://github.com/cs231n/cs231n.github.io/blob/master/convolutional-networks.md#conv>.
- [15] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [16] wikipedia, 2017. https://de.wikipedia.org/wiki/Convolutional_

Neural_Network#/media/File:
3D_Convolution_Animation.gif.

- [17] Zhou Y. and Chellappa R. Computation of optical flow using a neural network. *IEEE International Conference*, 71–78, 1988.