

SOFT-SUBSPACE CLUSTERING ON A HIGH-DIMENSIONAL MUSICAL DATASET

Emil Juzovitski

August 27, 2019

Master Thesis Presentation

Examiner: Pawel Herman

Supervisor: Johan Gustavsson

INTRODUCTION

BACKGROUND

PROBLEM STATEMENT

METHOD

RESULTS

DISCUSSION

What is Clustering analysis?

Finding **clusters** (groups) in a set of data objects

What is Clustering analysis?

Finding **clusters** (groups) in a set of data objects

- **feature similar** data objects partitioned into the same cluster
- **feature dissimilar** data objects partitioned into different cluster.

What is Clustering analysis?

Finding **clusters** (groups) in a set of data objects , based on **similarity**

Example Tasks:

- Fitting products into different aisles in a grocery store
- Grouping distributors based on the products they sell

What is Clustering analysis?

Finding **clusters** (groups) in a set of data objects , based on **similarity**

- Point/Object - vector of features describing the object
- How do we measure similarity?

$$D(X_1, X_2) = \sum_j^m d(x_{1j}, x_{2j})$$

Where **j** is a feature.

Task

Categorizing songs based on anonymous numerical audio features

-
-

Task

Categorizing songs based on anonymous numerical audio features

- Not specifically looking to categorize data based on genre
- Can we find usable themes, composers can not?

INTRODUCTION

BACKGROUND

PROBLEM STATEMENT

METHOD

RESULTS

DISCUSSION

k-means

Partition-based clustering measuring similarity through the Euclidean distance

-

***k*-means**

Partition-based clustering measuring similarity through the Euclidean distance

- Find *k* clusters by iteratively optimizing cluster- centers and members

k-means

Iteratively optimize one variable at a time for cost function $P(U, C)$ where U is the partition matrix, and C the cluster center.

P1. Assign each point to the most similar cluster.

Fix $C = \hat{C}$, optimize $P(U, \hat{C})$

P2. Update the mean for each cluster.

Fix $U = \hat{U}$ optimize $P(\hat{U}, C)$

k-means

Iteratively optimize one variable at a time for cost function $P(U, C)$ where U is the partition matrix, and C the cluster center.

P1. Assign each point to the most similar cluster.

Fix $C = \hat{C}$, optimize $P(U, \hat{C})$

$$u_{il} = \begin{cases} 1 & d(X_i, C_l) \leq d(X_i, C_t) \text{ for } 1 \leq t \leq k \\ 0 & \text{for } t \neq l \end{cases}$$

k-means

Iteratively optimize one variable at a time for cost function $P(U, C)$ where U is the partition matrix, and C the cluster center.

P2. Update the mean for each cluster.

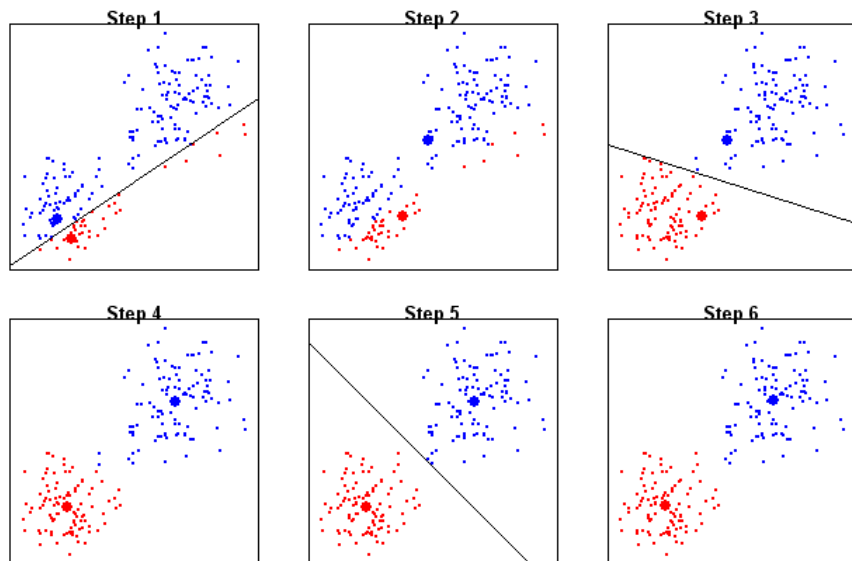
Fix $U = \hat{U}$ optimize $P(\hat{U}, C)$

$$c_{lj} = \frac{\sum_{i=1}^n u_{il} x_{ij}}{\sum_{i=1}^n u_{il}}$$

K-MEANS ITERATION STEPS

1. **CHOOSE INITIAL CLUSTER REPRESENTATIVES C^0**
2. **FIX $C^t = \hat{C}$, OPTIMIZE $P(U^t, \hat{C})$, OBTAIN $P(U^{t+1}, \hat{C})$**
IF $P(U^t, \hat{C}) = P(U^{t+1}, \hat{C})$
RETURN $P(U^t, \hat{C})$ (CONVERGENCE)
3. **FIX $U^{t+1} = \hat{U}$, OPTIMIZE $P(\hat{U}, C^t)$ OBTAIN $P(\hat{U}, C^{t+1})$**
IF $P(\hat{U}, C^t) = P(\hat{U}, C^{t+1})$
RETURN $P(\hat{U}, C^t)$ (CONVERGENCE)
4. **REPEAT STEPS 2. AND 3.**

K-MEANS ITERATION STEPS



Problems:

- Using the wrong features is destructive.
- Knowing what features to use is hard and time-intensive, sometimes impossible.
- Rhythm, Vocals? for theme₂?

Soft-subspace clustering

An extension of K-means with an additional variable to optimize W — allowing cluster unique feature weights.

- Automatic
- Embedded
- Cluster unique feature weights (Subspaces)
- Linear!

SSC ITERATION STEPS

1. **CHOOSE INITIAL CLUSTER REPRESENTATIVES C^0**
2. **FIX $C^t = \hat{C}$, $W^t = \hat{W}$, OPTIMIZE $P(U^t, \hat{C}, \hat{W})$. OBTAIN $P(U^{t+1}, \hat{C}, \hat{W})$**
IF $P(U^t, \hat{C}, \hat{W}) = P(U^{t+1}, \hat{C}, \hat{W})$
RETURN $P(U^t, \hat{C}, \hat{W})$ (CONVERGENCE)
3. **FIX $U^{t+1} = \hat{U}$, $W^t = \hat{W}$, OPTIMIZE $P(\hat{U}, C^t, \hat{W})$. OBTAIN $P(\hat{U}, C^{t+1}, \hat{W})$**
IF $P(\hat{U}, C^t, \hat{W}) = P(\hat{U}, C^{t+1}, \hat{W})$
RETURN $P(\hat{U}, C^t, \hat{W})$ (CONVERGENCE)
4. **FIX $\hat{U} = U^{t+1}$, $\hat{C} = C^{t+1}$, OPTIMIZE $P(\hat{U}^t, \hat{C}^t, W^t)$. OBTAIN $P(\hat{U}, \hat{C}, W^{t+1})$**
IF $P(\hat{U}, \hat{C}, W^t) = P(\hat{U}, \hat{C}, W^{t+1})$.
RETURN $P(\hat{U}, \hat{C}, W^t)$ (CONVERGENCE)

4. **FIX** $\hat{U} = U^{t+1}$, $\hat{C} = C^{t+1}$, **OPTIMIZE** $P(\hat{U}^t, \hat{C}^t, W^t)$. **OBTAIN** $P(\hat{U}, \hat{C}, W^{t+1})$
- IF** $P(\hat{U}, \hat{C}, W^t) = P(\hat{U}, \hat{C}, W^{t+1})$.
- RETURN** $P(\hat{U}, \hat{C}, W^t)$ **(CONVERGENCE)**

1. FSC

$$P(U, C, W) = \sum_{l=1}^k \left[\sum_{i=1}^n \sum_{j=1}^m u_{il} w_{lj}^{\beta} d(x_{ij}, c_{lj}) + \epsilon \sum_{j=1}^m w_{lj}^{\beta} \right] \quad (1)$$

$$w_{lj} = \frac{1}{\sum_{t=1}^m \left[\frac{D_{lj} + \epsilon}{D_{lt} + \epsilon} \right]^{\frac{1}{\beta-1}}} \quad (2)$$

$$D_{lj} = \sum_{x_i \in C_l} d(x_{ij}, c_{lj}) \quad (3)$$

2. EWKM

3. LEKM

1. FSC
2. EWKM

$$P(U, C, W) = \sum_{l=1}^k \left[\sum_{i=1}^n \sum_{j=1}^m u_{il} w_{lj} d(x_{ij}, c_{lj}) + \gamma \sum_{j=1}^m w_{lj} \log(w_{lj}) \right] \quad (1)$$

$$c_{lj} = \frac{\sum_{x_i \in C_l} x_{ij}}{\text{Count}_{x_i \in C_l}} \quad (2)$$

$$w_{lj} = \frac{\exp\left(\frac{-D_{lj}}{\gamma}\right)}{\sum_{t=1}^m \exp\left(\frac{-D_{lt}}{\gamma}\right)} \quad (3)$$

(4)

3. LEKM

1. FSC
2. EWKM
3. LEKM

$$P(U, C, W) = \sum_{l=1}^k \sum_{i=1}^n u_{il} \left[\sum_{j=1}^m w_{lj} \ln [1 + d(x_{ij}, c_{lj})] + \gamma \sum_{j=1}^m w_{lj} \ln(w_{lj}) \right] \quad (1)$$

$$c_{lj} = \frac{\sum_{x_i \in C_l} [1 + d(x_{ij}, c_{lj})]^{-1} x_{ij}}{\sum_{x_i \in C_l} [1 + d(x_{ij}, c_{lj})]^{-1}} \quad (2)$$

$$w_{lj} = \frac{\exp(\frac{-D_{lj}}{\gamma})}{\sum_{t=1}^m \exp(\frac{-D_{lt}}{\gamma})} \quad (3)$$

$$D_{lj} = \frac{\sum_{x_i \in C_l} \ln [1 + d(x_{ij}, c_{lj})]}{\text{Count}_{x_i \in C_l}} \quad (4)$$

INTRODUCTION

BACKGROUND

PROBLEM STATEMENT

METHOD

RESULTS

DISCUSSION

1. *What is a suitable soft-subspace algorithm for the given dataset?*
2. *How does the performance of the chosen soft-subspace clustering algorithm compare to K-means from the perspective of novelty and general quality?*

INTRODUCTION

BACKGROUND

PROBLEM STATEMENT

METHOD

RESULTS

DISCUSSION

Method PS 1. Tuning parameters on an external criteria — genre purity

Selecting the best algorithm.

Method PS 2. Asked a panel of music composers to rate clusters.

Hyperparameter	Description	Algorithm
k	Number of clusters	All
β	Weighting factor for fuzzy-weighting algorithms	FSC
γ	Entropy regularization factor	EWKM, LEKM

Criteria	Description
Audio Similarity	How well do songs in the cluster share audio patterns?
Cultural Similarity	How well do songs in the cluster share an audience?
Playlist Uniqueness	How easy would it be for a composer come up with a similar playlist?
General Quality	How useful is the cluster for the composers?

Cluster 0

Next >

Songs

Take This Chance - Anastacia



Called You Mine - Cape Lion



Sleeping Pad - Bendigo Fletcher



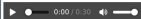
Happiness - 12" Version - The Pointer Sisters



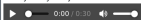
Invisible - Yung Sherman, Uli K, Bladee



Can't You See (feat. The Notorious B.I.G.) - Total, The Notorious B.I.G.



If Only - Maya Payne



Really - Al Kooper, Mike Bloomfield



Don't Look Back In Anger - Remastered - Oasis



I'm Sure - Two Plus Two



Genres

pop - 31.9%

country - 27.3%

rock - 25.6%

pop rock - 14.3%

soul - 9.1%

Common Artists

Chris Young

The Doobie Brothers

Dan + Shay

Evaluation

Cluster Description

Enter A Description

*Audio, Audience, or Whatever you think describes it

General Quality 1

Audio Similarity 1

Cultural Similarity 1

Playlist Uniqueness 1
How different would it be to our existing playlists?

Submit Cluster >

INTRODUCTION

BACKGROUND

PROBLEM STATEMENT

METHOD

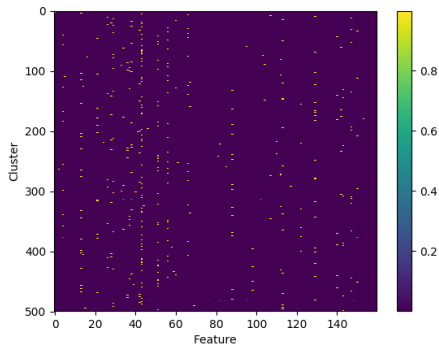
RESULTS

DISCUSSION

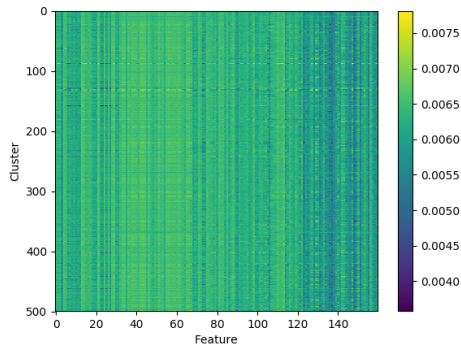
k	K-means	EWKM		LEKM		FSC	
	Purity	γ	Purity	γ	Purity	β	Purity
50	39.0%	0.005	28.2%	0.0	38.8%	30.0	38.0
100	40.8%	0.005	28.9%	2.1	40.6%	15	39.6
500	45.4%	0.001	28.6%	1.4	45.5%	30.0	43.8%

Best parameter- and purity-score (γ, β) given k , where scores represent the mean purity of three runs.

RESULTS PART I.



(a) EWKM



(b) LEKM

RESULTS PART 2.

Source	Mean And Standard Deviation			
	General Quality	Audio Similarity	Cultural Similarity	Playlist Uniqueness
<i>Playlist</i>	7.7 (1.2)	7.8 (1.5)	7.7 (0.6)	4.5 (3.5)
<i>K-means</i>	6.6 (2.5)	7.0 (2.3)	6.2 (2.9)	5.4 (3.2)
<i>LEKM</i>	6.4 (3.4)	6.4 (2.4)	6.0 (2.5)	3.8 (2.6)

Index	One-Way ANOVA			
	General Quality	Audio Similarity	Cultural Similarity	Playlist Uniqueness
$F_{0.05}$	0.21	0.45	0.47	0.37
F_{crit}	4.10	3.98	4.10	3.98
P_{-value}	0.81	0.65	0.64	0.70

1. Cannot differ k-means and LEKM in terms of purity
2. Cannot reject H_0 , $\mu_{playlist} = \mu_{kmeans} = \mu_{LEKM}$!

INTRODUCTION

BACKGROUND

PROBLEM STATEMENT

METHOD

RESULTS

DISCUSSION

- Limitations in testing purity, composer evaluation
- Underlying dataset problems, another limitation
- Dataset properties enjoy uniform weights?

1. Better Initialization
2. Mixed Data SSC.
3. Inter-cluster distances
4. Automize k-selection.

1. OPPONENT QUESTIONS.

2. OTHER QUESTIONS.