

Extraction of Main Dynamics in Platform Metrics

Evelyn Trautmann

Data 2018, Porto

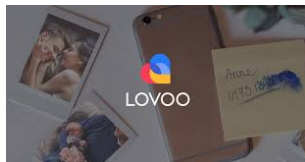
26-28 July, 2018



Lovoo



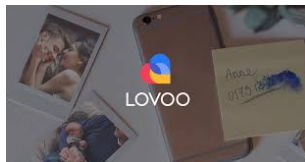
- dating platform with > 70 million users worldwide



Lovoo



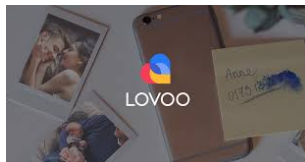
- dating platform with > 70 million users worldwide
- > 1.5 million daily active users, 4.5 million monthly active users



Lovoo



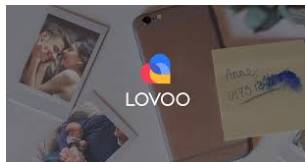
- dating platform with > 70 million users worldwide
- > 1.5 million daily active users, 4.5 million monthly active users
- 420 million chat messages, 4.2 billion match votes per month



Lovoo



- dating platform with > 70 million users worldwide
- > 1.5 million daily active users, 4.5 million monthly active users
- 420 million chat messages, 4.2 billion match votes per month
- $\approx 3.5TB$ analysis data per month



Introduction



- metrics might be influenced by several factors

Introduction



- metrics might be influenced by several factors
- decisions based on incomplete data can be misleading

Introduction



- metrics might be influenced by several factors
- decisions based on incomplete data can be misleading
- platform monitoring means keeping track of various metrics

Introduction



- metrics might be influenced by several factors
- decisions based on incomplete data can be misleading
- platform monitoring means keeping track of various metrics
- some effects are visible only in sub-dimensions like single countries, or device types

Introduction



- metrics might be influenced by several factors
- decisions based on incomplete data can be misleading
- platform monitoring means keeping track of various metrics
- some effects are visible only in sub-dimensions like single countries, or device types
- dimension reduction without losing important details

Clustering



- metrics considered as points in a vector space

Clustering



- metrics considered as points in a vector space
- similarity defines (inverse) distance measure

Clustering

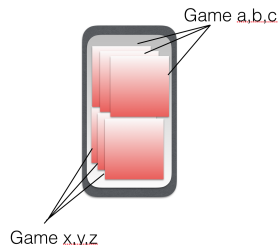
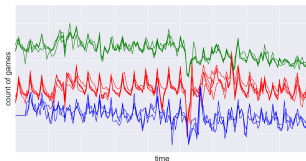


- metrics considered as points in a vector space
- similarity defines (inverse) distance measure
- find clusters of closely related points

Toy Example - Gaming App



- metric 0 game a: count of games played by paying users
- metric 1 game a: count of games played by non paying users
- metric 2 game b: count of games played by paying users
- metric 3 game b: count of games played by non paying users
- ...



Similarity Function



- Similarity Function $\phi : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, 1]$

Similarity Function



- Similarity Function $\phi : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, 1]$
- radial basis function [PVG⁺11]

$$\phi(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$$

Similarity Function



- Similarity Function $\phi : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, 1]$
- radial basis function [PVG⁺11]

$$\phi(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$$

- correlation based similarity

Similarity Function



- Similarity Function $\phi : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, 1]$
- radial basis function [PVG⁺11]

$$\phi(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$$

- correlation based similarity
- handle negative correlations

Similarity Graph

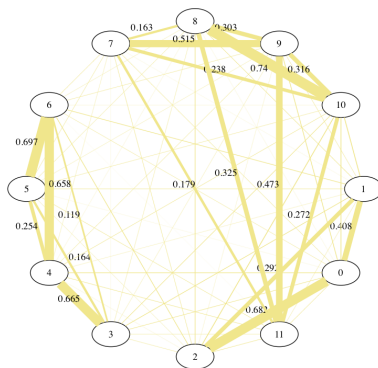


Figure: Similarity Graph

Matrix Representation

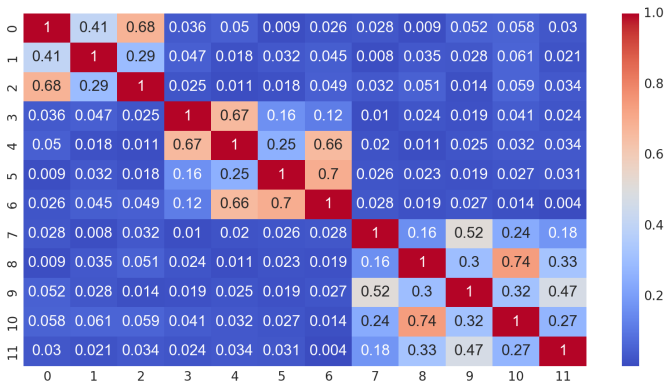


Figure: Matrix Representation

Matrix Representation

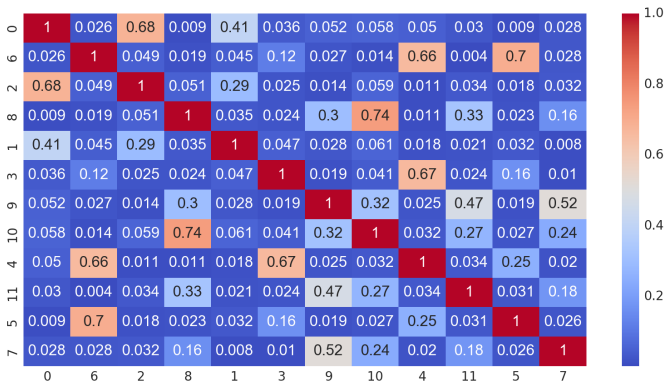


Figure: Matrix Representation

Graph Laplacian



- Graph Laplacian with similarities as off-diagonal entries

Graph Laplacian



- Graph Laplacian with similarities as off-diagonal entries
- negative row sums on diagonal

$$L = \begin{cases} l_{ij} \geq 0, \text{ for } i \neq j \\ l_{ii} = -\sum_{k \neq i} l_{ik}, \text{ for } i = j. \end{cases}$$

Graph Laplacian



- Graph Laplacian with similarities as off-diagonal entries
- negative row sums on diagonal

$$L = \begin{cases} l_{ij} \geq 0, \text{ for } i \neq j \\ l_{ii} = -\sum_{k \neq i} l_{ik}, \text{ for } i = j. \end{cases}$$

- largest eigenvalue is always 0

Graph Laplacian



- Graph Laplacian with similarities as off-diagonal entries
- negative row sums on diagonal

$$L = \begin{cases} l_{ij} \geq 0, \text{ for } i \neq j \\ l_{ii} = -\sum_{k \neq i} l_{ik}, \text{ for } i = j. \end{cases}$$

- largest eigenvalue is always 0
- and the respective eigenvalue is constant

Graph Laplacian



- Graph Laplacian with similarities as off-diagonal entries
- negative row sums on diagonal

$$L = \begin{cases} l_{ij} \geq 0, \text{ for } i \neq j \\ l_{ii} = -\sum_{k \neq i} l_{ik}, \text{ for } i = j. \end{cases}$$

- largest eigenvalue is always 0
- and the respective eigenvalue is constant
- block matrices have multiple largest eigenvalues and piecewise constant eigenvectors

Spectral Clustering [vL07]

- 1: **procedure** SPECTRALCLUSTERING(X, k)
- 2: compute similarity matrix with pairwise similarities
- 3: Transform to Graph Laplacian L
- 4: $v_1, \dots, v_k = \text{eig}(L, k)$ ▷ first k eigenvectors
- 5: $U = V^T$ ▷ k n -dimensional \rightarrow n k -dimensional vectors
- 6: $\text{clusterAssignment} = \text{kmeans}(U, k)$
- 7: **return** clusterAssignment
- 8: **end procedure**

https://github.com/metterlein/spectral_clustering

Number of Clusters



- with clear block structure eigenvalue gap is obvious

Number of Clusters



- with clear block structure eigenvalue gap is obvious
- if connectivity amongst blocks is too large, determination of cluster count is getting complex

Number of Clusters



- with clear block structure eigenvalue gap is obvious
- if connectivity amongst blocks is too large, determination of cluster count is getting complex
- possible approach:
 - PCA explained variance

Data Description



- metrics representing user activities, payments, etc

Data Description



- metrics representing user activities, payments, etc
- each metric is divided in several dimensions like country, gender, etc

Data Description



- metrics representing user activities, payments, etc
- each metric is divided in several dimensions like country, gender, etc
- combination of metrics and dimensions generates around 300 timeseries in our example.

Data Preparation



- aggregation per day

Data Preparation



- aggregation per day
- normalization $X = \frac{1}{\sigma} (X - \mu)$, with μ mean value and σ standard deviation

Data Preparation



- aggregation per day
- normalization $X = \frac{1}{\sigma} (X - \mu)$, with μ mean value and σ standard deviation
- rolling mean of 7 days to smooth weekly periodicity

Real Data Example [PVG⁺11]



Real Data Example.ipynb

https://github.com/metterlein/spectral_clustering

Conclusion



- clustering metrics reduces dimensionality of observation space

Conclusion



- clustering metrics reduces dimensionality of observation space
- small number of interpretable cluster representatives help keeping track of main platform dynamics

Conclusion



- clustering metrics reduces dimensionality of observation space
- small number of interpretable cluster representatives help keeping track of main platform dynamics
- by cluster assignments can be discovered unexpected relations between several metrics

References I



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** (2011), 2825–2830.



Ulrike von Luxburg, *A tutorial on spectral clustering*, CoRR **abs/0711.0189** (2007).

Thank You!

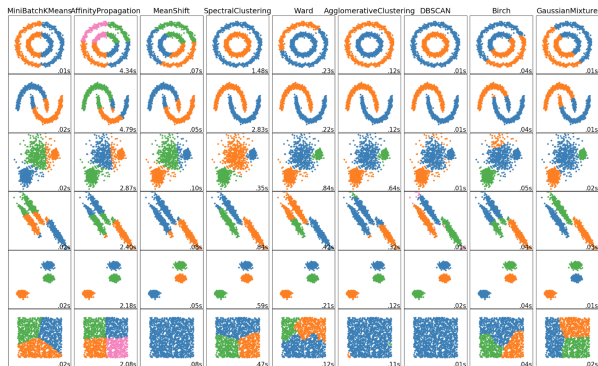
Questions?

`evelyn.trautmann@lovoo.com`

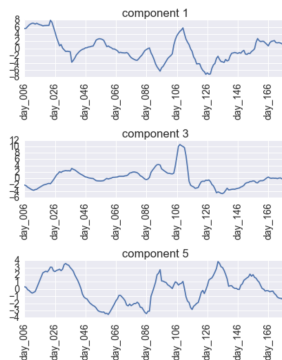
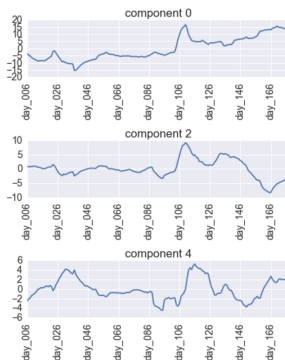
`https://github.com/metterlein/spectral_clustering`

Clustering Methods

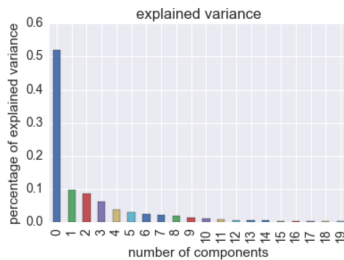
[PVG⁺11]



PCA

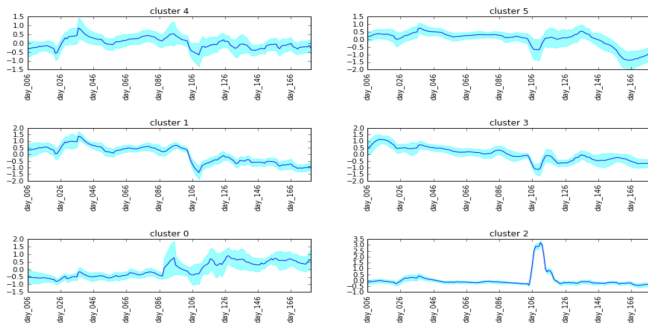


PCA



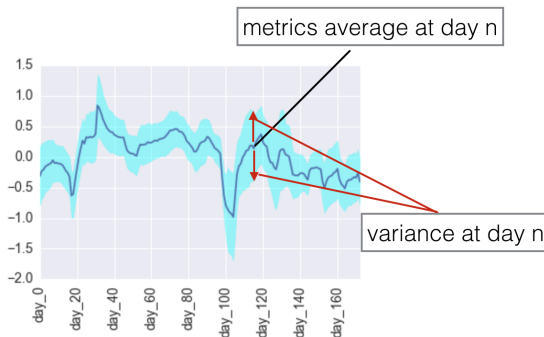
Extract Platform Dynamics by means of Cluster Centers

- compute average timeseries per cluster
- illustrate average cluster timeseries with variance corridor
- visualize metrics types and subdimensions entering respective clusters



Measure Clustering Quality

- compute variance for each timepoint over all cluster members
- optimal clustering minimizes variance for each cluster



Outlook



- investigate cluster assignment change over time

Outlook



- investigate cluster assignment change over time
- hierarchical approach: clustering sub-blocks

Outlook



- investigate cluster assignment change over time
- hierarchical approach: clustering sub-blocks
- add time shifted series to recognize Granger causalities