# Case Study of Data Mining Mutual Engagement

Nick Bryan-Kinns
Queen Mary University of London
London, E1 4NS. UK
n.bryan-kinns@qmul.ac.uk

**This paper presents a case study of using conventional Data Mining techniques to identify clusters of creative contributions made by users in a collaborative music making system. Results of the clustering suggest that when people mutually engage with each other they tend to converge on similar creative contributions, whereas smaller clusters may indicate higher quality contributions. We also show how we used Data Mining to discriminate between kinds of creative contributions including: complex melodic structures, simple melodic structures, musical motifs, and rhythmical contributions. Our use of Data Mining does not require direct interaction with users and so it may be useful in real-world study contexts.**

*Creativity; Creativity Support Tools; Mutual Engagement; Data Mining; Evaluation; Interaction Corpus*

## 1. INTRODUCTION

Collaborative music making is a rich domain in which to explore creativity and the design of Creativity Support Tools (CSTs) [24] as it is a fundamental form of human creativity which relies on mutual engagement between people to coproduce a joint creation [6]. CSTs for music making typically emerge from the field of New Interfaces for Musical Expression (NIME) [1] which has an interest in evaluating these creative and expressive systems [2]. Researchers in broader fields such as Sonic Interaction Design (SID) [22] highlight the need for evaluation methods for creative interactions involving sound as the primary modality. However, as noted by [16], evaluating creativity and CSTs is not straightforward. In addition to self-reporting by participants (e.g. [13] and [10]) there is the potential to use measures of people's interaction with CSTs and to assess features of the creative products themselves to evaluate creativity, collaboration, and user interfaces. In this paper we describe a case study of applying conventional Data Mining methods to analyzing logs of participants' interaction within a collaborative music making tool.

The field of Data Mining (DM) explores how large data sets can be modelled [20] for example by identifying recurring patterns and clusters within the data. This makes it an ideal candidate to analyze large sets of human-computer interaction data generated in creative activities such as music making. To date DM has been used to explore people's patterns of interaction with user interfaces in a 'data-driven approach to understanding user behavior' [26] for example predicting user behavior [11], categorizing users' interests [14] and identifying personality traits [4]. Users'

interaction with websites has also been used to identify online behavioral patterns [26], to identify emergent topics in online chat [9], and to model influence in social media [25]. To date there has been limited application of DM to clustering the content of the actual creative activities within CSTs. For example, [21] describe how DM is used to classify visual designs of web pages but it relies on extensive human labelling of content which is not feasible in the long term. In contrast [23] show how communications between designers can be mined to evaluate their collaboration, but the focus is on the communication about the design rather than the content of the designs themselves. We are interested in the feasibility of using Data Mining to investigate patterns of creative interaction between people in CSTs. In particular, how DM techniques can be used to identify clusters of creative contributions, and whether these clusters are affected by user interface features of CSTs.

## 2. DATA MINING MUTUAL ENGAGEMENT

An on-line collaborative music making system referred to as Daisyfield [6] allows participants to co-edit short loops (48 beats; 5 seconds) of music through a web based user interface. Daisyfield was designed to support research into *mutual engagement* – the points at which people creatively spark together [5]. In an empirical study [6] of 24 trios of participants (72 in total) over 100Mb of logged interaction data was collected as participants created music together over 18 hours. In the study Daisyfield's user interface configuration was manipulated in two ways to produce four versions: i) *authorship*: whether each participant's contributions had

a unique color or all participants' contributions were the same color; ii) *awareness*: whether participants could see each others' mouse pointers or not. The Mutual Engagement Interaction Corpus (MEIC) [8] contains a record of all 524,092 states of the collaboration and all musical interaction between participants in the studies, and [6] describes participants' creative interactions.

Figure 1 illustrates Daisyfield - the musical score is represented by the set of flower like circles referred to as *daisies*. Each daisy contains a circle of dots, and the currently played set of notes is indicated by the line radiating clockwise from the center. In this way multiple musical parts can be created and edited by multiple people at the same time. Each participant can create, edit, and delete their own and others' notes and daisies. The overall sound produced is one shared musical loop.
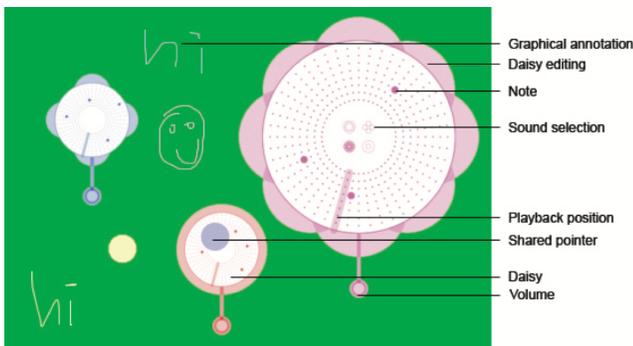


*Figure 1: Daisyfield User Interface from [6]*

## 2.1 Data Mining Workflow

We applied a conventional Data Mining Workflow to the MEIC in order to identify clusters of musical contributions by participants and to explore whether these clusters were somehow influenced by the user interface configuration. Our workflow involved: i) *Segmentation*; ii) *Clustering*; iii) *Visualization*.

### 2.1.1 Segmentation
The first step in Data Mining is typically to segment raw input data into a bag of words (BOW) – a collection of units of data. In text domains the bag of words contains text words extracted from the source documents. In the musical domain a different method of segmentation of the input data is required.

In Daisyfield a musical loop is usually composed of several shorter musical phrases. Therefore the loops need to be segmented into the short musical phrases which we refer to as patterns of notes (*PoN*).

Segmentation heuristics for identifying a PoN are based on [12] and proposed in [7]. A PoN has:

(i) At least three notes (not pauses), and

(ii) No more than two pauses between notes.

Segmentation code was written in MATLAB [18], and produced 225,097 PoNs from the MEIC (62,416 unique PoNs). Figure 2 shows an example PoN extracted using the segmentation heuristics. In the figure the vertical axis represents pitch of the notes in the sequence. The white squares indicate a musical note. The horizontal axis represents time from left to right. In this example there is a clear sequence of rising pitches of notes (white squares) followed by a sequence of descending notes. There are also some additional notes in the sequence which would create chorded notes in the sequence. This example is one of 345 unique PoNs extracted from 1432 PoNs produced by one participant in one fifteen minute creative session.



*Figure 2: An example PoN*

### 2.1.2 Clustering
There are two main DM approaches to identifying similar data in a BOW: i) *supervised* in which DM algorithms find more examples of researcher-specified patterns; and ii) *unsupervised*, which relies on bottom up data-driven algorithmic discovery of patterns. We are interested in discovering kinds of creative contributions and so we use unsupervised DM methods to discover *clusters* of PoNs. Clustering of data requires metrics of distance between elements in the BOW. Considering the PoNs as two dimensional binary matrices there are two common distance metrics: Binary Hamming Distance (BHD); Squared Euclidean distance of Local Binary Pattern (LBP).

Clustering code was written in MATLAB using standard MATLAB functions including *kmeans* with *hamming* and *sqeuclidean* distance metrics for the k-means clustering algorithm [17] which is a commonly used approach to clustering large sets of data.

Figure 3 illustrates six clusters generated for one participant in one session. The number of clusters is arbitrarily selected to provide an initial view into the data. The top row of the figure shows the averaged PoN in a cluster, and the columns show the top 15 PoNs in each cluster.



*Figure 3: BHD: PoN clusters for one participant*

2

In DM the clustering of large data sets needs to be optimized to avoid cluster over-saturation and outliers. An established method for optimizing the number of clusters is to use the silhouette [15] of data in clusters – a measure of how similar the data in a cluster is in comparison to data outside the cluster. Using the MATLAB *silhouette* function the optimal number of clusters of PoNs for the MEIC was found to be 26.

Figure 4 shows the results of the clustering of all unique PoNs for the whole data set for 26 clusters. Visually, these 26 clusters provide information about the *kinds* of PoNs that have been generated by participants, especially in the LBP clusters which illustrate the range of shapes of PoNs that have been generated by participants. For example, LBP generates several clusters of undulating PoNs, and some diagonal patterns similar to those see in Figure 3. In contrast, BHD highlights several PoN outliers in the clusters (indicated by the higher contrast average PoNs which indicate less spread of data), and clusters shorter PoNs together than LBP. The outliers highlighted by BHD could be used as indicators of novel contributions by participants, whereas the LBP clustering gives us a better spread of the typical kinds of PoNs created.
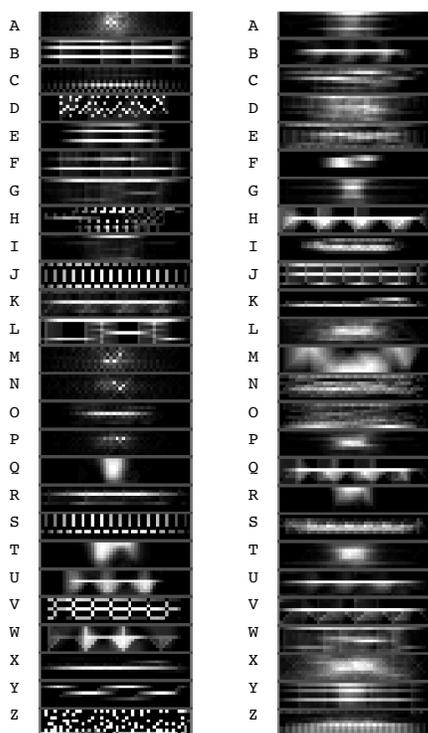


*Figure 4: a) BHD Clusters, b) LBP Clusters*

Examining the clusters of PoNs indicates that the following are typical PoNs created by participants as illustrated in Figure 4 where the clusters are referred to by the clustering method followed by the cluster number e.g. LBP.D is the fourth cluster from the top of the LBP clusters:

- **Complex melodic structure** e.g. BHD.U, BHD.W, LBP.H, LBP.Q, LBP.V. These clusters contain repeating, undulating, long PoNs which rise and fall several times, some of which extend to the whole length of the sequence. These PoNs would provide a consistent and coherent structure for the composition.

- **Simple melodic structure** e.g. BHD.L, BHD.Y, LBP.C. These PoNs alternate between two pitches over the length of the sequence, and would provide a simple (two tone) structure for the composition.

- **Musical motif** e.g. BHD.A, BHD.M, BDH.N, BHD.P, BHD.Q, LBP.F, LBP.G, LBP.L, LBP.P, LBP.R, LBP.T. These are short PoNs centered around the middle pitch, and would act as melodic motif or stand-out phrase in a composition.

- **Rhythmical** PoNs e.g. BHD.B, BHD.E, BHD.G, BHD.K, BHD.O, BHD.R, BHD.Y, LBP.B, LBP.K, LBP.W. These PoNs contain multiple notes of a single pitch – the timing of the repetition of the notes would provide a rhythmic structure for the composition.

### 2.2 Visualization

To visualize the effect of the user interface configurations of Daisyfield on the kinds of PoNs created, the clusters of PoNs were plotted against UI configurations. In the graphs on the left of Figure 5 and Figure 6 the x axis indicates whether or not Colors were present in the user interface (left: No Colors; right: Colors), and the y axis indicates whether or not Shared Pointers were visible (bottom: No Shared Pointers; top: Shared Pointers) e.g. a cluster of PoNs which only occur in conditions with Shared Pointers and Colors would be plotted in the top right corner of the graph. Figure 5 uses BHD and Figure 6 uses LBP to cluster the PoNs. In these figures the size of the circle indicates the number of PoNs in the cluster. The sizes of clusters are also detailed in the bar chart to the right of the plot, and for reference the cluster label corresponds to the labels in Figure 4.

The BHD.A and LBP.G clusters are the largest clusters identified (musical motifs from Figure 4). Furthermore, BHD.A and LBP.G are both plotted slightly off-center in the graphs with their centers being above left of the center of the graph. This indicates that BHD.A and LBP.G are slightly more frequent in conditions with Shared Pointers and No Colors. Indeed, examining the largest clusters (BHD.A, BHD.M, BHD.N, BHD.P, and LBP.G, LBP.T) shows that the **largest clusters are all musical motifs**, and are all plotted slightly up and left from center. This suggests that **having No Colors and Shared Pointers encourages participants to create more similar PoNs across groups** than other conditions, i.e. large clusters indicate that many people are creating similar PoNs which may indicate

convergent creativity between people. The experiment in [6] found that participants were more mutually engaged when they had No Colours and Shared Pointers, so production of **musical motifs may be an indicator of mutual engagement between participants**. For example, with No Colours and Shared Pointers, participants in [6] reported that they left more involved with the group significantly more frequently than in other conditions, and reported that they understood what was going on significantly more frequently than in other conditions.
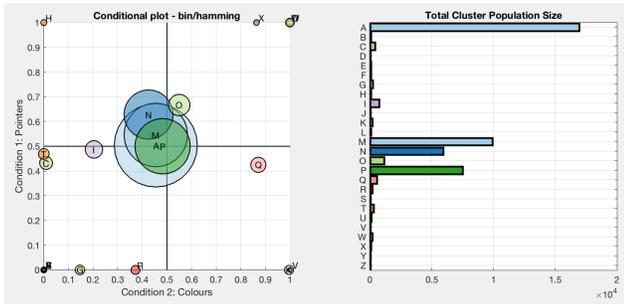


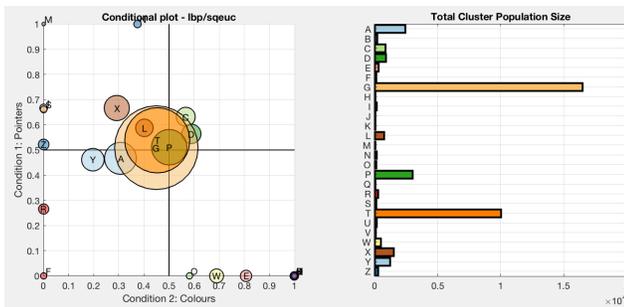**Figure 5**: *Plot of BHD Clusters against UI configurations*



**Figure 6**: *Plot of LBP Clusters against UI configurations*

In terms of outliers, both plots show a number of smaller clusters. However, unlike the larger clusters, there is no consistency between the spread of these smaller groups when using BHD and LBP which indicates that these **smaller clusters may be artefacts of the clustering method** rather than an indication of some participant behavior or an effect of experimental conditions. An alternative view is that **smaller clusters may indicate more novel and innovative creativity** of participants as they are not repeated across sessions or groups. Looking at the spread of the small clusters it appears that with LBP there are a large number of outliers with Colors and No Shared Pointers which the experiment in [6] suggests **may be better compositions**. However, in contrast, BHD plots the most outliers with Colors and Shared Pointers, and No Shared Pointers and No Colors. These differences make it difficult to draw concrete conclusions about the role of the outliers.

## 3. SUMMARY

We showed how Data Mining was used to identify clusters of creative contributions in logs of use of a collaborative music making system, Dasiyfield. We a segmentation method for musical contributions for DM and developed a method to visualize the effect of user interface features on clusters which helped us to explore how these features might affect creative collaborations.

The case study showed how DM could be used to identify different kinds of musical contribution: complex melodic structures; simple melodic structures; musical motifs; and rhythmical contributions. In terms of mutual engagement [6] we found that small clusters identified by LBP may be better compositions, and that production of musical motifs may be an indicator of mutual engagement between participants.

As with the use of DM in other domains, we found that it is important to consider: segmentation methods; optimization of cluster numbers; choice of distance metrics; choice of clustering algorithm; and interpretation of outliers. As we are using DM as a tool for exploring data we need to consider how the clusters are visualized for inspection. In Figure 4 we chose to visualize the clusters as the average of all PoNs in that cluster. This provides a rich feeling for the shape of the PoNs in the cluster, but can also produce quite ambiguous results, such as LBP.X which doesn't convey a great deal of salient information about the constituent PoNs. Alternative visualizations to be tested for utility include averaging a small number of PoNs for each cluster, or selecting the most central PoN in the cluster as the representative.

Our DM approach is non-intrusive – it does not require direct interaction with participants – so we believe it could be useful in analyzing behavior in other creative domains and contexts. Our future work will focus on analyzing people's creative contributions in live NIME performances [1] where logs of interaction would be readily available from the digital music instruments, and in public interactive art domains such as "Living Laboratories" [19] in which "the exhibition becomes a site for collaboration between curators, artists, and audiences" (ibid.). Such analysis of real-world data logs could also be used in conjunction with the increasing number of HCI techniques for studying interactive experiences in public places such as museums, galleries, libraries, festivals, and open spaces [3].

## 4. ACKNOWLEDGEMENTS

## 5. REFERENCES

1. Refsum Jensenius Alexander and Michael J. Lyons (eds.). 2017. *A NIME Reader: Fifteen Years of New Interfaces for Musical Expression*. Springer International Publishing AG, Switzerland. ISBN 9783319472133

2. Jeronimo Barbosa, Joseph Malloch, Marcelo M. Wanderley, and Stéphane Huot. 2015. What does "Evaluation" mean for the NIME community? In *Proceedings of the International Conference on New Interfaces for Musical Expression*, Baton Rouge, LA, USA, May 31-June 3, 2015.

3. Steve Benford, Chris Greenhalgh, Andy Crabtree, Martin Flintham, Brendan Walker, Joe Marshall, Boriana Koleva, Stefan Rennick Egglestone, Gabriella Giannachi, Matt Adams, Nick Tandavanitj, and Ju Row Farr. 2013. Performance-Led Research in the Wild. ACM Trans. Comput.-Hum. Interact. 20, 3, Article 14 (July 2013), 22 pages. DOI=http://dx.doi.org/10.1145/2491500.2491502

4. Eli T Brown, Alvitta Ottley, Helen Zhao, Quan Lin, Richard Souvenir, Alex Endert, Remco Chang. 2014. Finding waldo: Learning about users from their interactions. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):1663–1672.

5. Nick Bryan-Kinns and Fraser Hamilton. 2009. Identifying Mutual Engagement. *Behaviour & Information Technology*. DOI: 10.1080/01449290903377103

6. Nick Bryan-Kinns. 2013. Mutual Engagement and Collocation with Shared Representations. *International Journal of Human Computer Studies*. DOI: http://dx.doi.org/10.1016/j.ijhcs.2012.02.004

7. Nick Bryan-Kinns. 2014. Mutual Engagement in Digitally Mediated Public Art. In *Interactive Experience in the Digital Age*, Candy, L., & Ferguson, S. (Eds.), Springer.

8. Nick Bryan-Kinns. 2018. *Mutual Engagement Interaction Corpus*. Published May 2018 at https://c4dm.eecs.qmul.ac.uk/rdr/handle/123456789/43

9. Shuo Chang, Peng Dai, Jilin Chen, and Ed H. Chi. 2015. Got Many Labels?: Deriving Topic Labels from Multiple Sources for Social Media Posts using Crowdsourcing and Ensemble Learning. In *Proceedings of the 24th International Conference on World Wide Web (WWW '15 Companion)*. ACM, New York, NY, USA, 397-406. DOI: http://dx.doi.org/10.1145/2740908.2745401.

10. Erin Cherry and Celine Latulipe. 2014. Quantifying the Creativity Support of Digital Tools through the Creativity Support Index. *ACM Trans. Comput.-Hum. Interact. 21, 4, Article 21 (June 2014)*, 25 pages. DOI: http://dx.doi.org/10.1145/2617588.

11. Ed H. Chi, Peter Pirolli, and James Pitkow. 2000. The scent of a site: a system for analyzing and predicting information scent, usage, and usability of a Web site. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems (CHI '00)*. ACM, New York, NY, USA, 161-168. DOI=http://dx.doi.org/10.1145/332040.332423

12. Roger B. Dannenberg and Ning Hu. 2002. Pattern Discovery Techniques for Music Audio. In *Proceedings of 3rd International Conference on Music Information Retrieval (ISMIR), 2002.*

13. Nicholas Davis, Chih-PIn Hsiao, Kunwar Yashraj Singh, Lisa Li, and Brian Magerko. 2016. Empirically Studying Participatory Sense-Making in Abstract Drawing with a Co-Creative Cognitive Agent. In *Proceedings of the 21st International Conference on Intelligent User Interfaces (IUI '16)*. ACM, New York, NY, USA, 196-207. DOI: https://doi.org/10.1145/2856767.2856795

14. Jeffrey Heer and Ed H. Chi. 2002. Separating the swarm: categorization methods for user sessions on the web. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '02)*. ACM, New York, NY, USA, 243-250. DOI=http://dx.doi.org/10.1145/503376.503420

15. Leonard Kaufman and Peter J. Rouseeuw. 1990. *Finding Groups in Data: An Introduction to Cluster Analysis*. Hoboken, NJ: John Wiley & Sons, Inc.

16. Andruid Kerne, Andrew M. Webb, Steven M. Smith, Rhema Linder, Nic Lupfer, Yin Qu, Jon Moeller, and Sashikanth Damaraju. 2014. Using metrics of curation to evaluate information-based ideation. *ACM Trans. Comput.-Hum. Interact.* 21, 3, Article 14 (June 2014), 48 pages. DOI: http://dx.doi.org/10.1145/2591677

17. Stuart P. Lloyd. 1982. Least Squares Quantization in PCM. *IEEE Transactions on Information Theory*. 28, 129–137.

18. MATLAB. 2017. Retrieved August 2017 from https://uk.mathworks.com/products/matlab.html

19. Lizzie Muller and Ernest Edmonds. 2006. Living Laboratories: Making and Curating Interactive Art. In *ACM SIGGRAPH 2006 Art gallery (SIGGRAPH '06)*. ACM, New York, NY, USA, Article 160. DOI=http://dx.doi.org/10.1145/1178977.1179120

20. Anand Rajaraman, Jeffrey David Ullman. 2011. *Mining of Massive Datasets*. Cambridge University

Press. DOI:
https://doi.org/10.1017/CBO9781139058452

21. Arvind Satyanarayan, Maxine Lim, and Scott Klemmer. 2012. A platform for large-scale machine learning on web design. In *CHI '12 Extended Abstracts on Human Factors in Computing Systems (CHI EA '12).* ACM, New York, NY, USA, 1697-1702. DOI: http://dx.doi.org/10.1145/2212776.2223695

22. Stefania Serafin, Karmen Franinovic, Thomas Hermann, Guillaume Lemaitre, Michal Rinott, and Davide Rocchesso. 2011. Sonic interaction design. In Thomas Hermann, Andy Hunt, and John G. Neuhoff (Eds), *The Sonification Handbook*, Chapter 5, 87–110. Logos Publishing House, Berlin.

23. Simeon J Simoff, Mary Lou Maher. 2000. Analysing participation in collaborative design environments. Design Studies, 21(2), 119-144, Elsevier, ISSN 0142-694X

24. Ben Shneiderman. 2007. Creativity support tools: Accelerating discovery and innovation. *Communications of the ACM,* 50(12):20–32.

25. Jie Tang, Jimeng Sun, Chi Wang, and Zi Yang. 2009. Social influence analysis in large-scale networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '09)*. ACM, New York, NY, USA, 807-816. DOI: https://doi.org/10.1145/1557019.1557108

26. Gang Wang, Xinyi Zhang, Shiliang Tang, Haitao Zheng, and Ben Y. Zhao. 2016. Unsupervised Clickstream Clustering for User Behavior Analysis. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI '16)*. ACM, New York, NY, USA, 225-236. DOI: https://doi.org/10.1145/2858036.2858107