

On the genre-fication of Music: a percolation approach (long version)

R. Lambiotte* and M. Ausloos†

SUPRATECS, Université de Liège, B5 Sart-Tilman, B-4000 Liège, Belgium

(Dated: 09/07/2005)

In this paper, we analyze web-downloaded data on people sharing their music library. By attributing to each music group usual music genres (Rock, Pop...), and analysing correlations between music groups of different genres with percolation-idea based methods, we probe the reality of these subdivisions and construct a music genre cartography, with a tree representation. We also show the diversity of music genres with Shannon entropy arguments, and discuss an alternative objective way to classify music, that is based on the complex structure of the groups audience. Finally, a link is drawn with the theory of *hidden variables* in complex networks.

PACS numbers: 89.75.Fb, 89.65.Ef, 64.60.Ak

I. INTRODUCTION

Take a sample of people, and make them listen to a list of songs. If a majority of people should find an agreement on basic subdivisions, like *Rock/Jazz/Pop...*, a more refine description will lead to more and more disparate answers, even contradictions. These originate from the different background, taste, music knowledge, mood or *network of acquaintances* [1] of the listeners, i.e. in a statistical physics description, these processes correspond to ageing, internal fluctuations and neighbour-neighbour interactions. The more and more eclectic music offer, together with the constant mixing of old music genres into new ones make the problem still more complicated. Even artists seem to avoid the usual classifications by refusing to enter well-defined yokes, and prefer to characterise themselves as a unique mix-up of their old influences [2].

Obviously, categorising music, especially into finer genres or subgenres, is not an easy task, and is strongly subjective. This task is also complicated by the constant birth of newly emerging styles, and by the very large number of existing sub-divisions. For instance, the genre *Electronic music* is divided in *wikipedia* [3] into 9 sub-genres (*Ambient, Breakbeat...*), each of them being divided into several subsubgenres. This categorising is becoming more and more complex in the course of time.

This paper tries to find an answer to the above problems by showing in an "objective" way the existence of music trends that allow to classify music groups, as well as the relations between the usual genres and sub-genres. To do so, we use web-downloaded data from the web, and define classifications based on the listening habits of the groups audience. Thereby, we account for the fact that music perception is driven both by the people who make music (artists, Majors...), but also by the people who listen to it. Our analysis consists in characterising a large sample of individual musical habits from a statistical physics point of view, and in extracting collective

trends. In a previous work [4], we have shown that such collective listening habits may lead to the usual music subdivisions in some particular cases, but also to unexpected structures that do not fit the neat usual genres defined by the music industry. Those represent the non-conventional taste of listeners. Let us note that alternative music classifications based on signal analysis may also be considered [5, 6].

In section II, we describe the methodology, namely the analysis of empirical data from collaborative filtering websites, e.g. *audioscrobbler.com* and *musicmobs.com*. We will also give a short review of the statistical methods introduced in [4]. Mainly, these methods consist in evaluating the correlations between the groups, depending on their audience, and in using filtering methods, i.e. percolation idea-based (PIB) methods, in order to visualise the collective behaviours. In section III, we attribute lists of genres to a sample of music groups, by downloading data from the web. These data, that describe the different tags, i.e. genres, used by people to classify music groups, are analysed by using the Shannon entropy as a measure of the music group diversity. By examining correlations between these different music genres, we also use the statistical methods of section II in order to make a map of music genres (see [7] for an example from the social science). This cartography is justified by the fact that *alike* music genres are statistically correlated by their audience. It is shown that these correlations are *homophilic* [8], i.e. *alike* music genres tend to be listened to by the same audience. Homophily is known to occur in many *social* systems, including online communities [9], co-authorship networks [10, 11] and linking patterns between political bloggers [12].

Let us stress that the issues of this work are part of the intense ongoing physicist research activity on opinion formation [13, 14, 15, 16, 17], self organisation on networks [18, 19], including clique formation [20], percolation transitions [21], as well as on the identification of *a priori* unknown collective behaviours in complex networks [22], e.g. proteins [23], genes [24], linguistics [25, 26], industrial sectors [27], groups of people [28]...

*Electronic address: Renaud.Lambiotte@ulg.ac.be

†Electronic address: Marcel.Ausloos@ulg.ac.be

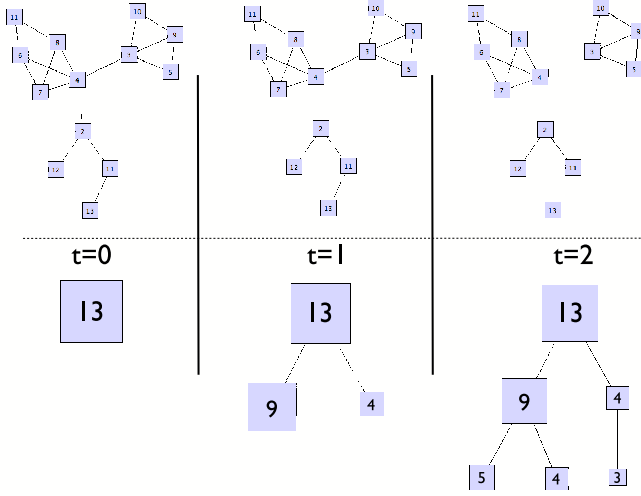


FIG. 1: Branching representation of a squared correlation matrix of 13 elements. At each increasing step ($t=0,1,2$) of the filter ϕ , links are removed, so that the network decomposes into isolated islands. These islands are represented by squares, whose size depends on the number of nodes in the island. Islands composed by only one music group are not depicted. Starting from the largest island, branches indicate a parent relation between the islands. The increasing filter method is applied until all links are removed.

II. METHODOLOGY

A. Data analysis

In this work, we analyze data retrieved from collaborative filtering websites (see [29] for a detailed definition). These sites propose people to share their profiles and experiences in order to help them discover new musics/books... that should (statistically) correspond to their own taste. In the present case, we focus on a database downloaded from *audioscrobbler.com* in January 2005. It consists of a listing of users (each represented by a number), together with the list of music groups the users own in their library. This structure directly leads to a bipartite network for the whole system. Namely, it is a network composed by two kinds of nodes, i.e. the persons, called users or listeners in the following, and the music groups. The network can be represented by a graph with edges running between a group i and a user μ , if μ owns i .

In the original data set, there are 617900 different music groups, although this value is skewed due to multiple (even erroneous) ways for a user to characterise an artist (e.g. *The Beatles*, *Beatles* and *The Beetles* count as three music groups) and 35916 users. On average, each user owns 140 music groups in his/her library, while each group is owned by 8 persons. For completeness, let us note that the listener with the most groups possesses 4072 groups (0.6% of the total music library) while the group with the largest audience, *Radiohead*, has 10194

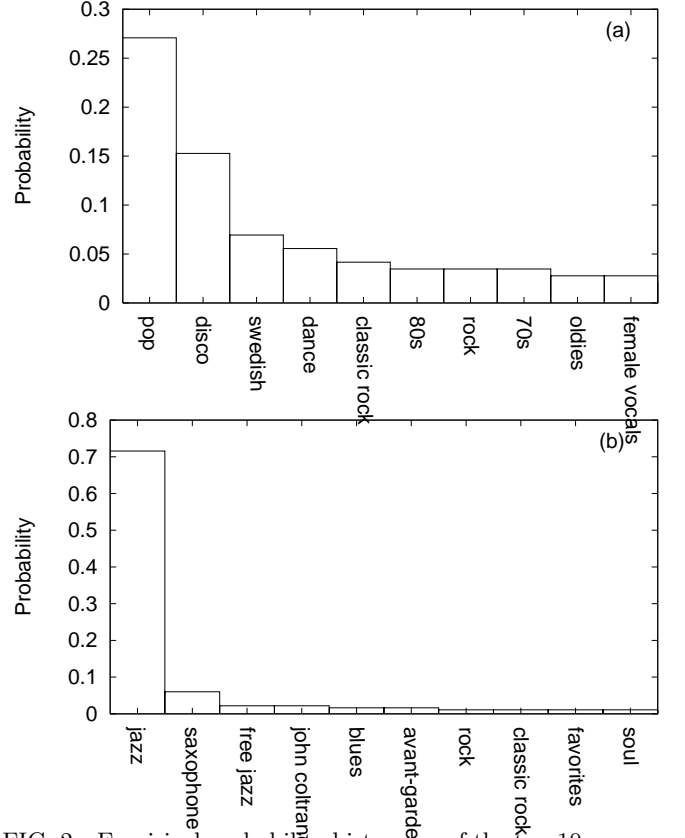


FIG. 2: Empirical probability histogram of the top 10 genres tagged by listeners to *ABBA* (a), and to *John Coltrane* (b). The data have been downloaded from <http://www.lastfm.com> in August 2005.

users (28% of the user community). This asymmetry in the bipartite network is expected as users have in general specific tastes that prevent them from listening to any kind of music, while there exist *mainstream* groups that are listened to by a very large audience. Let us stress that this asymmetry is also observable in the degree distributions for the people and for the groups.

In the following, we make a selection in the total number of groups for computational reasons, namely we have analysed a subset composed of the top 1000 most-owned groups. This limited choice was also motivated by the possibility to identify these groups at first sight.

B. Percolation idea-based filtering

In this section, we review the method introduced in [4] in order to extract collective structures from the data set. Each music group i is characterised by its signature, i.e. the vector:

$$\vec{\Gamma}^i = (..., 1, ..., 0, ..., 1, ...) \quad (1)$$

of n_L components, where $n_L = 35916$ is the total number of users in the system, and where $\Gamma_\mu^i = 1$ if the listener μ owns group i and $\Gamma_\mu^i = 0$ otherwise. By doing so, we

consider that the audience of a music group, i.e. the list of persons listening to it, identifies its signature.

In order to quantify the correlations between two music groups i and j , we calculate the symmetric correlation measure:

$$C^{ij} = \frac{\bar{\Gamma}^i \cdot \bar{\Gamma}^j}{|\bar{\Gamma}^i| |\bar{\Gamma}^j|} \equiv \cos \theta_{ij} \quad (2)$$

where $\bar{\Gamma}^i \cdot \bar{\Gamma}^j$ denotes the scalar product between the two n_L -vector, and $||$ its associated norm. This correlation measure, that corresponds to the cosine of the two vectors in the n_L -dimensional space, vanishes when the groups are owned by disconnected audiences, and is equal to 1 when their audiences strictly overlap.

In order to extract families of alike music groups from the correlation matrix C^{ij} , we use the PIB method [4]. We define the filter coefficient $\phi \in [0, 1[$, and filter the matrix elements so that $C_\phi^{ij} = 1$ if $C^{ij} > \phi$, and let $C_\phi^{ij} = 0$ otherwise. Starting from $\phi = 0.0$, namely a fully connected network, increasing values of the filtering coefficient remove less correlated links and lead to the shaping of well-defined islands, completely disconnected from the main island. Let us stress that this systematic removal of links is directly related to percolation theory, and that the internal correlations in the network displace and broaden the percolation transition [4, 32]. From a statistical physics point of view, the meaning of ϕ is that of the inverse of a temperature, i.e. high values of ϕ restrain the system to highly correlated islands; in the same way, low temperature restrains phase space exploration to low lying free energy wells. This observation suggests that PIB methods should be helpful in visualising free energy profiles and reaction coordinates between metastable states [30].

A branching representation of the community structuring is used to visualise the process (see Fig.1 for the sketch of three first steps of an arbitrary example). To do so, we start the procedure with the lowest value of $\phi = 0.0$, and we represent each isolated island by a square whose surface is proportional to its number of nodes (the music groups). Then, we increase slightly the value of ϕ , e.g. by 0.01, and we repeat the procedure. From one step to the next step, we draw a bond between emerging sub-islands and their parent island. The filter is increased until all bonds between nodes are eroded (that is, there is only one node left in each island). Let us note that islands composed by only one music group are not depicted, as these *lonely* music groups are self-excluded from the network structure, whence from any genre. Applied to the above correlation matrix C^{ij} , the tree structure gives some insight into the diversification process by following branches from their source (top of the figure) toward their extremity (bottom of the figure). The longer a given branch is followed, the more likely it is forming a well-defined music genre.

In [4], we have shown that the resulting tree representation exhibits long persisting branches, some of them

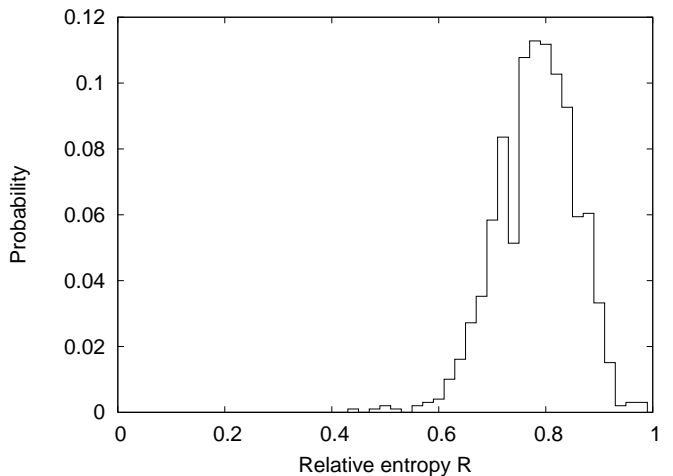


FIG. 3: Empirical probability histogram of the relative entropy R_i (see text for definition), obtained for the top 1000 music groups. The tagged genres have been downloaded from <http://www.lastfm.com> in August 2005.

leading to standard, homogenous style groupings, such as [Kylie Minogue, Dannii Minogue, Sophie Ellis Bextor] (dance pop), while many other islands are harder to explain from a standard genre-fication point of view and reveal evidence of unexpected collective listening habits.

III. GENRE CARTOGRAPHY

A. Measure of diversity

In view of the above analysis, attributing genres to music groups is a difficult problem. This complexity is made clearer by observing the different ways listeners characterise the same music group. To perform this analysis, we have downloaded from <http://www.lastfm.com> a list of the descriptions, i.e. genres, that people tag to music groups in their music library, together with the number of times this description occurred. For instance, from this site, one gets that *ABBA* (Fig.2a) is described by an eclectic range of different music sub-divisions. These sub-divisions are based on the group style (Pop, Rock...), on the time period (80s, 70s...) or on geographical grounds (Swedish) and their choice depends on the listener, i.e. his perception and subjective way to characterise music (see first paragraph of the introduction).

For this work, we have downloaded these lists of genres for the top 1000 groups, thereby empirically collecting a statistical genre-fication of the music groups. Let us stress that the data could not be downloaded for 5 of the groups, due to misprints in their name, e.g. *Bjadrk* instead of *Björk*. Consequently, we focus in the following on the $n_G = 995$ remaining music groups. One should also note that <http://www.lastfm.com> limits access to the top 25 genres of each group.

The statistical genre-fication of the sample may ex-

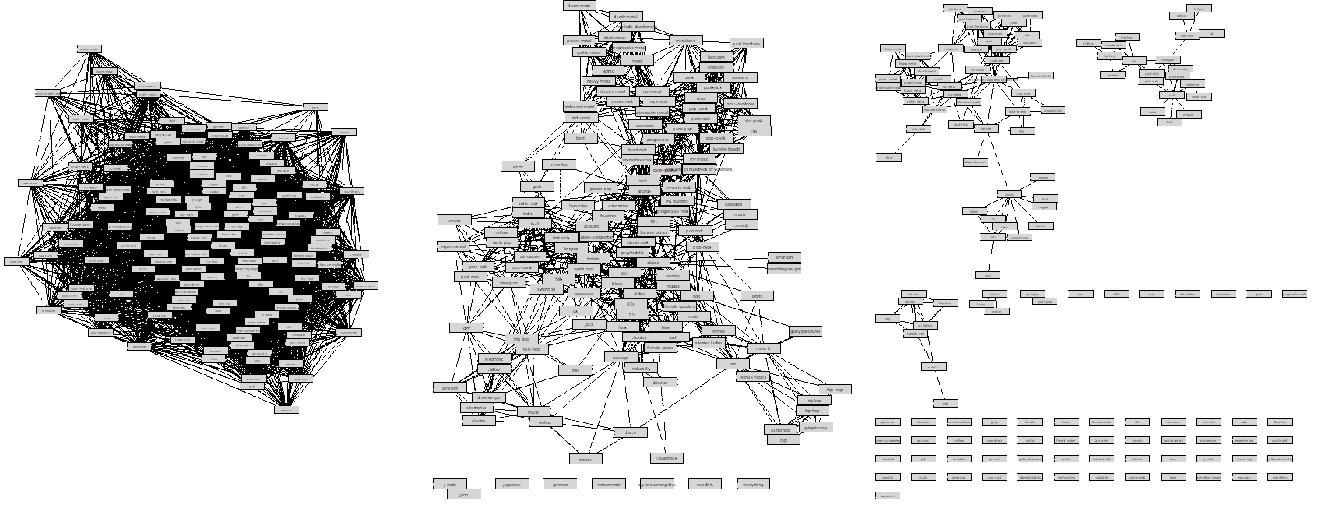


FIG. 4: Graph representation of the music genres filtered correlation matrix $M^{\mu_1 \mu_2}$ for 3 values of the filter parameter $\phi = 0.09$, $\phi = 0.12$ and $\phi = 0.15$, displayed from left to right. Rectangles represent the genres observed in the sample of 995 music groups. The action of filtering leads to a removal of less correlated links, thereby exhibiting the internal structure of the network. The graphs were plotted thanks to the *visone* graphical tools [31].

hibit quantitatively different behaviours. For instance, a music group like *John Coltrane* (Fig.2b) shows a peaked histogram, i.e. it is almost only described by the tag *jazz*, in contrast with ABBA that is described by a large variety of tags. In order to measure the complexity, or diversity of each music group i , we introduce the Shannon entropy [33]:

$$S_i = - \sum_g p_{i,g} \ln p_{i,g} \quad (3)$$

where $p_{i,g}$ is the probability for genre g to be tagged to the music group i , and the sum is performed over all possible genres (with, as said before, a maximum of 25). By construction, this quantity vanishes $S_i^{\min} = 0$ when the group i is wholly described by one tag g^* , i.e. $p_{i,g} = \delta_{gg^*}$ while it takes its maximum value $S_i^{\max} = \ln 25$ for the uniform distribution $p_{i,g} = \frac{1}{25}$. In order to restrain the problem to the interval $[0 : 1]$, we introduce the relative quantity $R_i = \frac{S_i}{\ln 25}$. This quantity is therefore representative of the number of different terms needed by listeners to describe the music group i , i.e. the diversity of the music group. In figure 3, we plot the empirical distribution of this relative entropy over the 995 considered groups. It shows clearly a high degree of diversity of the music groups, therefore requesting many different tags for characterisation.

B. Genres correlations

In this section, we use the methods of section IIB in order to analyse the correlations between genres attributed to each music group i . In the data set, we find 2394 different music genres. Nonetheless, in order to remove

irrelevant tags (due to misprints for instance) and to simplify our analysis, we restrict the scope to all music genres that have been attributed to at least 20 music groups. There are 142 such music genres, that we label with index $\gamma \in [1, 142]$. Let us note G_γ this list of genres, and $P_{i,\gamma}$ their probability for the music group i . For instance, these notations read as follows in the case of *John Coltrane*:

$$\begin{aligned} G &= [\dots, \text{jazz}, \dots, \text{saxophone}, \dots, \text{free jazz}, \dots] \\ P_{J.C.} &= [\dots, 0.72, \dots, 0.06, \dots, 0.02, \dots] \end{aligned} \quad (4)$$

In order to measure correlations between the 142 music genres, we define the 142×142 correlation matrix \mathbf{M} , based on the correlations \mathbf{C} between the music groups (see Eq.2):

$$M^{\gamma_1 \gamma_2} = \frac{\sum_i \sum_{j \neq i} C^{ij} P_{i,\gamma_1} P_{j,\gamma_2}}{N^{\gamma_1 \gamma_2}} \quad (5)$$

where \mathbf{N} is a normalisation matrix:

$$N^{\gamma_1 \gamma_2} = \sum_i \sum_{j \neq i} P_{i,\gamma_1} P_{j,\gamma_2} \quad (6)$$

Practically, we make a loop over the $n_G(n_G - 1)$ pairs of *different* music groups (i, j) , each pair being characterised by the correlation coefficient C^{ij} . For each of these pairs, we evaluate all pairs of music genres γ_1 and γ_2 such that $P_{i,\gamma_1} \neq 0$ and $P_{j,\gamma_2} \neq 0$, and increase the matrix element $M^{\gamma_1 \gamma_2}$ by the quantity $C^{ij} P_{i,\gamma_1} P_{j,\gamma_2}$. The normalisation matrix element $N^{\gamma_1 \gamma_2}$ is itself updated by $P_{i,\gamma_1} P_{j,\gamma_2}$. At the end of the loops, the correlation matrix is normalised: $M^{\gamma_1 \gamma_2} \rightarrow \frac{M^{\gamma_1 \gamma_2}}{N^{\gamma_1 \gamma_2}}$.

In order to reveal collective behaviours from the correlation matrix $M^{\gamma_1 \gamma_2}$, we apply PIB methods. Starting

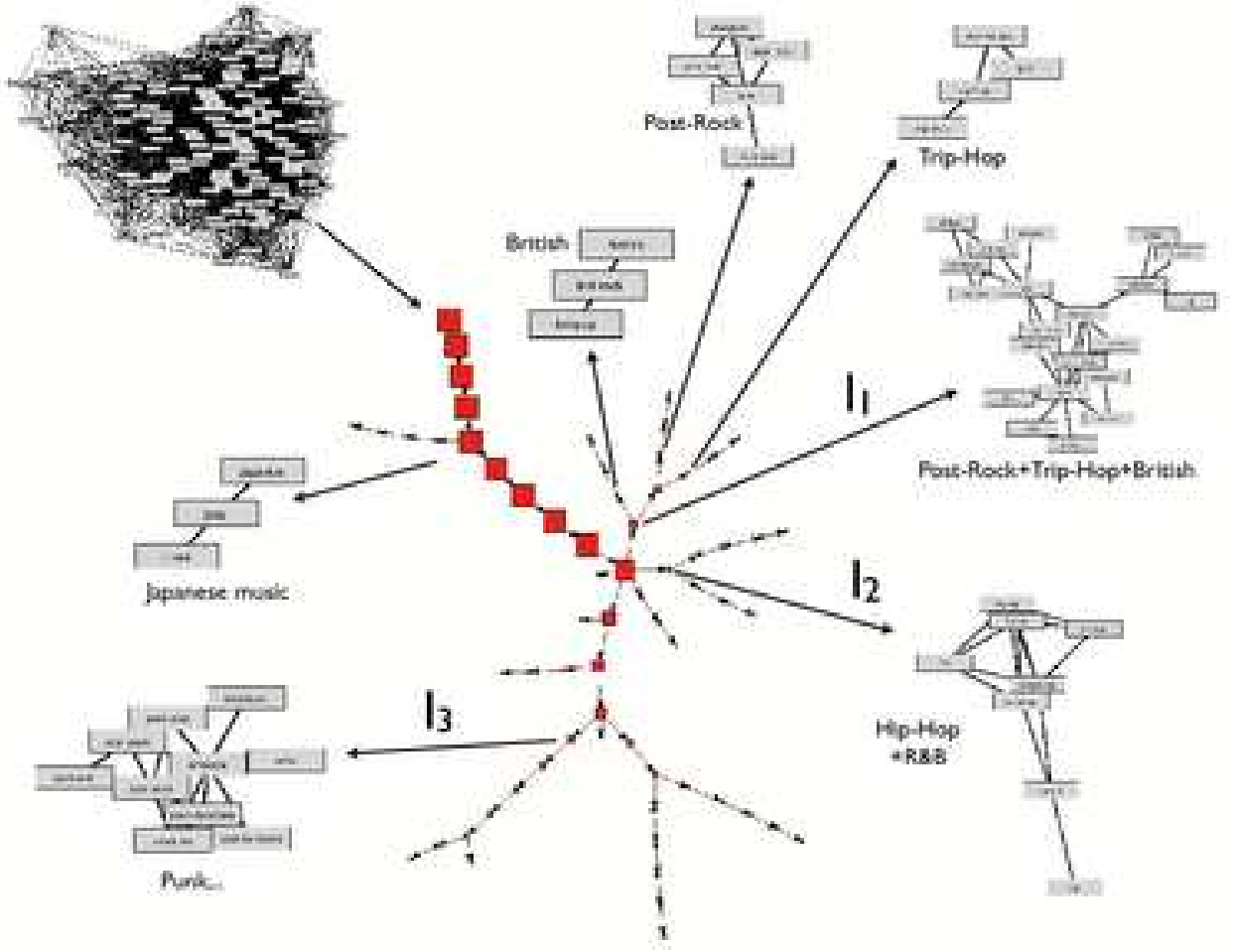


FIG. 5: Branching representation of the correlation matrix $M^{\gamma/2}$. The filtering parameter ϕ ranges from 0.05 to 0.25 (from top to bottom), and is increased at each step by 0.01 (the tree length is 20 steps). It induces a snake of squares at each filtering level. The shape of the snake as well as its direction are irrelevant. The tree obviously shows the emergence of homogeneous branches, that are composed of alike music-subdivisions, thereby showing evidence of genre families. The first island extraction occurs at $\phi = 0.1$, and corresponds to a family of genres related to Japanese music: [japanese, jpop, j-rock]. Among the different structures uncovered by the method, let us note the appearance of the islands I_1 ($\phi = 0.15$), I_2 ($\phi = 0.15$) and I_3 ($\phi = 0.18$) described in the main text.

at a very low value of the filtering coefficient (see Fig.4), say $\phi = 0.09$, the graph is fully connected. Increasing values of the filtering coefficient lead to the formation of cliques and to the emergence of disconnected islands, as those occurring in [4]. Finally, we plot in Fig.5 the tree representation of the filtering process. Poring over the branches of this tree is very instructive and confirms the existence of non-trivial correlations between the different music genres. These correlations shape the relations between genres, and give an objective definition to the notion of sub-genre, genre family....

For instance (see Fig.5), one observes at $\phi = 0.15$ the extraction of two large sub-islands, I_1 and I_2 . I_1 is composed of genres related to Post-Rock, Brit-Rock and Trip-Hop: [chillout, ambient, trip hop, downtempo,

trip-hop, idm, post-rock, post rock, shoegaze, alt-country, post-punk, indie pop, indie rock, lo-fi, emusic, indie, folk, brit rock, british, britpop, uk]. I_2 is itself composed of Hip-Hop and R&B genres: [hip hop, hiphop, hip-hop, gangsta rap, rap, us hiphop, r and b, rmb]. At $\phi = 0.16$, a small sub-island extracts from I_1 , composed of all British related tags, thereby defining a new sub-genre. Finally, at $\phi = 0.17$, I_1 breaks into two separated blocks, one related to Rock music, the other related to Trip-Hop music. Such a breaking also occurs for I_2 at $\phi = 0.16$, and leads to a Hip-Hop sub-genre and a R&B sub-genre. Finally, let us note the punk-related island I_3 emerging at $\phi = 0.18$.

Before concluding, one should insist on the *homogeneity* of the above sub-islands, i.e. their composition is ra-

tional given our a priori knowledge of music. This feature highlights the *homophily* [8] of the music groups, which means that similar groups, i.e. groups with similar tags, tend to be listened by the same audience.

IV. CONCLUSION

In this article, we study empirically the musical behaviours of a large sample of persons. Our analysis is based on web-downloaded data and uses complex network techniques in order to uncover collective trends from the data. To do so, we use percolation idea-based techniques [4] that consist in filtering correlation matrices, i.e. correlations between the music groups, and in visualising the resulting structures by a branching representation. Each of the music groups is characterised by a list of genres, that are tags used by the listeners in order to describe the music group. By studying correlations between these tags, we highlight non-trivial relations between the music genres. As a result, we draw a cartography of music, where large structures are statically uncovered and identified as a genre family. Let us stress that this work is closely related to the theory of hidden variables [34, 35], i.e. the hidden variables being here the music group tags. Consequently, this study should provide an empirical test for the theory.

This work has also many applications in marketing and show business, e.g. taste suggestions in online services, in publicity, libraries.... This kind of approach also opens the way to quantitative modelling of opinion/taste formation [36], and offers quantitative tools for sociologists and musicologists. For instance, G. d’Arcangelo [37] has recently used our analysis in order to discuss the emergence of a growing eclecticism of music listeners that is driven by curiosity and self-identification, in opposition to the uniform trends promoted by commercial radios and *Major* record labels [38]. Applications should also be considered in taxonomy [39], in scientometrics, i.e. how to classify scientific papers depending on their authors, journal, year, keywords..., and in linguistics [40], in order to highlight relations between a signifier (tag) and a signified (music group).

Acknowledgments

Figures 4 and 5 were plotted thanks to the *visone* graphical tools [31]. R.L. would like to thank especially G. D’Arcangelo for fruitful discussions. This work has been supported by European Commission Project CREEN FP6-2003-NEST-Path-012864 and COST P10 (Physics of Risks).

-
- [1] D. Sornette, F. Deschatres, T. Gilbert and Y. Ageon *Phys. Rev. Lett.*, **93** (2004) 228701
 - [2] For instance, see the interview of the Belgian band *dEUS*, on <http://www.metroactive.com/papers/metro/07.17.97/deus-9729.html>, whose diverse influences include *ABBA*, *Sonic Youth*, *Captain Beefheart*, *Franck Zappa*...
 - [3] http://en.wikipedia.org/wiki/List_of_electronic_music_genres
 - [4] R. Lambiotte and M. Ausloos, arXiv physics/0508233
 - [5] J.P. Boon and O. Decroly, *Chaos*, **5** (1995) 501
 - [6] H. D. Jennings, P.C. Ivanov, A. M. Martins, P.C. da Silva and G.M. Viswanathan, *Physica A*, **336** (2004) 585
 - [7] http://en.wikipedia.org/wiki/Image:Genealogy_cuban_music.png
 - [8] N. Masuda and N. Konno, *Physica A*, in press
 - [9] L.A. Adamic and E. Adar, *Social Networks*, **25** (2003) 211
 - [10] M. E. J. Newman, *Proc. Natl. Acad. Sci. USA*, **98** (2001) 404
 - [11] R. Lambiotte and M. Ausloos, physics/0508234
 - [12] L. Adamic and N. Glance, www.blogpulse.com/papers/2005/AdamicGlanceBlogWWW.pdf
 - [13] J.A. Holyst, K. Kacperski and F. Schweitzer, *Annual Review of Comput. Phys.* **9** (2001) 253
 - [14] S. Galam, *Eur. Phys. J. B*, **25** (2002) 403
 - [15] K. Sznajd and J. Sznajd, *Int. J. Mod. Phys. C*, **11** (2000) 1157
 - [16] C. Castellano, M. Marsili and A. Vespignani, *Phys. Rev. Lett.*, **85** (2000) 3536
 - [17] M.N. Kuperman and D.H. Zanette, *Eur. J. Phys. B*, **26** (2002) 387
 - [18] A. Aleksiejuk, J.A. Holyst and D. Stauffer, *Physica A*, **310** (2002) 260
 - [19] S. N. Dorogovtsev, A. V. Goltsev and J. F. F. Mendes, *Phys. Rev. E*, **66** (2002) 016104
 - [20] D. Fenn, O. Suleman, J. Efstathiou and N. F. Johnson, arXiv physics/0505071
 - [21] R. D’Hulst and G.J. Rodgers, *Physica A*, **308** (2002) 443
 - [22] G. Palla, I. Derenyi, I. Farkas and T. Vicsek, *Nature*, **435** (2005) 814
 - [23] E. Ravasz, A.L. Somera, D.A. Mongru, Z. Oltvai and A.L. Barabási, *Science*, **297** (2002) 1551
 - [24] J. Zivković, S. Thurner, N. Wick and B. Tadić, poster presented at the 2005 Next-SigmaPhi conference
 - [25] M. Thelwall and E. Price, *JASIST* (to appear)
 - [26] D. Stauffer and C. Schulze, *Phys. of Life Rev.*, **2** (2005) 89
 - [27] J.-P. Onnela, A. Chakraborti, K. Kaski, J. Kertész and A. Kanto, *Phys. Rev. E*, **68** (2003) 056110
 - [28] D.J. Watts, P.S. Dodds and M.E.J. Newman, *Science*, **296** (2002) 1302
 - [29] http://en.wikipedia.org/wiki/Collaborative_filtering
 - [30] T. S. van Erp, D. Moroni and P. G. Bolhuis, *J. Chem. Phys.*, **118** (2003) 7762
 - [31] <http://www.visone.de/>
 - [32] R. Lambiotte and M. Ausloos, arXiv physics/0507154
 - [33] C. Beck and F. Schlögl, *Thermodynamics of chaotic systems*, Cambridge Univ. Press (1993)
 - [34] A. Fronczak, P. Fronczak and J. A. Holyst *Phys. Rev. E*, **70** (2004) 056110

- [35] M. Boguñá and R. Pastor-Satorras *Phys. Rev. E*, **68** (2003) 036112
- [36] R. Lambiotte and M. Ausloos, in preparation
- [37] G. D'Arcangelo, in *Mobile Music: Ubiquity, Creativity and Community*, MIT Press, to be published (2006)
- [38] L. Margolis, *The Christ. Sc. Monitor* (April 11, 2003)
- [39] <http://en.wikipedia.org/wiki/Taxonomy>
- [40] F. de Saussure, *Cours de linguistique générale*, Ed. Payot (1964)