

A (dis)similarity index for comparing two character networks

François Bavaud

SLI / UNIL

francois.bavaud@unil.ch

Coline Métrailler

SLI / UNIL

coline.metrailler@unil.ch

1 Introduction

Very often, networks of fictional characters exist in two or more versions. For instance (section 4), network A is built from a novel, and network B from a movie adaptation. Besides the main characters common to both versions, there are characters proper to a single version only. Also, the importance of common characters (as e.g. measured by their relative occurrence), is bound to vary between the two versions, as is the strength of their mutual relations (as e.g. measured by their relative co-occurrence). Hence, the networks A and B differ both along character weights (nodes) and interaction weights (edges).

Our contribution proposes the definition of a single index measuring the overall similarity between A and B . In particular, $RV \in [0, 1]$, with $RV = 1$ iff A and B are identical, and $RV = 0$ iff A and B have no character in common. This similarity coefficient can be transformed into a dissimilarity coefficient, which can be additively decomposed into five interpretable components.

2 Formalism (RV index and MDS)

The proposed similarity index is a weighted version of the so-called RV-coefficient (Robert and Escoufier, 1976), initially defined as the *cosine similarity* between the covariances of two sets of multivariate features \mathbf{X}_A and \mathbf{X}_B .

In the dual, individual-centered perspective relevant for tackling the comparison of character networks, the proposed similarity index reads

$$RV = \frac{c_{AB}}{\sqrt{c_{AA} c_{BB}}} \quad c_{AB} = \text{trace}(\mathbf{K}_A \mathbf{K}_B) \quad (1)$$

where \mathbf{K}_A is the matrix of weighted scalar products between individuals or *kernel*, computed as

$$\mathbf{K}_A = \text{diag}(\sqrt{\mathbf{f}_A}) \mathbf{B}_A \text{diag}(\sqrt{\mathbf{f}_A}) \quad (2)$$

$$\mathbf{B}_A = -\frac{1}{2} \mathbf{H}_A \mathbf{D}_A \mathbf{H}_A^\top \quad (3)$$

where \mathbf{f}_A is the vector of individual relative weights, normalized to one, \mathbf{D}_A is the matrix of squared Euclidean dissimilarities between individuals (as computed from their features \mathbf{X}_A), $\mathbf{H}_A = \mathbf{I} - \mathbf{1} \mathbf{f}_A^\top$ is the centering matrix, and \mathbf{B}_A is the usual matrix of scalar products. Spectral decomposition of kernels \mathbf{K}_A (and \mathbf{K}_B) enables the *weighted classical multidimensional scaling* (MDS), permitting to extract individual coordinates reproducing the dissimilarities \mathbf{D}_A , and whose low-rank restriction expresses a maximum amount of the configuration dispersion or *inertia* $\Delta_A = \frac{1}{2} \mathbf{f}_A^\top \mathbf{D}_A \mathbf{f}_A$.

3 The case of character networks

Adapting the above formalism to the comparison of character networks raises interesting formal challenges:

- (1) character networks have to be converted into weighted configurations (\mathbf{f}, \mathbf{D}) , where \mathbf{D} is a squared Euclidean dissimilarity between characters; many available competing definitions and corresponding algorithms (commuting times, heat kernels, Sinkhorn, Metropolis) permit to compute the latter.
- (2) crucially, relative weights \mathbf{f}_A and \mathbf{f}_B (set to the uniform weights $1/n$ in most multivariate proposals) may *differ* to a spectacular extent: their supports $\text{supp}(\mathbf{f}_A)$ and $\text{supp}(\mathbf{f}_B)$ do not even coincide in general, since version A may contain characters absent in version B , and vice-versa.

Define the *compromise weight* \mathbf{h} as

$$h_i = \frac{\sqrt{f_i^A f_i^B}}{Z} \quad \text{where } Z = \sum_j \sqrt{f_j^A f_j^B} \in [0, 1] \quad (4)$$

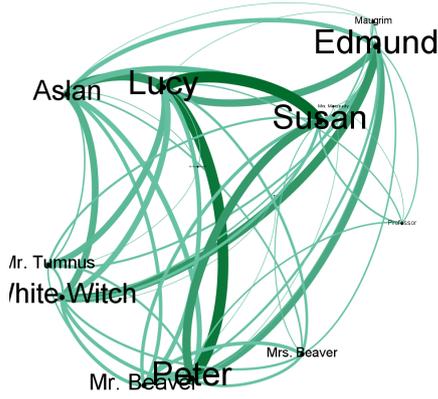


Figure 1: character network of the book *The Lion, the Witch and the Wardrobe*.

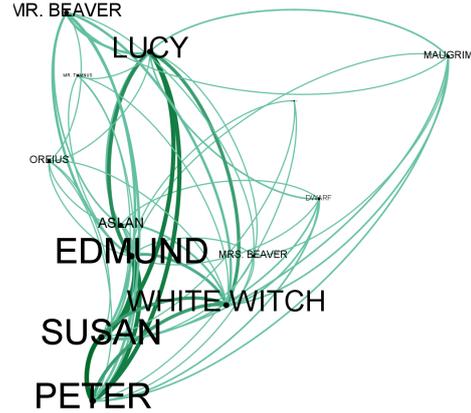


Figure 2: character network of the movie *The Chronicles of Narnia*.

By construction, $h_i = 0$ unless character i appears in both versions. A little algebra demonstrates the numerator of the similarity index (1) to express as

$$\frac{c_{AB}}{Z^2} = \text{trace}(\mathbf{K}_{h_A} \mathbf{K}_{h_B}) + \kappa_{AB} \quad (5)$$

$$\text{where } \kappa_{AB} = D_{h_{f_A}}^A D_{h_{f_B}}^B + 2 \sum_i h_i \ell_i^A \ell_i^B .$$

Here \mathbf{K}_{h_A} is the kernel associated to configuration $(\mathbf{h}, \mathbf{D}_A)$ and $D_{h_{f_A}}^A$ is the squared Euclidean distance between the gravity centers of $(\mathbf{h}, \mathbf{D}_A)$ and $(\mathbf{f}_A, \mathbf{D}_A)$ (terms involving B are defined analogously). Term κ_{AB} , which can be negative, represents a correction due to the non-coincidence of the \mathbf{f}_A - and \mathbf{h} -centroids in configuration A (respectively the \mathbf{f}_B - and \mathbf{h} -centroids in configuration B); its second component involves a weighted covariance between the \mathbf{h} -centered vectors

$$\ell_A = \mathbf{B}_{h_A} \mathbf{f}_A \quad \text{and} \quad \ell_B = \mathbf{B}_{h_B} \mathbf{f}_B$$

4 Illustration: The Chronicles of Narnia

The Lion, the Witch and the Wardrobe was the second of the seven novels of the *The Chronicles of Narnia*, written by C. S. Lewis in 1950, and adapted in a film directed by A. Adamson released in 2005. Among a total of 37 characters identified in the pre-treatments (manual annotation of all named entities throughout the book and the movie script, then gathered into groups of aliases to create a list of distinct characters), 16 are common to both the book and the movie (i.e. they constitute the support of \mathbf{h} in (4), with $Z = 0.890$), 8 occur in the book only, and 13 in the movie only.

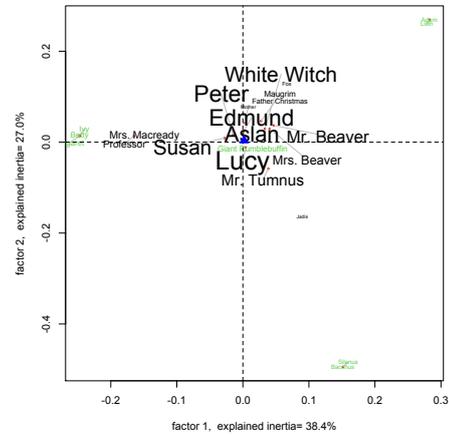


Figure 3: first MDS coordinates (x_i^A, y_i^A) of the 24 book characters (the 8 green ones appear in the book only), centered on $(\bar{x}_{f_A}^A, \bar{y}_{f_A}^A) = (0, 0)$. The blue dot depicts the centroid $(\bar{x}_h^A, \bar{y}_h^A) = (0.002, 0.004)$, fairly close to the origin: $D_{h_{f_A}}^A = 5.6 \cdot 10^{-5}$.

The symmetric *book cross-count matrix* $\mathbf{N}^A = (n_{ij}^A)$ counts the number of co-occurrences of characters i and j within a window of 5 paragraphs (each paragraph being delimited by a line break), with $n_{ii}^A = 0$. The associated character networks obtain by defining n_{ij}^A as the absolute weight of the unoriented edge ij , and $n_{i\bullet}^A = \sum_j n_{ij}^A$ as the absolute weight of node i . The relative book weights \mathbf{f}_A of the previous section are defined as $f_i^A = n_{i\bullet}^A / n_{\bullet\bullet}^A$. The same treatment yields the movie graph from the movie cross-count matrix \mathbf{N}^B .

The cross-counts matrices permit to define a so-called *instantaneous jump process*, permitting to navigate from node i to node j with a rate given by (minus) the components of the *normalized Lapla-*

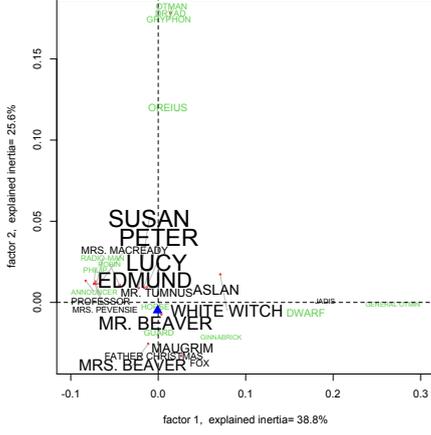


Figure 4: first MDS coordinates (x_i^B, y_i^B) of the 29 movie characters (the 13 green ones appear in the movie only), centered on $(\bar{x}_{f_B}^B, \bar{y}_{f_B}^B) = (0, 0)$. The blue dot depicts the centroid $(\bar{x}_h^B, \bar{y}_h^B) = (-0.0007, -0.005)$, fairly close to the origin again: $D_{h,f_B}^B = 4.0 \cdot 10^{-5}$.

can $L_{ij}^A = \delta_{ij} - n_{ij}^A / \sqrt{n_i^A n_j^A}$. They define in turn a positive semi-definite *weighted diffusion* or *heat kernel* by $\mathbf{K}_A(t) = \exp(-t\mathbf{L}_A) - \sqrt{\mathbf{f}_A} \sqrt{\mathbf{f}_A}^\top$, obeying $\mathbf{K}_A(t) \sqrt{\mathbf{f}_A} = 0$ as it must (see e.g. Bavaud, 2013). The free parameter $t > 0$ is the diffusion time, and spectral decomposition of $\mathbf{K}_A(t)$ (for $t = 10$) yields the MDS representation of the book characters in figure 3. Figure 4 depicts the analogous MDS representation of the movie characters constructed from \mathbf{N}^B .

At this stage, the book and the movie network of characters have been expressed into kernels, namely \mathbf{K}_A and \mathbf{K}_B ; alternative kernel constructions are possible, their discussion being beyond the scope of this paper. Applying the above formalism yields a coefficient of similarity (1) of $RV = 0.113$. Since \mathbf{h} centroids are close to the \mathbf{f}_A -, respectively \mathbf{f}_B -centroids, the relative importance of both terms in the r.h.s. of (5) is very contrasted, namely 100.11% and -0.11% (κ_{AB} being negative).

5 An exact decomposition formula

A similarity coefficient such as $RV \in [0, 1]$ can be converted into a *dissimilarity* coefficient $d \in [0, \infty)$ by $d = -\ln RV$. Applying the transformation to (1) yields the following exact decomposition for the dissimilarity between character net-

works A and B :

$$d_{AB} = -\ln RV = -\ln RV_h - 2 \ln Z - \frac{1}{2} \ln \Gamma_A - \frac{1}{2} \ln \Gamma_B - \ln(1 + \epsilon) \quad (6)$$

where

- $RV_h = \frac{\text{trace}(\mathbf{K}_{hA} \mathbf{K}_{hB})}{\sqrt{\text{trace}(\mathbf{K}_{hA}^2) \text{trace}(\mathbf{K}_{hB}^2)}} \in [0, 1]$ measures the similarity between dissimilarities \mathbf{D}_A and \mathbf{D}_B in the common compromise weighting \mathbf{h} .
- $Z \in [0, 1]$ in (4) is a measure of similarity between weights \mathbf{f}_A and \mathbf{f}_B .
- $\Gamma_A = \frac{\text{trace}(\mathbf{K}_{hA}^2)}{\text{trace}(\mathbf{K}_A^2)}$ is a measure of the ratio of the (quartic) dispersion of configuration \mathbf{D}_A in the compromise weighting \mathbf{h} to the dispersion of \mathbf{D}_A in the original weighting \mathbf{f}_A (Γ_B is defined analogously). One expects $\Gamma_A \in [0, 1]$ due to the loss of diversity entailed by the disappearance of A -specific characters in the compromise weighting.
- $\epsilon = \frac{\kappa_{AB}}{RV_h \sqrt{\Gamma_A \Gamma_B} \text{trace}(\mathbf{K}_A^2) \text{trace}(\mathbf{K}_B^2)}$ is a normalized measure of the centroid correction occurring in (5), whose magnitude is expected to be small since main characters are likely to occur in both versions A and B .

In the present study, (6) yields (in order)

$$2.1829 = 0.9385 + 0.2323 + 0.4349 + 0.5762 + 0.0011 \quad .$$

References

- François Bavaud. 2013. Testing spatial autocorrelation in weighted networks: the modes permutation test. *Journal of Geographical Systems*, 15(3):233–248.
- Paul Robert and Yves Escoufier. 1976. A unifying tool for linear multivariate statistical methods: the RV-coefficient. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 25(3):257–265.