# Fusion and Inference from Multiple Data Sources in a Commensurate Space

Z. Ma, D.J. Marchette, C.E. Priebe

**Abstract**

Given objects measured under multiple conditions—e.g., indoor lighting versus outdoor lighting for face recognition, multiple language translation for document matching, etc.—the challenging task is to perform data fusion and utilize all the available information for inferential purposes. We consider two exploitation tasks: (1) how to determine whether a set of feature vectors represent a single object measured under different conditions; and (2) how to create a classifier based on training data collected under one condition in order to classify objects measured in other conditions. The key to both problems is to transform all sets of feature vectors into one commensurate space, where the (transformed) feature vectors are comparable and would be treated as if they were collected under the same condition. Toward this end, we studied Procrustes analysis and developed a new approach, which uses the interpoint dissimilarities for each condition. We impute the dissimilarities between measurements of different conditions to create one omnibus dissimilarity matrix, which is then embedded into Euclidean space. We illustrate our methodology on English and French documents collected from Wikipedia, demonstrating superior performance compared to that obtained via standard Procrustes transformation.

**Key words:** fusion, inference, dissimilarity, multidimensional scaling, Procrustes transformation, embedding

## 1 Introduction

Information fusion techniques aim to merge information from multiple data sources in order to achieve more accurate inferences than using each single source alone. Information fusion has been a hot research field with various applications (Chengjun Liu, 2001; Ross and Jain, 2004; Sun et al., 2005; Kludas et al., 2008).

In general, the most often used information fusion approaches can be summarized into two categories: feature level fusion and decision level fusion. In feature level fusion, feature vectors extracted from different data sources are combined into the Cartesian product space, directly (Ma, Cardinal-Stakenas, Park, Trosset, and Priebe, 2010) or via some data transformation procedures (Sun et al., 2005). Decision level fusion involves combining results obtained separately from all data sources. An ensemble of classifiers is one such example, as is track fusion (Chang et al., 2002). The advantage

1

of these two types of information fusion stems from the fact that multiple sets of feature vectors extracted from the same set of objects usually reflect different characteristics of patterns. By fusing multiple disparate data sources, one generates a more complete representation of the space in which the objects live, and hence has more power for inferential tasks such as hypothesis testing, classification, etc.

In this work, we consider information fusion from a different perspective. Suppose that objects are measured under multiple conditions—e.g. indoor lighting versus outdoor lighting for face recognition, multiple language translation for document matching, etc. The challenging questions are: (1) how to determine whether a set of feature vectors represent a single object measured under different conditions? For example, whether pictures taken under different lighting conditions are the photos of the same individual or not; and (2) how to create a classifier based on training data measured under one condition, and use it to classify objects measured in other conditions? We refer the two problems as the implicit translation problem and the classification problem, respectively. A direct approach would involve finding the underlying mappings among all the spaces of measurements and transform all these measurements into one commensurate space through the derived mappings. In this commensurate space, all transformed feature vectors are treated equally as if they were from the same data source. The solutions to both questions will then be straightforward. In real applications, finding the mappings among all spaces of measurements is usually difficult. In fact, it is possible to fuse multiple spaces into one commensurate space without using the mappings among these spaces. (Generalized) Procrustes analysis is one potential solution. Consider a set of objects, each of which is measured under $K$ ($K \geq 2$) conditions, yielding $K$ feature vectors. Assuming all the feature vectors have been column centered, Procrustes solution rotates (possibly with dilation) the feature vectors to best match each other, and thereby defines a commensurate space.

The raw data in text or image analysis are usually high-dimensional. Dissimilarity analysis is one of the commonly applied approaches for finding low-dimensional representation of such data. Usually in dissimilarity analysis, one first calculates interpoint dissimilarities to obtain a dissimilarity matrix, and then embeds it into a low dimensional space via multidimensional scaling. We use a collection of Wikipedia documents to illustrate the two problems (implicit translation and classification) and the solutions. The two step approach, which we refer to as the P-approach, first embeds dissimilarity matrices derived from different data sources and then utilizes a Procrustes transformation on the embeddings to make them commensurate. We propose an approach that simultaneously embeds all dissimilarity matrices and finds the commensurate space. In this approach, dissimilarity matrices from different data sources are put onto the diagonal of an omnibus matrix, whose off-diagonal entries are imputed. Embedding this omnibus matrix results in feature vectors in one commensurate space. We refer this approach as the W-approach. Both approaches are studied in this work, and the results on Wikipedia example show that the W-approach leads to larger powers in testing and higher accuracy in classification, compared to the P-approach.

In Section 2, we describe the Wikipedia data set, the derivation of dissimilarity matrices, and the implicit transformation and classification problems. Section 3 details the traditional Procrustes

solution and the proposed W-approach. The results are given in Section 4. Section 5 provides conclusions.

## 2 Data

Wikipedia is an open-source Encyclopedia that is written by a large community of users (everyone who wants to, basically). There are versions in over 200 languages, with various amounts of content. The full data for the Wikipediae are freely available for download. A Wikipedia document has one or more of: title, unique ID number, text—the content of the document, images, internal links—links to other documents, external links—links to other content elsewhere on the web, and language links—links to "the same" document in other languages. Figure 1 shows an English Wikipedia document titled "Geometry". The multilingual Wikipediae provide a good testbed for developing methods for analysis of text, implicit translation, and fusion of text and graph information.
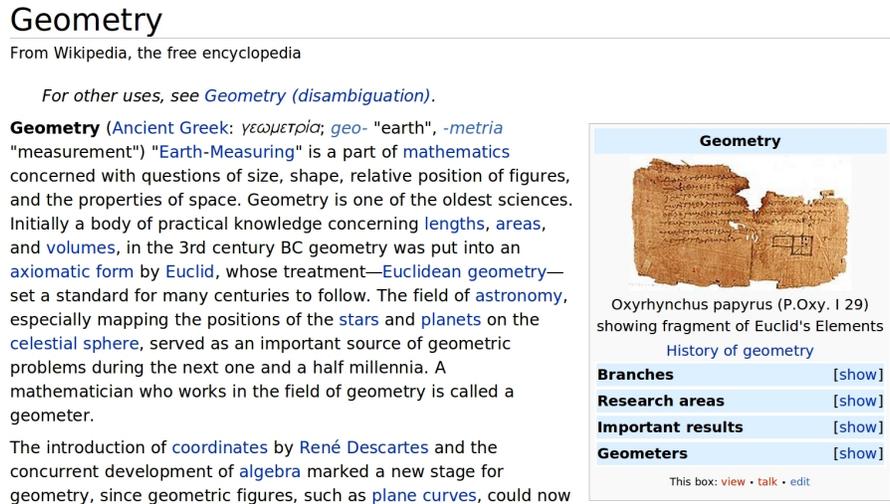
### Geometry
From Wikipedia, the free encyclopedia

*For other uses, see Geometry (disambiguation).*

**Geometry** (Ancient Greek: γεωμετρία; geo- "earth", -metria "measurement") "Earth-Measuring" is a part of mathematics concerned with questions of size, shape, relative position of figures, and the properties of space. Geometry is one of the oldest sciences. Initially a body of practical knowledge concerning lengths, areas, and volumes, in the 3rd century BC geometry was put into an axiomatic form by Euclid, whose treatment—Euclidean geometry—set a standard for many centuries to follow. The field of astronomy, especially mapping the positions of the stars and planets on the celestial sphere, served as an important source of geometric problems during the next one and a half millennia. A mathematician who works in the field of geometry is called a geometer.

The introduction of coordinates by René Descartes and the concurrent development of algebra marked a new stage for geometry, since geometric figures, such as plane curves, could now

**Geometry**

Oxyrhynchus papyrus (P.Oxy. I 29) showing fragment of Euclid's Elements
History of geometry

| Branches | [show] |
| Research areas | [show] |
| Important results | [show] |
| Geometers | [show] |

This box: view · talk · edit

**Figure 1:** "Geometry", an example of English Wikipedia documents. In general, a Wikipedia document has one or more of: title, unique ID number, text, images, internal links, external links, and language links.

### 2.1 Dissimilarities from Graph Structure and Textual Content

Let $G = (V, E)$ be a (directed) graph, where $V$ is the set of nodes—Wikipedia documents, and $E$ is the set of edges—the links within the documents. We consider two Wikipediae, English and French. A subset of the English and French Wikipediae is extracted such that there is an 1-1 correspondence between English documents and French documents. From the English subset, we take the (directed) 2-neighborhood of the document "Algebraic Geometry", yielding set $\boldsymbol{E} = \{\boldsymbol{x}_{1,0}, \ldots, \boldsymbol{x}_{n,0}\}$ ($n = 1382$). The associated documents in French constitute set $\boldsymbol{F} = \{\boldsymbol{x}_{1,1}, \ldots, \boldsymbol{x}_{n,1}\}$. Thus, the English graph with nodes in $\boldsymbol{E}$ is connected by construction, but the French graph with nodes in $\boldsymbol{F}$ need not be connected (and in fact it is not). We consider two types of data, both of which are given in the form of dissimilarity matrices denoted generically as $\mathbf{D}_0$ and $\mathbf{D}_1$: (1) dissimilarity matrices $\mathbf{G}_0$ and $\mathbf{G}_1$, developed from the graph structures of $\boldsymbol{E}$ and $\boldsymbol{F}$ respectively; (2) dissimilarity matrices

$\mathbf{T}_0$ and $\mathbf{T}_1$, obtained from the textual contents of $\boldsymbol{E}$ and $\boldsymbol{F}$ respectively.

To get dissimilarity matrices from graph structure, the adjacency matrices $\mathbf{A}_0$ and $\mathbf{A}_1$ are first created from $\boldsymbol{E}$ and $\boldsymbol{F}$. An adjacency matrix is a square binary matrix, with 1 in position $(i, j)$ only when the $i$th document contains an link to the $j$th document. Dissimilarity matrices $\mathbf{G}_0$ and $\mathbf{G}_1$ are then derived from $\mathbf{A}_0$ and $\mathbf{A}_1$, with $(i, j)$ entry as the number of steps it takes to reach node $j$ from node $i$. By the nature of the graphs, the elements of $\mathbf{G}_0$ take values in $\{0, \ldots, 4\}$, while the elements of $\mathbf{G}_1$ take values in $\{0, \ldots, 1384\}$, with 1384 meaning no directed path between two nodes. Because it is computationally expensive to develop $\mathbf{G}_1$, in practice we assign the value 6 to $\mathbf{G}_1(i, j)$ if it takes more than 4 steps to reach node $j$ from node $i$. Finally, $\mathbf{G}_0$ and $\mathbf{G}_1$ are symmetrized by averaging the corresponding lower- and upper-triangle entries, respectively.

For dissimilarity matrices of textual content, we use Lin & Pantel's approach (Lin and Pantel, 2002; Pantel and Lin, 2002) on Wikipedia documents $\boldsymbol{E}$ and $\boldsymbol{F}$ to obtain two mutual information feature matrices. Rare-word discounting (Lin and Pantel, 2002) is then applied to reduce the impact of infrequent words. Given feature vectors of two documents $\boldsymbol{a}, \boldsymbol{b} \in \mathbb{R}^d$, a dissimilarity function $\rho$ is defined as $\rho(\boldsymbol{a}, \boldsymbol{b}) = 1 - (\boldsymbol{a} \cdot \boldsymbol{b})/(\|\boldsymbol{a}\|_2 \|\boldsymbol{b}\|_2)$. Employing $\rho$ on the two feature matrices of $\boldsymbol{E}$ and $\boldsymbol{F}$ separately results in two dissimilarity matrices $\mathbf{T}_0$ and $\mathbf{T}_1$.

When a new English document $\boldsymbol{y}_0$ and a new French $\boldsymbol{y}_1$ are provided, we have access to the dissimilarities (for both graph structure and textual content) between $\boldsymbol{y}_0$ and $\boldsymbol{x}_{i,0}$, and those between $\boldsymbol{y}_1$ and $\boldsymbol{x}_{i,1}$, $i = 1, \ldots, n$. Therefore the Wikipeida data set contains four dissimilarity matrices $\mathbf{G}_0, \mathbf{G}_1, \mathbf{T}_0$ and $\mathbf{T}_1$, and each new document $\boldsymbol{y}_k$ will be represented by a dissimilarity vector $\{\delta(\boldsymbol{y}_k, \boldsymbol{x}_{i,k})\}_{i=1}^{n}$, $k = 0, 1$. ($\delta$ is a dissimilarity function.)

## 2.2  Implicit Translation and Classification

An implicit translation of a document, unlike a word-level or a real translation in any normal sense, is an association with another document in a different language that is on the same topic. In our framework, we treat each topic as an object with measurements (documents) taken under different conditions (languages). That is, topics are represented by documents of different languages. Consider the two collections of matched Wikipedia documents $\boldsymbol{E} = \{\boldsymbol{x}_{1,0}, \ldots, \boldsymbol{x}_{n,0}\}$ and $\boldsymbol{F} = \{\boldsymbol{x}_{1,1}, \ldots, \boldsymbol{x}_{n,1}\}$. Let $\boldsymbol{x}_{i,0} \sim \boldsymbol{x}_{i,1}$ denote that the English document $\boldsymbol{x}_{i,0}$ and the French document $\boldsymbol{x}_{i,1}$ are matched—they are the measurements (under $K = 2$ conditions) of the same topic. The goal of implicit translation is to determine whether a match is present between two new documents $\boldsymbol{y}_0$ and $\boldsymbol{y}_1$. That is, we consider the hypothesis testing:

$$H_0 : \boldsymbol{y}_0 \sim \boldsymbol{y}_1 \quad \text{versus} \quad H_A : \boldsymbol{y}_0 \nsim \boldsymbol{y}_1$$

Notice that we assume the two new documents represent a matched pair under $H_0$. This allows us to control the probability of missing a true match. This is practical in many applications where computer algorithms are used to eliminate easily rejected pairs and the remaining possibly matched pairs will be manually examined.

The second problem is to classify French documents by a classifier trained on English documents. Formally, consider two manifolds, $\Xi_0$ and $\Xi_1$. Let

$$(X, C, Z) \sim F_{X,C,Z},$$
$$C : \Omega \to J \cup \tilde{J},$$
$$Z : \Omega \to \{0, 1\},$$
$$X|Z = z : \Omega \to \Xi_z,$$

where $J$ and $\tilde{J}$ are two disjoint sets of class labels. Suppose the following training data are available

$$\mathcal{T}_0 = \{(x_i, c_i \in J, z_i = 0), \ i = 1, \ldots, n_0\},$$
$$\mathcal{T}_1 = \{(x_i, c_i \in J, z_i = 1), \ i = 1, \ldots, n_1\},$$
$$\tilde{\mathcal{T}}_0 = \{(x_i, c_i \in \tilde{J}, z_i = 0), \ i = 1, \ldots, m_0\}.$$

That is, there are training data from all the classes $J \cup \tilde{J}$ in space $\Xi_0$, but in space $\Xi_1$ only training data from classes $J$ are available. We are interested in creating a classifer $g$ based on the training data and use it to classify future observations in $\Xi_1$ into one of the classes in $\tilde{J}$. Figure 2 depicts the classification problem.
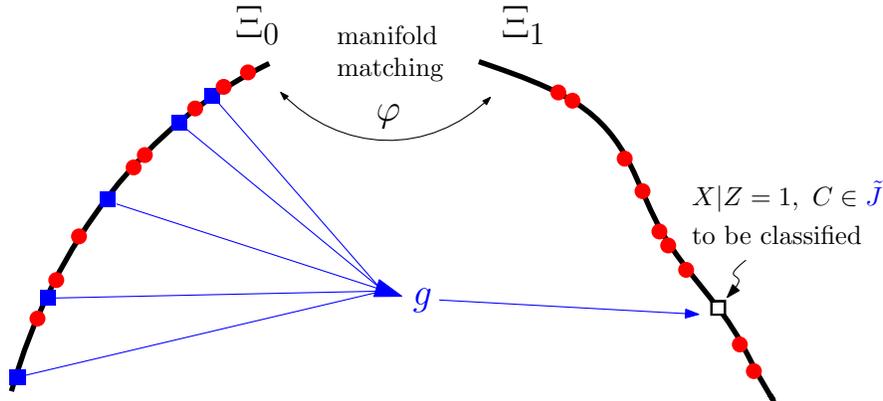


**Figure 2:** Classification problem. In space $\Xi_0$ training data from classes $J$ (red) and $\tilde{J}$ (blue) are available, while in space $\Xi_1$ only training data from classes $J$ are available. We are interested in training a rule $g$ to classify objects of classes $\tilde{J}$ in space $\Xi_1$. It is impossible to directly create such a classifier in $\Xi_1$ due to lack of training data.

We consider $\Xi_0$ and $\Xi_1$ to be the English and French Wikipedia document space, respectively. The 1382 Wikipedia documents are labeled into 5 groups. The two disjoint sets of class labels are $J = \{0, 1, 2\}$ and $\tilde{J} = \{3, 4\}$. We are interested in finding a way to create a classifier based on English documents and use it to classify French documents.

## 3   Methods

For the implicit translation problem, suppose that there is a way to embed $\boldsymbol{E} \in \Xi_0$ and $\boldsymbol{F} \in \Xi_1$ into a commensurate space $\Xi_c$, where the embeddings of English and French documents would be

treated equally, as if they were collected under the same condition. We can embed the two new documents $\boldsymbol{y}_0$ and $\boldsymbol{y}_1$, referred to as the out-of-sample documents, into the space $\Xi_c$. Whether a match is present is then determined by examining the distance between the embeddings of $\boldsymbol{y}_0$ and $\boldsymbol{y}_1$, with a large distance being evidence against $H_0$. There are two ways to determine critical values. The naïve way is to treat the distances between the embeddings of matched pairs in $\mathbf{E}$ and $\mathbf{F}$ as the ground truth, and use the $100(1-\alpha)$th percentile as the critical value for a level $\alpha$ test. However, this method dose not always lead to large powers, because the distribution of the distances between out-of-sample embeddings is usually slightly different from that of the original embeddings, even under the matched assumption $H_0$. Another way of obtaining critical values is by means of Monte Carlo simulation: (i) randomly choose a pair of matched documents $\boldsymbol{x}_{i,0}$ and $\boldsymbol{x}_{i,1}$ from $\boldsymbol{E}$ and $\boldsymbol{F}$, and treat them as out-of-sample documents; (ii) embed the selected documents into the space $\Xi_c$, and compute their distance; and (iii) repeat (i—ii) to obtain an empirical distribution of such distances. The critical value for a level $\alpha$ test is then calculated as the $100(1-\alpha)$th percentile of this empirical distribution. We use the latter method in this work and get larger powers than using the naïve approach.

For the classification problem, suppose a commensurate space $\Xi_c$ could be obtained through $\mathcal{T}_0$ and $\mathcal{T}_1$—English and French documents of classes $J$. We can embed the training English documents $\tilde{\mathcal{T}}_0$ and the new French documents into the space $\Xi_c$. In the commensurate space $\Xi_c$, building classifier $g$ based on English documents with labels in $\tilde{J}$ and using it to classify new French documents are then straightforward.

Therefore the key to both problems is: how shall we determine the commensurate space $\Xi_c$ and how shall we embed new documents into this space?

## 3.1   Procrustes Transformation

The Procrustes analysis (Sibson, 1978, and references contained therein) is to transform a configuration of points (source) to another (target) as closely as possible in the least-square sense. The permitted transformations are any combination of dilation (uniform scaling), rotation, reflection, and translation. We define the space where the target and the transformed source live as the commensurate space.

For the implicit translation problem, we embed $\mathbf{D}_0$ and $\mathbf{D}_1$ through multidimensional scaling to obtain $n \times d$ configurations $\mathbf{X}_0$ and $\mathbf{X}_1$ in the space $\mathbb{R}^d$ separately. The two new documents $\boldsymbol{y}_0$ and $\boldsymbol{y}_1$ are then embedded to $\tilde{\boldsymbol{y}}_0$ and $\tilde{\boldsymbol{y}}_1$ in $\mathbb{R}^d$ respectively via out-of-sample embedding (Trosset and Priebe, 2008). Notice that the coordinates in $\mathbf{X}_0$ and $\mathbf{X}_1$ may be given in different systems. Procrustes analysis is performed to transform one of the embeddings—e.g., $\mathbf{X}_1$—to best match the other one—e.g., $\mathbf{X}_0$. The resulting transformation function $t$ is then applied to the corresponding out-of-sample embedding $\tilde{\boldsymbol{y}}_1$ so that $t(\tilde{\boldsymbol{y}}_1)$ and $\tilde{\boldsymbol{y}}_0$ are commensurate.

For the classification problem, a similar procedure is performed. Let $\mathbf{D}_0^J$ and $\mathbf{D}_1^J$ denote the dissimilarity matrices among documents in $\mathcal{T}_0$ and $\mathcal{T}_1$. We embed $\mathbf{D}_0^J$ and $\mathbf{D}_1^J$ to $\mathbf{X}_0^J$ and $\mathbf{X}_1^J$ in $\mathbb{R}^d$

respectively. Then the English documents in $\tilde{\mathcal{T}}_0$ and the new French documents, whose class labels belong to $\tilde{J}$, are embedded to $\boldsymbol{X}_0^{\tilde{J}}$ and $\boldsymbol{X}_1^{\tilde{J}}$ in $\mathbb{R}^d$ respectively via out-of-sample embedding. Procrustes transformation function $t_J$ learned from $\mathbf{X}_0^J$ and $\mathbf{X}_1^J$ is then applied to $\boldsymbol{X}_1^{\tilde{J}}$ so that $t_J(\boldsymbol{X}_1^{\tilde{J}})$ and $\boldsymbol{X}_0^{\tilde{J}}$ are commensurate.

We refer this approach as the P-approach.

## 3.2 Our Approach

The P-approach creates commensurate space in two steps, namely embedding and Procrustes transformation. We develop a novel method, which defines commensurate space in one step. In implicit translation, we impute $\mathbf{W}$, the dissimilarities between $\boldsymbol{E}$ and $\boldsymbol{F}$, by the entrywise average of $\mathbf{D}_0$ and $\mathbf{D}_1$. An omnibus dissimilarity matrix $\mathbf{M}$ is then constructed by putting $\mathbf{D}_0$ and $\mathbf{D}_1$ on the diagonal, and putting $\mathbf{W}$ on the off-diagonal. We embed $\mathbf{M}$ to obtain a configuration of $2n$ points $\mathbf{X}$ in $\mathbb{R}^d$. We take the first $n$ points and the remaining $n$ points as embeddings of $\mathbf{D}_0$ and $\mathbf{D}_1$, respectively. Notice that $\mathbf{X}_0$ and $\mathbf{X}_1$ are already in the same space $\Xi_c$, because the distances between matched English and French document pairs have been taken into account when embedding $\mathbf{M}$—the imputed matrix $\mathbf{W}$ has all zeros on its diagonal. For any two additional documents $\boldsymbol{y}_0$ and $\boldsymbol{y}_1$, let $\boldsymbol{u}_0$ and $\boldsymbol{v}_1$ denote the dissimilarity vector between $\boldsymbol{y}_0$ and $\boldsymbol{E}$, $\boldsymbol{y}_1$ and $\boldsymbol{F}$, respectively. Under the null hypothesis that $\boldsymbol{y}_0$ and $\boldsymbol{y}_1$ are matched, we impute the dissimilarities between $\boldsymbol{y}_0$ and $\boldsymbol{F}$ (denoted by $\boldsymbol{v}_0$), and dissimilarities between $\boldsymbol{y}_1$ and $\boldsymbol{E}$ (denoted by $\boldsymbol{u}_1$) by entrywise average of $\boldsymbol{u}_0$ and $\boldsymbol{v}_1$. That is, $\boldsymbol{v}_0 = \boldsymbol{u}_1 = (\boldsymbol{u}_0 + \boldsymbol{v}_1)/2$. Out-of-sample embedding is used to embed $(\boldsymbol{u}_0^t, \boldsymbol{v}_0^t)^t$ and $(\boldsymbol{u}_1^t, \boldsymbol{v}_1^t)^t$ into $\Xi_c$. Figure 3 depicts the construction of the omnibus dissimilarity matrix $\mathbf{M}$.

$$
\overset{2n \times 2n}{\mathbf{M}} = \begin{bmatrix} \overset{n \times n}{\mathbf{D}_0} & \overset{n \times n}{\mathbf{W}} \\ \mathbf{W}^T & \overset{n \times n}{\mathbf{D}_1} \end{bmatrix} \begin{matrix} \overset{n \times 1}{\boldsymbol{u}_0} \overset{n \times 1}{\boldsymbol{u}_1} \\ \overset{n \times 1}{\boldsymbol{v}_0} \overset{n \times 1}{\boldsymbol{v}_1} \end{matrix}
$$
$$
\begin{matrix} \boldsymbol{y}_0 & \boldsymbol{u}_0^t & \boldsymbol{v}_0^t \\ \boldsymbol{y}_1 & \boldsymbol{u}_1^t & \boldsymbol{v}_1^t \end{matrix}
$$

**Figure 3:** We impute $\mathbf{W}$, the dissimilarities between $\boldsymbol{E}$ and $\boldsymbol{F}$, by $(\mathbf{D}_0 + \mathbf{D}_1)/2$ to construct $\mathbf{M}$, which is then embedded into the space $\Xi_c$. We impute $\boldsymbol{u}_1$ and $\boldsymbol{v}_0$ by $(\boldsymbol{u}_0 + \boldsymbol{v}_1)/2$. Finally, out-of-sample embedding is used to embed $(\boldsymbol{u}_0^t, \boldsymbol{v}_0^t)^t$ and $(\boldsymbol{u}_1^t, \boldsymbol{v}_1^t)^t$ into $\Xi_c$.

In the classification problem, similarly we create omnibus matrix $\mathbf{M}^J$ from $\mathbf{D}_0^J$, $\mathbf{D}_1^J$ and the imputed matrix $\mathbf{W}^J = (\mathbf{D}_0^J + \mathbf{D}_1^J)/2$. The omnibus matrix $\mathbf{M}^J$ is then embedded into a commensurate space $\Xi_c$. To embed out-of-sample English documents in $\tilde{\mathcal{T}}_0$, we first impute the dissimilarity between $\boldsymbol{x}_{i,0} \in \tilde{\mathcal{T}}_0$ and $\boldsymbol{x}_{j,1} \in \mathcal{T}_1$ by the average of the dissimilarities between $\boldsymbol{x}_{j,1}$ and $\boldsymbol{x}_{i,0}$'s 3 nearest neighbors in $\mathcal{T}_0$. (These dissimilarities can be found in $\mathbf{W}^J$.) All the imputed dissimilarities are stored in $\mathbf{D}_{01}^{\tilde{J}J}$. The dissimilarities between documents in $\tilde{\mathcal{T}}_0$ and $\mathcal{T}_0$ are given by $\mathbf{D}_0^{\tilde{J}J}$, and the dissimilarities among $\tilde{\mathcal{T}}_0$ are given by $\mathbf{D}_0^{\tilde{J}}$. Trosset and Priebe (2008)'s out-of-sample embedding approach is then used to embed $\tilde{\mathcal{T}}_0$ into the space $\Xi_c$. Similarly, new French documents of classes $\tilde{J}$ are embedded into $\Xi_c$. Figure 4 depicts the construction of the omnibus dissimilarity matrix $\mathbf{M}^J$

and how to out-of-sample embed documents in $\tilde{\mathcal{T}}_0$.

$$\mathbf{M}^J \quad = \quad \begin{array}{|c|c|c|} \hline \mathbf{D}_0^J & \mathbf{W}^J & \mathbf{D}_0^{J\tilde{J}} \\ \hline \mathbf{W}^J & \mathbf{D}_1^J & \mathbf{D}_{10}^{J\tilde{J}} \\ \hline \mathbf{D}_0^{\tilde{J}J} & \mathbf{D}_{01}^{\tilde{J}J} & \mathbf{D}_0^{\tilde{J}} \\ \hline \end{array}$$

**Figure 4:** We impute $\mathbf{W}^J$, the dissimilarities between documents in $\mathcal{T}_0$ and $\mathcal{T}_1$, by $(\mathbf{D}_0^J + \mathbf{D}_1^J)/2$ to construct $\mathbf{M}^J$, which is then embedded into the space $\Xi_c$. The dissimilarities between documents in $\tilde{\mathcal{T}}_0$ and $\mathcal{T}_0$ are given by $\mathbf{D}_0^{\tilde{J}J}$ ($\mathbf{D}_0^{J\tilde{J}}$ is the transpose of $\mathbf{D}_0^{\tilde{J}J}$). The dissimilarity between $\boldsymbol{x}_{i,0} \in \tilde{\mathcal{T}}_0$ and $\boldsymbol{x}_{j,1} \in \mathcal{T}_1$ are imputed by the average of the $\mathbf{W}^J$ entries that are corresponding to $\boldsymbol{x}_{i,1}$ and $\boldsymbol{x}_{i,0}$'s 3 nearest neighbors in $\mathcal{T}_0$. All the imputed dissimilarities are stored in $\mathbf{D}_{01}^{\tilde{J}J}$ ($\mathbf{D}_{10}^{J\tilde{J}}$ is the transpose of $\mathbf{D}_{01}^{\tilde{J}J}$).

We refer this approach as the W-approach.

## 3.3  Fusion

We consider one additional step, to combine the data of textual content and graph structure. Ideally both sources of data contain complementary information so that their fusion leads to larger power in testing and higher accuracy in classification than using either textual content data or graph structure data alone. We achieve the fusion by combining the embeddings obtained in the P- or W-approach via the Cartesian product (Ma et al., 2010).

# 4  Results

To compute critical values and estimate powers in hypothesis testing, we randomly select two pairs of matched documents from $\boldsymbol{E}$ and $\boldsymbol{F}$. That is, we leave out four documents, two from each language, and they result in two matched pairs and two non-matched pairs. (Notice that in a real problem we only need to leave one matched pair out to get critical values; leaving two matched pairs out makes it also possible to estimate testing powers.) The approaches introduced in Section 3 are then applied to obtain the distances between the two matched pairs (denoted by $d_0$), and the distances between the two non-matched pairs (denoted by $d_A$). We use Classical Multidimensional Scaling (CMDS) (Torgerson, 1952; Cox and Cox, 2001) in the embedding. Embedding dimension $d = 6$ is determined by Zhu and Ghodsi's automatic dimensionality selection (2006). We use ranks of the distances $d_A$ based on 200 Monte Carlo simulations to estimate the powers for different levels of $\alpha$, where the power $\beta_\alpha$ is the probability of rejecting the null hypothesis when rejection is in fact the correct decision and $\alpha$ is the probability of missing a true match. That is, for each $\alpha \in [0, 1]$, the critical value $c_\alpha$ is defined as the $(100\alpha)$th percentile of $d_0$, and the corresponding power is the percentage of distances in $d_A$ that are larger than the critical value $c_\alpha$. The power at level $\alpha$ is our

performance in determining that a non-match is in fact a non-match. The $\beta$ against $\alpha$ ROC curves are shown in Figure 5. For example, at $\alpha = 0.05$ (missing 5% of the true matches), we obtain a power of $\hat{\beta}_{W\text{-}fusion} = 0.560$ (correctly eliminating 56% of the false matches) via W-fusion. This is a statistical significant improvement over the results obtained sans fusion ($\hat{\beta}_{P\text{-}G} = 0.135$, $\hat{\beta}_{P\text{-}T} = 0.379$, $\hat{\beta}_{W\text{-}G} = 0.403$, $\hat{\beta}_{W\text{-}T} = 0.468$. See Figure 5).
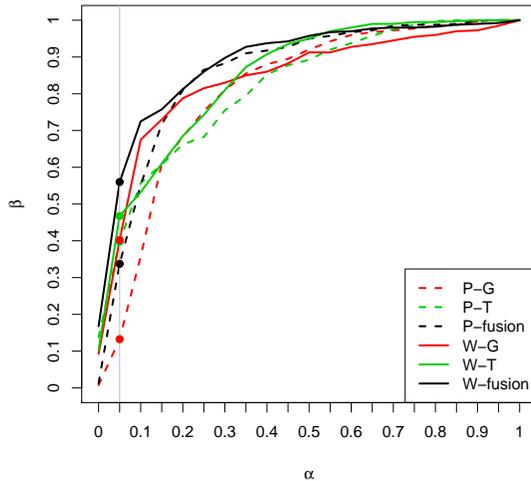


**Figure 5:** The ROC curve depicts that W-approach is generally superior to P-approach; T is generally superior to G; Fusion is generally superior to either G or T alone.

As mentioned in Section 3, commensurate space $\Xi_c$ in the classification problem is determined by $\mathbf{D}_0^J$ and $\mathbf{D}_1^J$. Training English documents in $\tilde{\mathcal{T}}_0$ and new French documents are then embedded into $\Xi_c$. We consider two types of association relations between $\mathcal{T}_0$ and $\mathcal{T}_1$, 1-to-1 association and group association. When assuming 1-to-1 association, we use the information of 1-to-1 correspondence between the training English and French documents with classes in $J$; while for group association, we use only the class label information between English and French documents, but not use the 1-to-1 relationship between them. Introducing group association between $\mathcal{T}_0$ and $\mathcal{T}_1$ makes it possible to define a commensurate space through non-matched English and French documents. When assuming group association, in P-approach we learn transformation matrix through the group means of embeddings, while in W-approach we impute the dissimilarities among same group by 0s and those between different groups by the dissimilarities between group means.

In the commensurate space, we train a linear classifier $g$ based on the embedding of $\tilde{\mathcal{T}}_0$. We then apply $g$ to the embeddings of new French documents. Classification errors are given in Table 1. It is clear that W-approach results in smaller classification errors than P-approach. But combining data from graph structure and text content does not, in general, improve performance.

| Assocation | P-G | P-T | P-fusion | W-G | W-T | W-fusion |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| $1 - 1$ | 0.417 | 0.496 | 0.493 | 0.300 | 0.285 | 0.282 |
| Group | 0.404 | 0.470 | 0.470 | 0.301 | 0.069 | 0.122 |

**Table 1:** Given the association between the training data $\mathcal{T}_0$ and $\mathcal{T}_1$, one-to-one or group-to-group, we transform $\Xi_0$ and $\Xi_1$ into one commensurate space by P- or W-approach. A linear disscriminant classifier is then created based on $\tilde{\mathcal{T}}_0$ and then tested on $\tilde{\mathcal{T}}_1$. The symbols G and T indicate that the Graph and Text data, respectively.

## 5    Conclusion

We have discussed two problems regarding fusion from multiple data sources in a commensurate space:

1. how to determine whether a set of feature vectors represent a single object measured under different conditions?

2. how to create a classifier based on training data measured under one condition in order to classify objects measured in other conditions?

The key to both problems is to construct a commensurate space, where the (transformed) feature vectors of different sources are comparable and would be treated as if they were collected under the same condition. Two approaches were studied. In P-approach, embedding dissimilarity matrices and defining commensurate space are performed separately. W-approach achieves the two procedures simultaneously, by constructing an omnibus dissimilarity matrix. Applying both approaches on Wikipedia data set showed that W-approach leads to higher hypothesis testing powers in the implicit translation problem and smaller errors in the classification problem, compared to P-approach.

## References

Chang, K. C., T. Zhi, and R. K. Saha (2002). Performance evaluation of track fusion with information matrix filter. *IEEE Transactions on Aerospace and Electronic Systems 38*, 455–466.

Chengjun Liu, Wechsler, H. (2001, April). A shape- and texture-based enhanced fisher classifier for face recognition. *Image Processing, IEEE Transactions on 10*(4), 598–608.

Cox, T. F. and M. A. A. Cox (2001). *Multidimensional scaling*. Boca Raton: Chapman & Hall/CRC.

Kludas, J., E. Bruno, and S. M. Maillet (2008). *Information Fusion in Multimedia Information Retrieval*. Berlin, Heidelberg: Springer-Verlag.

Lin, D. and P. Pantel (2002). Concept discovery from text. In *Proceedings of the 19th international conference on Computational linguistics*, Morristown, NJ, USA, pp. 1–7. Association for Computational Linguistics.

Ma, Z., A. Cardinal-Stakenas, Y. Park, M. W. Trosset, and C. E. Priebe (2010, February). Combining dissimilarity representations in embedding product space. *Journal of Classification*, accepted for publication.

Pantel, P. and D. Lin (2002). Discovering word senses from text. In *In Proceedings of ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 613–619.

Ross, A. and A. K. Jain (2004). Multimodal biometrics: An overview. In *Proceedings of 12th Signal Processing Conference (EUSIPCO)*, pp. 1221–1224.

Sibson, R. (1978). Studies in the robustness of multidimensional scaling: Procrustes statistics. *Journal of the Royal Statistical Society. Series B (Methodological) 40*(2), 234–238.

Sun, Q.-S., S.-G. Zeng, Y. Liu, P.-A. Heng, and D.-S. Xia (2005). A new method of feature fusion and its application in image recognition. *Pattern Recognition 38*(12), 2437 – 2448.

Torgerson, W. (1952, December). Multidimensional scaling: I. theory and method. *Psychometrika 17*(4), 401–419.

Trosset, M. W. and C. E. Priebe (2008). The out-of-sample problem for classical multidimensional scaling. *Computational Statistics & Data Analysis 52*(10), 4635–4642.

Zhu, M. and A. Ghodsi (2006, November). Automatic dimensionality selection from the scree plot via the use of profile likelihood. *Computational Statistics & Data Analysis 51*(2), 918–930.