

# Enhancing the Stability of Spectral Ordering with Sparsification and Partial Supervision: Application to Paleontological Data

Dimitrios Mavroeidis  
Department of Informatics  
Athens University of Economics and Business  
Greece  
dmavr@aueb.gr

Ella Bingham  
Helsinki Institute for Information Technology  
University of Helsinki  
Finland  
ella@iki.fi

## Abstract

*Recent studies have demonstrated the prospects of data mining algorithms for addressing the task of seriation in paleontological data (i.e. the age-based ordering of the sites of excavation). A prominent approach is spectral ordering that computes a similarity measure between the sites and orders them such that similar sites become adjacent and dissimilar sites are placed far apart. In the paleontological domain, the similarity measure is based on the mammal genera whose remains are retrieved at each site of excavation. Although spectral ordering achieves good performance in the seriation task, it ignores the background knowledge that is naturally present in the domain, as paleontologists can derive the ages of the sites of excavation within some accuracy. On the other hand, the age information is uncertain, so the best approach would be to combine the background knowledge with the information on mammal co-occurrences. Motivated by this kind of partial supervision we propose a novel semi-supervised spectral ordering algorithm. Our algorithm modifies the Laplacian matrix used in spectral ordering, such that domain knowledge of the ordering is taken into account. Also, it performs feature selection (sparsification) by discarding features that contribute most to the unwanted variability of the data in bootstrap sampling. The theoretical properties of the proposed algorithm are thoroughly analyzed and it is demonstrated that the proposed framework enhances the stability of the spectral ordering output and induces computational gains.*

## 1. Introduction

In this paper we consider the task of ordering the observations in the data, accompanied by partial supervision and sparsification<sup>1</sup>, aiming at a more stable ordering. Although

<sup>1</sup>The terms “sparsification” and “feature selection” will be used interchangeably

it may initially seem surprising that we employ partial supervision and feature selection in a common framework, it is analytically demonstrated in the paper that each component addresses a different cause of instability of the results. In our context, stability refers to the variation of the end result with respect to small changes in the data; in practice we will measure this by bootstrap sampling.

In distance based ordering, the task is to find a permutation of objects such that similar objects become adjacent; in addition, the more dissimilar the objects are, the larger the order distance between them. The standard optimization problem formulation used for deriving the distance based ordering is known to be *NP-hard* [8], and spectral ordering [2, 4] presents a popular, algorithmically feasible approach for deriving approximate solutions. Despite the name “ordering”, the aim is not to rank the objects into any preference ranking, and the first and last object in the ordering are merely those that are maximally dissimilar to each other. Algorithmically, the order solution is derived by the eigenvector corresponding to the second eigenvalue of the data Laplacian matrix.

Our main application area is paleontology: our observations (objects, instances) are sites of excavation and our features (attributes, variables) are mammal genera whose remains are found at these sites. In addition, we have auxiliary information on the estimated ages of the sites; this information is uncertain to some degree. Spectral ordering of the sites of excavation can be based solely on the co-occurrences of the mammal genera, irrespective of the ages of the sites. It has been shown [6] that this kind of plain spectral ordering is a fast and standardized way of biochronological ordering of the sites. Albeit the favorable results in the biochronological ordering task, the spectral ordering does not take into account the background knowledge that naturally exists in the domain. The successful incorporation of domain knowledge is expected to increase the quality of the results.

In the current study, we take advantage of the domain knowledge of the ages of the sites and combine that with the spectral ordering, ending up with a semi-supervised spectral ordering<sup>2</sup>. In addition, we consider feature selection. Towards this target the features that contribute most to the unwanted variation of the data (measured by bootstrap sampling) will be removed. These features correspond to mammals whose observations are noisy. The paleontological data is noisy in many respects [7]: the preservation, recovery and identification of fossils are all random to some extent. These uncertainties are, however, hard to quantify, and a systematic way of characterizing the uncertainty would be most welcome — the behaviour of the features in bootstrap sampling is here chosen for this task.

The two components of the proposed framework, namely partial supervision and feature selection will make the resulting ordering more stable with respect to small variations in the data. As it is analyzed in detail in Section 6, each component of the framework addresses a different cause of instability of the spectral ordering results.

The theoretical analysis suggests and the experiments verify that the main advantages of the proposed framework as induced by the enhancement of stability are twofold: Firstly, results become more resilient to perturbations of the input, thus the reliability of the results is increased. Secondly, the power method [17] computes the ordering result more efficiently than in the original setting.

## 2. Spectral ordering

Given a set of  $n$  objects and a pairwise similarity measure between them, the task of distance based ordering is to derive the order indexes of the objects such that similar objects are placed in adjacent orders while dissimilar objects are placed far apart. More formally, distance sensitive ordering considers the following optimization problem:

$$\min_r \sum_{i,j} (r(i) - r(j))^2 w_{ij}$$

where  $w_{ij}$  is the similarity between objects  $i$  and  $j$  and vector  $r$  is the permutation of  $\{1, 2, \dots, n\}$  that optimizes the objective function. The values of the elements  $r(i)$  of vector  $r$  reflect the ordering of object  $i$ .

It is known that the general optimization problem related to distance based ordering is *NP*-hard [8], and thus approximate solutions should be considered. A popular approach is spectral ordering [2, 4] that performs a continuous relaxation on the solution vector  $r$ , and reduces the optimization problem to a standard eigenvalue-eigenvector problem. In the context of this work we rely on a slight

<sup>2</sup>In the context of this work we will use the terms “semi-supervised” and “partial supervision” to refer to the domain knowledge interchangeably

modification of the standard spectral ordering formulation as derived by [4], where the authors derive the ordering solution as the second eigenvector<sup>3</sup> of the normalized Laplacian matrix  $L = D^{-1/2}WD^{-1/2}$ . Here,  $W$  is the object-object similarity matrix  $W = X^T X$ ,  $D$  is the diagonal degree matrix containing the row sums of  $W$ , and the data matrix  $X$  contains the objects as its columns and the features as its rows. Other choices of  $W$  are also possible:  $W$  can essentially be any object-object similarity matrix. The use of the normalized Laplacian facilitates the theoretical analysis of the proposed semi-supervised spectral ordering framework and also presents theoretical advantages [16] over the unnormalized Laplacian that is commonly used for spectral ordering.

It should be noted that in the spectral graph theory literature the normalized Laplacian matrix is commonly referred to as  $L = I - D^{-1/2}WD^{-1/2}$ , however in the context of this paper we will employ the aforementioned notation and consider the normalized Laplacian as  $L = D^{-1/2}WD^{-1/2}$ . This matrix is well studied in the context of spectral graph theory (e.g. [15] and references within) and it is known to have 1 as its largest eigenvalue. Moreover, by defining the object-similarity matrix  $W = X^T X$ ,  $L = D^{-1/2}WD^{-1/2}$  becomes positive semi-definite.

## 3. Two factors that determine the stability of spectral ordering

A common approach for measuring the stability of spectral algorithms requires the quantification, in the form of an error perturbation matrix  $E$ , of the uncertainty associated with the input matrix (eg. [11]). Based on matrix  $E$  the stability of spectral ordering is determined by the similarity of the ordering solution as derived by the original Laplacian matrix  $L$  versus the perturbed Laplacian matrix  $L + E$ . Further details on the computation of  $E$  in the domain of interest will be provided in Section 6.3.

Based on this formulation, the stability of the ordering solution can be derived by Matrix Perturbation Theory, and more precisely Stewart’s theorem on the perturbation of invariant subspaces [14]. Based on Stewart’s theorem we can derive an upper bound on the difference between the ordering solution of  $L$  versus  $L + E$ . The upper bound applies when the smallest eigengap between the second eigenvalue of  $L$  and the rest is larger than four times the spectral norm of matrix  $E$ . In the case of spectral ordering the smallest eigengap is determined by the eigengap between the first and the second eigenvalue of the Laplacian matrix and the eigengap between the second and the third.

<sup>3</sup>We consider the eigenvalues ordered in decreasing order, i.e. the first eigenvalue is the largest eigenvalue and so on. The first eigenvector is the eigenvector that corresponds to the largest eigenvalue and so on.

The upper bound gets smaller as the eigengap enlarges and the norm of the perturbation matrix  $E$  decreases. Thus, the stability depends on two factors: *the size of the eigengap and the norm of the perturbation*.

As we analyze further in the subsequent section, these eigengaps are not a mere theoretical artifact but are associated with the data-structure as well as computational issues related to the derivation of the spectral ordering solution.

## 4. Semantics of the eigengaps

### 4.1. Eigengap $\lambda_1 - \lambda_2$

The eigengap between the first and the second eigenvalue of the Laplacian matrix is associated with the level of data connectivity. More precisely, if we consider the Laplacian  $D^{-1/2}WD^{-1/2}$  and the associated graph (i.e. a graph with edge weights  $W(i, j)$ ), then the size of the second eigenvalue is associated with the cost of producing two separated clusters [5, 15]. In fact when the eigengap is 0, i.e. the algebraic multiplicity of first eigenvalue is larger than 1, then the graph is disconnected and the clusters can be produced with zero cost. The following theorem illustrates this relation (note that we have appropriately changed the theorem statement from [15] to take into account that we consider the Laplacian  $D^{-1/2}WD^{-1/2}$  instead of  $I - D^{-1/2}WD^{-1/2}$ ):

**Theorem 4.1** (can be found in [15]). *Let  $G$  be an undirected graph with non-negative weights  $W$ . Then the multiplicity  $k$  of the eigenvalue 1 of matrix  $L = D^{-1/2}WD^{-1/2}$  equals the number of connected components in the graph. The eigenspace of 1 is spanned by the vectors  $D^{1/2}e_{A_i}$  of those components, where  $e_{A_i}$  is such that  $e(j)_{A_i} = 1$  for all vertices  $j$  that belong to the connected component  $A_i$ .*

Theorem 4.1 signifies that when the second eigenvalue is close to the first, a small amount of perturbation can make the graph disconnected, thus significantly affecting the second eigenvector. Thus, spectral graph theory provides us with the necessary tools for understanding the source of instability when the eigengap between the first and the second eigengap is small.

### 4.2. Eigengap $\lambda_2 - \lambda_3$

In order to study the eigengap between the second and the third eigenvalue of the Laplacian matrix  $L$ , we assume that the data is adequately connected (i.e. the algebraic multiplicity of the largest eigenvalue is 1) and consider the following transformation:  $L' = L - vv^T$ , where  $v$  is the first eigenvector of Laplacian  $L$  (i.e.  $v = \frac{D^{1/2}e}{\|D^{1/2}e\|}$  with  $D$  being the degree matrix of the Laplacian  $L$  and  $e$  a unit vector,  $e(i) = 1$  for all  $i$ ). With this definition the matrix  $L'$ , apart from  $v$ , has

exactly the same eigenvectors and eigenvalues as  $L$ . Thus, the second eigenvalue of  $L$  is the largest eigenvalue of  $L'$ . This transformation is always possible and requires solely the computation of the degree matrix  $D$ .

The transformation of matrix  $L$  makes apparent the relevance of the power method [17] for computing the spectral ordering solution. Recall that the power method does not derive the full eigen-decomposition of a matrix and can compute solely the dominant eigenvalue and corresponding eigenvector. It starts with an initial vector  $b_0$ , and then computes iteratively  $b_{k+1} = \frac{Ab_k}{\|Ab_k\|}$ . If matrix  $A$  has an eigenvalue that is strictly larger than the rest and if the initial vector  $b_0$  has a non-zero component in the direction of the dominant eigenvector, then the rate of convergence of  $b_k$  will be determined by  $\frac{|\lambda_2|}{|\lambda_1|}$ , where  $\lambda_1$  is the dominant in magnitude eigenvalue of  $A$  and  $\lambda_2$  is the second in magnitude eigenvalue. The larger the eigengap between  $|\lambda_2|$  and  $|\lambda_1|$  the faster the convergence.

Based on  $L'$ , the power method can be used to derive the ordering solution. The power method will converge with rate  $\frac{\lambda_3}{\lambda_2}$ , where  $\lambda_2$  is the second eigenvalue of  $L$  (and thus the dominant eigenvalue of  $L'$ ) and  $\lambda_3$  is the third eigenvalue of  $L$  (and thus the second eigenvalue of  $L'$ ).

This analysis illustrates that the convergence of the power method for computing the ordering solution depends on the eigengap between the second and the third eigenvalue of the Laplacian matrix. A method that successfully enlarges this eigengap will increase the efficiency of the power method.

## 5. Elements of linear algebra

In order to study the behavior and the properties of the proposed spectral ordering framework, we need to recall certain elements of linear algebra. Firstly, we recall Weyl's theorem on the perturbation of eigenvalues.

**Theorem 5.1** (Weyl, can be found in [14]). *Let  $A$  be a symmetric matrix with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  and  $E$  a symmetric perturbation with eigenvalues  $\epsilon_1 \geq \epsilon_2 \geq \dots \geq \epsilon_n$ . Then for  $i = 1, \dots, n$  the eigenvalues  $\bar{\lambda}_i$  of  $A + E$  will lie in the interval  $[\lambda_i + \epsilon_n, \lambda_i + \epsilon_1]$ .*

Another theorem we will employ is concerned with the affect of rank- $k$  updates to matrix eigenvalues.

**Theorem 5.2** (Wilkinson [17], can also be found in [12]). *Suppose  $B = A + \tau \cdot uu^T$  where  $A \in \mathbb{R}^{n \times n}$  is symmetric,  $u \in \mathbb{R}^n$  has unit Euclidean norm and  $\tau \in \mathbb{R}$ . Then, there exist  $m_1, \dots, m_n \geq 0$ ,  $\sum_{i=1}^n m_i = 1$ , such that*

$$\lambda_i(B) = \lambda_i(A) + m_i\tau, \quad i = 1, \dots, n$$

*Moreover, concerning rank- $k$  updates  $B = A + \sum_{i=1}^k \tau_i \cdot uu^T$ , there exist  $m_{ij} \geq 0$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, k$  with  $\sum_{i=1}^n m_{ij} = 1$ ,*

such that

$$\lambda_i(B) = \lambda_i(A) + \sum_{j=1}^k m_{ij} \tau_j, \quad i = 1, \dots, n.$$

## 6. Proposed spectral ordering framework

As we have mentioned in the introductory section, the proposed framework considers partial supervision and feature selection with the general aim of stabilizing the spectral ordering results. In this section we will present each component of the framework and demonstrate their contribution to the stability of the results. Recall that in Section 3 we have stated that stability essentially depends on two factors, namely the size of the relevant eigengaps as well as the uncertainty associated with the Laplacian matrix estimates. In the subsequent sections it is analytically demonstrated that the semi-supervised component is associated with the enlargement of the eigengaps, while the feature selection is concerned with the reduction of uncertainty.

### 6.1. Semi-supervised framework

The semi-supervised component assumes that an initial ordering of the objects is provided and aims at adjusting the original object similarities such that the input ordering is taken into account. Recall that the original object similarities are used for computing the Laplacian matrix  $D^{-1/2}WD^{-1/2}$  (i.e.  $W(i, j)$  is the similarity between object  $i$  and  $j$ ) that derives the ordering solution. The proposed method essentially aims at adjusting the values of the  $W$  matrix based on the input ordering.

In order to achieve this goal, we initially consider the definition of a Laplacian matrix that produces the initial input ordering, i.e. the second eigenvector derives the same results as the input order. If we consider the initial input ordering  $r$  (i.e.  $r(i)$  is the order of object  $i$ ) as a permutation of  $\{1, 2, \dots, n\}$ , and a degree matrix  $D$ , we can define the initial input Laplacian as:

$$L_{input} = v_0 v_0^T + \frac{1}{2} v_1 v_1^T$$

where  $v_0 = \frac{D^{1/2}e}{\|(D^{1/2}e)\|}$ , with  $e$  being the unit vector (i.e.  $e_i = 1$  for all  $i$ ) and

$$v_1 = \frac{r(i) - (\sum_i r(i) \sqrt{d_i}) / (\sum_i \sqrt{d_i})}{\|r(i) - (\sum_i r(i) \sqrt{d_i}) / (\sum_i \sqrt{d_i})\|}, \quad (1)$$

with  $d_i$  being the  $i^{th}$  diagonal element of the Degree matrix.

In order to understand the definition of the  $L_{input}$  matrix, one should initially observe that vector  $v_0$  is essentially the largest eigenvector of any Laplacian matrix with degree matrix  $D$  (if there are no disconnected components). Moreover,

vector  $v_1$  is by construction orthogonal to  $v_0$  and produces exactly the same ordering as  $r$ . Based on the above we can write  $L_{input}$  in the form of a Laplacian with degree matrix  $D$ , i.e.  $L_{input} = D^{-1/2}W_{input}D^{-1/2}$ , which has exactly two eigenvectors  $v_0$  and  $v_1$ , with corresponding eigenvalues 1 and  $\frac{1}{2}$ . The  $W_{input}$  matrix will contain the object similarities that generate the input ordering. Notice that this construction is possible for any degree matrix  $D$ .

It should also be noted that there exist different possible definitions of the  $v_1$  eigenvector that are orthogonal to  $v_0$  and also preserve the initial input order. However, the specific choice of  $v_1$  imposes equal distances between the elements of the eigenvector  $v_1$  and thus also on the ‘‘continuous’’ ordering solution between the objects. In the absence of further knowledge on the initial input ordering it would not be reasonable to impose the additional bias of unequal distances between the objects.

Based on the definition of  $L_{input}$  we derive the final Laplacian as a linear combination of the original data Laplacian (thereafter referred to as  $L_{data}$ ) and  $L_{input}$  as:

$$L_{semi} = cL_{data} + (1 - c)L_{input}$$

where  $0 \leq c \leq 1$  is a confidence factor associated with each component of the summation. The behavior of  $L_{semi}$  can be understood if we write  $L_{semi}$  as:

$$L_{semi} = D^{-1/2}(cW_{data} + (1 - c)W_{input})D^{-1/2}$$

which is possible since  $L_{input}$  is defined with the same degree matrix as  $L_{data}$ . This illustrates the main intuition of the semi-supervised framework that essentially adjusts the similarities of the original Laplacian such that the ordering is taken into account.

Intuitively one would expect that the use of supervision increases the reliability of the ordering results. This intuition is reflected in the eigengaps of  $L_{semi}$ . As demonstrated in the subsequent analysis, they can be enlarged with the appropriate choice of the  $c$  parameter, as compared to  $L_{data}$ .

### 6.2. Theoretical analysis of the semi-supervised framework

We will now analyze theoretically the behavior of the eigenvalues of  $L_{semi}$  with respect to the parameter  $c$ , the eigenstructure of  $L_{data}$  as well as the ordering solutions of  $L_{data}$  and  $L_{input}$ . In most theorems we derive the required amount of supervision (i.e. required value for  $(1 - c)$  or  $c$ ) such that the desired eigenvalue bounds, or eigengaps are achieved. We can summarize the theoretical results as follows:

- Theorem 6.1 demonstrates that parameter  $c$  can fully control the eigenvalues of  $L_{semi}$ , almost independent of the structure of the Laplacians  $L_{data}$  and  $L_{input}$ .

- Theorem 6.2 demonstrates that if the eigenvalues of  $L_{data}$  are close to the bounds we wish to derive for the eigenvalues of  $L_{semi}$ , then these can be achieved with little supervision (i.e. small values for  $(1 - c)$ ).
- Theorems 6.3,6.4,6.5 demonstrate that the behavior of the eigenvalues depends also on the ordering solutions as derived by  $L_{data}$  and  $L_{input}$ . When the ordering solutions conform to a high degree, then the eigengaps are enlarged even with little supervision (i.e. small values for  $(1 - c)$ ).

We will start with the dependence of the eigenvalues of  $L_{semi}$  with respect to the parameter  $c$ . The following theorem demonstrates that with the appropriate choice of parameter  $c$ , large eigengaps can be achieved.

**Theorem 6.1.** *Let  $L_{data}$  be an  $n \times n$  normalized Graph Laplacian,  $c$  a real number such that  $0 \leq c \leq 1$  and  $L_{input}$  be the Laplacian as derived by an initial input ordering. Define the matrix  $L_{semi} = cL_{data} + (1 - c)L_{input}$ . Its largest eigenvalue will be  $\lambda_1(L_{semi}) = 1$ , its second eigenvalue will reside in the interval  $\frac{1}{2} - \frac{c}{2} + c\lambda_n(L_{data}) \leq \lambda_2(L_{semi}) \leq \frac{1}{2} + \frac{c}{2}$ , where  $\lambda_n(L_{data})$  is the smallest eigenvalue of matrix  $L_{data}$ , and its third eigenvalue will be smaller than  $c$ ,  $\lambda_3(L_{semi}) \leq c$ .*

*Proof.* In order to compute the appropriate bounds for the eigenvalues of  $L_{semi}$  we can employ Weyl's theorem on the matrices  $cL_{data}$ ,  $(1 - c)L_{input}$  and  $L_{semi} = cL_{data} + (1 - c)L_{input}$  and derive for the largest eigenvalue of  $L_{semi}$ ,  $\lambda_1(L_{semi})$ :

$$\lambda_1(L_{semi}) \leq \lambda_1(cL_{data}) + \lambda_1((1 - c)L_{input})$$

Based on the fact that  $\lambda_1(cL_{data}) = c \cdot 1 = c$  (since the largest eigenvalue of  $L_{data}$  is 1) and  $\lambda_1((1 - c)L_{input}) = (1 - c) \cdot 1$  (since the largest eigenvalue of  $L_{input}$  is 1) we can derive:

$$\lambda_1(L_{semi}) \leq 1.$$

Moreover for the first Laplacian eigenvector  $v_0$ , we have that  $L_{semi}v_0 = [cL_{data} + (1 - c)L_{input}]v_0 = cL_{data}v_0 + (1 - c)L_{input}v_0 = c \cdot v_0 + (1 - c) \cdot v_0 = v_0$ . Thus  $v_0$  is an eigenvector of  $L_{semi}$  with corresponding eigenvalue 1. Thus

$$\lambda_1(L_{semi}) = 1.$$

Concerning the second eigenvalue of  $L_{semi}$  we can employ Weyl's theorem and state:

$$\lambda_2((1 - c)L_{input}) + \lambda_n(cL_{data}) \leq \lambda_2(L_{semi}) \leq \lambda_2((1 - c)L_{input}) + \lambda_1(cL_{data}).$$

It holds  $\lambda_2((1 - c)L_{input}) = (1 - c)\frac{1}{2}$ ,  $\lambda_n(cL_{data}) = c\lambda_n(L_{data})$  and  $\lambda_1(cL_{data}) = c$ . Thus,

$$(1 - c)\frac{1}{2} + c\lambda_n(L_{data}) \leq \lambda_2(L_{semi}) \leq (1 - c)\frac{1}{2} + c \Leftrightarrow \frac{1}{2} - \frac{c}{2} + c\lambda_n(L_{data}) \leq \lambda_2(L_{semi}) \leq \frac{1}{2} + \frac{c}{2}$$

Concerning the third eigenvalue of  $L_{semi}$  we can employ Weyl's theorem and state:

$$\lambda_3(L_{semi}) \leq \lambda_3((1 - c)L_{input}) + \lambda_1(cL_{data}).$$

We have  $\lambda_3((1 - c)L_{input}) = (1 - c) \cdot 0 = 0$  (since  $L_{input}$  has only two non-zero eigenvalues) and  $\lambda_1(cL_{data}) = c$ . Thus

$$\lambda_3(L_{semi}) \leq c$$

□

The bounds derived in the theorem above depend solely on the parameter  $c$  and illustrate that with the appropriate choice of parameter  $c$ , large eigengaps can be achieved. However, if the eigengaps of matrix  $L_{data}$  are already large, then little supervision (i.e. smaller values of  $(1 - c)$ ) is required. The subsequent theorem illustrates this connection.

**Theorem 6.2.** *Let  $L_{data}$  be an  $n \times n$  normalized Graph Laplacian, and  $L_{input}$  be the the Laplacian as derived by an initial input ordering. Define the matrix  $L_{semi} = [cL_{data} + (1 - c)L_{input}]$ . In order to derive an upper bound  $\bar{\lambda}_2 \geq \frac{1}{2}$  on the second eigenvalue of  $L_{semi}$ ,  $\lambda_2(L_{semi}) \leq \bar{\lambda}_2$ , we must set  $c = \frac{\bar{\lambda}_2 - \frac{1}{2}}{\lambda_2(L_{data}) - \frac{1}{2}}$ . In order to derive an upper bound on the third eigenvalue of  $L_{semi}$ ,  $\lambda_3(L_{semi}) \leq \bar{\lambda}_3$ , we must set  $c \leq \frac{\bar{\lambda}_3 + \lambda_2 - \frac{1}{2}}{\lambda_3(L_{data}) + \lambda_2(L_{data}) - \frac{1}{2}}$ .*

*Proof.* In order to apply Wilkinson's theorem, we consider that matrix  $L_{semi}$  is composed by a rank-2 update on matrix  $cL_{data}$ . We can write for the three largest eigenvalues of  $L_{semi}$ :

$$\begin{aligned} \lambda_1(L_{semi}) &= c\lambda_1(L_{data}) + m_{11}(1 - c) + m_{12}\frac{1-c}{2} \\ \lambda_2(L_{semi}) &= c\lambda_2(L_{data}) + m_{21}(1 - c) + m_{22}\frac{1-c}{2} \\ \lambda_3(L_{semi}) &= c\lambda_3(L_{data}) + m_{31}(1 - c) + m_{32}\frac{1-c}{2} \end{aligned}$$

Since the largest eigenvalue of  $L_{semi}$  is equal to 1, we have:  $\lambda_1(L_{semi}) = 1 \Rightarrow c\lambda_1(L_{data}) + m_{11}(1 - c) + m_{12}\frac{1-c}{2} = 1 \Rightarrow c + (m_{11} + \frac{m_{12}}{2})(1 - c) = 1 \Rightarrow m_{11} + \frac{m_{12}}{2} = 1$ .

Moreover, we have  $\sum_{i=1}^n (m_{i1} + \frac{m_{i2}}{2}) = 1 + \frac{1}{2} \Rightarrow m_{11} + \frac{m_{12}}{2} + \sum_{i=2}^n (m_{i1} + \frac{m_{i2}}{2}) = 1 + \frac{1}{2} \Rightarrow \sum_{i=2}^n (m_{i1} + \frac{m_{i2}}{2}) = \frac{1}{2}$ .

Thus  $m_{21} + \frac{m_{22}}{2} \leq \frac{1}{2}$ .

Now for the second eigenvalue we can write:

$$\lambda_2(L_{semi}) = c\lambda_2(L_{data}) + m_{21}(1 - c) + m_{22}\frac{1-c}{2} \leq c\lambda_2(L_{data}) + \frac{1-c}{2}.$$

Recall that we aim at determining the appropriate  $c$  such that the upper bound  $\bar{\lambda}_2$  is achieved. Thus we have:

$$c\lambda_2(L_{data}) + \frac{1-c}{2} = \bar{\lambda}_2 \Rightarrow c = \frac{\bar{\lambda}_2 - \frac{1}{2}}{\lambda_2(L_{data}) - \frac{1}{2}}$$

In order to derive the appropriate bound for the third eigenvalue we should initially observe that  $m_{21} + \frac{m_{22}}{2} = \frac{\lambda_2(L_{semi}) - c\lambda_2(L_{data})}{1-c}$ .

Thus,  $\sum_{i=3}^n (m_{i1} + \frac{m_{i2}}{2}) = \frac{1}{2} - \frac{\lambda_2(L_{semi}) - c\lambda_2(L_{data})}{1-c} \Rightarrow m_{31} + \frac{m_{32}}{2} \leq$

$$\frac{1}{2} - \frac{\lambda_2(L_{semi}) - c\lambda_2(L_{data})}{1-c}.$$

Now for the third eigenvalue we can write:

$$\lambda_3(L_{semi}) = c\lambda_3(L_{data}) + m_{31}(1-c) + m_{32}\frac{1-c}{2} \leq c\lambda_3(L_{data}) + \frac{1-c}{2} - \lambda_2(L_{semi}) + c\lambda_2(L_{data}).$$

Recall that we aim at determining the appropriate  $c$  such that the upper bound  $\bar{\lambda}_3$  is achieved. Thus we have:

$$c\lambda_3(L_{data}) + \frac{1-c}{2} - \lambda_2(L_{semi}) + c\lambda_2(L_{data}) = \bar{\lambda}_3 \Rightarrow c = \frac{\bar{\lambda}_3 + \lambda_2(L_{semi}) - \frac{1}{2}}{\lambda_2(L_{data}) + \lambda_3(L_{data}) - \frac{1}{2}} \leq \frac{\bar{\lambda}_3 + \lambda_2 - \frac{1}{2}}{\lambda_2(L_{data}) + \lambda_3(L_{data}) - \frac{1}{2}}$$

□

The derived  $c$  for the second eigenvalue is meaningful when the desired upper bound  $\bar{\lambda}_2(L_{semi})$  is smaller than  $\lambda_2(L_{data})$ , and when both are larger than  $\frac{1}{2}$ , as this ensures that  $c \in [0, 1]$ . This is a natural setup because in order to achieve stability one should lower the second eigenvalue, as this will enlarge the eigengap between the first eigenvalue (which is always equal to 1) and the second. Concerning the derived  $c$  for the third eigenvalue, it is meaningful (i.e.  $c \in [0, 1]$ ), when  $\bar{\lambda}_3(L_{semi})$  is smaller than  $\lambda_3(L_{data})$ .

One would generally expect the behavior of the  $L_{semi} = cL_{data} + (1-c)L_{input}$  matrix to also depend on the eigenvectors of  $L_{data}$  and  $L_{input}$  and not solely on the eigenvalues. It would be intuitive to consider that when the ordering solutions as derived by  $L_{data}$  and  $L_{input}$  conform to a high degree, then even with little supervision (i.e. small values of  $(1-c)$ ), the reliability of the ordering results is rapidly increased. This is demonstrated in the following theorems.

**Theorem 6.3** (Best Case Scenario). *Let  $L_{data} = v_0v_0^T + \lambda_2v_2v_2^T + \dots + \lambda_nv_nv_n^T$  be the data Laplacian matrix and  $L_{input} = v_0v_0^T + \frac{1}{2}v_1v_1^T$ . If the ordering solution as derived by the second eigenvector of  $L_{data}$  is equal to the provided supervision  $v_2 = v_1$ , then the eigenvalues of matrix  $L_{semi} = cL_{data} + (1-c)L_{input}$  will be  $\lambda_1(L_{semi}) = 1$ ,  $\lambda_2(L_{semi}) = c\lambda_2(L_{data}) + \frac{1-c}{2}$ , and  $\lambda_i(L_{semi}) = c\lambda_i(L)$ , for  $i = 3, \dots, n$ . Moreover, the required supervision for achieving the eigengap  $\lambda_1(L_{semi}) - \lambda_2(L_{semi}) = gap$ , is  $c = \frac{1/2 - gap}{\lambda_2(L_{data}) - 1/2}$ , and the required supervision for achieving the eigengap  $\lambda_2(L_{semi}) - \lambda_3(L_{semi}) = gap$ , is  $c = \frac{1/2 - gap}{1/2 - (\lambda_2(L_{data}) - \lambda_3(L_{data}))}$ .*

*Proof.* We have that the original data Laplacian is decomposed as  $L_{data} = v_0v_0^T + \lambda_2v_2v_2^T + \dots + \lambda_nv_nv_n^T$  and  $L_{input} = v_0v_0^T + \frac{1}{2}v_2v_2^T$  (since the two matrices induce the same order solution, i.e.  $v_2 = v_1$ ). Thus:

$$L_{semi} = cL_{data} + (1-c)L_{input} = v_0v_0^T + (c\lambda_2(L_{data}) + \frac{1-c}{2})v_2v_2^T + c\lambda_3(L_{data})v_3v_3^T + \dots + c\lambda_n(L_{data})v_nv_n^T.$$

Based on the above, we can derive the required  $c$  value as:

$$\lambda_1(L_{semi}) - \lambda_2(L_{semi}) = gap \Rightarrow 1 - c\lambda_2(L) - \frac{1-c}{2} = gap \Rightarrow$$

$$c = \frac{1/2 - gap}{\lambda_2(L_{data}) - 1/2}$$

$$\lambda_2(L_{semi}) - \lambda_3(L_{semi}) = gap \Rightarrow c\lambda_2(L_{data}) + \frac{1-c}{2} - c\lambda_3(L_{data}) =$$

$$gap \Rightarrow c = \frac{1/2 - gap}{1/2 - (\lambda_2(L_{data}) - \lambda_3(L_{data}))}. \quad \square$$

On the other hand, when the initial input ordering solution corresponds to the eigenvector of  $L_{data}$  that is associ-

ated with the smallest eigenvector, then more supervision (i.e. larger values of  $(1-c)$ ) is required.

**Theorem 6.4** (Worst Case Scenario). *Let  $L_{data} = v_0v_0^T + \lambda_2v_2v_2^T + \dots + \lambda_nv_nv_n^T$  be the data Laplacian matrix and  $L_{input} = v_0v_0^T + \frac{1}{2}v_1v_1^T$ . If the provided supervision is equal to the last eigenvector of the Laplacian matrix  $v_1 = v_n$ , then the eigenvalues of matrix  $L_{semi} = cL_{data} + (1-c)L_{input}$ , will be  $\lambda_1(L_{semi}) = 1$ ,  $\lambda_n(L_{semi}) = c\lambda_n(L_{data}) + \frac{1-c}{2}$  and the rest will be  $\lambda_i(L_{semi}) = c\lambda_i(L_{data})$ , for  $i = 2, \dots, n-1$ . Moreover, the required supervision for achieving the eigengap  $\lambda_2(L_{semi}) - \lambda_3(L_{semi}) = gap$ , is  $c = \frac{1/2 - gap}{1/2 + (\lambda_2(L_{data}) - \lambda_n(L_{data}))}$ .*

*Proof.* We have that  $L_{semi} = cL_{data} + (1-c)L_{input} = v_0v_0^T + c\lambda_2(L_{data})v_2v_2^T + c\lambda_3(L_{data})v_3v_3^T + \dots + (c\lambda_n(L_{data}) + \frac{1-c}{2})v_nv_n^T$ . Thus the eigengap between the second and the third eigenvalue will steadily become smaller as supervision increases (i.e.  $(1-c)$  increases), until the eigenvalue corresponding to the eigenvector  $v_n$  gets larger than  $c\lambda_2(L_{data})$ . Based on the above, we can derive the required  $c$  value as:

$$\lambda_2(L_{semi}) - \lambda_3(L_{semi}) = gap \Rightarrow c\lambda_n(L_{data}) + \frac{1-c}{2} - c\lambda_2(L_{data}) = gap \Rightarrow c = \frac{1/2 - gap}{1/2 - (\lambda_n(L_{data}) - \lambda_2(L_{data}))} \quad \square$$

In general, we can express the initial input ordering solution (i.e. the second eigenvector of  $L_{input}$ ) as a linear combination of the eigenvectors of  $L_{data}$ . Based on this decomposition, it would be intuitive to expect that the eigenvectors that do not participate in the input ranking solutions are downgraded in importance. This is demonstrated in the subsequent theorem.

**Theorem 6.5.** *Let  $L_{data} = v_0v_0^T + \lambda_2v_2v_2^T + \dots + \lambda_nv_nv_n^T$  be the data Laplacian matrix and  $L_{input} = v_0v_0^T + \frac{1}{2}v_1v_1^T$ . Write  $v_1$  as a linear combination of the eigenvectors of the data Laplacian matrix<sup>4</sup>,  $v_1 = w_2v_2 + w_3v_3 + \dots + w_nv_n$ . Then the eigenvalues of  $L_{semi} = cL_{data} + (1-c)L_{input}$  are  $\lambda_1(L_{semi}) = 1$ , and  $\lambda_i(L_{semi}) = c\lambda_i(L_{data})$  for all  $i$  such that  $w_i = 0$ .*

*Proof.* We have that  $L_{input} = v_0v_0^T + \frac{1}{2}v_1v_1^T = v_0v_0^T + \frac{1}{2}(w_2v_2 + w_3v_3 + \dots + w_nv_n)(w_2v_2 + w_3v_3 + \dots + w_nv_n)$ . It is evident that for those  $v_i$  such that  $w_i = 0$  we will have  $L_{input}v_i = 0$ . Thus,  $L_{semi}v_i = c\lambda_i(L_{data})v_i$ . □

This theorem signifies that the eigenvectors that do not participate in the input ranking solution will be quickly downgraded in importance (through the shrinkage of their eigenvalues), while the rest will finally converge to  $v_1$ . The same effect will take place concerning the eigenvectors that have small significance in the solution (i.e.  $w_i \approx 0$ ).

<sup>4</sup>This is always possible since  $\{u_2, \dots, u_n\}$  are orthogonal to  $v_0$  and to each other, thus forming a basis for every vector that is orthogonal to  $v_0$ .

### 6.3. Quantification of uncertainty

As we have analyzed in section 3, an integral component of stability assessment is the quantification of uncertainty in the form of an error-perturbation matrix  $E$ . Since, we have already defined three matrices in the previous section ( $L_{data}$ ,  $L_{input}$  and  $L_{semi}$ ), we will need to define an appropriate perturbation matrix  $E$  for each. We will begin with  $L_{input}$  that is associated with the initial input ordering. Suppose we can characterize the degree of reliability in the supervision by comparing two rankings produced by the domain knowledge: if these are close to each other, then the domain knowledge is reliable. For both rankings we generate a corresponding eigenvector as was described in Section 6.1, and the difference between these vectors will be denoted as  $v = u_1 - u_2$  where  $u_1$  and  $u_2$  are the two ranking eigenvectors. The element  $v(i)$  gives the uncertainty related to object  $i$ . The input perturbation matrix  $E_{input}$  is a rank-1 matrix

$$E_{input} = 1/2vv^T. \quad (2)$$

We will define the error-perturbation matrix for the  $L_{data}$  matrix in a way that will enable feature selection for uncertainty reduction. We initially observe that the order solution of  $L_{data} = D^{-1/2}WD^{-1/2} = D^{-1/2}X^T XD^{-1/2}$  can be derived by  $D^{-1/2}X^T u_2$ , where  $u_2$  is the second eigenvector of the ‘‘feature Laplacian’’  $L_{feat} = XD^{-1}X^T$ . Notice that  $L_{data}$  and  $L_{feat}$  have the same eigenvalues and if  $u$  is an eigenvector of  $L_{feat}$ , then  $D^{-1/2}X^T u$  is an eigenvector of  $L_{data}$ . Thus, the stability of the ordering solution can be derived by the stability of the  $L_{feat}$  matrix. In order to quantify the uncertainty associated with the elements of  $L_{feat}$ , we bootstrap the observations (here, excavation sites) and produce bootstrap confidence intervals for the elements of the  $L_{feat}$  matrix (pair-wise feature similarities). Consequently, we define matrix  $E_{data}$  such that  $E_{data}(i, j)$  is the maximum difference between  $L_{feat}(i, j)$  and the endpoints of the respective confidence interval.

The error-perturbation matrix of  $L_{semi}$  is derived by the norms of the matrices that take part in the summation. More precisely, we define

$$\|E_{semi}\|_2 = c\|E_{data}\|_2 + (1 - c)\|E_{input}\|_2. \quad (3)$$

Having defined all the appropriate error-perturbation matrices, we can move on to evaluate the stability of the spectral ordering framework and explore possible approaches for uncertainty reduction.

### 6.4. Feature selection for uncertainty reduction

Based on the definition of  $E_{data}$  as the perturbation of a  $feature \times feature$  matrix, we can consider feature selection

for uncertainty reduction. The proposed framework is similar in spirit with [11], where the features that contribute maximally to the norm of  $E_{data}$  matrix are sequentially removed. More precisely, at each step of the algorithm, the feature that corresponds to the column (or row) of matrix  $E_{data}$  that has the highest norm is removed. Although, we employ feature selection in the same manner as in [11] we should stress that there are some important differences. The main difference is concerned with the fact that the new perturbation matrix  $E'_{data}$ , as induced by the removal of a feature, will not be a principal submatrix of  $E_{data}$ . This is because the removal of a feature will influence the values of the degree matrix  $D$ , thus affecting the confidence intervals of all the feature-pairs. In order to address this issue, we recompute the confidence intervals and  $E_{data}$  matrix after each feature is removed. However, it should be noted that when there is a large number of features, we can expect that the degree matrix is not severely affected and thus we can consider the principal submatrix of  $E_{data}$  (after removing the row and column  $i$  that corresponds to the removed feature) as an accurate approximation of the new perturbation matrix  $E'_{data}$ . When this is the case, it is guaranteed that the the uncertainty, as expressed by the norm  $\|E_{data}\|_2$  will be reduced.

## 7. Related work

Concerning the semi-supervised component, our work is conceptually related to personalized Pagerank [9]. Personalized Pagerank derives the stationary probability of the random walk based on a weighted linear combination of the transition matrix and a prior distribution, in the form of  $A = [cP + (1 - c)S]^T$ , where  $P$  is the row-stochastic transition matrix and  $S = eu^T$ , where  $u$  contains the prior distribution. Apart from the intuitive probabilistic interpretation of the  $A$  matrix, it has been shown that parameter  $c$  can control the eigengap between the largest and the second eigenvalue.

**Theorem 7.1** (Haveliwala and Kamvar [10]). *Let  $P$  be an  $n \times n$  row-stochastic matrix. Let  $c$  be a real number such that  $0 \leq c \leq 1$ . Let  $S$  be the  $n \times n$  rank-one row-stochastic matrix  $S = eu^T$ , where  $e$  is the  $n$ -vector whose elements are all  $e_i = 1$  and  $u$  is an  $n$ -vector that represents a probability distribution. Define the matrix  $A = [cP + (1 - c)S]^T$ . Its second eigenvalue is  $|\lambda_2| \leq c$ .*

Concerning the feature selection component our work is conceptually related to Stability based Sparse PCA [11]. In this work the authors consider the use of feature selection for uncertainty reduction in the context of PCA, and demonstrate empirically that feature selection can stabilize the PCA results in several real-world UCI datasets.

We use results from matrix perturbation theory [14], stating that the rank- $k$  approximation of a matrix  $A$  is close to a

rank- $k$  approximation of  $A + E$ , if  $E$  has weak spectral properties compared to those of  $A$ . Somewhat similar properties have been used in a different setting, namely speeding up SVD and kernel PCA: Achlioptas [1] shows how to choose the perturbation  $E$  based on the elements of the  $A$  matrix, such that the matrix  $A + E$  is either a quantized or sampled version of  $A$ , making eigenvalue algorithms work faster.

The prospects of spectral ordering in the paleontological domain have been demonstrated by Fortelius et al [6]. In this work, plain spectral ordering of the sites, based on mammal co-occurrences and discarding the age information of the sites, was considered. In addition, Puolamäki et al [13] present a full probabilistic model that again only considers the co-occurrences in the data.

## 8 Empirical results

In the experiments we aim at verifying that the proposed framework enhances the stability of spectral ordering and increases the relevant eigengaps in the paleontological data. Recall that this will increase the reliability of the ordering results and improve on the convergence rate of the power method. The experiments indeed verify the anticipated behavior and increase the relevant eigengaps (demonstrated in Sections 8.2,8.3) and also the stability of the ordering results (demonstrated in Section 8.3).

### 8.1 Data

The paleontological data we are considering consists of findings of land mammal genera in Europe and Asia 25 to 2 million years ago. The data set is stored and coordinated in the University of Helsinki, Department of Geology [3]. Our version of the data set was downloaded on June 26, 2007.

The observations in our data are the sites of excavation and the features are mammal genera whose remains are found at these sites. In total we have 1887 observations and 823 features. The data matrix is 0-1 valued: an entry  $x_{ij} = 1$  means that mammal  $i$  was found at site  $j$ , and 0 otherwise. The data is very sparse: about 1 per cent of the entries are nonzero. We will also work with a small subset of data containing 1123 observations and 18 features; this subset is more dense, having 12 per cent of its entries nonzero. Thereafter we will refer to the sparse dataset as *paleo<sub>sp</sub>* and the dense dataset as *paleo<sub>d</sub>*.

In addition, we have auxiliary information on the estimated ages of the sites: an approximate age for each site, and also a more precise age for some sites; the methods available for estimating the ages vary from site to site, and thus at some sites the information is more certain than at others. The approximate ages will be used to construct an initial ranking  $r_{input}$  of the sites, and this will be used as an input in the semi-supervised setting, results of which will

be presented in Section 8.2. Both the precise and approximate ages will be needed when quantifying our belief in the initial ranking, that is, defining the perturbation matrix  $E_{input}$  of  $L_{input}$  as discussed in Section 6.3; empirical results on this will be shown in Section 8.3.

We will assume that the data is fully connected in that the algebraic multiplicity of the first eigenvalue of the data Laplacian is 1: if this is not the case, the removal of disconnected observations will be a preprocessing step. In addition, we will preprocess the data such that almost-disconnected components are removed too: these correspond to objects that are very weakly connected to the rest of the objects. For such objects  $i$ , the value  $r(i)$  in the order vector  $r$  (obtained by sorting the second eigenvector) is very large compared to other  $r(j)$ .

### 8.2 Effect of supervision on the stability

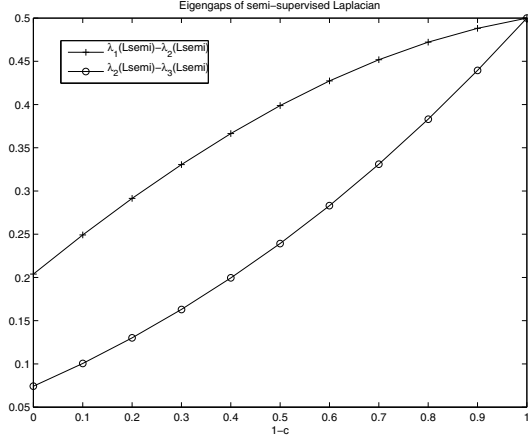
Let us first demonstrate that the eigengaps of the data Laplacian increase when domain knowledge is taken into account. These experiments were performed on the sparse and large *paleo<sub>sp</sub>* dataset, where the initial eigengaps are small. Recall that the stability of spectral ordering essentially depends on two factors, one of which are the eigengaps between the first and second eigenvalue and the second and third eigenvalue of the Laplacian. Figure 1 shows the behavior of the eigengaps of the semi-supervised Laplacian  $L_{semi} = cL_{data} + (1 - c)L_{input}$  at a varying level of supervision. Choosing  $1 - c = 0$  corresponds to no supervision, in which domain knowledge is not taken into account and the spectral ordering is done based on feature co-occurrences only; the eigengaps at  $c = 1$  thus show the eigengaps of the data Laplacian. In contrast,  $1 - c = 1$  corresponds to the trivial case of full supervision of the ranking, in which co-occurrences in the data are not taken into account but only the domain knowledge ranking is used. We observe that both eigengaps increase rapidly when the level of supervision increases. Thus the spectral ordering becomes more stable as more emphasis is put on the domain knowledge.

### 8.3 Effect of feature selection on the stability

We will then demonstrate that the stability of the spectral ordering increases as features are removed step by step. The removed features will be chosen based on their contribution on the variability of the feature-feature similarity matrix, measured as matrix  $E_{data}$  discussed in Section 6.4. It should be noted that after each feature is removed,  $L_{semi}$  is reevaluated based on  $L_{input}$  and  $L_{data}$  which are appropriately recomputed.

We will measure the stability of the spectral ordering by a “stability factor”  $sf$  that depends on the eigengaps and the





**Figure 1. Eigengaps  $\lambda_1(L_{semi}) - \lambda_2(L_{semi})$  (+) and  $\lambda_2(L_{semi}) - \lambda_3(L_{semi})$  (o) at varying level of supervision. Horizontal axis:  $1 - c$ , confidence in domain knowledge.  $1 - c = 0$ : no supervision;  $1 - c = 1$ : full supervision.**

norm of the perturbation matrix:

$$sf_{semi} = \frac{\min(\lambda_1(L_{semi}) - \lambda_2(L_{semi}), \lambda_2(L_{semi}) - \lambda_3(L_{semi}))}{\|E_{semi}\|_2} \quad (4)$$

For the stability factor we will need appropriate values for  $L_{semi}$  and  $\|E_{semi}\|_2$ .

Let us first construct the semi-supervised Laplacian  $L_{semi} = cL_{data} + (1 - c)L_{input}$ . For this we need to carefully choose the confidence factor  $c$  reflecting our belief in the observed data versus the initial ranking. We can either rely on a domain expert or better still, derive  $c$  from the body of domain knowledge: we will choose to define

$$c = \|E_{input}\|_2 / (\|E_{data}\|_2 + \|E_{input}\|_2) \quad (5)$$

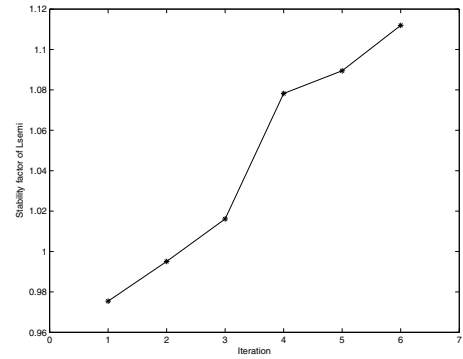
which naturally characterizes the confidence such that a large perturbation in the initial input ranking leads to a high confidence in the observed data, and vice versa. In the definition (5), the data perturbation  $E_{data}$  will be obtained by bootstrap sampling as discussed in section 6.3. The input perturbation  $E_{input}$  for paleontological data will be derived based on the availability of approximate or precise ages for each site: in addition to the initial ranking  $r_{input}$  based on approximate ages of the sites, we construct another initial ranking  $r_s$  using the precise ages available for some of the sites. (The sites for which a precise age is not available will get an average ranking in  $r_s$ .) For both rankings  $r_{input}$  and  $r_s$  we generate a corresponding eigenvector,  $v_{input}$  and  $v_s$ , using Equation (1). We then take the difference between these eigenvectors as  $v = v_{input} - v_s$  and use that in place of  $v$  in Equation (2), to measure the difference between the two

orderings. This gives us the perturbation  $E_{input}$  associated with the domain knowledge. Having now collected all the necessary components, we can construct the matrix  $L_{semi}$ .

For the stability factor in Equation (4) we also need the value for  $\|E_{semi}\|_2$ . Based on the definition for  $c$ , the equation (3) now simplifies

$$\|E_{semi}\|_2 = c\|E_{data}\|_2 + (1 - c)\|E_{input}\|_2 = \frac{2\|E_{data}\|_2\|E_{input}\|_2}{\|E_{data}\|_2 + \|E_{input}\|_2} \quad (6)$$

Having now defined all components of the stability factor (4) let us then see how it behaves when features are iteratively removed. Figure 2 shows the results. We have employed subset *paleo<sub>d</sub>* of 1123 observations and 18 features in this experiment. At each iteration, one feature is removed based on its contribution to the data perturbation. Simultaneously, a few observations typically get removed, as they have become disconnected with the other observations due to the removal of the feature in question. The stability factor of the semi-supervised Laplacian increases during feature selection, showing that feature selection enhances the stability of spectral ordering.

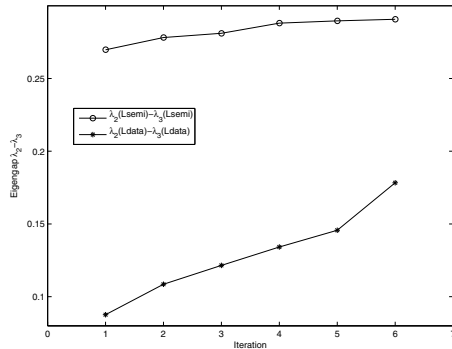


**Figure 2. Stability factor during feature selection. Horizontal axis: iteration. One feature is removed at each iteration.**

In addition, recall that the eigengap between the second and third eigenvalue affects the convergence of the power method, as discussed in Section 4.2. Figure 3 shows that this eigengap increases during feature selection, both in the original data and in the semi-supervised setting. Thus feature selection enhances the behaviour of the power method in both the original spectral ordering framework and the semi-supervised framework.

## 9 Discussion

In this paper we have shown how to increase the stability of spectral ordering using two separate tools: par-



**Figure 3. Eigengap between the 2nd and 3rd eigenvalue during feature selection. Semi-supervised spectral ordering (○), original spectral ordering (\*). One feature is removed at each iteration (horizontal axis).**

tial supervision in the form of a (possibly uncertain) domain knowledge ordering, and sparsification in the form of feature selection. We have presented a detailed theoretical analysis showing how the eigengaps between the first and second eigenvalue, and the second and third eigenvalue, of the Laplacian affect the stability, and how partial supervision will increase the eigengaps. Feature selection in turn will decrease the norm of the perturbation matrix  $E$  that quantifies the uncertainty associated with the observed data.

Our main application area is paleontology: we have considered the ordering of the sites of excavation in paleontological data, by complementing spectral ordering with domain knowledge of the approximate ages of the sites. The paleontological data is noisy in that many observations are missing, and prone to small changes when the findings are more carefully examined. Also, we never have access to the exact ages of the sites. Thus when ordering the sites, the best we can aim at is an ordering that is as stable as possible with respect to small variations in the data. This motivates our task of optimizing the stability of spectral ordering. We have shown that in the paleontological data, the eigengaps quickly increase as semi-supervision is used. Also, feature selection, by removing the mammals that contribute most to the variation of the results in bootstrap sampling, is demonstrated to increase the stability of spectral ordering.

In future work we aim at exploring the potentials of our framework in different application domains, where partial supervision is naturally present. Moreover, we aim at extending the proposed framework to spectral clustering.

**Acknowledgements.** The authors wish to thank professor Mikael Fortelius for fruitful discussions regarding paleontological data and professor Heikki Mannila for his insights in

spectral ordering. E. Bingham was supported by Academy of Finland grant 118653 (ALGODAN).

## References

- [1] D. Achlioptas. Random matrices in data analysis. In J.-F. Boulicaut, F. Esposito, F. Giannotti, and D. Pedreschi, editors, *Machine Learning: ECML 2004*, number 3201 in LNAI, pages 1–7. Springer, 2004.
- [2] J. E. Atkins, E. G. Boman, and B. Hendrickson. A spectral algorithm for seriation and the consecutive ones problem. *SIAM Journal on Computing*, 28(1):297–310, 1999.
- [3] M. F. (coordinator). Neogene of the Old World database of fossil mammals (NOW). University of Helsinki. <http://www.helsinki.fi/science/now/>, 2007.
- [4] C. H. Q. Ding and X. He. Linearized cluster assignment via spectral ordering. In C. E. Brodley, editor, *Proc. 21st Intl Conf on Machine Learning (ICML)*. ACM, 2004.
- [5] C. H. Q. Ding, X. He, and H. Zha. A spectral method to separate disconnected and nearly-disconnected web graph components. In *Proc. 7th Intl Conf on Knowledge Discovery and Data Mining (KDD)*, pages 275–280, 2001.
- [6] M. Fortelius, A. Gionis, J. Jernvall, and H. Mannila. Spectral ordering and biochronology of European fossil mammals. *Paleobiology*, 32(2):206–214, 2006.
- [7] M. Fortelius, L. Werdelin, P. Andrews, R. L. Bernor, A. Gentry, L. Humphrey, W. Mittmann, and S. Viranta. Provinciality, diversity, turnover and paleoecology in land mammal faunas of the later Miocene of western Eurasia. In R. Bernor, V. Fahlbusch, and W. Mittmann, editors, *The Evolution of Western Eurasian Neogene Mammal Faunas*, pages 414–448. Columbia University Press, 1996.
- [8] A. George and A. Pothen. An analysis of spectral envelope reduction via quadratic assignment problems. *SIAM Journal on Matrix Analysis and Applications*, 18(3):706–732, 1997.
- [9] T. Haveliwala, S. Kamvar, and G. Jeh. An analytical comparison of approaches to personalizing PageRank. In *Technical Report: <http://dbpubs.stanford.edu:8090/pub/2003-35>*, 2003.
- [10] T. Haveliwala, S. Kamvar, and G. Jeh. The second eigenvalue of the Google matrix. In *Technical Report: <http://dbpubs.stanford.edu:8090/pub/2003-35>*, 2003.
- [11] D. Mavroudis and M. Vazirgiannis. Stability based sparse LSI/PCA: Incorporating feature selection in LSI and PCA. In *Machine Learning: ECML 2007*, pages 226–237, 2007.
- [12] S. Mika. *Kernel Fisher Discriminants*. PhD thesis, University of Technology, Berlin, 2002.
- [13] K. Puolamäki, M. Fortelius, and H. Mannila. Seriation in paleontological data using Markov Chain Monte Carlo methods. *PLoS Computational Biology*, 2(2):e6, 2006.
- [14] G. W. Stewart and G.-J. Sun. *Matrix Perturbation Theory*. Academic Press, 1990.
- [15] U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- [16] U. von Luxburg, M. Belkin, and O. Bousquet. Consistency of spectral clustering. *Annals of Statistics*, 36(2):555–586, 2008.
- [17] J. H. Wilkinson. *The Algebraic Eigenvalue Problem*. Oxford University Press, 2004.