

QUANTIFYING THE THEORY VS. PROGRAMMING DISPARITY USING SPECTRAL BIPARTIVITY ANALYSIS AND PRINCIPAL COMPONENT ANALYSIS

Natarajan Meghanathan

Department of Electrical & Computer Engineering and Computer Science,
Jackson State University, Jackson, MS, USA

ABSTRACT

Some students in the Computer Science and related majors excel very well in programming-related assignments, but not equally well in the theoretical assignments (that are not programming-based) and vice-versa. We refer to this as the "Theory vs. Programming Disparity (TPD)". In this paper, we propose a spectral bipartivity analysis-based approach to quantify the TPD metric for any student in a course based on the percentage scores (considered as decimal values in the range of 0 to 1) of the student in the course assignments (that involves both theoretical and programming-based assignments). We also propose a principal component analysis (PCA)-based approach to quantify the TPD metric for the entire class based on the percentage scores (in a scale of 0 to 100) of the students in the theoretical and programming assignments. The spectral analysis approach partitions the set of theoretical and programming assignments to two disjoint sets whose constituents are closer to each other within each set and relatively more different from each across the two sets. The TPD metric for a student is computed on the basis of the Euclidean distance between the tuples representing the actual numbers of theoretical and programming assignments vis-a-vis the number of theoretical and programming assignments in each of the two disjoint sets. The PCA-based analysis identifies the dominating principal components within the sets of theoretical and programming assignments and computes the TPD metric for the entire class as a weighted average of the correlation coefficients between the dominating principal components representing these two sets.

KEYWORDS

Spectral Analysis, Principal Component Analysis, Correlation Coefficient, Theory vs. Programming Disparity, Eigenvector, Bipartivity

1. INTRODUCTION

Spectral analysis of a complex network has been observed to reveal significant structural details that would have been hitherto unknown in the scientific community [1]. Spectral analysis of a network graph typically involves the computation of the Eigenvectors and their corresponding Eigenvalues using one of the symmetric matrices such as the adjacency matrix, Laplacian matrix [7] and etc that reflect the topology of the network [2]. The indexes of the entries in an Eigenvector correspond to the sorted order (increasing order) of the node ids. Spectral analysis of an $n \times n$ matrix results in ' n ' Eigenvalues and the corresponding ' n ' Eigenvectors (one Eigenvector per Eigenvalue). Any two Eigenvectors are orthogonal to each other (i.e., the dot product of any two Eigenvectors is zero). The largest Eigenvalue is called the "Principal Eigenvalue" and the Eigenvector corresponding to the principal Eigenvalue is called the "Principal Eigenvector" [3]. If all the entries in the underlying matrix used for spectral analysis are positive (≥ 0), then the principal Eigenvalue is guaranteed to be positive and all the entries in the principal Eigenvector are also positive. Hence, when computed based on matrices with positive entries, in order to still

satisfy the *mutually orthogonal property*, (unless all the entries in an Eigenvector are zero) at least one non-zero entry in every Eigenvector other than the Principal Eigenvector is guaranteed to be negative so that the dot product of any two Eigenvectors evaluates to zero.

In a classical work [4], Estrada et al proposed that the extent of bipartivity (in the form of a bipartivity index) among the vertices in a network could be quantified using the Eigenvalues of the adjacency matrix of the network graph. A graph is said to be bipartite if the vertices of the graph could be grouped into two disjoint sets such that the end vertices of any edge are in the two different partitions and not in the same partition. Estrada et al observed that if the underlying graph is bipartite, the vertices with positive and negative entries in the Eigenvector (hereafter referred to as the 'bipartite Eigenvector') corresponding to the smallest Eigenvalue represent the two disjoint partitions of vertices in the graph. Estrada et al also observed that if the underlying graph is not bipartite, the vertices with positive and negative entries in the bipartite Eigenvector could still be construed to form the two disjoint partitions of the vertices of the graph such that there are exist a minimal number of edges (referred to as the 'frustrated edges') between vertices in the same partition and a majority of the edges are between vertices across the two partitions.

In the first half of the paper, we conduct spectral bipartivity analysis of the scores earned by a student in theoretical and programming assignments of a Computer Science course and seek to quantify the extent of disparity in the scores earned by the student in the theoretical assignments vs. programming assignments. Some students in the Computer Science and related majors excel very well in programming-related assignments, but not equally well in the theoretical assignments (that are not programming-based) and vice-versa. We refer to this as the "Theory vs. Programming Disparity (*TPD*)". Our methodology is briefly described here (more details are in Section 2): The student score in each assignment is considered in a decimal percentage scale of 0 to 1 (i.e., each assignment is evaluated for 100% and the decimal percentage score for a student in the assignment is the percentage score divided by 100: for example, if an assignment score is 81%, the decimal percentage score is $81/100 = 0.81$). We first determine the Difference Matrix (*DM*) of the student scores in the assignments wherein an entry DM_{uv} is the absolute difference in the decimal percentage scores of the two assignments u and v . We then determine the bipartite Eigenvector of the *DM* by subjecting it to spectral analysis. We identify the index entries with positive signs and negative signs, and the corresponding assignment IDs are grouped into two separate (disjoint) sets. We observe that any two assignments with similar (closer) values for the decimal percentage scores are more likely to be grouped into the same set and two assignments with appreciably different decimal percentage scores are more likely to be grouped in separate sets. That is any two vertices u and v whose DM_{uv} entry is closer to 0 are more likely to end up in the same set of vertices and vertices u and v whose DM_{uv} entry is much greater than 0 are more likely to end up in different sets of vertices. Such an observation is consistent with the observations made by Estrada et al for bipartivity analysis using an adjacency matrix A (i.e., vertices u and v whose A_{uv} entries were 0 are more likely to be in the same partition and vertices u and v whose A_{uv} entries were 1 are more likely in different partitions). We quantify the *TPD* on the basis of the Euclidean distance between the actual number of theoretical and programming assignments vs. the number of theoretical and programming assignments in the two sets of disjoint assignments identified through spectral bipartivity analysis.

In the second half of the paper, we present a principal component analysis (PCA)-based approach to quantify the *TPD* for an entire class. PCA is a dimensionality reduction technique that maximally captures the variances across a dataset of m features to n dominating principal components such that $n \ll m$. The dominating principal components are considered to cumulatively capture at least 80% of the variances among the feature values in the dataset. In this context, we propose an approach wherein we will first separately identify the dominating principal components of the set of theoretical assignments and the set of programming

assignments for the entire class. We will then determine the pair-wise correlation coefficients (Pearson's) between the principal components across these two categories of assignments. If there is disparity among the students between the two categories of assignments, we expect these Pearson's correlation coefficients to be negative: indicating that students who scored high in the theoretical assignments did not score as high in the programming assignments and vice-versa. If the correlation coefficients are positive, it indicates the students who scored high (low) in the theoretical (programming) assignments scored high (low) in the other category as well. We propose to quantify the TPD for an entire class as a weighted average of these Pearson's correlation coefficients, with the weights being the product of the variances of the pair-wise dominating principal components of the theoretical and programming sets of assignments.

The rest of the paper is organized as follows: In Section 2, we present our proposed spectral bipartivity analysis-based methodology to quantify the *TPD* metric per student using a running example. Section 3 evaluates the effectiveness of the proposed *TPD* per student approach with two of the well-known metrics (Bipartivity index: *BPI* [4] and Hausdorff Distance: *HD* [5]) that exist in the literature to study the effectiveness of partitioning of a data set to two clusters. Section 3 also highlights the uniqueness of the proposed *TPD* per student approach. Section 4 presents the PCA-based approach proposed to quantify the *TPD* metric for an entire class. Section 5 reviews related work and Section 6 concludes the paper. Throughout the paper, the terms 'set' and 'partition', 'network' and 'graph', 'edge' and 'link' are used interchangeably. They mean the same.

2. SPECTRAL BIPARTIVITY ANALYSIS-BASED APPROACH TO QUANTIFY THEORY VS. PROGRAMMING DISPARITY PER STUDENT

Let P and T be respectively the set of scores (represented in decimal percentage format) earned by a student in programming and theoretical assignments. The indexes for the assignments in the set P range from 0 to $|P| - 1$, where $|P|$ is the cardinality of the set P (i.e., the number of programming assignments). The indexes for the assignments in the set T range from $|P|$ to $|P| + |T| - 1$, where $|T|$ is the cardinality of the set T (i.e., the number of theoretical assignments). Let S be the union of the two sets P and T . That is, the set S comprises of the scores earned by the student in the programming assignments, followed by the theoretical assignments. The indexes for the assignments in S are the same as their indexes in the sets P or T , whichever they come from. Let DM be a symmetric/square matrix whose dimensions correspond to the cardinality of the set S . An entry DM_{ij} for row index i and column index j represents the absolute difference in the decimal percentage scores for the i^{th} and j^{th} element/assignment in the set S . Figure 1 presents the sets P , T and S (of size 8, 4 and 12 respectively) as well as the DM matrix (of dimensions 12 x 12) for a sample data set that is used as a running example to explain the proposed methodology in this section. The indexes for the assignments in the sets P and T range from 0...7 and 8...11 respectively; the indexes of these assignments are retained in the set S that is the amalgamation of the assignments in the sets P and T (in the same order). Likewise, the indexes in DM correspond to the indexes in the set S .

Figure 2 presents the 12 Eigenvalues and the entries of the Bipartite Eigenvector, which is the Eigenvector corresponding to the smallest Eigenvalue of -2.2285, of the Difference Matrix (DM). The indexes (colored in blue) whose entries are positive (≥ 0) are grouped to partition (set) X and the indexes (colored in red) whose entries are negative (< 0) are grouped to partition (set) Y . These are the two disjoint partitions of the assignments in the set S predicted based on spectral bipartivity analysis. In Figure 2, we also display the submatrices of the Difference Matrix that show the difference in the decimal percentage scores for any two assignments within the sets X and Y as well as between the two sets X and Y . We notice the entries in the submatrices

corresponding to the differences in the assignment scores within the sets X or Y are relatively much smaller (closer to 0) compared to the entries in the submatrix corresponding to the differences in the assignment scores between an assignment in set X and an assignment in set Y.

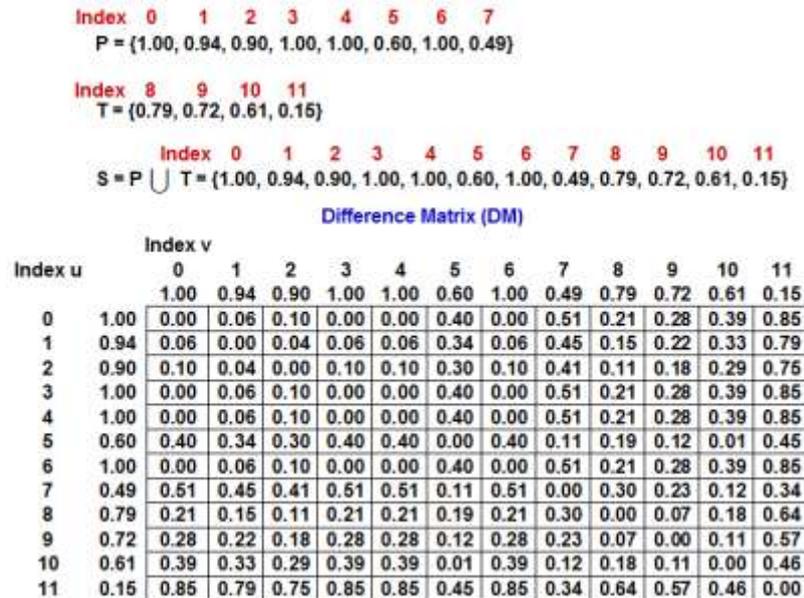


Figure 1. Sample Data set to Illustrate the Spectral Bipartivity-based Analysis for Theory vs. Programming Disparity

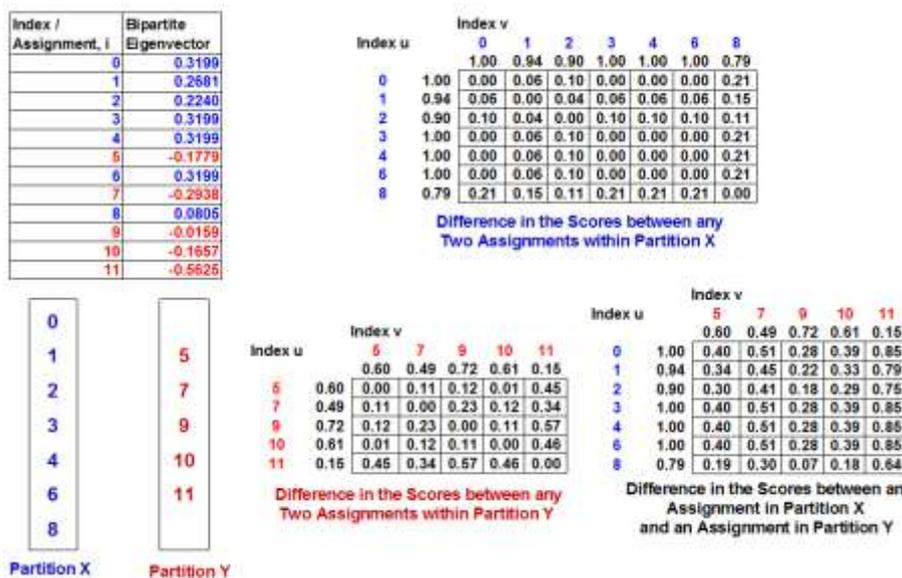


Figure 2. Bipartite Eigenvector for the Difference Matrix of Figure 1 as well as the Submatrices Representing the Difference Values between Assignments within the two Partitions and across the Two Partitions

We now seek to quantify the extent to which the grouping of the assignments in sets X and Y are closer to the grouping of the programming and theoretical assignments in the sets T and P. In order for our hypothesis (that there exists a disparity in the scores earned by the students in the programming vs. theoretical assignments) to be true, we would like the two sets X and Y to be the

same as the two sets T and P (or P and T); that is, we would prefer the set X to be all theoretical assignments and the set Y to be all programming assignments or the set X to be all programming assignments and the set Y to be all theoretical assignments. In this pursuit, we first determine the number of theoretical assignments and the number of programming assignments in each of the sets X and Y and let these be indicated using symbols $|X_T|$, $|X_P|$, $|Y_T|$ and $|Y_P|$. We then determine the Euclidean distance between the tuples $(|X_T|, |Y_P|)$ and $(|T|, |P|)$ as well as between the tuples $(|Y_T|, |X_P|)$ and $(|T|, |P|)$ and refer to the minimum of these two Euclidean distances as the Theory vs. Programming Tuple Proximity ($TPTP$) distance for the given data set.

The maximum value for the $TPTP$ distance is $\sqrt{|T|^2 + |P|^2}$ and we will incur it when all the four values $|X_T|$, $|X_P|$, $|Y_T|$ and $|Y_P|$ are zero each (i.e., the assignment IDs in the partitions X and Y identified through spectral bipartivity analysis have no overlap with the assignment IDs in the partitions P and T). Such a scenario occurs when there is minimal or no disparity among the scores in the programming vs. theoretical assignments and the distribution of the assignment IDs in the partitions X and Y are random. On the other hand, if there is maximum disparity, the assignment IDs in partitions X and Y will overlap with those of P and T , and either $|X_P| \approx |P|$ and $|Y_T| \approx |T|$ or $|Y_P| \approx |P|$ and $|X_T| \approx |T|$. This would make the $TPTP$ distance much smaller than the maximum value of $\sqrt{|T|^2 + |P|^2}$. Considering the above interpretation of the $TPTP$ distance, we formulate the TPD (Theory vs. Programming Disparity) metric as follows:

$$TPD = 1 - \frac{TPTP}{\sqrt{|T|^2 + |P|^2}}$$

If there is maximum disparity, then the $TPTP$ distance will be either closer to 0 or the ratio $TPTP / \sqrt{|T|^2 + |P|^2}$ be closer to 0, making the TPD metric score to be closer to 1. On the other hand, if there is no disparity, the $TPTP$ distance will be closer to the maximum value of $\sqrt{|T|^2 + |P|^2}$ and as a result the ratio $TPTP / \sqrt{|T|^2 + |P|^2}$ be closer to 1, making the TPD metric score to be closer to 0.

Figure 3 presents the computation of the $|X_T|$, $|X_P|$, $|Y_T|$ and $|Y_P|$ numbers for the running example of Figures 1-2 and the computation of the TPD metric for the sample data set. The TPD metric value for this data set is 0.75, indicating that there is an appreciable disparity in the theoretical vs. programming scores in this data set. Such a conclusion could also be justified by visually looking at the proximity of the tuple $(|Y_T| = 3, |X_P| = 6)$ corresponding to the $TPTP$ distance to the tuple $(|T| = 4, |P| = 8)$ in Figure 3 as well as the raw data set values $\{0.79, 0.72, 0.61, 0.15\}$ and $\{1.00, 0.94, 0.90, 1.00, 1.00, 0.60, 1.00, 0.49\}$ for the sets T and P respectively.

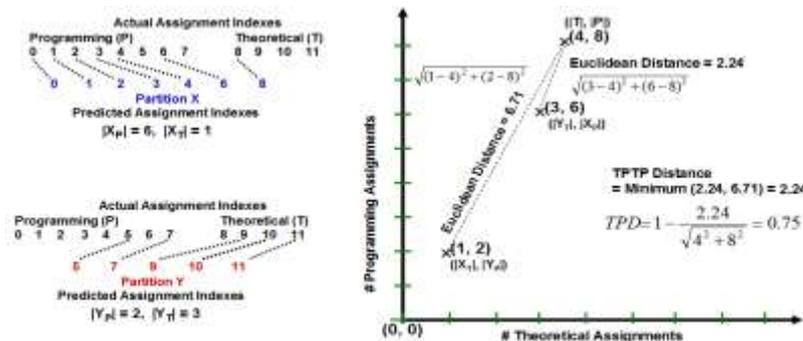


Figure 3. Computation of the TPD Metric for the Sample Data Set of Figure 1

3. EVALUATION OF THE SPECTRAL BIPARTIVITY ANALYSIS APPROACH

In this section, we apply the proposed spectral bipartivity analysis-based approach to quantify the theory vs. programming disparity per student in the CSC 228 Data Structures and Algorithms course taught in Spring 2020 at Jackson State University, MS, USA. In addition to the *TPD* metric, we consider two other metrics that appear to be potentially applicable to quantify the extent of disparity in a data set with respect to two different categories (in this case, theoretical vs. programming assignments). These are:

(i) Bipartivity Index (*BPI*): The *BPI* was originally proposed by Estrada et al [4] to quantify the extent of bipartivity between the two partitions of vertices identified using the Eigenvector (referred to as the Bipartite Eigenvector) corresponding to the smallest Eigenvalue. The input matrix for Estrada et al's spectral bipartivity analysis is a 0-1 adjacency matrix. If the underlying graph is not bipartite, the two partitions of vertices identified using the Bipartite Eigenvector have as few edges as possible between vertices within the same partition and a majority of the edges are between vertices across the two partitions. If the underlying graph is indeed bipartite, the two partitions of vertices identified using the Bipartite Eigenvector will have no edges between vertices within the same partition and all the edges in the graph will be between vertices across the two partitions. The *BPI* of a graph of '*n*' vertices is computed using the following formulation based on the '*n*' Eigenvalues (λ) of its 0-1 adjacency matrix. If the underlying graph is bipartite, the sum of the sinh values of the '*n*' Eigenvalues will be zero and the *BPI* of the graph will be 1.0. If the underlying graph is not bipartite, then the sum of the sinh values of the '*n*' Eigenvalues will be greater than 0, and the *BPI* will be less than 1.0.

$$BPI = \frac{\sum_{i=0}^{n-1} \cosh(\lambda_i)}{\sum_{i=0}^{n-1} \cosh(\lambda_i) + \sum_{i=0}^{n-1} \sinh(\lambda_i)}$$

In this section, we will explore how well the Eigenvalues of the Difference Matrix (*DM*) of an assignment score data set capture the Theory vs. Programming Disparity such that the *BPI* values are closer to 1.0 for data sets with larger values for the *TPD* metric and vice-versa. Figure 4 displays the 12 Eigenvalues of the running example assignment scores data set of Section 2: the sums of the cosh and sinh functions of the Eigenvalues are 33.4068 and 12.2509 respectively, leading to a *BPI* of $33.4068 / (33.4068 + 12.2509) = 0.73$.

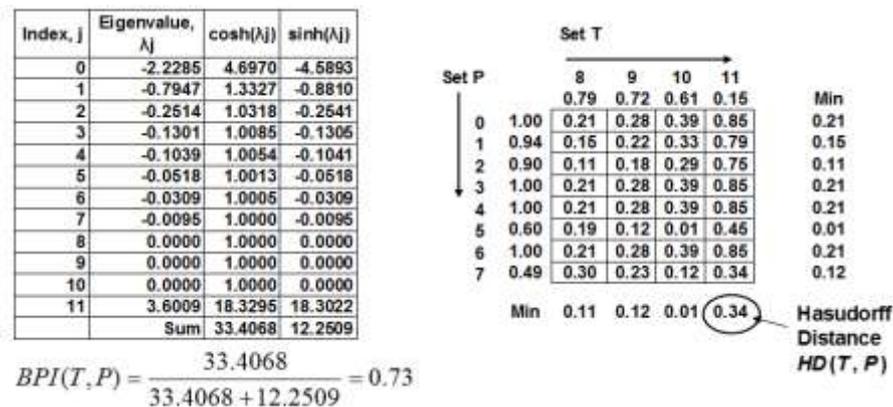


Figure 4. Computation of the Bipartivity Index (*BPI*) and Hausdorff Distance (*HD*) Metric Values for the Sample Data Set of Figure 1

(ii) Hausdorff Distance: The Hausdorff Distance (*HD*) metric [5] has been traditionally used to quantify how far are two data sets in a particular metric space. In the context of quantifying the Theoretical vs. Programming Disparity, we propose to compute the Hausdorff Distance (see below for the formulation) between the decimal percentage scores (referred to as data points) in the sets of theoretical assignments (*T*) vs. programming assignments (*P*). For every data point in data set *T* (and likewise, *P*), we determine the closest distance (in our case, the absolute difference) to a data point in the other data set *P* (*T*). The Hausdorff Distance for the *P* vs. *T* scores for a student data set is the maximum of the closest distances determined as mentioned above. The Hausdorff Distance will be thus smaller if every data point in one data set is closer to some data point in the other data set.

$$HD(T, P) = \text{Max} \left\{ \text{Min}_{i \in T} \left[|T(i) - P(j)| \right]_{j \in P}, \text{Min}_{j \in P} \left[|P(j) - T(i)| \right]_{i \in T} \right\}$$

For *T* vs. *P* data sets that exhibit larger disparity, we expect several data points (assignment scores) in one data set to be appreciably different from those of the other data set. However, even if there exists one outlier data point in either of the two data sets (that is farther away from the other data points in the two data sets), the Hausdorff Distance metric has the vulnerability to get larger and not be an accurate reflection of the closeness or the extent of disparity among the assignment scores in the two sets *T* and *P*. Figure 4 displays the computation of the Hausdorff Distance metric values for the running example *P* vs. *T* data set of Section 2: we observe the presence of a lower theoretical assignment score (0.15 corresponding to index 11 in set *T*) contributes to a relatively larger $HD(T, P)$ value of 0.34 (note that the next largest value among all the minimum values computed across the two data sets is 0.21).

Figure 5 presents a real-time data set comprising of the assignment scores (13 programming assignments and 4 theoretical assignments) for 17 students of the CSC 228 Data Structures and Algorithms course taught in Spring 2020 at Jackson State University, MS. Figure 6 presents the values for the *TPD*, *BPI* and *HD* metrics (also visually compared using a heat map [6]) obtained for the data set of 17 students as well as plots the *TPD* vs. *BPI* and *TPD* vs. *HD* distributions. In the heat map shown in Figure 6, the red, yellow/orange and green colors are respectively indicators of high, moderate and lower values for the *TPD*, *BPI* and *HD* metrics (all of which can be represented in a scale of 0 to 1). We observe the *BPI* and *HD* metrics to have a tendency of over rating (too much red cells for the *BPI* metric) and under rating (too much green cells for the *HD* metric) the theoretical vs. programming disparity. On the other hand, the values for the *TPD* metric are within a moderate range of 0.54 to 0.77 that is sufficient enough to distinguish students with respect to the theoretical vs. programming disparity.

The tendency of the *BPI* and *HD* metrics to respectively over rate and under rate the theoretical vs. programming disparity can be observed in the case of students 10, 12 and 13, who all have just one or two submissions in one of the two categories (programming or theoretical). For such students, the majority of the entries are 0.0. The *BPI* metric tends to cluster the assignment scores for each of these students to two sets: a set of assignments for which submission has been made and a set for which no submission has been made, with the edge weights capturing the difference between these two sets. On the other hand, the *HD* metric tends to give more weight to the minimal difference in the scores between any two assignments. Both the *BPI* and *HD* metrics tend to treat all the 17 assignments as one set (i.e., one dimension) and tends to partition to two clusters (*BPI*) or find the minimum score between any two assignments; whereas, the *TPD* approach considers the problem in a two-dimensional perspective of theoretical vs. programming assignments; the clustering in the two-dimensional space gives leverage to consider four possible combinations for two clusters: the number of theoretical assignments in the sets *X* and *Y* as well

as the number of programming assignments in the sets X and Y identified using spectral analysis as well as makes use of the actual number of theoretical and programming assignments to compute the Euclidean distances as formulated in Section 2. Due to the different approaches taken, we observe only a weak-moderate correlation between the TPD vs. BPI scores and the TPD vs. HD scores for the data set analyzed in Figures 5 and 6.

Student / Index	Programming Assignments												Theoretical Assignments				
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
0	0.00	0.00	0.00	0.00	0.80	0.65	0.00	0.00	0.30	0.00	1.00	0.00	0.25	0.00	0.25	0.00	0.56
1	1.00	1.00	0.80	0.95	1.00	1.00	1.00	1.00	1.00	0.98	1.00	0.92	1.00	0.96	1.00	0.96	1.00
2	0.00	0.00	0.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	1.00	0.80	0.50	0.75	0.00	0.82	0.78
3	0.25	0.30	0.25	1.00	0.75	0.30	0.25	0.35	0.55	0.00	0.00	0.00	0.00	0.30	0.00	0.00	0.00
4	0.00	0.00	0.00	0.45	0.35	0.60	0.00	0.90	0.70	0.00	0.25	0.00	0.00	0.00	0.25	0.00	0.71
5	0.00	0.20	1.00	0.80	0.30	0.65	0.00	0.00	0.00	0.90	0.85	0.50	0.00	0.00	0.00	0.48	0.59
6	0.75	0.00	0.60	0.82	0.75	0.75	0.78	0.90	0.55	0.00	1.00	1.00	0.30	0.93	0.35	0.43	0.90
7	0.00	0.85	0.50	0.85	1.00	0.90	1.00	0.90	0.85	0.65	1.00	1.00	0.40	0.62	0.50	0.94	0.98
8	0.00	0.00	0.40	0.80	0.00	0.10	0.10	0.25	0.70	0.10	0.30	0.50	0.30	0.28	0.20	0.33	0.85
9	0.60	0.00	0.15	0.85	0.50	0.50	1.00	0.90	0.30	0.10	0.85	1.00	0.50	0.76	0.28	0.42	0.65
10	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.30	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
11	0.00	0.00	0.00	0.80	0.00	0.00	0.75	0.90	0.50	0.60	0.00	1.00	0.00	0.00	0.00	0.41	0.51
12	0.00	0.00	0.00	0.00	0.00	0.00	0.75	0.00	0.00	0.00	0.00	0.00	0.00	0.05	0.00	0.00	0.00
13	0.00	0.00	0.00	0.40	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.54	0.00	0.00	0.04
14	0.25	0.15	0.00	0.20	0.00	0.10	0.00	0.90	0.30	0.10	0.00	1.00	0.00	0.25	0.30	0.19	0.29
15	1.00	0.50	0.60	0.90	1.00	1.00	1.00	0.90	0.85	0.85	1.00	1.00	0.55	0.52	0.70	1.00	0.97
16	0.25	1.00	1.00	0.95	1.00	0.78	0.60	1.00	1.00	0.65	1.00	1.00	0.75	0.98	1.00	0.82	0.74

Figure 5. Real-time Data Set used to Apply and Evaluate the Proposed Approach

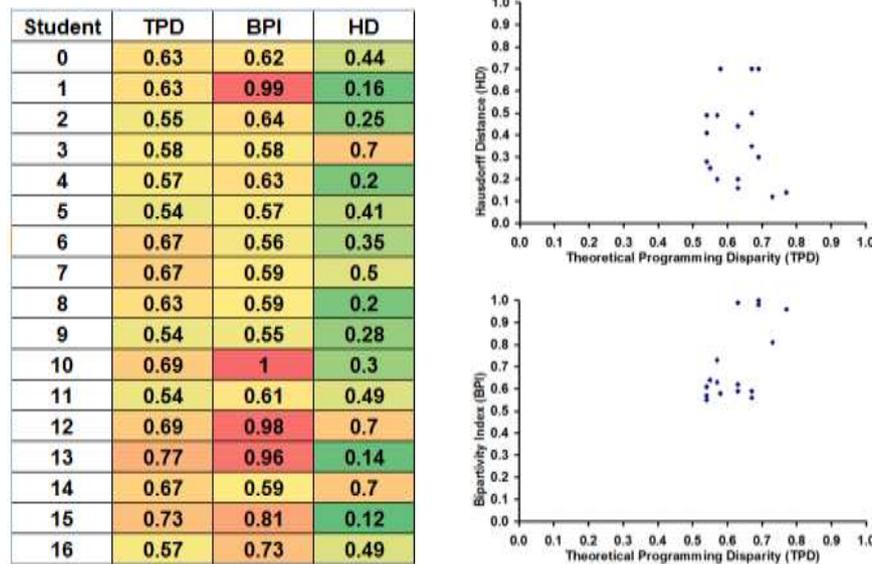


Figure 6. Comparison the Quantitative Assessments of Theory vs. Programming Disparity using the Proposed Theoretical Programming Disparity (TPD) Metric vs. the Bipartivity Index (BPI) and Hausdorff Distance (HD) Metrics

4. PRINCIPAL COMPONENT ANALYSIS-BASED APPROACH TO QUANTIFY THEORY VS. PROGRAMMING DISPARITY FOR AN ENTIRE CLASS

In this section, we propose a principal component analysis (PCA) [22]-based approach to quantify the theoretical vs. programming disparity for an entire class. Such an analysis will be useful to come up with measures for continuous improvement in the assessment cycles of the degree programs in academic institutions. PCA is a classical dimensionality reduction technique used in data analytics research wherein the dataset comprises of a number of features that exhibit some level of correlation between each other. Our analysis is based on the hypothesis that the

performance of the students in the assignments within each category (theoretical vs. programming) have some level of correlation among themselves. We seek to quantify whether there is correlation in student performance (scores) across the two categories of assignments for the entire class. A negative correlation in the student scores across the categories of assignments is an indication of theoretical vs. programming disparity for the entire class. We will first give an outline of the PCA procedure, then explain the procedure to quantify the TPD value for an entire class and finally illustrate the execution of these procedures using two example datasets.

4.1. PCA Procedure

Given a raw dataset with m features and S records, we execute the following steps in this order:

(1) We first standardize the dataset using the mean and standard deviation of the values for each feature in the raw dataset so that the mean and standard deviation of the values for each feature in the standardized dataset (denoted 'Z' in this discussion) are 0 and 1 respectively.

(2) We determine the Covariance Matrix of Z by multiplying the transpose of Z with Z (i.e., $Z^T * Z$).

(3) We determine the Eigenvalues and Eigenvectors of the Covariance Matrix of Z. There will be m Eigenvalues and m Eigenvectors (an Eigenvector corresponding to each Eigenvalue). Each Eigenvector will have m entries, corresponding to the m features of the dataset. The Eigenvector corresponding to the largest Eigenvalue is called the principal Eigenvector.

(4) We determine the m principal components by multiplying the standardized dataset Z of dimension $S \times m$ with the Eigenvectors of dimension $m \times 1$. The principal components will be of dimensions $S \times 1$, with the entries corresponding to the S records.

(5) We determine the variances of the entries in each principal component and rank the principal components in the decreasing order of these variances. We refer to a subset (the top n principal components, in the decreasing order of their variances) of the m principal components as dominating principal components such that the sum of their variances is at least 80% of the sum of the variances of all the m principal components.

The n dominating principal components are considered to maximally capture the variances among the values of the m features in the original dataset with a reduced number of dimensions (we expect $n \ll m$).

4.2. Quantifying TPD using the Principal Components

Let there be a standardized dataset with m_T and m_P theoretical and programming assignments respectively. Let T-PC₁, T-PC₂, ..., T-PC_X be the X dominating principal components generated based on the scores in the theoretical assignments ($X \ll m_T$) and P-PC₁, P-PC₂, ..., P-PC_Y be the Y dominating principal components generated based on the scores in the programming assignments ($Y \ll m_P$) in the dataset. Note that X need not be equal to Y . Let σ -T₁, σ -T₂, ..., σ -T_X be the fractions of the variances of the X dominating principal components for the theoretical assignments vis-a-vis the sum of the variances of the m_T principal components corresponding to the theoretical assignments. Likewise, let σ -P₁, σ -P₂, ..., σ -P_Y be the fractions of the variances of the Y dominating principal components for the programming assignments vis-a-vis the sum of the variances of the m_P principal components for the programming assignments. We determine a $X \times Y$ correlation matrix wherein an entry $(i, j; 1 \leq i \leq X$ and $1 \leq j \leq Y)$ corresponds the Pearson's correlation coefficient between the i^{th} dominating principal component for the theoretical assignments (i.e., T-PC _{i}) and the j^{th} dominating principal component for the programming assignments (i.e., P-PC _{j}). The TPD metric for the entire class is then computed as the weighted average of the $X * Y$ Pearson's correlation coefficients, wherein the weight for an entry $(i, j; 1 \leq i \leq X$ and $1 \leq j \leq Y)$ is the product of the fractions of the variances σ -T _{i} and σ -P _{j} .

4.3. Example 1

In this section, we use a dataset wherein there is a positive correlation (i.e., there is no disparity) among student scores across the theoretical and programming assignments. The dataset (see Figure 7) represents students scores (percentages) in 13 programming assignments (identified as P-1, P-2, ..., P-13 in Figure 7) and 4 theoretical assignments (identified as T-1, T-2, T-3 and T-4 in Figure 7) in the CSC 228 Data Structures and Algorithms class offered in Spring 2020 at Jackson State University, MS, USA. Figure 8 presents the dominating principal components for the theoretical assignments (T-PC1, T-PC2) and programming assignments (P-PC1, P-PC2, P-P-PC3 and P-PC4), their fractions of the variances vis-a-vis the sum of the variances of the principal components for the theoretical and programming assignments as well as the Pearson's correlation matrix between the entries in the dominating principal components for the theoretical vs. programming assignments. Figure 8 also displays the weights for the entries in the Pearson's correlation matrix that are computed as the products of the fractions of the variances of the corresponding dominating principal components of the theoretical vs. programming assignments.

Student	P-1	P-2	P-3	P-4	P-5	P-6	P-7	P-8	P-9	P-10	P-11	P-12	P-13	T-1	T-2	T-3	T-4
0	96.3	100	100	90	100	100	100	80	35	100	100	100	73.8	74	77	92	79
1	100.0	100	100	100	100	100	100	100	100	100	100	100	86.2	97.5	95	94	99
2	41.3	50	0	40	65	80	0	90	82	100	75	100	27.7	19	52	89	82
3	12.5	0	0	40	65	80	100	90	82	100	50	100	40.0	10	50	61	77
4	35.0	0	0	80	0	100	0	0	0	0	0	0	0.0	20	0	0	0
5	37.5	0	0	75	0	100	0	100	78	90	15	0	21.5	49	7	0	0
6	7.5	0	0	85	25	100	5	40	85	100	25	100	29.2	35	44	54	18
7	40.0	100	100	100	88	100	75	90	35	100	100	100	63.1	66	92	94	61
8	83.8	100	100	75	45	100	82	90	80	100	100	100	83.1	66	64	73	26
9	62.5	100	25	100	80	100	50	98	85	100	50	100	43.1	76	76	97	30
10	65.0	0	0	80	0	100	0	95	85	0	5	0	29.2	38	0	22	0
11	42.5	87	15	75	0	60	0	88	85	0	50	100	29.2	41	47	63	29
12	53.8	80	50	80	25	100	50	100	20	100	25	100	43.1	25	62	31	53
13	31.3	100	100	100	0	100	65	90	100	100	100	100	86.2	79	83	95	100
14	100.0	100	100	100	100	60	90	100	100	100	100	100	86.2	100	100	100	100
15	55.0	100	80	85	95	60	90	80	82	90	50	100	49.2	44	62	83	52
16	91.3	100	85	100	0	60	50	80	100	100	60	100	61.5	68	73	100	61
17	81.3	100	80	90	92	100	95	100	100	100	100	100	70.8	57	87	93	100

Figure 7. Raw Dataset 1 used for Principal Component Analysis (PCA)

Based on the results presented in Figure 8, the TPD for the entire class is computed as the weighted average of the entries in the Pearson's correlation matrix of the dominating theoretical (T) vs. programming (P) principal components. The calculation is shown below.

$$TPD = \frac{(0.3699 * 0.8638) + (0.0905 * 0.3961) + (0.0747 * 0.2862) + (0.0548 * 0.0520) + (0.0678 * -0.2280) + (0.0163 * 0.4601) + (0.0134 * 0.0707) + (0.0099 * 0.0980)}{(0.3699 + 0.0905 + 0.0747 + 0.0548 + 0.0678 + 0.0163 + 0.0134 + 0.0099)} = 0.5357$$

We can observe the TPD value of 0.5357 to be more than half of the Pearson's correlation coefficient (0.8638) between the topmost dominating principal components of the theoretical and programming assignments (i.e., T-PC1 and P-PC1). The positive TPD value of 0.5357 indicates that there is no disparity in student performance in the theoretical vs. programming assignments and indicates that students who score high (low) in the theoretical assignments are likely to score high (low) in the programming assignments as well.

Dominating Principal Components for the Theoretical (T) and Programming (P) Assignments

Student	P-PC1	P-PC2	P-PC3	P-PC4	T-PC1	T-PC2
0	2.4071	1.5440	-1.0642	-0.5320	1.7545	-0.5638
1	3.2003	0.6424	0.1118	0.7880	2.0600	-0.1721
2	-1.2737	-2.4055	-0.8369	-0.2570	-2.1509	-1.5271
3	-1.1729	-2.6344	-1.8444	-0.2202	-3.2044	-0.8096
4	-5.4890	2.7009	-0.3412	-1.5644	-1.8571	0.6241
5	-3.2150	-0.2510	0.1697	1.9476	-1.9394	0.3533
6	-2.6534	-0.3999	-1.0232	-0.0492	-2.1804	1.3260
7	1.6700	1.2051	-1.2855	-0.3640	1.0868	0.8625
8	2.0825	0.3204	-0.2659	0.4500	1.1833	-1.0313
9	0.5866	0.0754	-0.1477	0.5715	0.5525	0.5828
10	-3.5313	0.5500	1.6400	1.6316	-1.3705	0.1068
11	-1.7648	-1.1543	2.1424	-1.2176	-0.6581	-0.1390
12	-0.5335	0.6261	-0.9517	-0.0141	-0.0145	-0.1257
13	1.5868	0.3339	0.2544	0.9141	0.9457	1.0125
14	3.2493	-0.3402	1.2515	-0.5914	2.0600	-0.1721
15	1.1859	-0.7151	0.3203	-1.5490	0.7220	-0.0504
16	1.1940	-0.1653	2.1051	-0.6486	1.7391	-0.0005
17	2.4709	0.0677	-0.2345	0.7045	1.2714	-0.2762

Fractions of the Variances of the Dominating Principal Components for the Theoretical (T) and Programming (P) Assignments

σ .P1	σ .P2	σ .P3	σ .P4	σ .T1	σ .T2
0.5063	0.1239	0.1022	0.0750	0.7306	0.1314

Weights for the T vs. P Pearson's Correlation Matrix

	σ .P1	σ .P2	σ .P3	σ .P4
σ .T1	0.3699	0.0905	0.0747	0.0548
σ .T2	0.0678	0.0163	0.0134	0.0099

Pearson's Correlation Matrix between the Dominating Principal Components

	P-PC1	P-PC2	P-PC3	P-PC4
T-PC1	0.8638	0.3961	0.2862	-0.0520
T-PC2	-0.228	0.4601	0.0707	0.0980

Figure 8. Dominating Principal Components for the Theoretical (T) and Programming (P) Assignments of Dataset 1, Fractions of the Variances, Pearson's Correlation Matrix of the Dominating T vs. P Principal Components and the Weights of the Entries in the Matrix

4.4. Example 2

We now present the TPD calculations for another dataset (CSC 323 Algorithm Design and Analysis course, Fall 2021 at Jackson State University, MS, USA) wherein we indeed observe disparity in student performance between the theoretical (T) and programming (P) assignments. Figure 9 presents the raw data and Figure 10 presents the dominating principal components of the T vs. P assignments, their fractions of the variances and the Pearson's correlation matrix as well as the weights of the entries in this matrix.

Student	P-1	P-2	P-3	P-4	P-5	P-6	P-7	P-8	T-1	T-2	T-3	T-4
0	100	100	92	75	25	100	100	78	85	34.5	97	68.5
1	90	100	95	75	25	85	85	76	78	42.5	51	55
2	85	100	95	75	25	85	85	50	81	78	67	53.5
3	100	100	95	75	25	85	100	68	90	85	99	79
4	88	100	95	75	25	85	85	72	69	71	48.5	60
5	55	100	72	75	25	100	20	50	55.5	41	14	29
6	100	100	100	88	100	85	85	64	75	74.5	43.5	50.5
7	80	25	100	88	25	100	100	58	96	70.5	95	94
8	15	100	77	45	25	85	20	76	47	32	8	41
9	50	100	95	75	25	85	60	56	83	87.5	28	84
10	90	100	91	75	25	85	85	50	81	53	33	44.5

Figure 9. Raw Dataset 2 used for Principal Component Analysis (PCA)

Based on the results presented in Figure 10, the TPD for the entire class is computed as follows:

$$TPD = \frac{(0.3139 * -0.7820) + (0.1570 * -0.1470) + (0.1029 * 0.2573) + (0.0735 * -0.2630) + (0.0789 * 0.0616) + (0.0394 * 0.2502) + (0.0259 * -0.6632) + (0.0185 * -0.3013)}{(0.3139 + 0.1570 + 0.1029 + 0.0735 + 0.0789 + 0.0394 + 0.0259 + 0.0185)} = -0.3326$$

We can observe the TPD value of -0.3326 to be less than (but reasonably close to) half of the Pearson's correlation coefficient (-0.7820) between the topmost dominating principal components of the theoretical and programming assignments (i.e., T-PC1 and P-PC1). However, a majority of the entries in the Pearson's correlation coefficients matrix are negative and hence the weighted average of these coefficients (i.e., the TPD value) remains negative as well. The negative TPD value of -0.3326 indicates that there is noticeable disparity in student performance in the theoretical vs. programming assignments and indicates that students who score high (low) in the theoretical assignments are likely to score low (high) in the programming assignments.

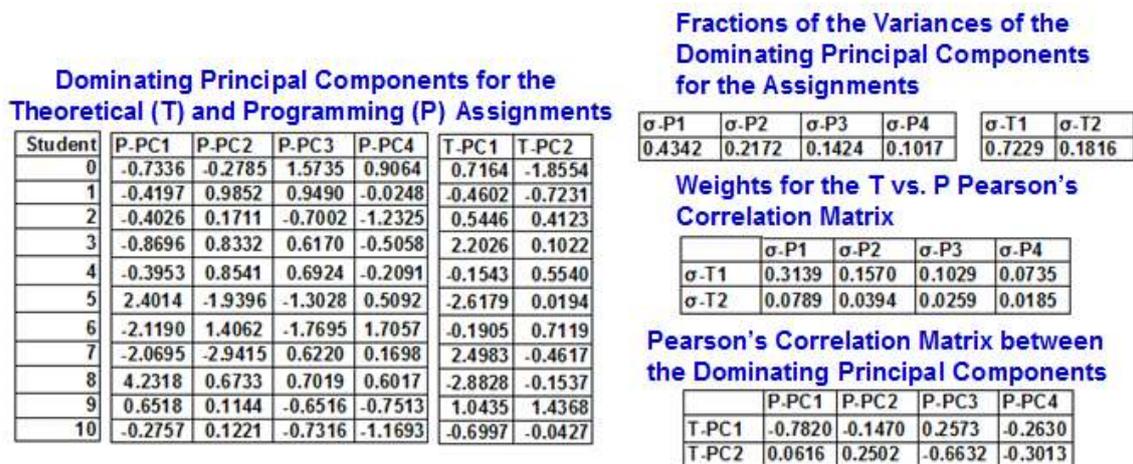


Figure 10. Dominating Principal Components for the Theoretical (T) and Programming (P) Assignments of Dataset 2, Fractions of the Variances, Pearson's Correlation Matrix of the Dominating T vs. P Principal Components and the Weights of the Entries in the Matrix

5. RELATED WORK

To the best of our knowledge, we have not come across any work in the literature that focuses on assessing and quantifying the disparity found between two different categories of assignments on a per-student basis. Disparity studies in academic settings have been so far mainly focused on gender [12] and race [13] as well as on the class, as a whole (e.g., [9-11]), and not on a per-student basis. The closest work we have come across related to our topic is the work of [14] wherein the authors apply principles from phenomenography [15] and variation theory [16] to explore the practices that are needed to bridge the gap in learning Computer Science theory and learning Computer Science practice (programming). But, there are no efforts to quantify the extent of the gap (or the disparity), as is done in our paper.

Below is a review of the works that we came across in the literature that focus on studies conducted to assess the contributing factors to the success or hardship for students majoring in Computer Science and the programming component of it. In [17], the authors did a survey to find out that students think Computer Science-ability is something both innate as well as extensible through effort. In [18], the authors surveyed Computer Science (CS) student performance data for 10 years and conclude that a successful CS student needs to be strong both in critical thinking skills and the core CS skills. They observed the critical thinking skills for a CS student typically come from the Math and Physics courses and concluded that these courses need to be enforced as pre-requisites for CS courses early on in the curriculum instead of being taken along with CS courses. However, no analysis has been reported in [18] regarding the skills that influence the theory vs. programming disparity found among CS students. In [19], the authors report there is no statistically significant influence of the assessment mode (programming in a computer vs. writing

a program in pen by hand) on the performance of students in a programming course. In [20], the authors observed that novice programmers tend to program using problem-solving skills obtained from domains familiar to them. In [21], the authors used reflective essays-based Attribution Theory to elucidate the internal and external causes that influence the performance of students in their first programming course.

The following works focus on analyzing the impact of one examination format on another. In [9], the authors build a model to predict the performance of students on the basis of examination formats: whether or not the performance of students in practical examinations can be predicted using their performance in the standardized examinations? The answer reported in [9] is No, as the two examination formats are observed to test different skill sets. Likewise, Haberyan [10] found no correlation between performance in weekly quizzes and examinations among students majoring in Biology. However, in [11], the authors observed that psychology undergraduates performed well in the examinations when they were also given weekly reading assignment-based quizzes throughout the course.

6. CONCLUSIONS

The high-level contribution of this paper is a spectral bipartivity analysis-based approach to quantify the disparity (the proposed metric is referred to as Theoretical vs. Programming Disparity: *TPD* metric) in the scores earned by students in two categories of assignments: theoretical and programming as well as a principal component analysis (PCA)-based approach to quantify the *TPD* metric for an entire class. The uniqueness of the spectral bipartivity analysis-based approach is that it quantifies the disparity on a per-student basis; the typical approach in the academic community so far is to use quantitative metrics that capture the disparity for an entire class as a whole [8]. Also, traditionally for such problems, the correlation measures [3, 8, 9] are used to quantify the extent of influence on one category of assignments over another. But, to use correlation measures, we need to have the same number of assignments under both the categories. With both the spectral bipartivity analysis and principal component analysis-based approaches, there can be a different number of assignments for the two categories. The pre-requisites for the spectral analysis-based approach are just the need to input the classification of the assignments ids (as either theoretical or programming) and the actual number of assignments under the two categories. We have demonstrated the characteristics, uniqueness and effectiveness of both the spectral bipartivity analysis and the PCA-based approaches through exhaustive evaluations with real-time datasets as well as through comparisons with quantitative metrics that appear to be compelling enough to be suitable to capture the disparity in the student scores in the assignment categories. The disparity problem on a per-student basis among two data sets of uneven size, especially in an academic setting, has not been so far considered in the literature for quantitative evaluation; likewise, there is no prior work that does dimensionality reduction and uses the principal components for the scores in theoretical and programming assignments to quantify the *TPD* metric for an entire class. We expect the proposed computation approaches to quantify *TPD* for a student and for an entire class would be valuable to both academicians and researchers.

ACKNOWLEDGEMENTS

The author would like to thank everyone, just everyone!

REFERENCES

- [1] Estrada, E, (2010) "Structural Patterns in Complex Networks through Spectral Analysis," *Proceedings of the 2010 Joint IAPR International Conference on Structural, Syntactic, and Statistical Pattern Recognition*, pp. 45-59. Springer-Verlag, Cesme Izmir, Turkey.
- [2] Sarkar, C., and Jalan, (2018) "Spectral Properties of Complex Networks," *Chaos: An Interdisciplinary Journal of Nonlinear Science*, vol. 28, no. 10, 102101.
- [3] Strang, G (2019) *Linear Algebra and Learning from Data*, 1st edition, Wellesley-Cambridge Press, Wellesley, MA, USA.
- [4] Ernada, E., and Rodriguez-Velazquez, J. A (2005) "Spectral Measures of Bipartivity in Complex Network," *Physical Review E*, vol. 72, no. 4, 2, 046105.
- [5] Birsan, T., and Tiba, D (2006) "One Hundred Years since the Introduction of the Set Distance by Dimitrie Pompeiu," *Proceedings of the IFIP Conference on System Modeling and Optimization*, vol. 199, pp. 35-39. Springer, Turin, Italy.
- [6] Wilkinson, L., and Friendly, M (2009) "The History of the Cluster Heat Map," *The American Statistician*, vol. 63, no. 2, pp. 179-184.
- [7] Godsil, C., and Royle, G. F (2013) *Algebraic Graph Theory*, 1st edition, Springer, Berlin, Germany.
- [8] Caven, M (2019) "Quantification, Inequality, and the Contestation of School Closures in Philadelphia," *Sociology of Education*, vol. 92, no. 1, pp. 21-40.
- [9] Davison, C. B., and Dustova, G (2017) "A Quantitative Assessment of Student Performance and Examination Format," *Journal of Instructional Pedagogies*, vol. 18, pp. 1-10.
- [10] Haberyan, K. (2003) "Do Weekly Quizzes Improve Student Performance on General Biology Exams?" *The American Biology Teacher*, vol. 65, pp. 110-114.
- [11] Johnson, B. C., and Kiviniemi, M. T (2009) "The Effect of Online Chapter Quizzes on Exam Performance in an Undergraduate Social Psychology Course," *Teaching of Psychology*, vol. 36, no. 1, pp. 33-37.
- [12] Master, A., Meltzoff, A. N., and Cheryan, S (2021) "Gender Stereotypes about Interests Start Early and Cause Gender Disparities in Computer Science and Engineering," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 118, no. 48, e2100030118.
- [13] Kozlowski, D., Lariviere, V., Sugimoto, C. R., and Monroe-White, T (2002) "Intersectional Inequalities in Science," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 119, no. 2, e2113067119.
- [14] Thune, M., and Eckerdal, A (2019) "Analysis of Students' Learning of Computer Programming in a Computer Laboratory Context," *European Journal of Engineering Education*, vol. 44, no. 5, pp. 769-786.
- [15] Farton, M (1986) "Phenomenography - A Research Approach Investigating Different Understandings of Reality," *Journal of Thought*, vol. 21, no. 2, pp. 28-49.
- [16] Bussey, T. J., Orgill, M., and Crippen, K. J (2013) "Variation Theory: A Theory of Learning and a Useful Theoretical Framework for Chemical Education Research," *Chemical Education Research Practice*, vol. 14, pp. 9-22.
- [17] Lewis, C. M., Yasuhara, K., and Anderson, R. E (2011). "Deciding to Major in Computer Science: A Grounded Theory of Students' Self-Assessment of Ability," *Proceedings of the 7th International Workshop on Computing Education Research*, pp. 3-10, ACM, Providence, RI, USA.
- [18] Yang, H., Olson, T. W., and Puder, A (2021) "Analyzing Computer Science Students' Performance Data to Identify Impactful Curricular Changes," *Proceedings of the IEEE Frontiers in Education Conference*, pp. 1-9, IEEE, Lincoln, NE, USA.
- [19] Oqvist, M., and Nouri, J (2018) "Coding by Hand or on the Computer? Evaluating the Effect of Assessment Mode on Performance of Students Learning Programming," *Journal of Computers in Education*, vol. 5, pp. 199-219.
- [20] Wellons, J., and Johnson, J (2011) "A Grounded Theory Analysis of Introductory Computer Science Pedagogy," *Systemics, Cybernetics and Informatics*, vol. 9, no. 6, pp. 9-14.
- [21] Vivian, R., Falkner, K., and Falkner, N (2013) "Computer Science Students' Casual Attributions for Successful and Unsuccessful Outcomes in Programming Assignments," *Proceedings of the 13th Koli Calling International Conference on Computing Education Research*, pp. 125-134, ACM, Koli, Finland.
- [22] Jolliffe, I. T (2002) *Principal Component Analysis*, Springer Series in Statistics, New York, USA.

AUTHORS

Natarajan Meghanathan is a tenured Full Professor of Computer Science at Jackson State University, MS, USA. He graduated with a PhD in Computer Science from The University of Texas at Dallas in 2005. His primary areas of research interests are Network Science, Machine Learning and Cyber Security. He has authored more than 175 peer-reviewed research publications in these areas.

