# Data-Driven Peer Group Selection for Salary Comparison in Higher Education: An Applied Analytics Approach to Building Trust

Eric Breimer
ebreimer@siena.edu

Kim Sangahn
skim@siena.edu

Seung Jin Wang
swang@siena.edu

## Abstract

This study introduces a data-driven method to select peer institutions in higher education for faculty salary comparison. Given a target institution, the goal is to form a peer group of similar colleges using stakeholder-identified variables/features like enrollment, finances, and student outcomes, but excluding salary data. An effective peer group places the target near the median salary level. Previous work raised equity concerns because the methodology generated separate peer groups, including one for base salaries and several for high-demand accredited disciplines. Concerns about objectivity and fairness emerged due to the use of subjective filters and post-hoc adjustments, such as including aspirational institutions. We seek a more consistent, data-driven approach that uses principal component analysis and nearest neighbor search to create a single unified peer group that can be used for all salary comparisons. By employing a more transparent, analytics-based method, we aim to enhance trust in the process to promote acceptance of the peer group among faculty and administrative stakeholders.

**Keywords:** Data-driven methodology, peer institution selection, salary benchmarking, higher education compensation, principal component analysis, nearest neighbor analysis.

# Data-Driven Peer Group Selection for Salary Comparison in Higher Education: An Applied Analytics Approach to Building Trust

*Eric Breimer. Lo, Samgahn and Seung Jim Wang*

## 1. INTRODUCTION

In higher education, equitable compensation for fulltime faculty is a cornerstone of institutional stability and morale. Compensation models often rely on benchmarking against peer institutions to ensure comparability, accounting for variations across disciplines where market forces may drive salary differentials. Traditionally, institutions have employed multiple peer groups tailored to specific purposes: a base group for general salaries and specialized groups for high-demand accredited fields that often have significantly higher salaries.

A fragmented approach to peer group selection can create perceptions of inequity among faculty and administrative stakeholders. For example, when peer groups for different disciplines vary significantly in composition, disparities in financial metrics may raise concerns about fairness. Additionally, common practices such as ad hoc filtering to exclude problematic peers and the subjective inclusion of aspirational institutions often lack consistent, rigorous criteria, undermining trust in the process. When faculty and administrators lack confidence in the peer group selection, it becomes challenging to accept compensation decisions based on comparisons with those groups.

This paper presents a data-driven methodology termed the Unified Peer Group (UPG), designed to streamline peer selection into a single, consistent framework. The UPG combines the top overall nearest neighbors with the most similar peers that share discipline-specific accreditations with the target institution. This approach enables the entire UPG to guide base salary decisions, while subsets can inform adjustments for accredited high-demand disciplines.

By leveraging Principal Component Analysis (PCA) and nearest neighbor analysis, the UPG identifies institutions most similar to the target institution across a multidimensional space. Salary data are excluded from the peer selection process to ensure impartiality. The objective is to form a peer group where the target institution aligns near the median for multiple variables/features. This approach assumes that the target's salaries will similarly approximate the peer group's median. If this assumption fails, it indicates that the target's salaries deviate from those of comparable peers, potentially justifying compensation adjustments.

Our new methodology improves upon prior work from 2021 and 2024 where subjective filters on Carnegie classification (American Council on Education, 2025), public/private status, and geography were deemed essential.

Empirical evidence shows that the refined, data-driven selection process efficiently excludes unsuitable peers, minimizing the need for ad hoc filters. Although administrative stakeholders recommended retaining one filter, this methodology has significantly strengthened trust in the peer group and decision-making process at the authors' institution.

## 2. BACKGROUND

Peer institution selection is a critical process in institutional research, evolving from subjective, bias-prone methods to sophisticated, data-driven approaches. Early peer selection relied on subjective criteria like geographic proximity or mission alignment, which often introduced inconsistencies (D'Allegro, 2017; D'Allegro & Zhou, 2013). These studies highlight the limitations of such approaches, advocating for objective methodologies using Integrated Postsecondary Education Data System (IPEDS) data. For instance, McLaughlin et al. (2011) proposed nearest neighbor algorithms to form peer groups based on key institutional metrics such as enrollment, finances, and student outcomes, offering a reproducible framework that minimizes bias.

To enhance the precision of peer selection, advanced analytical techniques like Principal Component Analysis (PCA) have gained prominence. PCA reduces correlated variables into uncorrelated principal components, capturing essential data variance while simplifying complex datasets (Jolliffe & Cadima, 2016). When integrated with nearest neighbor algorithms, PCA improves classification accuracy, as demonstrated in educational and non-educational

contexts (Lubis et al., 2020; McLaughlin et al., 2011). This synergy of PCA and nearest neighbor methods provides a robust foundation for equitable peer comparisons, particularly in contexts like salary benchmarking, where fairness and transparency are paramount.

In higher education, peer benchmarking informs critical policy decisions, such as funding and compensation strategies (Kelchen et al., 2024). However, existing multi-group models often fail to account for contextual factors like accreditation, which can significantly influence institutional profiles (AACSB, 2025; CCNE, 2025). Recent analytics applications in institutional research, such as big data in healthcare and campus crime analysis, underscore the importance of transparent, data-driven frameworks to build trust and ensure equitable comparisons (Mohammed & Lind, 2024; Kline et al., 2020). Despite these advances, gaps remain in integrating categorical factors like accreditation into unified peer selection models.

This study addresses these gaps by proposing a novel framework that combines PCA, nearest neighbor algorithms, and accreditation as a categorical factor. By synthesizing data-driven methodologies (McLaughlin et al., 2011; Lubis et al., 2020) with contextual considerations (AACSB, 2025; CCNE, 2025), this approach aims to enhance the accuracy and equity of salary benchmarking, contributing to more informed resource allocation in higher education.

## 3. Methods

### Data Sources and Variables
In our work, data were downloaded from IPEDS (NCES, 2023) focusing on 2,605 institutions with sufficient reported data (at least 11 of 14 key columns). Missing values were imputed using nearest neighbors following the approach of Troyanskaya et al. (2001). The authors' home institution has two key accreditations that impact salaries, Business (AACSB, 2025) and Nursing (CCNE, 2025). An important stakeholder goal was to include a balanced mix of peers with AACSB accreditation, CCNE accreditation, both accreditations, and neither. The accreditation status of institutions in not included in IPEDS and was scraped directly from the AACSB and CCNE websites.

Stakeholders expressed concerns about deviating significantly from past approaches, emphasizing the need for year-to-year consistency. Our goal was to introduce a new peer selection methodology that would gain broad acceptance without significantly altering the core variable set, which could raise additional concerns. Once a methodology is adopted, future work can explore refinements to variable selection.

The key variables shown below were derived from 14 IPEDS columns (NCES, 2023) and direct accreditation data. These variables seek to capture institutional size (student and faculty counts), financial health, and student success metrics, aligning with past practices. We introduced one new variable to capture key accreditations. We excluded Carnegie classification, geographic location, and institutional type (public vs. private) which were the subjects of subjective ad hoc filtering that previously complicated peer selection.

**FTEGD**: Full-time equivalent graduate students.

**FTEUG**: Full-time equivalent undergraduates

**Revenue**: Total operating revenue

**Endowment**: Value of the endowment

**Net Assets**: Total assets including endowment

**Ret Rate**: Retention rate from year 1 to 2

**Grad Rate**: Four-year graduation rate.

**Adm Rate**: Students enrolled divided by students admitted

**Faculty FTE**: Full-time equivalent faculty, weighted by full/part-time

**Net Price**: Average price after discounts, weighted by graduate/undergraduate

**ACCRED**: Accreditation

An accreditation (ACCRED) value of 1 indicates a school that shares all the key accreditations of the target institution, the value 0 indicates no shared accreditations, and the intermediate values indicate the percentage of shared accreditations. For example, given a target institution with five key accreditations, a peer that shares 4 out of the 5 accreditations would have an ACCRED value of 0.8. Representing all accreditations as one column helps to avoid the over-weighting that might occur when considering many key accreditations stored as separate variables. While this paper focuses on the authors' home institution with two key accreditations, our methodology can scale for institutions with many key accreditations.

### Analytical Process
PCA was applied after normalizing the variables with standard scaling (Jolliffe, 2002). The top 9

components, explaining 99% of the variance, were selected to ensure the ACCRED variable, which had low weight in components 1–8, influences peer selection. Figure 1 shows the cumulative explained variance of the principal components. Figure 2 shows the loadings or weighting of each of the direct variables on the first five principal components (see Appendix A for the full table).
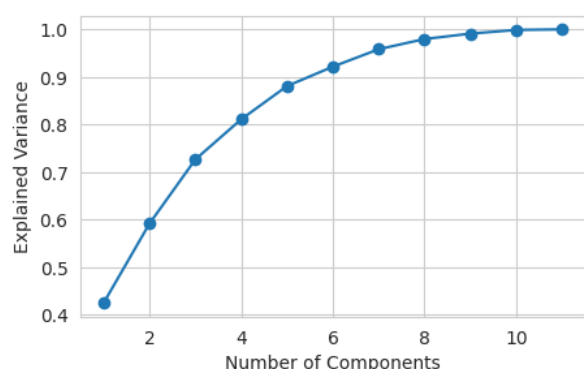


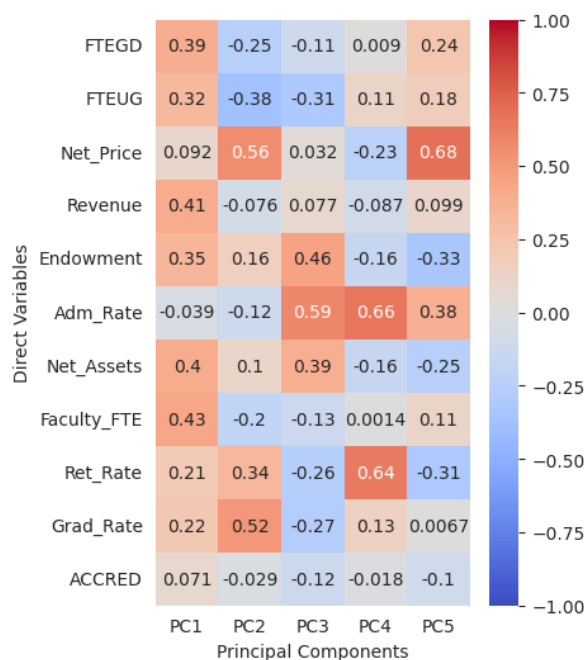**Figure 1: Cumulative explained variance by PCA components**



**Figure 2: Loadings (weighting) of direct variables on first 5 principal components**

Nearest neighbors were computed using Gower distance (Gower, 1927) rather than Euclidean distance, due the presence of the ACCRED variable. Gower distance is considered a good choice given categorical variables. While the ACCRED variable could be considered the percentage of matched accreditations, it exhibits

characteristics of an ordinal categorical value when considering on a few key accreditations. Additionally, IPEDS includes some variables with small value ranges that are rounded, which also have categorical characteristics. Thus, we felt Gower distance was the best choice given the composition of variables.

The Unified Peer Group (UPG) concept combines the top overall nearest neighbors with the top accredited neighbors for each key accreditation. This method ensures sufficient accredited peers to provide robust salary data for high-demand disciplines while maintaining a manageable UPG size for stakeholder review. At the authors' institution, stakeholders deemed a UPG exceeding 50 institutions too large and fewer than 15 insufficient for reliable salary data. Historically, peer groups ranged from 15 to 30 institutions, and significant deviations from this range raised concerns about reduced stakeholder acceptance.

To account for stakeholder concerns, the authors' institution defined the UPG as the union of:

1. Top 30 overall nearest neighbors.
2. Top 15 AACSB-accredited neighbors.
3. Top 15 CCNE-accredited neighbors.
4. Additional non-accredited neighbors to ensure at least 33% of the UPG lacks either accreditation.

Although the specific number of schools in the UPG is not determined through data-driven methods, the size selection aims to align with historical peer groups to enhance stakeholder acceptance. In general, the UPG definition should vary based on stakeholder concerns and constraints at the target institution.

In the general case, it important to note that the intersection of top overall nearest neighbors and accredited neighbors may vary. In one extreme, the top overall peers may include no accredited institutions, but the UPG definition ensures a minimum number of peers for each key accreditation. Conversely, if the top overall peers hold all key accreditations, stakeholders may raise concerns about their over-representation. However, the UPG definition can be adjusted to ensure a minimum number of non-accredited peers.

At the authors' institution, stakeholders recommended that one-third of the peer group consist of schools without either accreditation to reflect historical peer group composition. This proportion can be set to zero for institutions whose peer groups historically consisted entirely

of accredited institutions.

Although the selection of UPG size and composition involves inherent subjectivity, the UPG framework establishes overarching goals rather than ad hoc procedures for specific institution selection. For instance, past practices—such as excluding schools based on Carnegie classification—served as tactical steps to refine the peer list, not as strategic aims for achieving a particular Carnegie composition. Our work emphasizes developing an unbiased, data-driven process for selecting peers within stakeholder-defined parameters, leaving size determination to consensus. While our goal was to implement a fully data-driven approach, stakeholder feedback at the author's institution necessitated one subjective post-hoc filter: the exclusion of doctoral institutions.

Finally, the entire process is implemented in Python and documented in a Google Colab notebook (Google, n.d.). The notebook includes data acquisition (downloading and scraping), cleaning, principal component analysis, nearest neighbor calculations, and supporting visualizations. The notebook serves as an audit trail, enabling stakeholders to thoroughly examine the process.

### 4. Results & Analysis

Applying PCA and nearest neighbor yielded a ranking of the 2,605 institutions that we considered. The distribution of the distances to the authors' home institution (the target) are shown in Appendix B. Appendix C shows the correlation of key variables and the targets position in the distribution of key variables. We repeated the analysis for 20 random targets. This section includes a summary focused on geographic proximity, public vs private, Carnegie class and accreditation.

**Geographic Proximity**
In the past, geographic filters were an important element in peer group selection. Stakeholders felt that schools in distant regions would be poor matches and should be filtered out. However, such filtering suffers from subjectivity, especially bias in defining the boundaries of the target region. For instance, restricting a peer group to Northeastern states might exclude Ohio, even though Ohio may exhibit substantial similarity to the target region. Thus, the exclusion of Ohio-based schools might represent an ineffective filter.

Our results indicated that geographic filtering

may not be necessary. Figure 3 shows the geographic clustering of peer groups from three randomly selected target institutions. The other randomly selected targets exhibited similar clustering where peers tend to be geographically closer to the target institution.
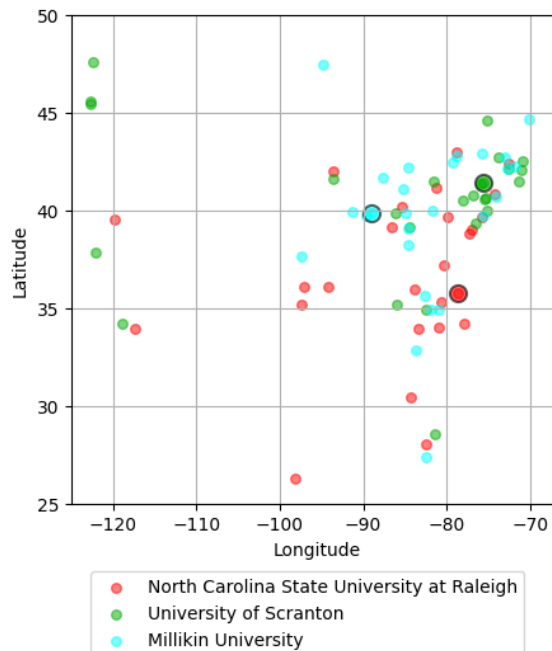


**Figure 3: Geographic distribution of three example peer groups**

Figure 4 shows a breakdown of Scranton University in Pennsylvania (PA) and Millikin University in Illinois (IL). The 10 geographically closest states to the targets are shown in green highlighting the tendency for peers to be in the nearest states.
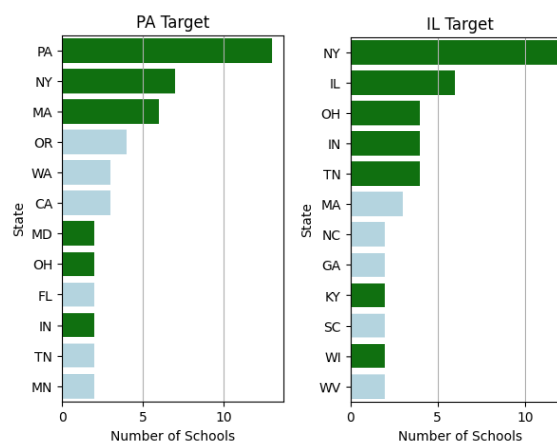


**Figure 4: State distribution of peer groups for target schools in PA and IL**

There are many reasons institutions that are

similar in key variables might be geographically close. Economically prosperous regions can support types of institutions that other regions cannot support. Since many students attend colleges near home (Turley, 2009; Acton, 2024), local schools compete to attract the same cohorts. Thus, nearby schools may converge in student quality similarity.

At the authors' home institution in the Northeast, stakeholders were initially skeptical when schools far outside the region were selected as nearest neighbors. However, after researching these geographic outliers, stakeholder agreed that they were good selections. Stakeholder were also comfortable including a few geographic outliers as long as the majority of peers were within the region. Using Nearest Neighbor without any geographic filtering is an opportunity to discover excellent matches outside of the region. And, the tendency to select peers in the region gives stakeholder confidence in the overall process.

### Public vs Private

The authors' home institution is a private 4-year college and stakeholders felt that filtering out public institutions was essential. However, only 3 public institutions ranked among the top 200 nearest neighbors (ranked #185, 187 and 198, respectively). When key financial variables are included, public and private institutions demonstrate significant difference. This gave stakeholder further confidence in the nearest neighbor approach in selecting appropriate peers.



**Figure 5: Visualization of nearest neighbors of a target public 4-year and private 4-year.**

Figure 5 visualizes the 100 nearest neighbors of a public 4-year institution and a private 4-year institution. For the public 4-year target, 32 out of the top 50 peers were also public 4-year institutions. Other public 4-year targets exhibited similar distributions. Private four-year institutions are the most common among the

2,605 schools considered. As a result, more private schools are available for selection, and public institutions rarely rank among the top peers of a private target institution.

### Carnegie Class

In 2021, the authors' home institution in the Northeast region was reclassified from *Baccalaureate Colleges: Arts & Sciences Focus* to *Master's Colleges & Universities: Smaller Programs*. This reclassification stemmed from earning AACSB accreditation in 2007, CCNE accreditation in 2017, launching a Master of Science in Accounting in 2009, and introducing an MBA program in 2018. This context is very important because some institutions may on the edge between two Carnegie classes or may overlap with two or more classes.

In the past, peer groups were selected by considering institutions that matched the target's current and most recent previous Carnegie classification. This decision was made without investigating the similarity of the target to schools in other Carnegie classes. As a result, this filter yielded very few peers with either AACSB or CCNE accreditation. Afterwards, stakeholder would advocate for the inclusion of accredited aspirants, which was an inherently subjective process.
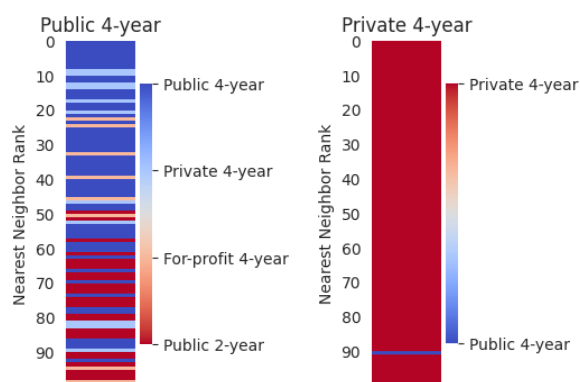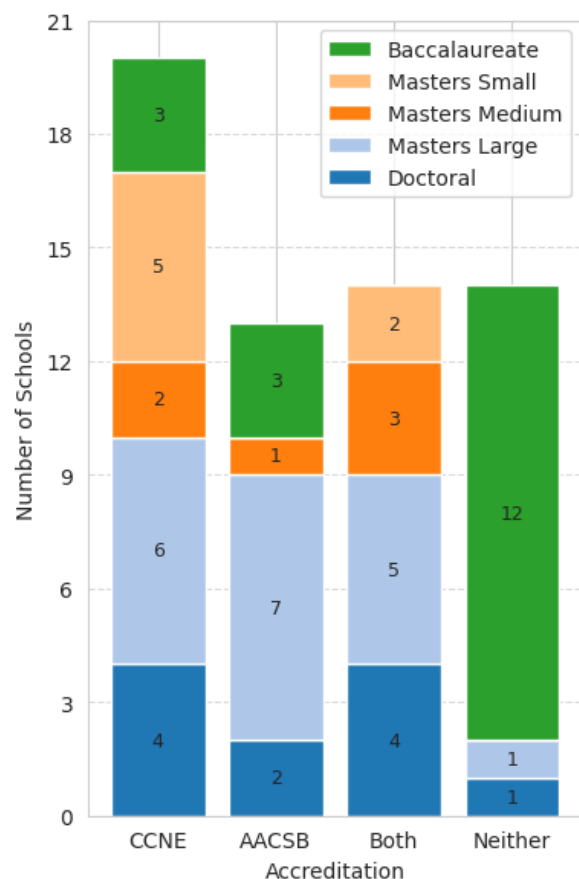
**Figure 6: Carnegie breakdown of the top 60 nearest neighbors to the target.**

Figure 6 shows the top 60 nearest neighbors to the authors' home institution, which includes (a) 20 schools with only CCNE accreditation, (b) 13 schools with only AACSB accreditation, (c) 14 schools with both accreditations, and (d) 14 with neither.

In examining these four groups, it became clear that peers in Carnegie class *Large and Medium Master's Colleges & Universities* should be selected in order to form a group with a sufficient number of accredited peers (at least 15 of each).

Nearest neighbor revealed that many of the closest matched accredited peers were outside of the target's Carnegie class. Examining these peers revealed significant similarities, particularly in overall enrollment and financial variables. Although the Carnegie size classification reflects graduate program size—and the target institution has relatively smaller graduate enrollment—there was strong alignment in tuition-based revenue and weighted cost. Stakeholders ultimately concluded that including medium and large master's institutions was appropriate, especially

given the institution's strategic goal of increasing graduate enrollment.

Carnegie *Doctoral/Professional Universities* (R3), ranked among the top 60 at positions 6, 9, 20, 21, 24, 26, 38, 42, 52, 54, and 60. Further investigation revealed that these nearest neighbors had only a few small doctoral programs. However, since the target institution has no plans to start any doctoral programs, stakeholders felt that peers with doctoral programs should be excluded from the UPG. This was the only post-hoc filter applied, and it was widely accepted by stakeholders.

**The Unified Peer Group**
After excluding doctoral institutions and applying the union criteria, the UPG consisted of 36 institutions with the following characteristics:

- 7 with both accreditations.
- 12 with neither (33% threshold met).
- 16 CCNE total (9 CCNE-only + 7 both).
- 15 AACSB total (8 AACSB-only + 7 both).

Note that the 15th selected AACSB peer is also CCNE-accredited.



**Figure 7: Unified Peer Group (UPG) overlap**

Figure 7 illustrates the overlap among the 24 schools with one or both accreditations. This overlap—including the 7 schools with both accreditations—results in a relatively small peer group that captures sufficient CCNE and AACSB peers for discipline-specific salary comparisons.

To understand the influence of accreditation on selection, we re-ran PCA and nearest neighbor without the ACCRED variable and it yielded a nearly identical UPG of 36 institutions. In this new UPG, the lowest ranked peer with both accreditations was dropped in favor of a better match with only AACSB accreditation. With ACCRED excluded, the 15th-ranked CCNE and AACSB peers were ranked 39th and 43rd respective out of 2,605 total institutions (874 with CCNE and 554 with AACSB). Among the top 30

nearest neighbors were 9 AACSB peers and 13 CCNE peers.

## 5. Discussion & Conclusions

The overall goal of this work was to establish a peer group methodology that could be widely accepted by stakeholders. To achieve this goal, we considered two related objectives. First, we aimed to establish a single unified peer group that could be used to determine base salaries for all faculty, as well as discipline-specific salaries for accredited programs. Second, we sought to develop a data-driven process that eliminated as many subjective filters, ad hoc processes, and post-hoc decisions as possible.

In the past, independent peer groups were developed for discipline-specific salary comparisons, dividing the institution. For instance, the base group and Nursing group were selected with Carnegie filters that excluded most institutions with graduate programs, whereas the Business group, out of necessity, primarily included institutions with medium and large master's programs.

The UPG mitigates institutional division, equity concerns and stakeholder mistrust in two key ways. First, the selection process is consistent for all accredited discipline-specific groups. Second, base salaries are influenced by the inclusion of accredited peers in the UPG.

When a target institution holds multiple key accreditations, a key challenge is preventing the UPG from becoming excessively large. Stakeholders often seek to investigate all peers more deeply to understand each selection, which is not practical with a very large UPG. However, the UPG must include enough peers per accredited discipline to enable robust salary comparisons. The UPG framework allows adjustments to achieve a manageable size, such as modifying the number of peers per accreditation or the proportion of non-accredited peers. Thus, the UPG provides a flexible, scalable framework that can be adjusted to the institution's accreditations and stakeholder needs.

Our approach demonstrates empirically that several contentious subjective filters may be unnecessary. Among the top 200 nearest neighbors, only three were public institutions, and the closest Carnegie R1 university ranked 136th, indicating that nearest neighbor analysis naturally excludes unsuitable peers at top matches. The data-driven process revealed that

geographic and Carnegie classification filters excluded some of the best-matched peers, as stakeholders observed. This approach enabled stakeholders to recognize that excellent matches frequently included institutions beyond the target's geographic region or Carnegie classification.

At the authors' institution, the new methodology increased trust, consistent with literature on analytics adoption (Mohammed & Lind, 2024). Specifically, a unified, data-driven peer selection process—using PCA and nearest neighbors—produced a peer group that was quickly accepted with only one post-hoc adjustment: filtering out doctoral universities. Previously, post-hoc changes, such as including aspirational institutions, required subjective and time-consuming negotiation among faculty and administrative stakeholders.

## 6. Future Work

Our analysis revealed that accredited peers are often selected as nearest neighbors, even when accreditation is excluded as an input variable. However, this finding is based on a single target case, and further analysis is needed to assess accreditation's direct impact on peer selection. Accreditation likely correlates with other variables, leading to the selection of accredited peers due to their similarity in other features.

With confidence in the new methodology, a broader range of variables can be considered for updating the peer group in future years. Sources like IPEDS provide hundreds of variables, making comprehensive feature analysis a logical next step for refining the process. Past methodologies raised concerns about overfitting and weighting bias when using numerous correlated variables. By employing PCA, our approach enables the inclusion of a broader range of variables to assess similarity with the target institution.

At the time of this writing, administrative stakeholders are examining the salaries of institutions in the UPG. If the target institution is near the median of the UPG then our methodology will be further validated.
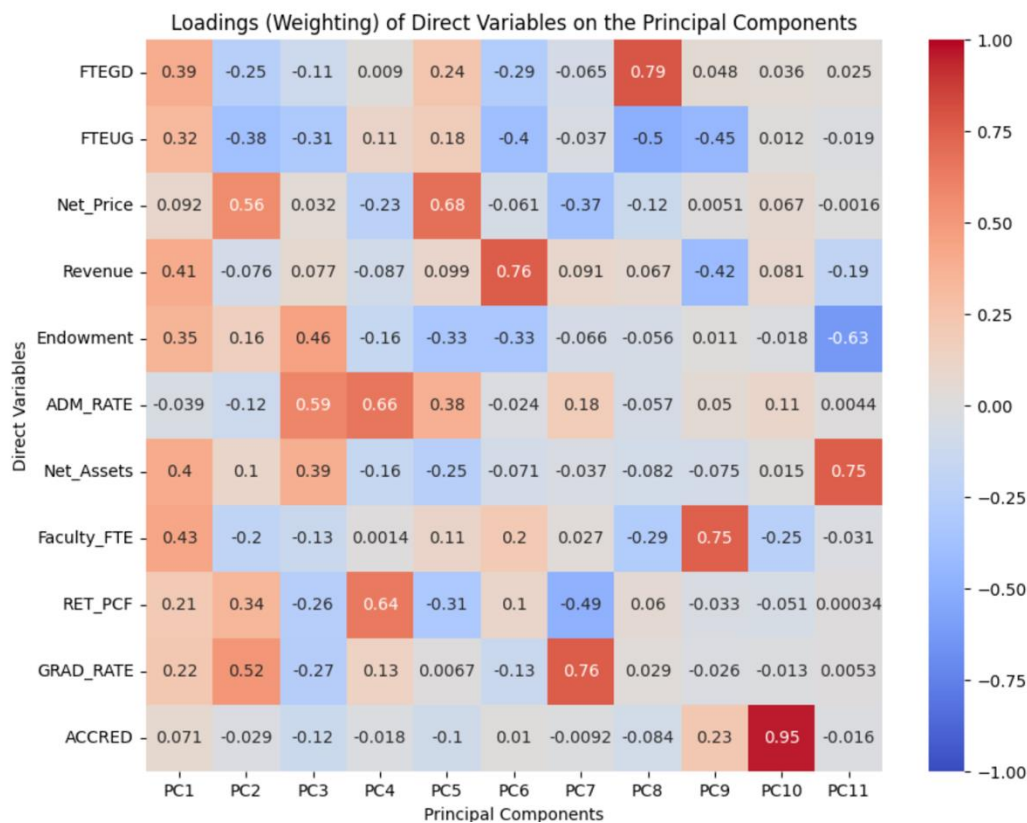
## 9. REFERENCES

AACSB (2025). Association to Advance Collegiate Schools of Business. (n.d.). AACSB International. https://www.aacsb.edu/

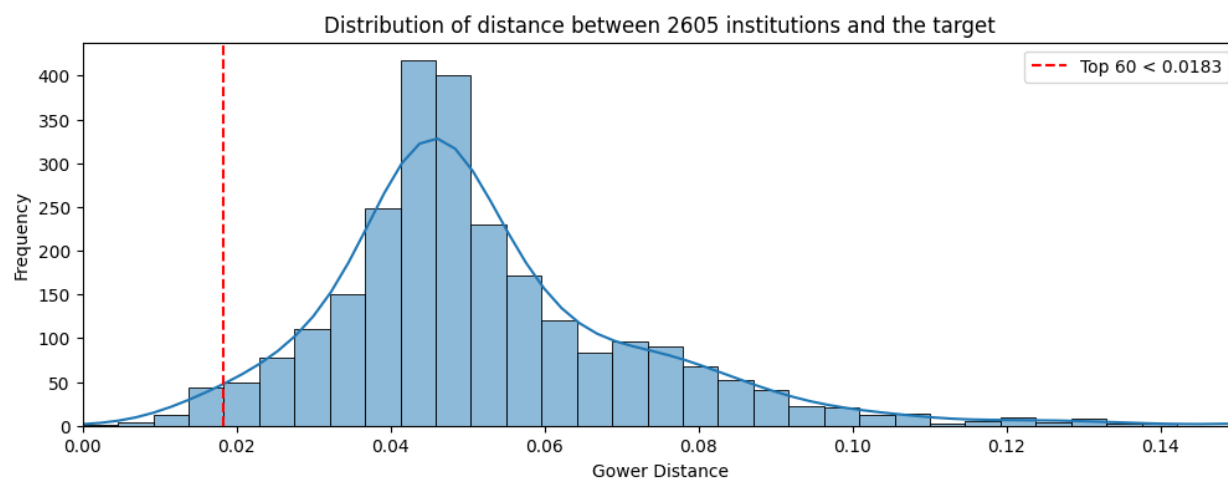Acton, R., Cortes, K. E., Miller, L., & Morales, C. (2024). Distance to degrees: How college

proximity shapes students' enrollment choices and attainment across race-ethnicity and socioeconomic status (EdWorkingPaper No. 24-1055). Annenberg Institute at Brown University. https://doi.org/10.26300/vjyg-ta27

American Council on Education (2025). Carnegie Classifications. https://carnegieclassifications.acenet.edu/

CCNE (2025). Commission on Collegiate Nursing Education. (n.d.). CCNE accreditation. American Association of Colleges of Nursing. https://www.aacnnursing.org/ccne-accreditation

D'Allegro, M. L. (2017). A case study to examine three peer grouping methodologies. The AIR Professional File, (142). https://doi.org/10.34315/apf1422017

D'Allegro, M. L., & Zhou, K. (2013). A case study to examine peer grouping and aspirant selection. The AIR Professional File, (132). https://doi.org/10.34315/apf1322013

Google. (n.d.). Google Colab (Version 3.0) [Software]. https://colab.research.google.com

Gower, J. C. (1971). A general coefficient of similarity and some of its properties. Biometrics, 27(4), 857–871.

Jolliffe, I. T. (2002). Principal component analysis (2nd ed.). Springer.

Jolliffe, I. T., & Cadima, J. (2016). Principal component analysis: A review and recent developments. Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 374(2065). https://doi.org/10.1098/rsta.2015.0202

Kelchen, R., Ortagus, J., Rosinger, K., Baker, D., & Lingo, M. (2024). The relationships between state higher education funding strategies and college access and success. Educational Researcher, 53(2), 100–110. https://doi.org/10.3102/0013189X231208964

Kline, D., Vetter, R., Clark, U., (2020). Understanding Campus Crime with a Multi-University Analytics System. *Journal of Information Systems Applied Research* 13(3) pp 21-28.

Lubis, A. H., Sihombing, P., & Nababan, E. B. (2020). Analysis of accuracy improvement in K-nearest neighbor using principal component analysis (PCA). Journal of Physics: Conference Series, 1566, 012062. https://doi.org/10.1088/1742-6596/1566/1/012062

McLaughlin, G. W., Howard, R. D., & McLaughlin, J. (2011, May 21–25). Forming and using peer groups based on nearest neighbors with IPEDS data [Paper presentation]. 51st Annual Forum of the Association for Institutional Research, Toronto, Ontario, Canada. https://eric.ed.gov/?id=ED504414

Mohammed, A., Lind, M., (2024). Examining Factors Influencing the Acceptance of Big Data Analytics in Healthcare. Journal of Information Systems Applied Research 17(2) pp 31-44. https://doi.org/10.62273/QNDU3179

NCES (2023). National Center for Education Statistics. Integrated Postsecondary Education Data System (IPEDS), 2022–23 final data. U.S. Department of Education. https://nces.ed.gov/ipeds/use-the-data/download-access-database

Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie, T., Tibshirani, R., Botstein, D., Altman, R. B. (2001). Missing value estimation methods for DNA microarrays. Bioinformatics, 17(6), 520–525. https://doi.org/10.1093/bioinformatics/17.6.520

Turley, R. N. L. (2009). College proximity and the urban isolation of urban youth. Social Science Research, 38(3), 628–646. https://doi.org/10.1016/j.ssresearch.2009.02.002

# Appendices and Annexures

## APPENDIX A

### Loadings (Weighting) of Direct Variables on the Principal Components

| Direct Variables | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | PC11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FTEGD | 0.39 | -0.25 | -0.11 | 0.009 | 0.24 | -0.29 | -0.065 | 0.79 | 0.048 | 0.036 | 0.025 |
| FTEUG | 0.32 | -0.38 | -0.31 | 0.11 | 0.18 | -0.4 | -0.037 | -0.5 | -0.45 | 0.012 | -0.019 |
| Net_Price | 0.092 | 0.56 | 0.032 | -0.23 | 0.68 | -0.061 | -0.37 | -0.12 | 0.0051 | 0.067 | -0.0016 |
| Revenue | 0.41 | -0.076 | 0.077 | -0.087 | 0.099 | 0.76 | 0.091 | 0.067 | -0.42 | 0.081 | -0.19 |
| Endowment | 0.35 | 0.16 | 0.46 | -0.16 | -0.33 | -0.33 | -0.066 | -0.056 | 0.011 | -0.018 | -0.63 |
| ADM_RATE | -0.039 | -0.12 | 0.59 | 0.66 | 0.38 | -0.024 | 0.18 | -0.057 | 0.05 | 0.11 | 0.0044 |
| Net_Assets | 0.4 | 0.1 | 0.39 | -0.16 | -0.25 | -0.071 | -0.037 | -0.082 | -0.075 | 0.015 | 0.75 |
| Faculty_FTE | 0.43 | -0.2 | -0.13 | 0.0014 | 0.11 | 0.2 | 0.027 | -0.29 | 0.75 | -0.25 | -0.031 |
| RET_PCF | 0.21 | 0.34 | -0.26 | 0.64 | -0.31 | 0.1 | -0.49 | 0.06 | -0.033 | -0.051 | 0.00034 |
| GRAD_RATE | 0.22 | 0.52 | -0.27 | 0.13 | 0.0067 | -0.13 | 0.76 | 0.029 | -0.026 | -0.013 | 0.0053 |
| ACCRED | 0.071 | -0.029 | -0.12 | -0.018 | -0.1 | 0.01 | -0.0092 | -0.084 | 0.23 | 0.95 | -0.016 |

## APPENDIX B

### Distribution of distance between 2605 institutions and the target



Top 60 < 0.0183

**APPENDIX C**



Correlation Matrix of Key Variables



Distribution of FTEGD



Distribution of FTEUG